

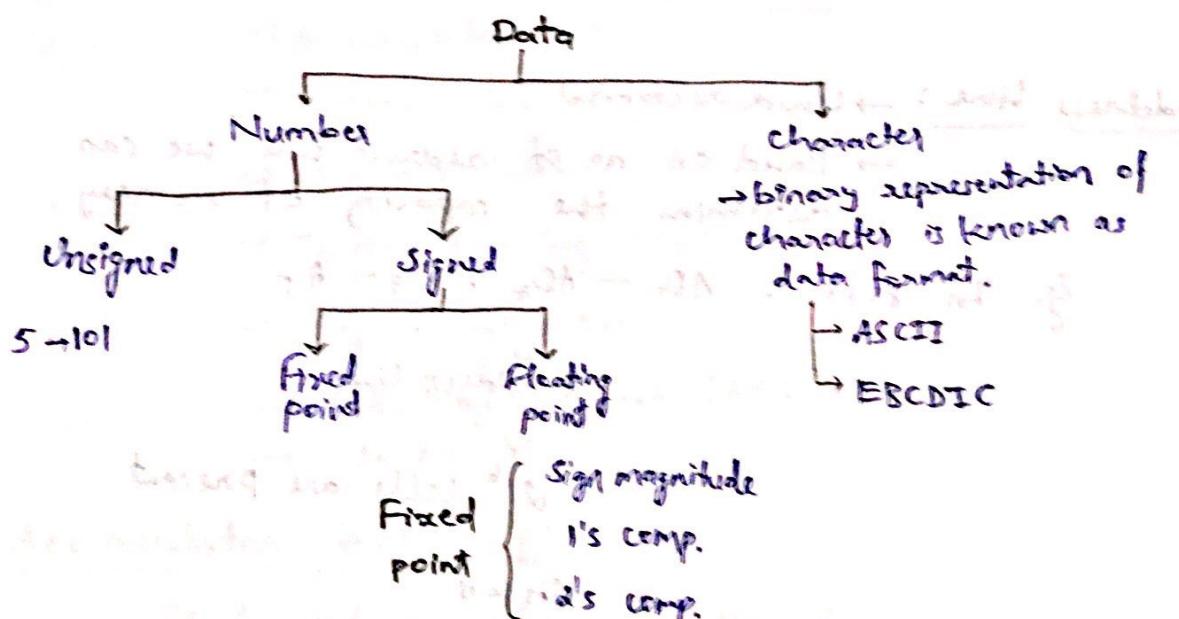
* COMPUTER ORGANIZATION & ARCHITECTURE *

* Computer Architecture :

Conceptual design and fundamental operational structure of a computer is known as computer architecture.
→ CPU design, Instructions, Addressing mode, Data format

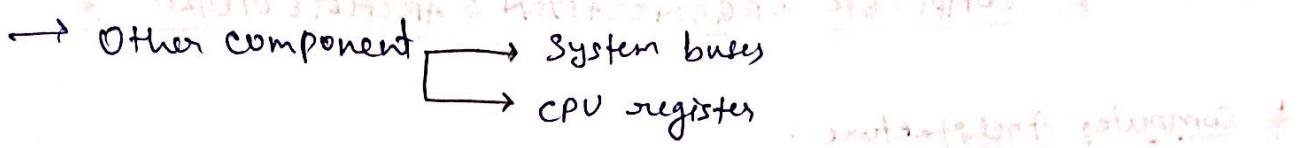
* Computer Organization : It deals with physical devices and their interconnection with a perspective of improving the performance.

- I/O organization
- Memory organization
- Performance improvement
 - ↳ Pipelining

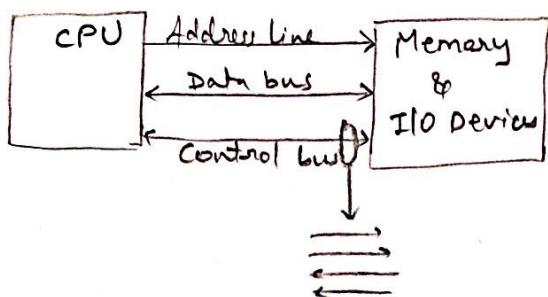


* Component of a computer

- CPU
 - ↳ Control Unit
 - ↳ ALU
- Memory
 - ↳ Primary or Main memory
 - ↳ Secondary Memory or Auxiliary
- I/O
 - ↳ Input
 - ↳ Output



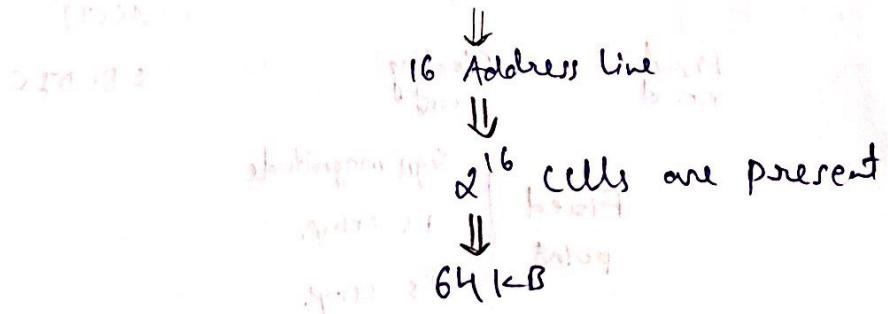
System Bus: Collection of communication line within the computer system.



Address line: → Unidirectional

→ Based on no. of address line we can determine the capacity of memory.

Eg. In 8085: $A_0 - A_7$, $A_8 - A_{15}$



Data line: These lines are used to carry the binary sequence of between CPU memory and input output.

→ These lines are bidirectional.

→ Based on no. of data line we can determine word length of processor.

Eg. 8085: $A_0 - A_7$, $A_8 - A_{15}$

↓
8-bit processor

Control Lines: These lines are used to carry the control signals and timing signals.

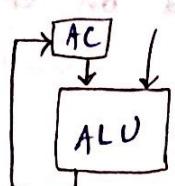
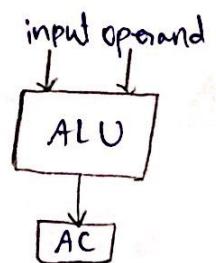
→ Control signals used to indicate the type of operation and timing signals are used to synchronise memory, and Input output operations with the processor clock.

* **Registers:** Small memories within the CPU.

- General purpose registers → Used to store general content like intermediate result.
- Special purpose registers
 - Accumulator (AC)
 - Program counter (PC)
 - Stack pointer (SP)
 - Instruction Reg. (IR)
 - Flag or status Reg.
 - Addressing Register (AR)
 - Data Reg.

* **Accumulator (AC)**

It is used to store results of ALU and sometimes one of the inputs operand for the ALU also.



Accumulator based
Architecture

Input Operand

Both the inputs taken from CPU registers called as register-based architecture

Stack-based Architecture

→ Both the inputs taken from stack

Complex system Architecture

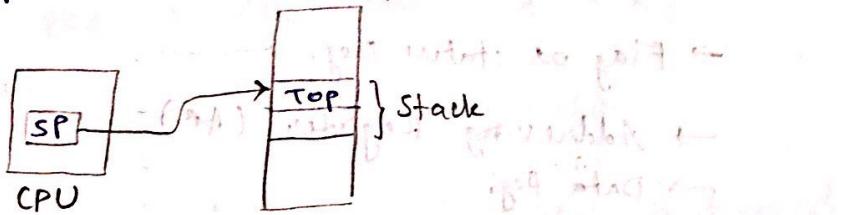
→ Both the inputs taken either from reg. or from memory

Program Counter (PC)

It is used to store the address of the next instruction to be executed.

Stack Pointer Reg.

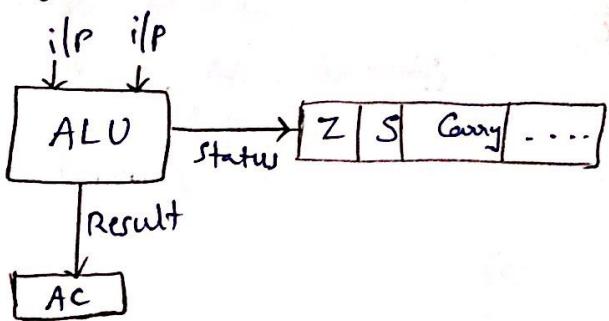
It is used to store address of top of process stack.



Instruction Register

It is used to store the current instruction which is being executed on CPU.

Flag or status Reg.



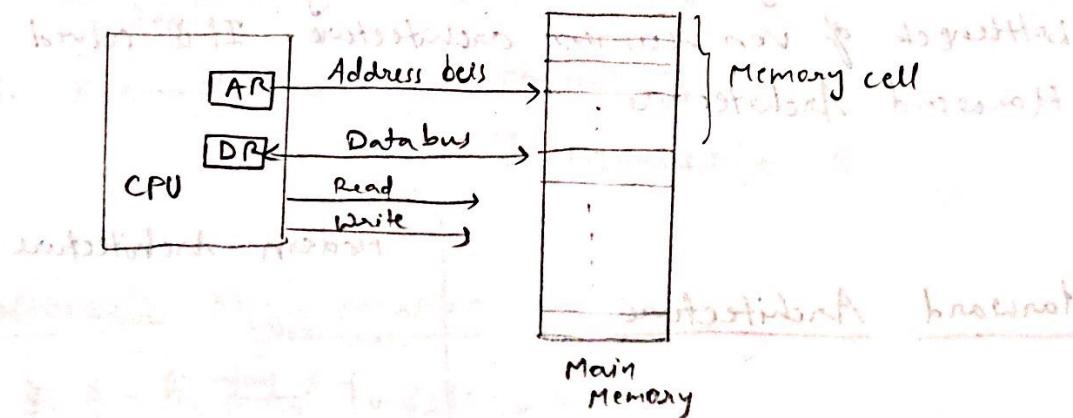
- It is used to store the status of ALU result.
 - These flags are used to evaluate conditions.
Eg. if $(a > b) \rightarrow a - b = 0$, we ignore memory
and go to branching off the first path.
- Program Status Word
- only status register } depends on implementation
AC + status register }

Address Register (AR) (Memory Address Register)

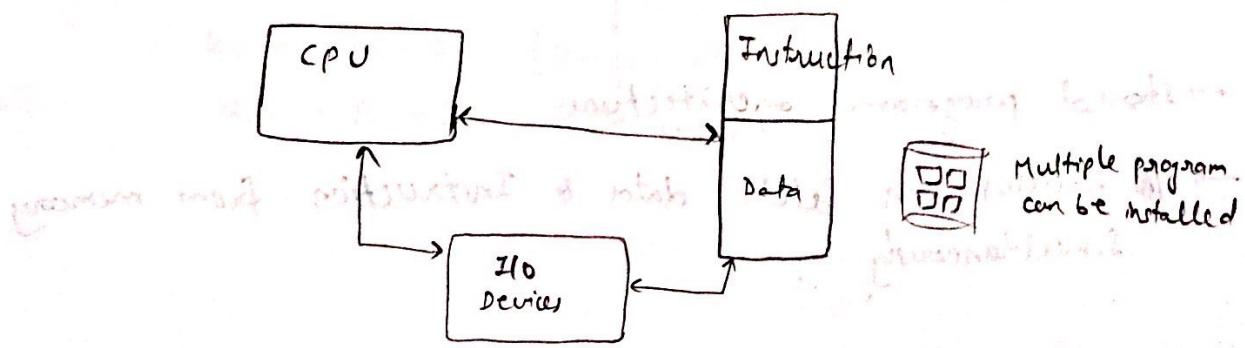
It is used to send address to memory.

Data Register (Memory Data Registers) (Mem. Buffer Reg.)

It is used to send data to memory (Mem. Write) and to receive data from memory (Mem. Read).



Von Neumann Architecture



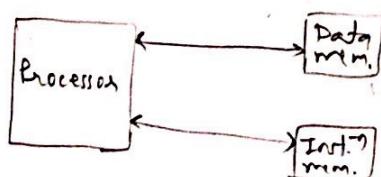
→ Before von neuman architecture, there were ~~fixed~~ program computers, their function is very specific and they could not be programmed.

→ Stored Program Computers : These can be programmed to carry out many & different tasks, applications are stored on them.
This is introduced by Von Neumann.

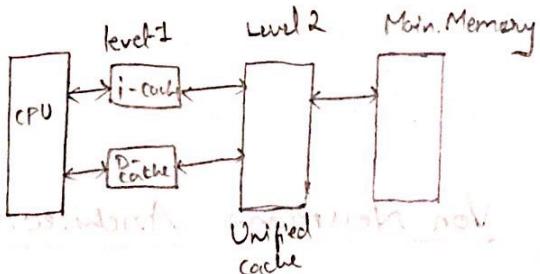
→ In von Neumann architecture, programs and data are stored in a separate storage unit called memory and are treated the same.

→ Processor can not fetch instruction & data both simultaneously from the memory this is known as bottleneck of von neumann architecture. It is solved by Harvard Architecture

Harward Architecture



Modern Architecture



→ stored program architecture

→ process can fetch data & Instruction from memory simultaneously

Micro Operations : The operations performed on data stored in the registers.

- CPU performs all the operations by performing one or multiple micro operations.
- Atomic operation of processor.

Types of Micro Operation

(i) Data Transfer / Register Transfer

(a) Register to Register

$$\text{e.g. } R_1 \leftarrow R_2 \\ R_2 \leftarrow R_1$$

(b) Memory operations

Read

Write

$$\text{e.g. } R_1 \leftarrow M[\text{Address}] \quad \text{Memory} \leftarrow \text{CPU}$$

$$\text{e.g. } M[\text{Address}] \leftarrow R$$

(ii) Arithmetic Microoperations

$$R_1 \leftarrow R_2 + R_3 \quad (\text{Addition})$$

$$R_1 \leftarrow R_2 - R_3 \quad (\text{Subtraction})$$

$$R_1 \leftarrow \bar{R}_2 \quad (\text{1's complement})$$

$$R_1 \leftarrow \bar{R}_2 + 1 \quad (\text{2's complement})$$

$$R_1 \leftarrow R_1 + 1 \quad (\text{Increment})$$

$$R_1 \leftarrow R_1 - 1 \quad (\text{Decrement})$$

(iii) Logical Microoperation

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

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

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

$$R_1 \leftarrow \overline{R_2 \oplus R_3} \quad \text{logical XNOR}$$

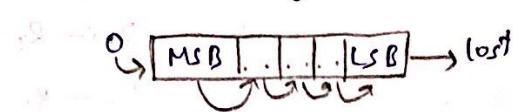
(iv) Shift Microoperation

Logical → all bits are moving left or right along with sign bit

Shift left

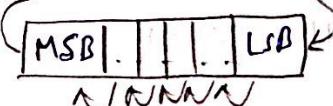


Shift right

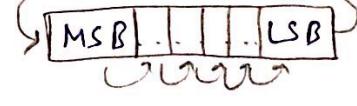


Circular

left

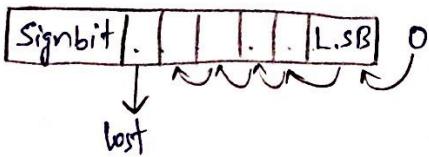


right

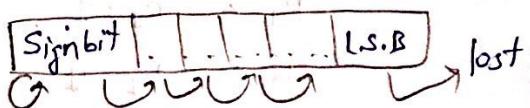


Arithmetic → All the bits are moving left or right except sign bit

Shift left

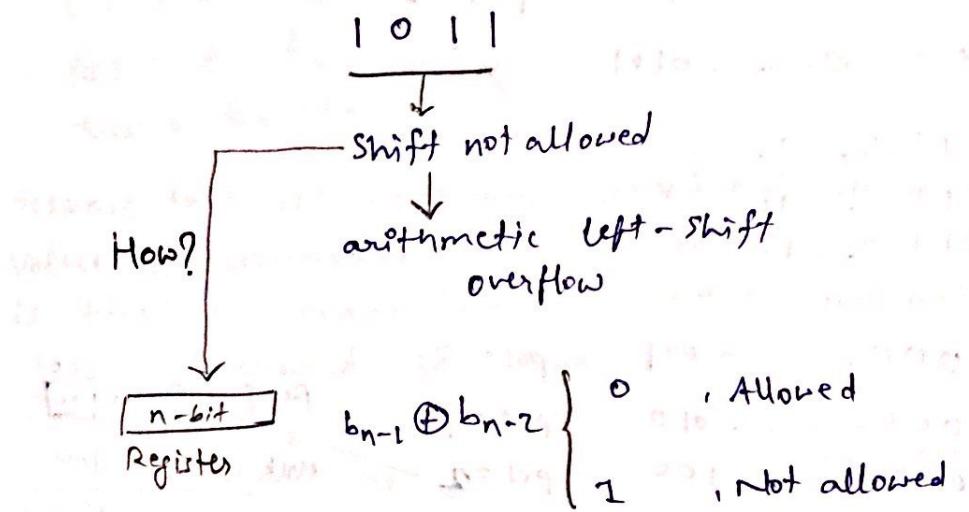


Shift Right



→ perform only when CPU ensure that sign bit will not change

Eg. Arithmetic Shift Left



Q- Consider a new instruction named Branch ^{on} Bit set the instruction "bbs" reg, pos, label" jumps to label if bit in position pos of register operand ~~reg~~ reg is 1. A register is 32-bits wide & bits are numbered 0 to 31, bit in position 0 being the LSB. Consider the following simulation of this instruction on process that does not have bbs implemented.

temp \leftarrow reg & mask

Branch to label if temp is non-zero.

Variable temp is a temporary reg. for the correct immaturity the variable mask must be generated by

(A) mask $\leftarrow 0x1 << pos$

(B) mask $\leftarrow 0xFFFFFFFF >> pos$

(C) mask $\leftarrow pos$

(D) mask $\leftarrow 0xF$

Rg.

...	bit	...
-----	-----	-----

80 79 78 77 76 75 74 73 72

Pos

Msk

0...	0	1	0....	0
------	---	---	-------	---

Temp 0...0 1 0...0

$$t_1 = a + b$$

$$t_2 = c + d$$

$$t_3 = e - t_2$$

$$t_4 = t_1 - t_3$$

Assume that all operands are initially in memory. the final value of computation should be in the memory. what is the minimum no. of move instruction in the code generated for basic block.



MOV M[1000], R₁
ADD M[2000], R₁
MOV M[3000], R₂
ADD M[4000], R₂
SUB M[5000], R₂
SUB R₁, R₂
MOV R₂, M[6000]

MOV M[1000], R₁
ADD M[2000], R₁
MOV M[3000], R₂
ADD M[4000], R₂
SUB M[5000], R₂
SUB R₁, R₂
MOV R₂, M[6000]

minimum no. of move instⁿ = 3

* Instruction

A group of bits which instructs computer to perform some operation.

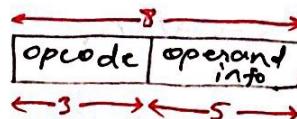
e.g. Add a, b
opcode

Instruction Format

opcode | operand info.

Instruction

Example: For a CPU, instruction is of 8-bits.



Max. operation possible = $\alpha^3 = 8$

Max. type of instructions supported by CPU

000	- ADD
001	- SUB
010	- MOV
011	- :
100	- :
101	- :
110	- :
111	- :

⇒ Types of instruction

4-address instruction

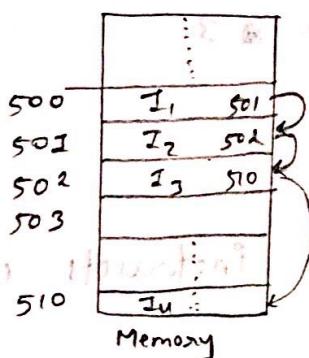
Maximum four addresses can be specified within an instruction.

Such instructions are used in those computer which do not have program counter.

Opcode	Add 1	Add 2	Add 3	Add 4
For operand			Next instruction address	

Such instructions are used in those computer which do not have program counter.

- All processor support the program counter register as a mandatory register, therefore this format is not used.



- (i) Larger sized instructions, hence more memory is required to store the program.
- (ii) If program is relocated then the fourth address part of each instruction should be updated and update of instruction is costly operation.

3-address Instruction

Maximum three addresses can be specified within an instruction.

opcode	Add. 1	Add. 2	Add. 3
For operand			

- Those computer which are using 3-add. instruction have to use program counter.

Example : Add r_0, r_1, r_2

- This 3-add. instruction format is used in modern computer system.
- Smaller size compared to 4-add. instⁿ, hence required less memory to store program.
- If relocation is required then no need to update instructions.

2-address Instruction

Max. Two addresses can be specified in instruction.

Opcode	Add. 1	Add. 2
For operand		

opcode	Add. 1.	Add. 2.	Source	Destination
			reg.	reg.

- One of the operand is use as source & destination both.



- the value of common operand is overwritten by the result of operation
- More no. of instruction required for program as compared to 3-add. instⁿ.

1-address Instruction

Max. one address can be possible within instruction

opcode	Add I.
--------	--------

→ Accumulator is used as 2nd operand implicitly.

Example: ADD R₁



$$AC \leftarrow AC + R_1$$

1) Load M[5000]

2) STORE M[4000]



(i) More number of instructions for program as compare to 2-address instruction.

→ 1-addr. Instruction is used in Accumulator based Arch.

0-address instruction

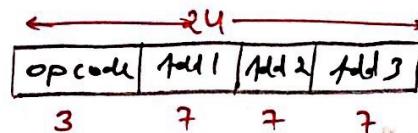
No any address specified in instruction.

opcode

→ 0-address instruction is used in stack based arch.



Q- Digital computer has 3 Byte long instruction, 7 bit address. Min & Max 3-address instⁿ, supported by this system.



$$\text{Max.} = 2^3 = 8$$

$$\text{Min.} = 1$$

Q- Computer supports 64 2-add. instⁿ in which each add. is of 8-bit. Size of an instⁿ is

#64	8bit	8 bit

$$\therefore 2^6 = 64$$

$$\therefore 6 + 8 + 8 = 22 \text{ bit} \approx 3 \text{ Byte}$$

Q- In above Question if we assume that instⁿ are stored in memory in Byte aligned fashion, the amount of memory required to store the program stack with contains 200 instⁿ is

$$22 \times 200 = \frac{4400}{8} = 550 \text{ Byte}$$

For 1 instⁿ to 22 bit ≈ 3 Byte

$$\therefore 200 \text{ inst}^n \rightarrow 200 \times 3 \\ = 600 \text{ Byte}$$

Instruction Set Architecture

Collection of all the instruction supported by CPU is known as inst. set. arch. of that CPU.

Q-17 40 distinct instⁿ

24 gen. reg.

32-bit instⁿ

opcode, 2 reg. operand and an immediate operand.

No. of bits available for immediate operand field

i).

32 Bit			
opcode	reg	reg.	imm. ope.
6	5	5	16

$$32 - 16 = 16 \text{ bit}$$

#instⁿ 40 reg. = 24 reg = 24

$$\log_2^{40} = \log_2^{24} + \log_2^{24} = 6 + 6 = 12$$

$$= 6 + 5 + 5$$

$$= 16$$

to determine no. of bits required for 32-bit word

Q-15 32-bit instⁿ. with 1 word long instⁿ

- 64-registers each of 32 bit long. words

45 instⁿ have immediate operand in addition to two register operands.

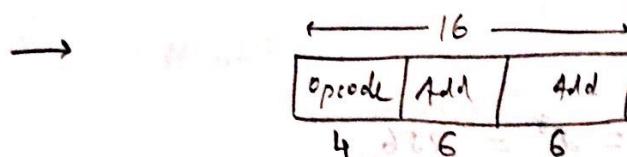
Assume immediate operand is an unsigned integer max. value of the immediate operand is.

6	6	6	14
Opcode	Reg.	Reg.	Imm. ope.
$\log 45$	$\log 64$	$\log 64$	
= 6	= 6	= 6	

∴ Max. value of Immediate operand

$$\Rightarrow 2^{14} = 16384 - 1 = 16383$$

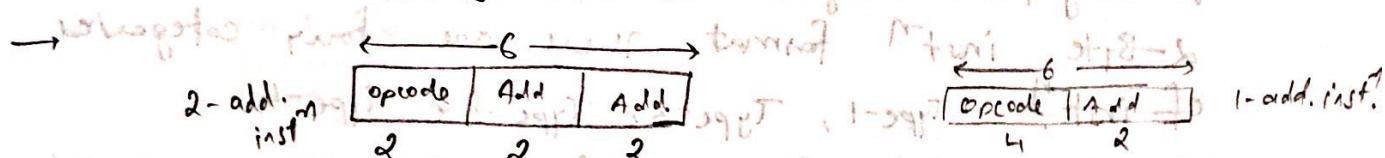
Q-16. If CPU used expanding op-code. If has 16 bit Inst? with 6-bit add. It supports one add & 2 add. inst? only. If there are 'n' two add. inst? the max. no. of one add. inst? is



$$\# \text{ 2 add. inst} = 2^4$$

~~Note~~ → $\# \text{ 1 add. inst} = (2^4 - n) \times 2^6$

Q- A digital computer has 6-bit instructions and 2-bit Add. It supports 2-add. & 1-add. inst's both. If there are 3 2-add. inst. then min & max. how many 1-add. inst's it can support.



$$\# \text{ 2-add. inst} = 2^2 = 4$$

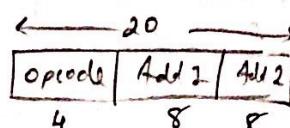
~~# 1-add. inst~~ $= (4 - 3) \times 2^2$

$$(2^2) \text{ binary numbers} = 4$$

$$\text{Min} = 1 \text{ to binary } 1-000$$

$$\text{Max} = 4 \text{ to binary } 100$$

Q- Consider a computer which have 20-bit inst & a 8-bit add. if there are 12 2-add. inst then max. how many 1-add. inst possible.



$$\# \text{ 2-add} = 2^4 = 16$$

$$\therefore \text{Max} = 1024$$

$$\# \text{ 1-add} = 16 - 12 = 4 \times 2^8$$

Q. - Digital computers has 32-bit instⁿ & 32-bit add.
 If there are 254 2-add. instⁿ & 8000 0-add. instⁿ. possible
 instⁿ? Then Max. how many 0-add. instⁿ possible?

32-bit for one instruction		
Opcode	Add1	Add2
8	12	12

$$\# \text{ 2-add. inst}^n = 2^8 = 256$$

$$\begin{aligned}\# \text{ 1-add. inst}^n &= 256 - 254 \\ &= 2 \times 2^1 = 2^2\end{aligned}$$

$$\begin{aligned}\# \text{ 0-add. inst}^n &= \frac{256 - 254}{2} = \frac{2}{2} = 1 \\ &= 2^0\end{aligned}$$

Ans. 1 instruction has 32 bits = 192×2^{12} bits. (192 bits for 12 bits)

Processor has 16 integer registers (R₀, R₁, ..., R₁₅) & 64 floating point registers (F₀, F₁, ..., F₆₃). It uses 2-Byte instⁿ format. There are four categories of instⁿ: Type-1, Type-2, Type-3 & Type-4.

Type-1 consists of 4 instⁿs. each with 3 integer register operands. (3Rs)

Type-2 category consists of 8 instⁿs. each with 2-floating point register operand (2Fs)

Type-3 consists of 14 instⁿs each with one integer reg. operand and 1-floating point reg. operand (1R + 1F)

Type-4 consists of N instⁿs each with a floating point reg. operand (1F).
 The max. value of N is.

$$\text{Type 1: } 4 \times 2^4 \times 2^4 \times 2^4 = 2^{14}$$

$$\text{Type 2: } 8 \times 2^4 \times 2^4 = 2^{15}$$

$$\text{Type 3: } 14 \times 2^4 \times 2^4 = 2^{15}$$



j¹⁶

Max. value of j = $2^{16} - 1 = 65535$

Total instruction = 2^{16}

Type 1: at (3A₃)

opcode	R	R/R
4-bit	4	4

$$\#inst = 4$$

$$\therefore \cancel{4} \times \cancel{4} \times \cancel{4} = 4^3$$

$$Max = 2^4 = 16$$

$$Used = 4$$

$$\underline{\text{Unused}} = 12$$

Type 2: (2F₃)

opcode	F	F
4-bit	6	6

~~Unused~~ unused

$$Max = 12 \times 2^0$$

$$\therefore \cancel{12} = 12$$

$$Used = 8$$

$$\underline{\text{Unused}} = 4$$

Type 3: (1R + (F))

opcode	R	F
6	4	6

$$\therefore \underline{\text{Unused}} = 4$$

$$\text{at first Max} = 4 \times 2^2$$

$$\therefore \underline{\text{Unused}} = 16$$

$$Used = 14$$

$$\underline{\text{Unused}} = 2$$

Type 4: (IF)

opcode	F
6	6

Brackets at middle of 10 in 6

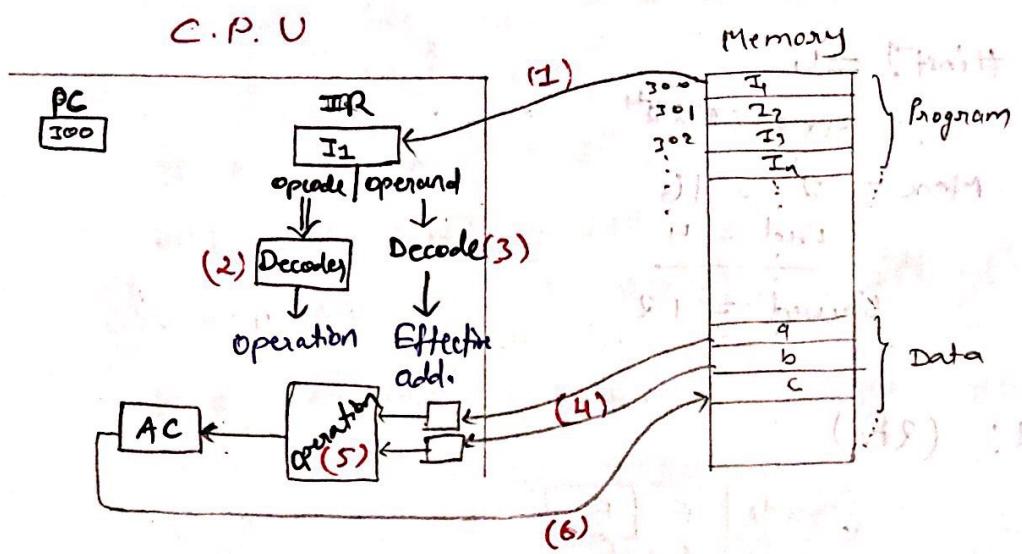
$$\therefore \underline{\text{Unused}} = 2^{2+4} = 2^6$$

$$Max = 2 \times 2^4$$

$$= 32$$

* Instruction Cycle

Six phases are required to execute the instruction.



(1) Instruction Fetch: Using the program counter value the next instruction is copied from memory to Instruction Register. During the IF only, the value of PC is incremented by the size of instruction.

(2) Instruction Decode: CPU decodes the opcode part of instruction to obtain the instruction.

(3) Effective Address Calculation: CPU decodes the operand information part of the instruction to obtain the effective address.

Effective Address: Address of operand in computation type instruction and target address in branch type instruction is effective address.

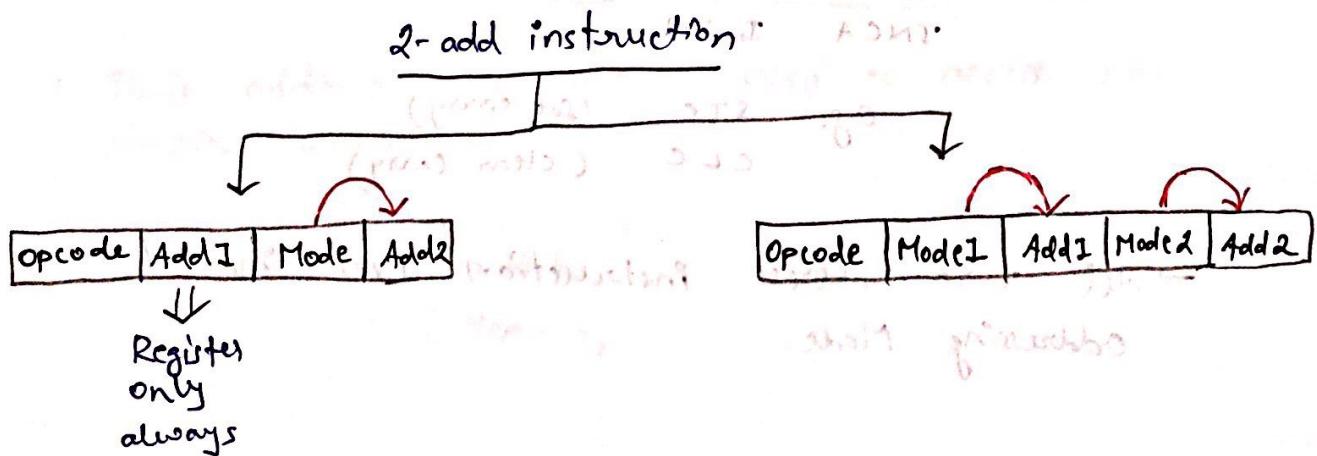
(4) Operand Fetch: The operands are fetched from the source using the effective address.

(5) Execution: CPU performs the operation and generates the result.

(6) Write Back: The generated result is copied back to the destination.

NOTE: All type of instructions do not require all six phases.

* Addressing Mode: It specifies how and from where the operand is obtained using the address field of instructions.



I-add. instr

Principle: In I-add. instruction, the operand address is implicit.

opcode	Mode	Add
--------	------	-----

Ex. ↓
00 ⇒ Operand value

01 ⇒ Register reference

10 ⇒ Memory Address

Memory reference

→ Addressing mode helps to obtain the Effective address.

* Types of Addressing Mode

(1) Implied Mode: The opcode definition itself specifies or (Implicit) how the operand, hence no need to specify any explicit address for the operand.

→ There is no effective address in implied mode.

Ex :

opcode	Mode	Add
--------	------	-----

INCA Implied

e.g. STC (Set carry)

CLC (Clear carry)

→ All zero address instruction uses implied addressing mode.

(ii) Immediate AM: The address field of instruction specifies the operand value.
→ This mode is used to initialize register with constant value.

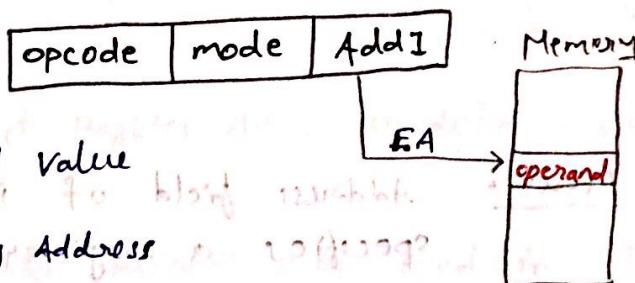
opcode	Add.
	↓ Data (Operand Value)

e.g. ADDI R10 23H : $r_{10} \leftarrow r_{10} + 23$

→ This addressing mode is used to access the constant.

(iii) Direct / Absolute AM:

The address field of instruction specifies the effective address.



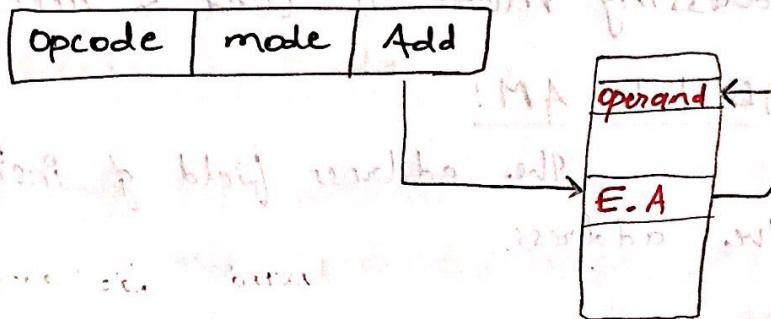
→ This addressing mode is used to access the static variable.

$$\begin{aligned} \text{Data} &= [\text{E.A}] \\ &= [\text{Memory}] \end{aligned}$$

(iv) Indirect Addressing Mode:

The address field of instruction specifies the address of E.A.

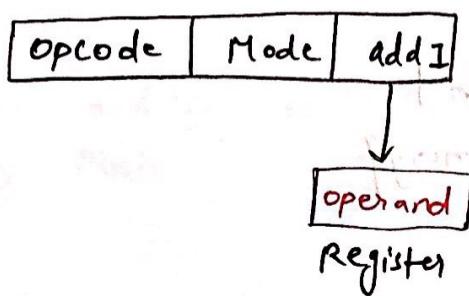
- This addressing mode is used to implement the pointer.
- In this addressing mode the data is present in the memory, that memory address may be available in register or memory.



(v) Register Mode:

Address field of instruction specifies a register which holds the operand.

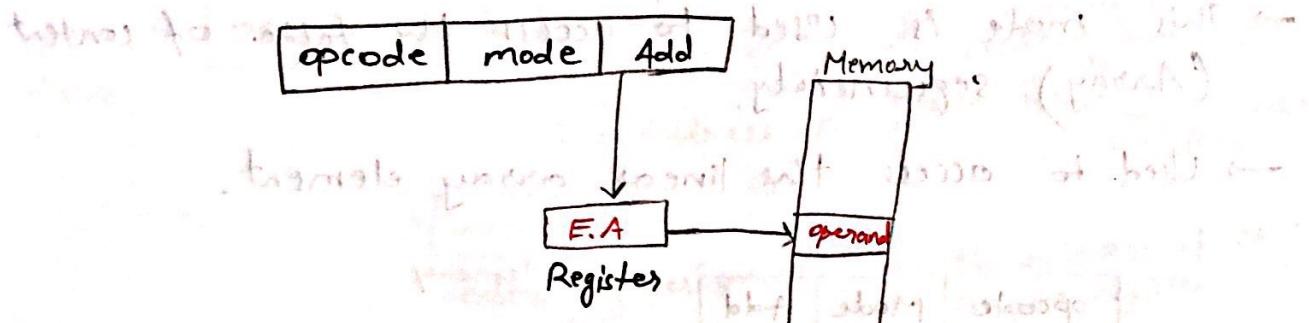
- E.A. is always the register name.



- this mode is used to access the local variable.

(vi) Register Indirect: The address field of instruction specifies a register which holds the effective address.

→ This mode is used to shorten the instruction length.



How ?

4 - GB RAM = Memory address = 32-bit

and

64 registers (G.P.Rs) = Reg. reference = 6-bit

For Direct ~~Register~~ A.M = $x + y + 32$ bit

For Indirect Register A.M = $x + y + 6$ bit

Hence, Register Indirect A.M used to shorten the instruction length

Access time?

Time to access operand in ~~Register mode~~ =

Direct Addressing mode = 1 memory access time

Time to access operand

in Indirect register A.M = 1 reg. access time + 1 Mem. access time

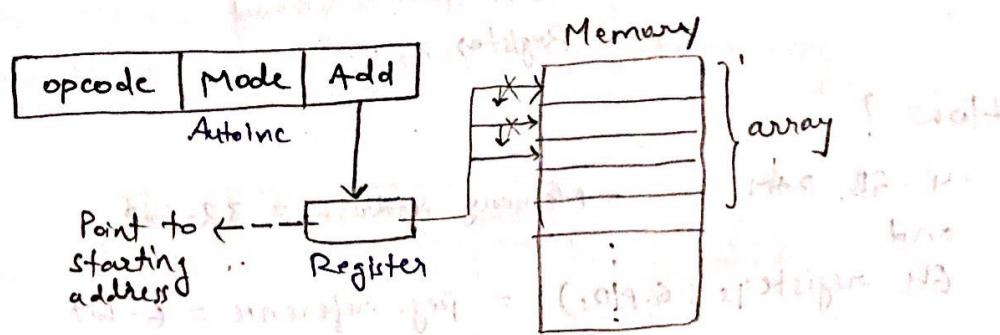
\approx 1 memory access time

(vii) Autoincrement / Autodecrement

If it is a variant of register indirect mode

in which the content of register (EA) is automatically incremented or decremented.

- This mode is used to access the table of content (Array) sequentially.
- Used to access the linear array element.



- EA can be calculated either by increasing or decreasing the step size with the base address.

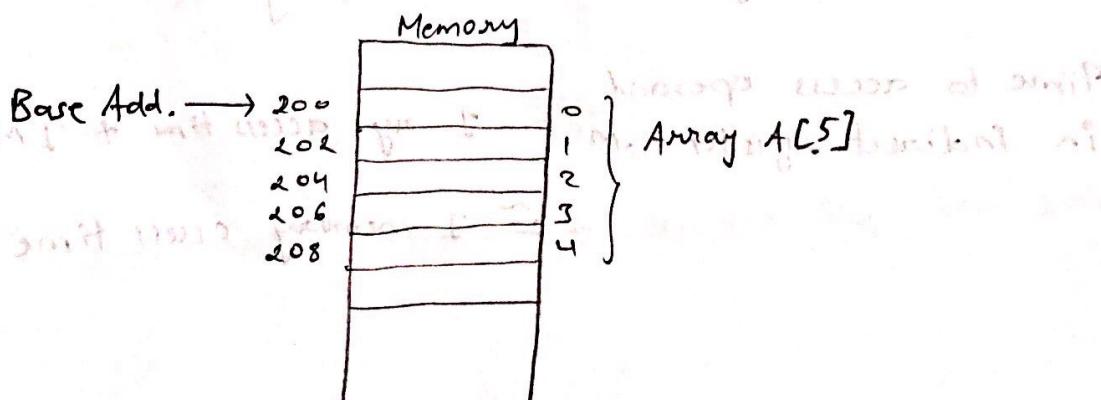
Autoincrement \Rightarrow post inc. offset treated as

Autodecrement \Rightarrow Pre dec.

(viii) Index Addressing Mode or Index Reg. Mode

This addressing mode is used to implement or access the array.

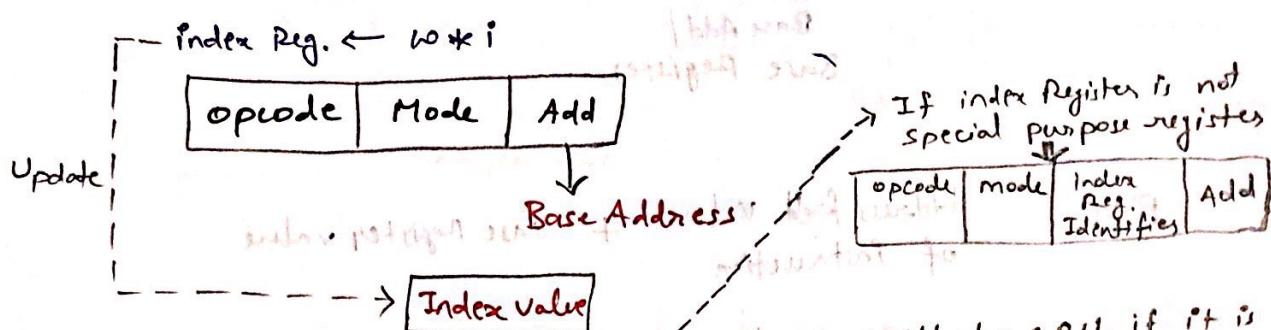
- This API is also called based index add. mode.



$$\text{address } (A[i]) = \text{Base Address} + (w * i)$$

↓
Size of an element

Index Value



Index Register → Access implicitly by C.P.U if it is special purpose registers.

$\therefore E.A = [\text{address fixed value}] + [\text{Index Register value}]$

$$= \text{Base Address} + \text{Index value}$$

(xii) ~~Indexed Mode~~

$$\begin{aligned} \therefore \text{Data} &= [E.A] \\ &= [\text{Base. Add.} + \underbrace{[\text{Index Pg.}]}_{\substack{\text{I Reg. Ref.} \\ \text{I Arithmetic computation}}} \Big] \end{aligned}$$

I memory reference

NOTE: If the program data (array) is relocated then the new base address should be updated in the address part of the instruction. Updating of instruction is costly operation.

To solve this problem we have Base Register Add. Mode

(ix) Base Addressing Mode

opcode	mode	Add
--------	------	-----

↓
Index

Base Add
Base Register

$$E.A = \text{Address field value of instruction} + \text{Base Register value}$$

NOTE: If program data (array) is relocated then the new Base address should be updated in the Base register, hence no need to update the instruction

(x) Base + Index Register A.M :

opcode	Mode	Add.
Base Add. Base Reg.	Index Index Reg.	

Add. field is of no use as we are using registers for Base Add and index.

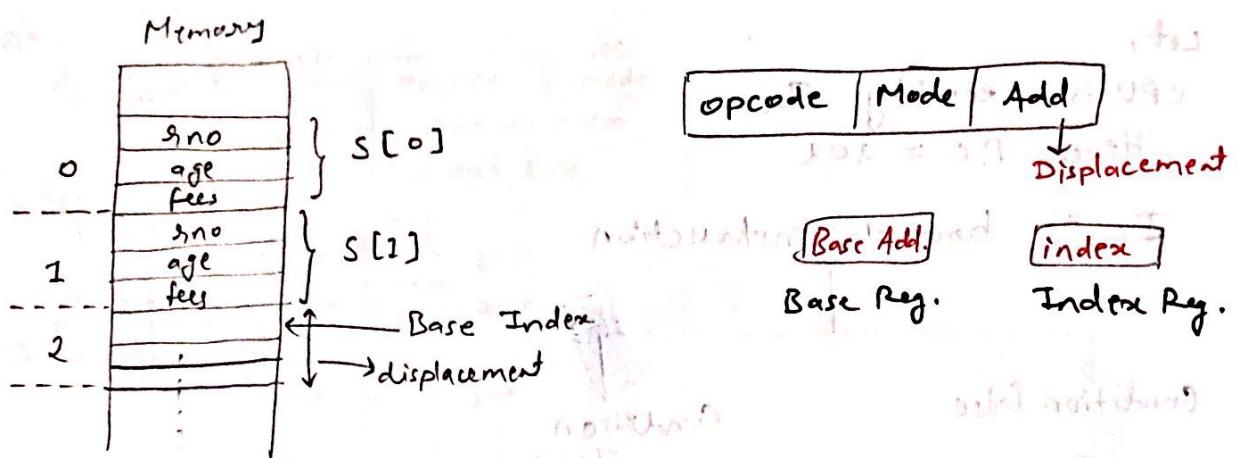
(xi) Base + Index Reg. + Displacement Mode

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

Example: struct Student S [10];

```
int rollno;  
int age;  
float fee;
```

}

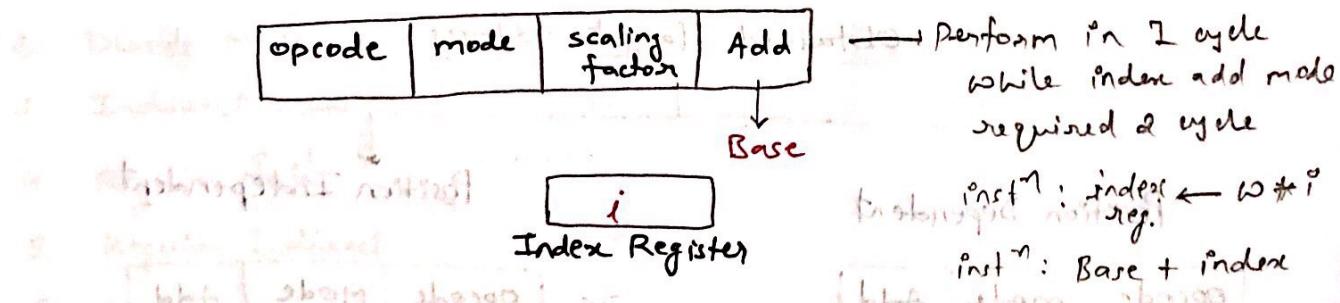


$$E.A. = [\text{Base Register value}] + [\text{Index Register value}] + [\text{Add. field value of inst.}]$$

↓
Base & Index
↓
Index addition
↓
Final Address

$$= \text{Base Address} + \text{Index} + \text{Displacement}$$

(xii) Scaled Mode:



$$E.A. = \text{Address field of Instruction} + \text{Scaling factor} * \text{Index Register value}$$

(xiii) PC - Relative Mode / position independent Mode

This mode is used for branch type of instructions.

Let,

CPU is executing I_2
Hence $P.C = 202$

I_2 is branch instruction

Target address = 206

Condition false

Condition True

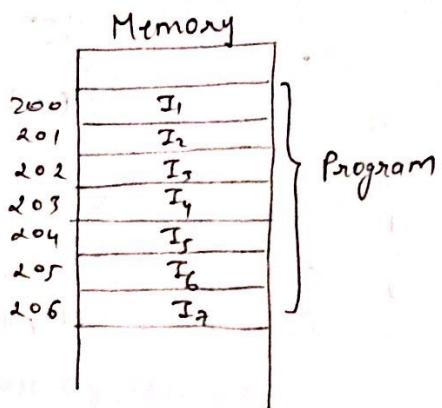
Next instruction
in the sequence
is executed
i.e., I_3

↓
Next instruction
executed is target
instruction (Let I_7)

\therefore Target Add. = 206

→ We have to
change P.C value

$P.C = \text{Target Address}$



→ No change in P.C.

Position Dependent

opcode	mode	Add
--------	------	-----

Target
Address

Position Independent

opcode	Mode	Add
--------	------	-----

$$\text{Target Address} = P.C + \text{Relative Value}$$

Instruction

Address

Value

Location

Address

Value

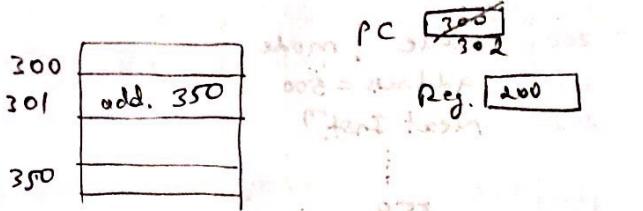
Example:

PC	200 202	opcode mode	
Reg.	400 399	address = 500	
XR	100	Next Inst. ^m	
Index Reg.			
		399	850
		400	700 + 7
		500	800
		600	900
		702	Target inst. ^m
		350	2957
		800	350

Program to locate target instruction (350) from address 2957. Address field has 350. No of bytes in Next inst. is 002. Hence, target address is 350. Address field has 350. Target inst. is 350.

Addressing Modes	Effective Add.	operand
1. Immediate	201	500
2. Direct Mode	500	800
3. Indirect Mode	800	350
4. Register Mode	-	400
5. Register Indirect	400	700
6. Auto-decrement	399	850
7. Index Mode	500 + 100 = 600	900
8. PC-relative	202 + 500 = 702	-

- Q- An instruction is stored at location 300 where its add. field is at location 202. The add. field value has 350. A processor reg. contains the no. 200. Evaluate the E.A if the add. mode is
- Direct Mode
 - Immediate Mode
 - Reg. Indirect
 - PC-relative



- (i) Direct Mode : E.A = 350
- (ii) Immediate Mode : E.A = 301
- (iii) Reg. Indirect : E.A = 200
- (iv) PC relative : EA = 302 + 350 = 652

Q. Relative Branch Mode Type instⁿ is stored at Memory 300. The Branch is made to an add. 450.

- (i) What should be the value of relative add. field of the instⁿ?
- (ii) Determine the value of PC before instⁿ fetch, after instⁿ fetch & after Execution phase of instⁿ.

→ P.C [300] (i) ~~450 - 300 = 150~~ → ~~150 = 300 + x~~ → ~~x = 150 - 300~~ → ~~x = -150~~

 301 (ii) P.C before instⁿ fetch = 300
 P.C after instⁿ fetch = 301
 P.C after instⁿ Execute = 450

(i) E.A = PC + Displacement (Relative add. field value)

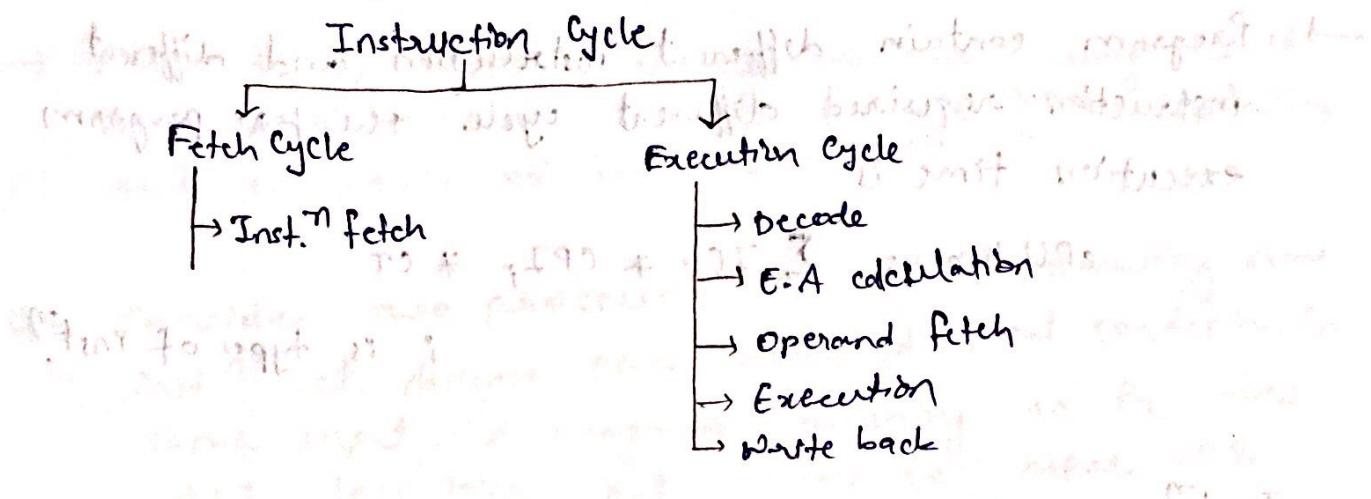
$$450 = 300 + x$$

$$450 - 300 = x$$

$$\therefore 150 = x$$

Q-2 For computers based on 3-field instⁿ format, each add. field can be used to specify reference to:

- S1: memory operand ✓
 - S2: processor register ✓
 - S3: An implied accumulator reg. → X
- Any special purpose reg. will not be accessed explicitly, it can be accessed implicitly only.



Q-10 ADD A [R₀], @B

No. of Memory cycle

$$M[A + [R_0]] \leftarrow M[A + [R_0]], M[M[B]]$$

1 mem. cycle 1 mem. cycle 2 Mem. cycle

Total = 4 Memory cycle.

* CPU Design

CPU cycle time: the amount of time in which CPU can perform one micro operation.

Cycle Per Inst? (CPI): No. of CPU cycle required to execute an instruction.

$$\text{Clock Rate} = \frac{1}{\text{CPU cycle time}}$$

$$\text{CPU time} = \frac{\text{Execution Time}}{\text{(E.T)}} = \# \text{Inst}^n * \text{CPI} * \text{cycle time}$$

$$\text{E.T} = \frac{\# \text{Inst}^n * \text{CPI}}{\text{clock rate}}$$

→ Program contain different instruction and different instruction required different cycle therefore program execution time is

$$c\text{PU.time} = \sum_i IC_i * CPI_i * CT$$

i is types of Inst?

Inst?	CPI	no. of inst?	Time
Load & store	6	40	240
ALU	4	20	80
Branch	5	30	150
other	4	10	40
			Total = 510
			Total = 100

$$\therefore CPI_{avg} = \frac{510}{100} = 5.1$$

510 divided by 100 for branch instruction steps 492

$$\therefore CPI_{avg} = \frac{\sum_{i=1}^m CPI_i * n_i}{\sum_{i=1}^m n_i}$$

at branching step ALU instruction 492 clock cycles
instruction 100 instruction 492 clock cycles

Performance of C.P.U

Performance of C.P.U is given by MIPS (Million Instⁿ per Sec.)

$$\text{MIPS} = \frac{\#inst.}{\text{Execution time} * 10^6} \quad (T3)$$

$$\text{MIPS} = \frac{\#inst \text{ H.s. do}}{CT + CPI * Cydetime * 10^6} = \frac{\text{Clock rate}}{CPI_{avg} * 10^6}$$

→ Sometimes MIPS may provide the wrong result because it does not include (consider) the capability and complexity of instruction.

- Q- Consider two processor P_1 and P_2 executing same instⁿ set. Assume that under identical condition for same input. A program running on P_2 takes 25% less time but incurs 20% more CPI. compared to program running on P_1 . If the clock freq. of P_1 is 1 GHz. then what is the clock freq. of P_2 .



$$1 \alpha = \frac{IC * CPI}{1 \text{ GHz}} * \cancel{\text{clock rate}}$$

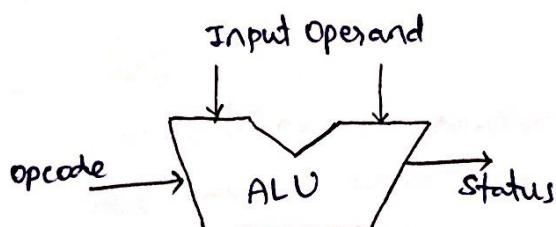
$$0.75 = \frac{IC * 1.2 CPI}{\text{clock rate}} * \cancel{10^9 \text{ ns/C clock cycle}}$$

$$\frac{1 \times 1 \text{ GHz}}{CPI} = \frac{0.75 * \cancel{10^9 \text{ ns/C clock cycle}}}{1.2 CPI} \rightarrow \cancel{10^9 \text{ ns/C clock cycle}}$$

$$\frac{1 \text{ GHz} * 1.2}{0.75} = \text{clock rate}$$

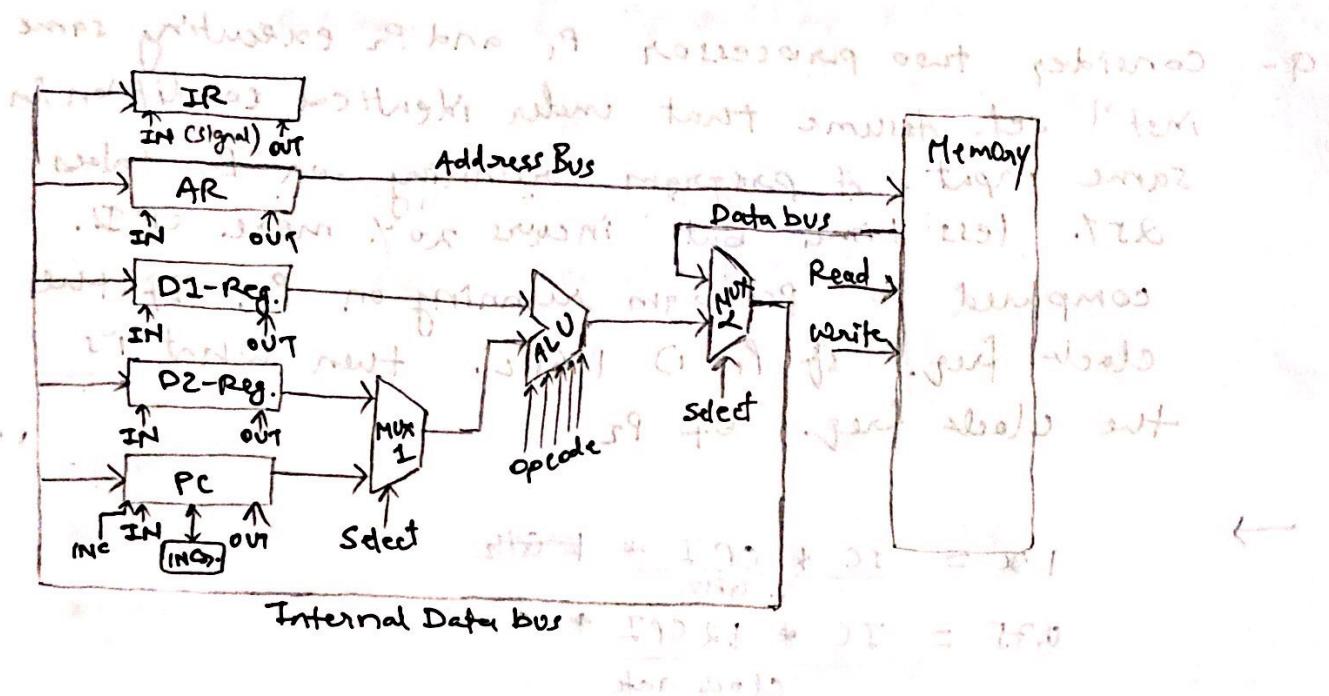
$$1.6 \text{ GHz} = \text{clock rate}$$

* ALU (Arithmetic & Logical Unit)



→ provide arithmetic & logic operation.

Data Path collection of functional unit ALU, mux etc. to perform dataprocessing operation B known as data path



Instruction fetch

$IR \leftarrow M[PC]$ --> Memory can not be access directly, AF can access memory. So we have to copy PC value to AR & IR can't access memory.

Then opⁿ should be performed in sequence if both are perform simultaneously then when internal data bus carrying data and, IAT signal of IP and AR both enable which leads to inconsistency in IR and AR.

These both operation can be perform simultaneously as they are mutually exclusive

Control Unit

It is part of CPU which generates control signal, and sends those signals to different functional unit so that the units can perform their respective operations.

Control Variable: The name of the control signal.
Ex. Read, Write

Control Word: Collection of all control signals generated by the control unit

IR	AR	D1	D2	PC	MUX1 Sele.	MUX2 Sele.	ALU	Memory Read/Write
IN	OUT	IN	OUT	IN	OUT INC		11111	

$AR \leftarrow PC$

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

1 0 0 1 0 0 0 0 0 0 0 1	X	0 0 0 0 0 0 0 0 0 0 0 1	IR
-------------------------	---	-------------------------	----

(word)

Control Unit Organization

- Hardwired Control Unit
- Microprogrammed Control Unit

 faster mode of operation (faster control unit)

→ Updation in control logic is difficult because rearranging the wire among circuits is difficult

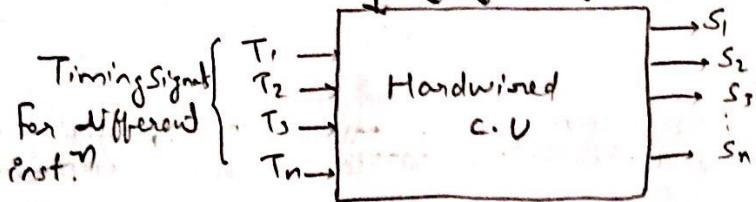
→ It cannot handle complex inst. as circuit is complex

→ More costlier

→ RISC is hardwired C.U.

Instruction $I_1, I_2, I_3, \dots, I_n$ \equiv set of $\{\equiv\}$

microop? \equiv set of microoperations to perform Inst? $I_1, I_2, I_3, \dots, I_n$



at T_1 ,
To execute
 $I_1 \wedge T_1$
will enable
these set
of signals

	I_1	I_2	I_3	\dots	I_n
T_1	S_1, S_2, S_4			\dots	
T_2	S_2, S_3	\dots	\dots	\dots	\dots
T_3	\dots	\dots	\dots	\dots	\dots
T_n	\dots	\dots	\dots	\dots	\dots

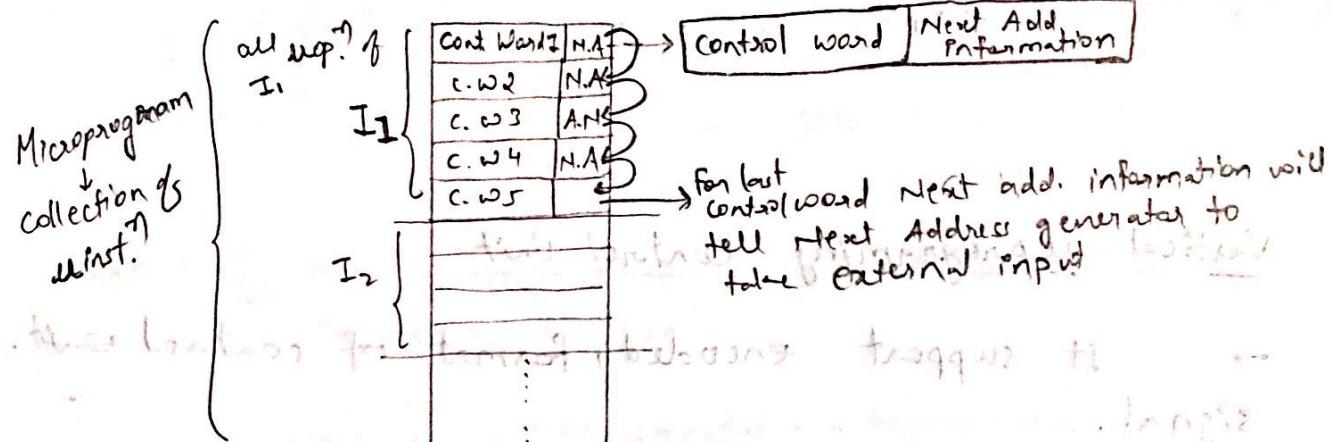
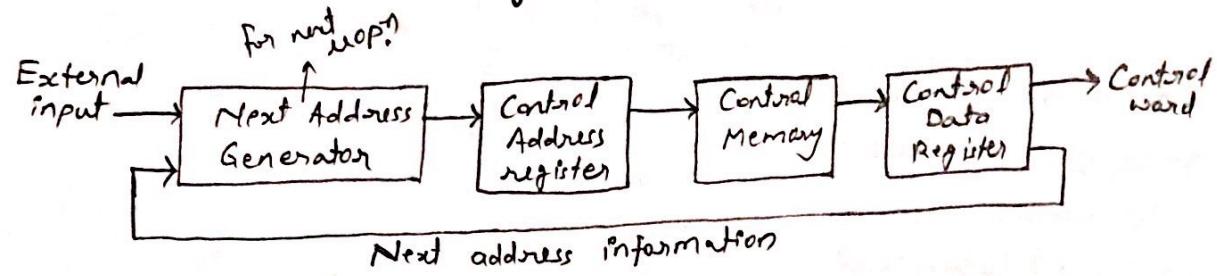
Microprogrammed control Unit

All possible control words are stored in a memory (control memory) and based on the requirement, the control word is read from the memory and the signals are sent to the respective units.

- 😊 (i) Updation in control logic is easy 😊
- 😊 (ii) Can handle complex microops (for storage control)
- 😢 (iii) slower than hardwired control unit as memory access is slower

→ Used in CISC (Complex Instruction Set Comp.)

Control word Sequencing.



Types of micro. control Unit

→ Horizontal ~~micro~~ control unit triggers all control signals at same time.

→ Vertical

stacking of various blocks of logic in width of one bit width of control word.

Horizontal micro. Control Unit

< Decoded format of control signal >

→ It supports decoded format of control signal means 2 bits / control signal.

Eg. If processor supports N control signal then this format required N bits in control field.

→ It supports longer control word.

→ There is no need of external decoder to generate the control signal. Therefore it is faster than the vertical micro. control unit.

Vertical reprogramming control unit

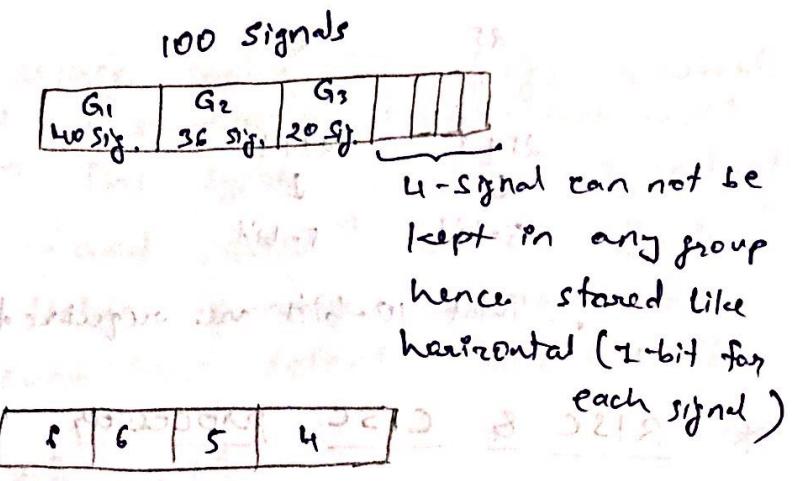
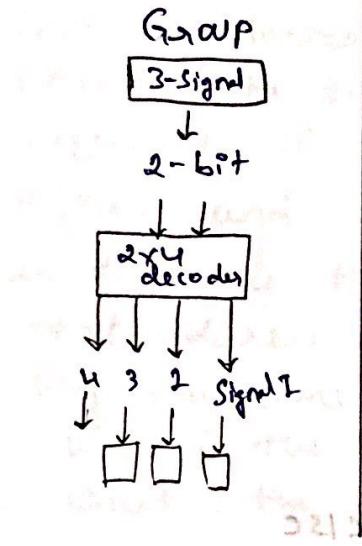
- It supports encoded format of control word signal.
- 'N' encoded control signal requires $\log_2 N$ bits.
- It supports shorter control word.
- There is need of external decoder to generate the control signal. therefore, it is slower than horizontal.

Longer format of 30 signals (horizontal format + 1 extra bit)
(default)
G1 | G2
16 | 14

point: min. 1 signal should be enabled from each group
min 0 signal [No any signal enabled from group]

4-bit | 4-bit
G1 | G2

point: 8-bit control word



NOTE: (1) The ascending order of control unit design in terms of speed is

Vertical → Horizontal → Hardwired logic
 low $\xrightarrow{\text{fast}} \text{High}$

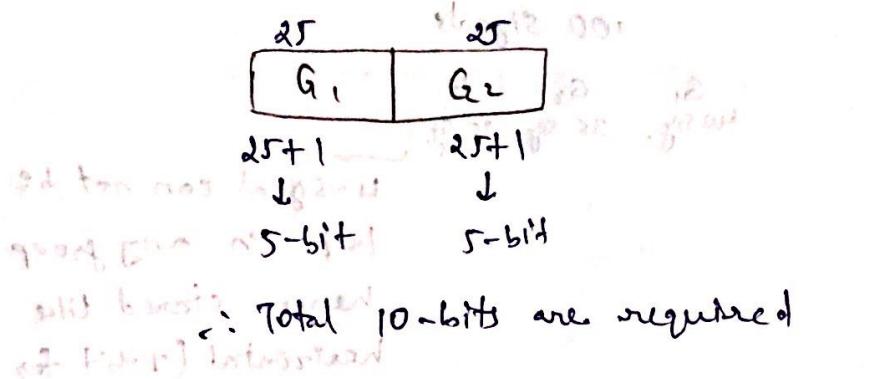
(2) The ascending order of control unit design in terms of flexibility.

Hardwired → Horizontal → Vertical
 low $\xrightarrow{\text{flexible}} \text{High}$

A progr. control unit is required to generate total of 25 signals. Assume that during any instⁿ at most 2 control signals are active. Min no. of bits required in the control word to generate the required control signal will be?

G ₁	G ₂
4	2 ¹

if
 all combination of 2 signals cannot be active at a time (if we require too signal & which are from G₂ group which is not possible)



* RISC & CISC processor.

RISC	CISC
→ Reduce instr ⁿ set computer	→ Complex instruction set computer
→ Support simple type of instr ⁿ	→ support complex instr ⁿ type
→ It support less no. of instruction	→ support more number of instruction
→ It has fixed size instruction	→ Variable size instruction
→ less no. of addressing mode	→ More no. of addressing mode
→ Support hardwired control unit	→ Support both hardwired & programmed control
→ Register to Register arithmetic op. ⁿ only	→ Reg. to Mem. or mem. to Reg. Arith. op ⁿ possible.
→ Support more no. of register	→ Support less no. of registers
Eg. Motorola, ARM powerPC	Eg. Pentium

Q - Consider Horizontal micro. control unit design in which CPU has to support 128 distinct instⁿ. Each instⁿ required 16 opⁿ. The system is using 63 control signals and has 1-add. minstⁿ format. Each minstⁿ has three fields namely control word, next address & condition select. There are 16 flag conditions.

What is the size of minstⁿ.

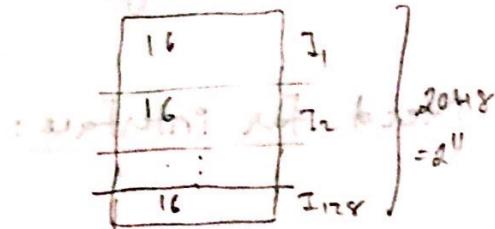
What is the size of control memory in bits.

$$\rightarrow \# \text{inst}^n = 128$$

$$\text{Each inst}^n = 16 - \text{micro-op}^n$$

$$\begin{aligned} \text{Total minst}^n &= 128 \times 16 \\ &= 2048 \\ &= 2^{11} \end{aligned}$$

$$\text{Control memory address} = 11 \text{ bits}$$



control word	Next Add	condition select
(63 - Cont.Sig.) 63-bit [Horizontal]	11-bit	4-bit $\lceil \log_2 16 \rceil$ flag

$$\text{Control memory size} = 2^{11} * 78 \text{ bits}$$

$$= 2^{11} * 78 * 8 \text{ bytes}$$

which will be reduced to 63 bits through OR operation. This is because each instruction has 63 control signals.

Condition select field for condition selection of 16 conditions will be converted to 4 bits. These 4 bits will be converted to 4 bits through OR operation.

Condition select field for condition selection of 16 conditions will be converted to 4 bits. These 4 bits will be converted to 4 bits through OR operation.

Condition select field for condition selection of 16 conditions will be converted to 4 bits. These 4 bits will be converted to 4 bits through OR operation.

Condition select field for condition selection of 16 conditions will be converted to 4 bits. These 4 bits will be converted to 4 bits through OR operation.

* I/O Organization

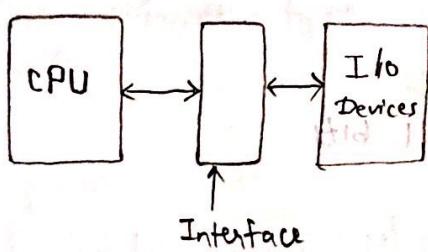
↳ I/O organization based upon hardware (hardware) → R
↳ I/O organization based upon software (software) → S
↳ I/O organization based upon both hardware & software.

Peripheral Device: The devices which are connected to CPU externally apart from memory.

↳ Peripheral Devices

- Input Device
- Output Device
- Storage Device

Need for interface:



- To communicate between CPU and I/O device.
- Interface provide synchronisation between them, as CPU is faster device & I/O devices are slower.
- Conversion of signals.
- I/O devices should not interfere to other device & CPU which is managed by interface.
- Conversion of data formats.
- Interfaces are H/w + S/w

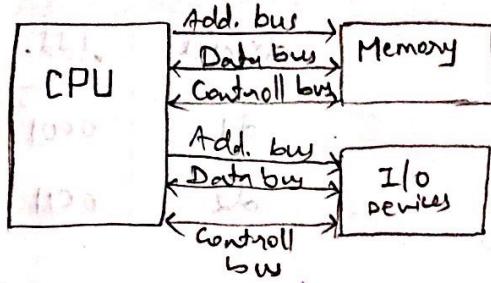
I/O processor

- Interface
- DMAE (DMA controller)
- I/O instruction execution

Communication of CPU with Memory and I/O

In three ways CPU can communicate with memory and I/O.

(i) Separate buses for both memory and I/O.



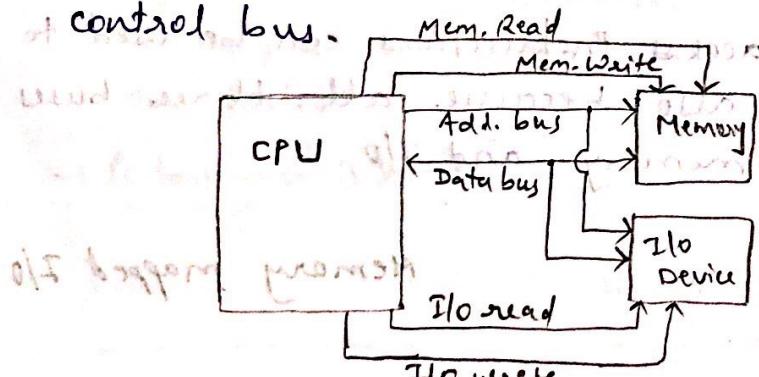
Adv.

Easy to implement

Disadv.

Cost will be more

(ii) Address bus & data bus is common but separate control bus. Mem. Read



Isolated I/O

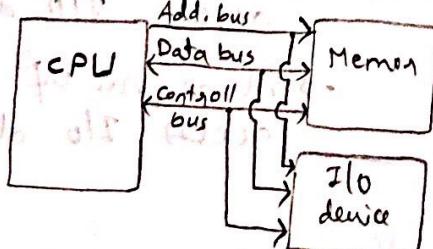
Ilo-mapped Ilo

Part-mapped I/O

(iii) All the bows

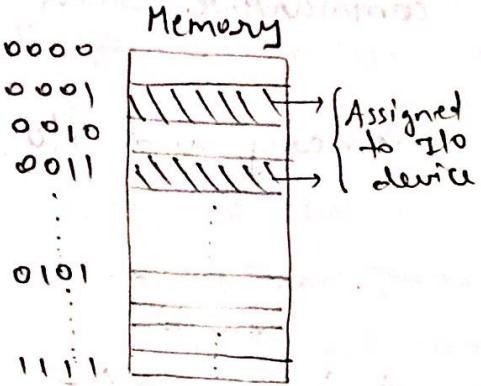
are common

for flo and memory



Memory-mapped I/O

Assume system with 16 Byte memory (Byte Addressable) and 2 I/O devices (d_1, d_2)



Few addresses are assigned to I/O devices

device	Add.
d_1	0001
d_2	0011

Whenever CPU generates address, it checks whether it is assigned to I/O devices, if its assigned then I/O devices is accessed, otherwise memory will access.

→ Some memory are wasted.

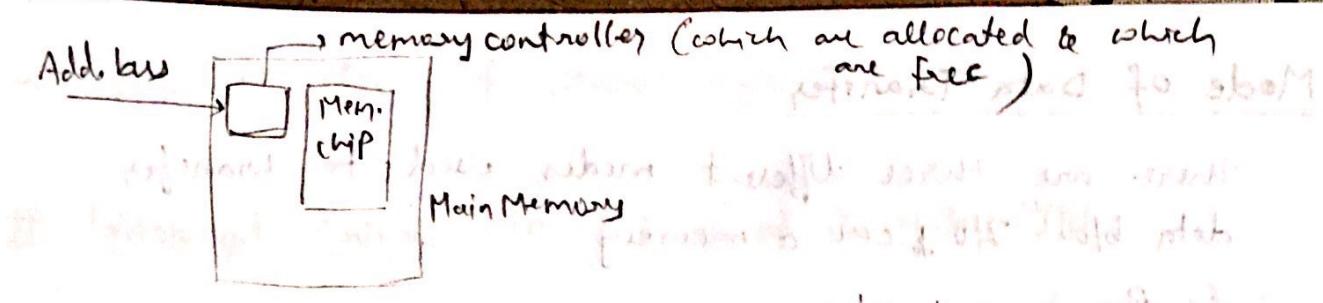
→ All the memory access instructions can be used to access I/O devices also because all three buses are common for memory and I/O.

I/O-mapped I/O

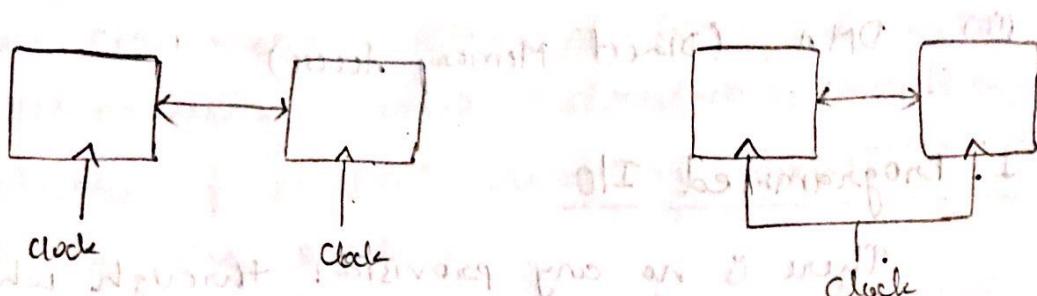
- No memory wastage
- less no. of instructions to access I/O devices compared to Mem. mapped I/O
- Less no. of addressing modes to access I/O devices
- Separate address space for both memory & I/O

Memory mapped I/O

- Some memory are wasted
- More no. of insts. to access I/O devices.
- More no. of add. modes to access I/O devices.
- Memory address space is shared between both.

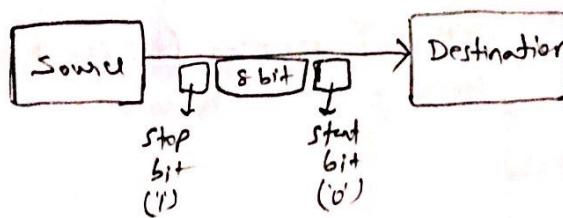


Serial Asynchronous Data Transfer



ATM should Against Slavery pros are in sync
Serial asynchronous transfer mode (ATM)
other data transfer. These ATM Plesio Synchronous

- separate clock generator for both device
- common clock generator for both device
- required external clock for both device, because one may be slower & one faster
- if any device is slower then transfer will be synchronized based on slower device
- performance is degraded



Apparent state transition of start bit = wait for next state
start bit

Apparent state transition of beginning with a start bit = state loss
wait working principle

Mode of Data Transfer

There are three different modes used to transfer data b/w I/O & CPU & memory

- (i) Programmed I/O
- (ii) Interrupt driven I/O
- (iii) DMA (Direct Memory Access)

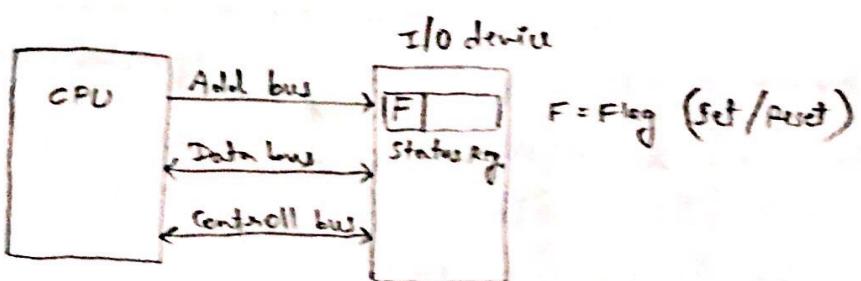
I. Programmed I/O

There is no any provision through which I/O device can inform CPU about data transfer. In such case, the device sets its own status and waits.

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

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

→ CPU time will be wasted because CPU undergoes waiting until the completion of the I/O operation. This waiting time is depending on the speed of I/O device.



$$\begin{aligned}\text{Total Data Transfer Time} &= \frac{\text{Time to check status}}{} + \frac{\text{Actual data transfer time}}{} \\ &= \frac{\text{Time required to read status register from device}}{} + \frac{\text{actual data transfer time}}{}\end{aligned}$$

→ Default size of status register is 1 Byte.

II. Interrupt Driven I/O / Interrupt Initiated I/O

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

When CPU receives interrupt from the device

(i) CPU executes current instruction completely

(ii) Status of current running process is stored back to stack

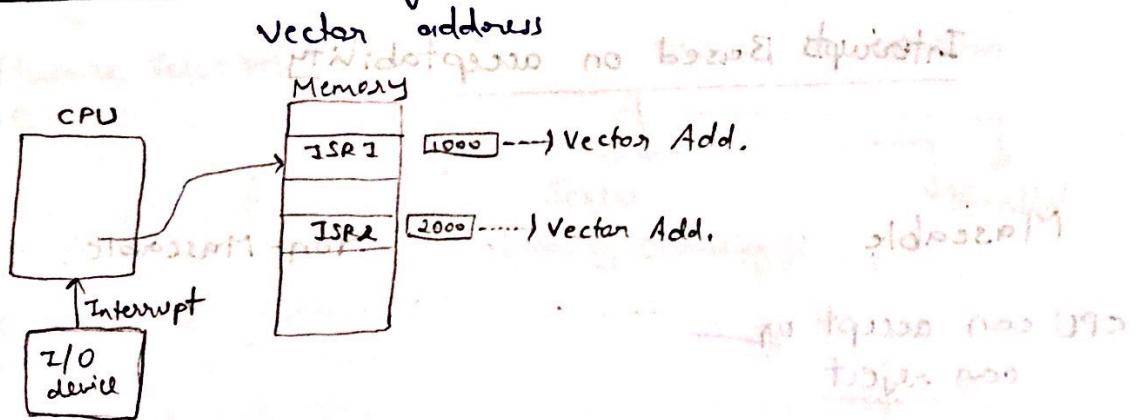
(iii) CPU takes a branch to service the interrupt

(iv) after the interrupt service, the previous program is resumed using the status from the stack.

ISR :- interrupt Service Routine

A function or routine execution of which services the interrupt.

Vector Address : Starting address of ISR is known as



Interrupts based on vectored Add.

Vectored Interrupt

device sends starting add. of ISR along with the interrupt signal. In such case CPU jumps to that add. directly and executes ISR.

Non-vectored interrupt (Scalar)

Device sends only interrupt signal to CPU in such case. CPU executes a default service routine first. which helps CPU to resolve the scenario and helps CPU to get actual ISR Address.

Non-vectored interrupt

- Every time an interrupt occurs CPU takes branch to a specific address only, this type of interrupt is non-vectored interrupt.

Interrupts Based on acceptability

Mascable

CPU can accept or can reject

Non-Mascable

Abort the interrupt

Keep in pending

Interrupts based on Source

External
(H/w)

Internal
(S/w)

If any I/O device generates interrupt

During interrupt instruction execution, because of any unacceptable error CPU interrupts itself.

Eg. (i) page fault

(ii) Addressing error in segmentation

(iii) Divide by zero.

→ If two device generates interrupt simultaneously, then the interrupt from highest priority device should be interrupt first.

Interrupts based on priority handling

Software Solution
(Polling)

Hardware Solution

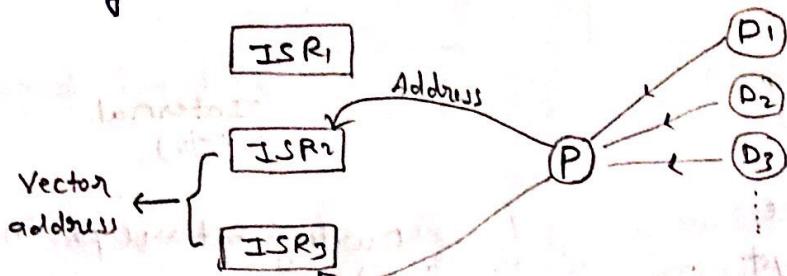
Serial

Parallel

all interrupts are serviced by branching to the same service program. This program then checks with each device if it is the one generating the interrupt. The order of checking is determined by the priority that has to be set. → this method is quite slow.

(Daisy Chaining)

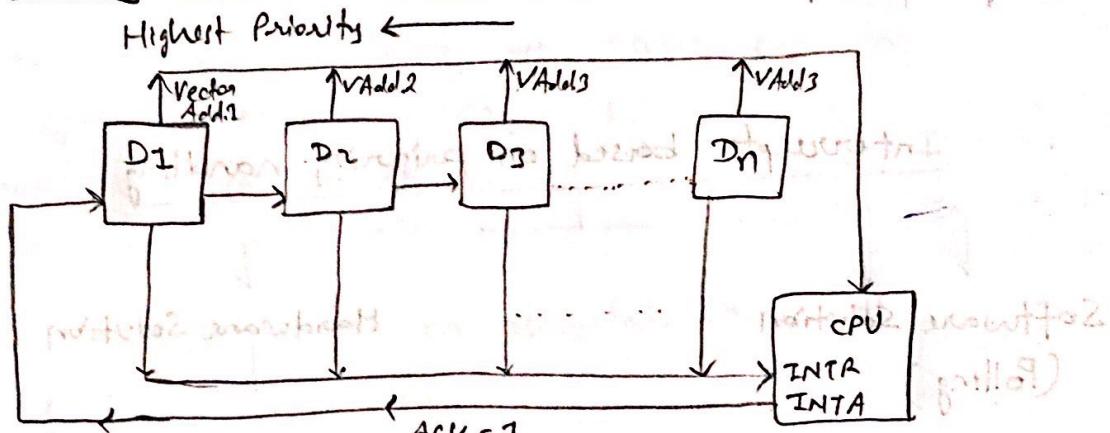
Polling



→ Device having highest priority is checked first and then devices are checked in descending order of priority.

If device P is checked to be generating the interrupt, another service program is called which works specifically for that particular device.

Daisy Chain (H/W polling)



→ ACK sent to device one, if it doesn't need it, it send to another device like wise up to the end.

→ If two device raise the interrupt, then priority is given to one which is closer to line

→ Priority is based on the physical location of device

→ Starvation is theoretically possible but not practically.

Time Required for

$$\text{data transfer} = \text{interrupt overhead} + \text{actual interrupt service time}$$

(interrupt service) [Execution of ISR]

Q-6 A device with data transfer rate 10 kB/sec is connected to CPU. Data is transferred byte-wise. Interrupt overhead be 4 μsec. What is the min. performance gain of operating device under interrupt mode over operating it under program controlled mode?

→ Programmed

In ~~Interrupt~~ mode:

$$10 \text{ kB/sec} \rightarrow 1 \text{ sec.}$$

$$10 \text{ B} \rightarrow \frac{1 \text{ B} \times 1 \text{ sec.}}{10 \times 10^3 \text{ B}} = 0.1 \text{ ms} \quad [\text{To check status}]$$

In interrupt mode = 4 μsec [Overhead only]

$$\text{Performance gain} = \frac{0.1 \times 10^{-3}}{4 \times 10^{-6}} = 25$$

Performance gain = Performance of interrupt mode

Speedup = Performance of prog. mode

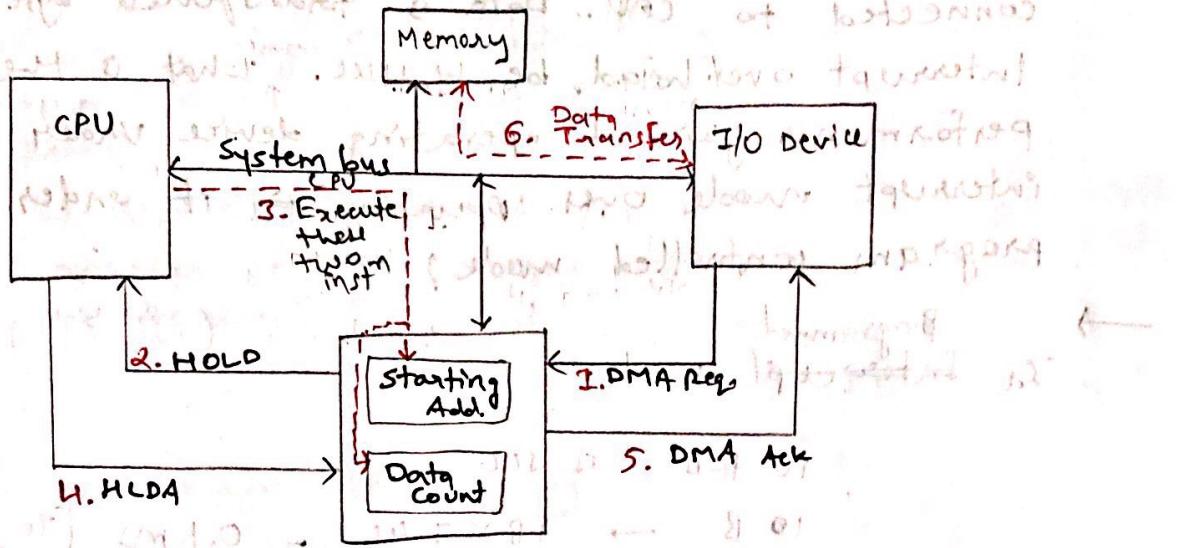
$$= \frac{1}{t_{\text{interrupt}}} = \frac{t_{\text{prog.}}}{t_{\text{interrupt}}} = \frac{0.1 \text{ ms}}{4 \mu\text{s}} = 25$$

$t_{\text{prog.}}$

DMA

It enables direct data transfer between memory and I/O device without (any) interference of CPU.

To implement DMA technique a specific hardware is used known as DMA controller (DMAC), which A



Starting Add. : The address starting from where the data transfer should be performed in memory.

Data Count : Size of data to be transferred.

In Byte \Rightarrow Byte addressable memory

In words \Rightarrow Word addressable memory

- 1. I/O device will send DMA req. to DMAC
- 2. DMAC will send HOLD signal to CPU
- 3. CPU execute two instruction to load DMA registers IOAR and Data Count
- 4. CPU wait for next DMA breakpoint (complete execution of current instn) and CPU will send HLDA

to DMA, it means microprocessor has released control of address bus, data bus to DMA.

- 5. Now DMA will send one ACK to IO device.
- 6. Now with the help of IP-OP read and memory write signal the data is transferred from IO device to memory.
After word is transfer, IOAR & Data count updated.

- During DMA transfer, CPU can perform only those operations which do not required System Bus, which means CPU will be idle mostly.
- During data transfer, if Data count not yet reached zero but I/O device is not ready to send or receive next data, then DMA controller release the system bus by deactivating bus request.
- If Data Count is zero, then DMA controller interrupts CPU and CPU will check status of DMA.

DMA Transfer Mode

Once DMA controller gets the control of the System bus, after how much time (how much data is transferred) The CPU takes the control of the buses back.

(1) Burst Mode: The entire amount of data is transferred before CPU takes back the control of buses.

(2) Block Mode:

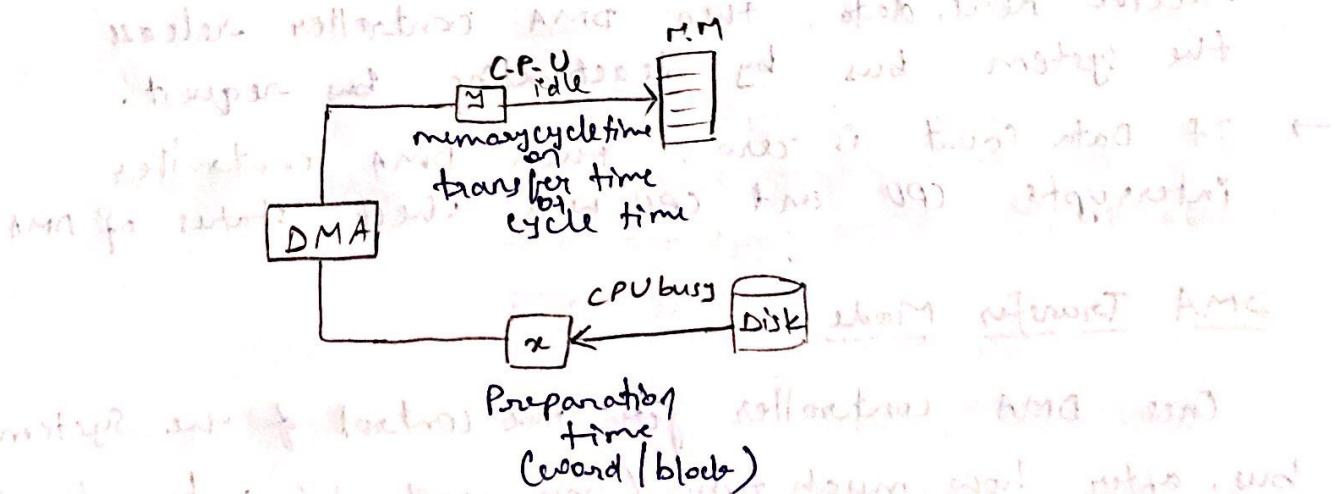
DMA returns the control of buses after a transfer of block.

Size of Block : 512 B - 1024 B
(General)

Batching long read & write operation of memory

(3) Cycle Stealing Mode:

The slow IO device takes some time to prepare the data, during that time CPU keeps the control of buses. Once the data is ready the control of the buses will given to DMA controller for one cycle in which the prepared data can be transferred to memory bus.



Time required to prepare the data = t_x

Time required to transfer the data = t_y

$$\% \text{ of CPU is blocked} = \frac{t_y}{t_x + t_y} * 100$$

$$\% \text{ of CPU is Busy} = \frac{t_x}{t_x + t_y} * 100$$

The data transfer preparation can be performed during data transfer.

$$\% \text{ of time CPU is blocked} = \frac{\text{tby}}{\text{ttx}} * 100\%,$$

; $t_{tx} \geq t_{by}$

Notes: If the interrupt occurs $\frac{t_{by}}{t_{tx} + t_{by}} * 100\%$, otherwise it needs to abort private steps and gets back to main program. interrupt of each interrupt happens at different part of code.

Transfer Type	Burst Mode	Block Transfer Mode	Cycle Stealing Mode
1. Used to Transfer	Small data size	Moderate	Large size data
2. In one time data transfer	Large	Moderate	Small
3. CPU Waiting	More	Moderate	Less
4. DMA Waiting	Less	Moderate	More

→ DMA controller is a special purpose processor which is used for data transfer b/w memory & I/O. It generates addresses & control signals like read & write for memory.

Q-2 The size of data count register of DMA controller is 16 bits. The processor needs to transfer file of 29,154 KB from disk to m.m. Mem. is byte addressed. The min. no. of times the DMA controller need to get control of system bus from processor to transfer the file from disk to m.m. is.

$$\frac{\text{Filesize}}{\text{Reg.Count}} = \frac{29154 \times 2^{10} \text{ Byte}}{2^{16}-1} = 455.5$$

↓
at count 0
we are not transferring
any byte

Q - Consider a device with 1MBPS transfer rate which is operating in cycle stealing mode of DMA. It requires 2us to transfer the data from device to memory when it is ready or prepare.

The data preparation starts only after the previous data is transferred. What is the % of time processor is blocked due to DMA. If data size is 16 Bytes.



$$1 \text{ MB} = 1 \text{ sec}$$

$$16 \text{ B} = \frac{16 \text{ B}}{1 \text{ MB}} = 16 \text{ } \cancel{\mu\text{s}}$$

$$\therefore \text{preparation time} = 16 \text{ } \mu\text{s}$$

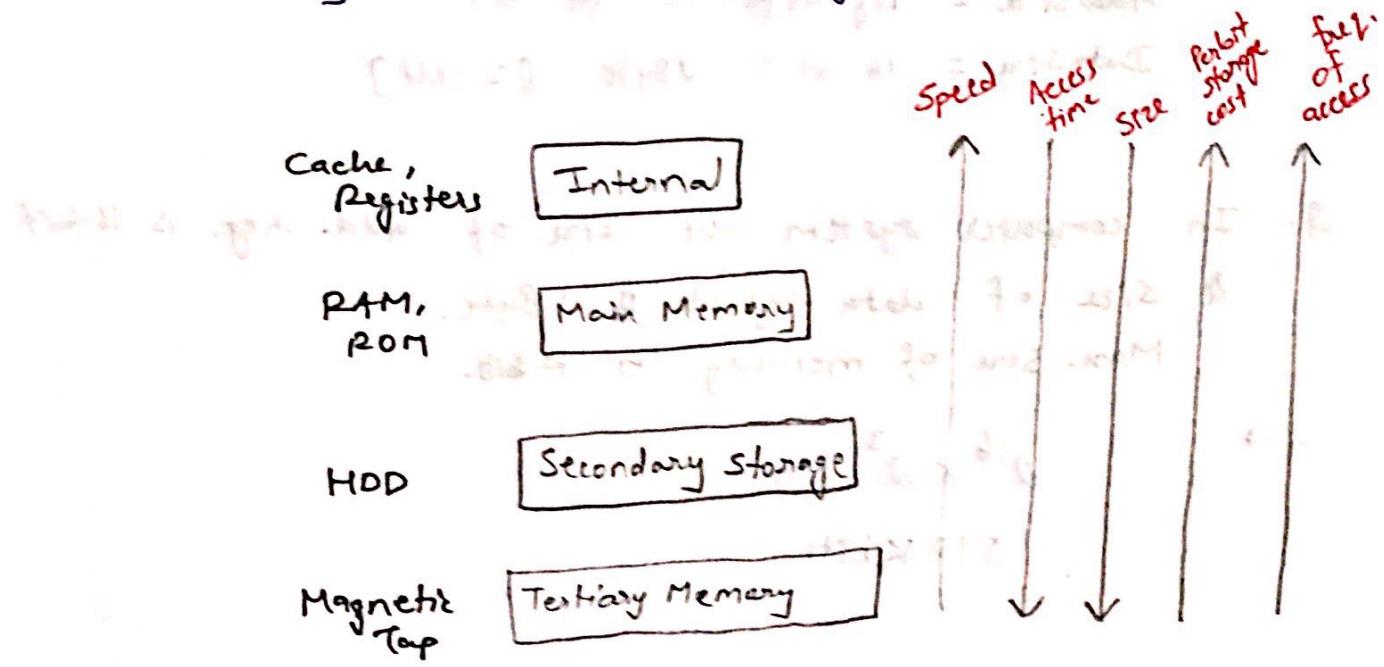
$$\% \text{ CPU block} = \frac{2}{16+2} = \frac{2}{18} = 11.11\%$$

Following are advantages of direct memory access (DMA):
 1) Transfer of data is faster than CPU.
 2) DMA can transfer data directly between memory and peripheral devices without involving CPU.
 3) DMA can transfer data in parallel with CPU.
 4) DMA can transfer data in burst mode.

* Memory Organization

Memory: A physical device to store the content in the system.

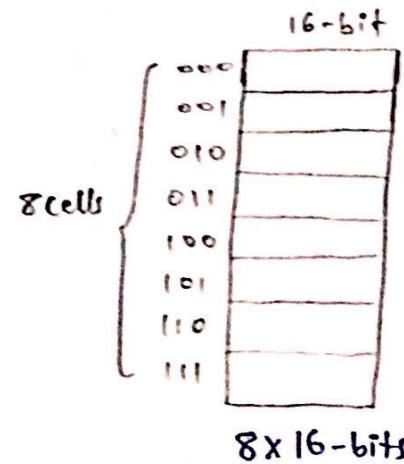
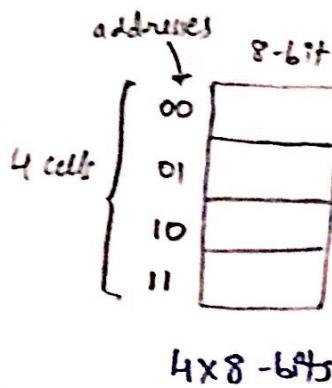
→ Instead of using single memory a memory hierarchy are used in the system.



Goal:
1. To maximize memory access speed.
2. To minimize per bit storage cost.

Memory Representation

Memory is represented by no. of cells \times 1 cell capacity.



$$\text{Memory Representation} = \# \text{cell} * 1 \text{ cell capacity}$$

Q- Memory: $128 \text{ K} \times 16$ -bit

$$\text{Add. size} = \log(28 \text{ K}) = 17\text{-bit}$$

$$\text{Data size} = 16\text{-bit} = 1 \text{ Byte} [1 \text{ cell}]$$

Q- In computer system the size of add. reg. is 16-bit & size of data register is 1 Byte.
Max. size of memory is 1 KB.



$$2^{16} \times 2^3 \text{ bits}$$

$$512 \text{ Kbit}$$

Main Memory

It is used to store current running processes (Instruction) and their data.

ROM: non-volatile memory

RAM: Volatile Memory

Types of ROM

(i) EROM: (Erasable ROM)

With the help of UV rays or laser rays the content of such ROM can be erased.

(ii) PROM: (Programmable ROM)

One time writing can be performed in such ROM and after that only read.

(iii) **EPROM**: Erasable Programmable ROM

(iv) **EEPROM**: Electrically Erasable Programmable ROM

With the help of electric current the content can be erased.

(v) **EAPROM**: Electrically Alterable Programmable ROM

With the help of electric current, content can be update or alter.

- Q- A ROM is used to store table of multiplication of two 4-bit unsigned integer. The size of ROM is,

→

A	B	$A * B$	$A + B$	Mult. Table	Add. Table
4-bits	4-bits	8 bits (max)	5 bits	$2^8 \times 2^8$ bits	$2^8 \times 5$ bits
n-bits	n-bits	2^n bits	$(n+1)$ bits	$2^n \times 2^n$ bits	$2^n \times (n+1)$ bits

$$\text{Memory} = 2^8 \times 8 \text{ bits}$$
$$= 2 \text{ K bits}$$

RAM

Types of RAM

Static (SRAM)

Built : Flip-flop

No recharge or refresh

Dynamic (DRAM)

Built : Capacitors (Content is stored in the form of electrical charges)

Electrical charges tend to discharge
Periodical refresh or recharge is required

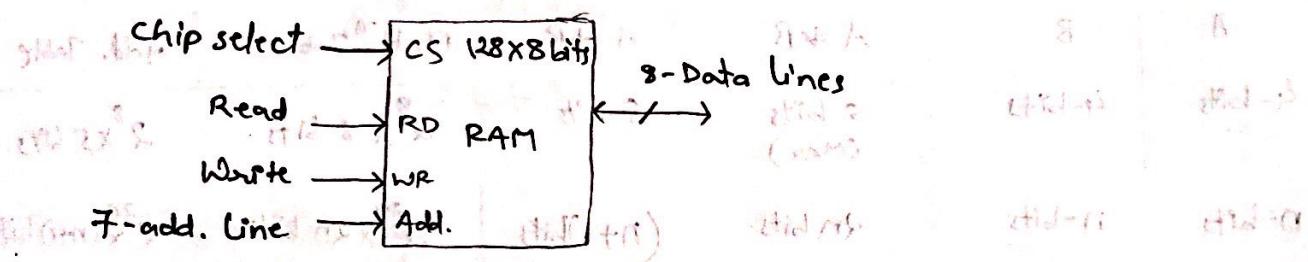
SRAM

- Faster
- Costly
- More power consumption
- Used for cache memory

DRAM

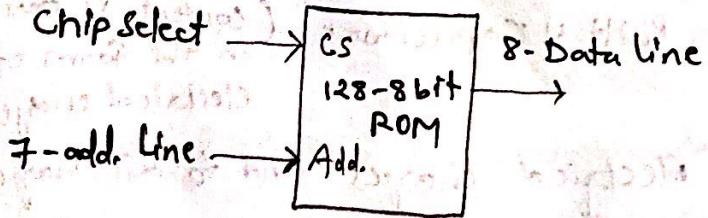
- Slower [during refresh, read/write can not be performed]
- less costly
- less power consumption
- Uses for main memory of computer systems

RAM Chip



CS	RD	WR	operation
0	X	X	No op. ⁿ
1	0	0	No op. ⁿ
1	0	1	Write op. ⁿ
1	1	X	Read op. ⁿ

Chip Select



CS	operation
0	No op. ⁿ
1	Read

$$\text{No. of pins} = 1 + 7 + 8 \\ = 16 - \text{pins}$$

- Q- A ROM chip of 1024×8 bits has 4 select line and operates from a 5 watt power supply. How many pins are needed.

→ Select line = 4

add. line = 10

data line = 8

Power supply = 1

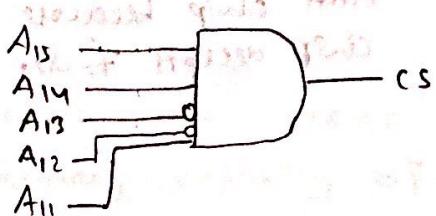
Ground = 1

Total = 24

- Q- In memory select design system

- Q- Chip select logic for certain DRAM chip in memory system design is shown below.

Assume that the memory system has 16-address line A_{15} to A_0 . what is the range of add. (In Hex.) of the memory system that can get enable by chip select (cs) signal.



- (A) C800 to CFFF
 (B) C800 to C8FF
 (C) DA00 to DFFF
 (D) IC400 to UFFF

A_{15}	A_{14}	A_{13}	A_{12}	A_{11}	$A_{10} \dots A_0$
1	1	0	0	1	0 0 . . . 0
1	1	0	1	1	0 0 . . . 1
1	1	1	1	1	0 0 . . . 1 0
1	1	0	0	1	1 1 . . . 1 1

Using Multiple Chips (M.M)

Total capacity = #chips * 1 chip capacity

Assume,

in memory system 2 chips of size 128×8 bits are used

$\therefore 128 \times 8$ - bits

7-bit address \Rightarrow RAM 0

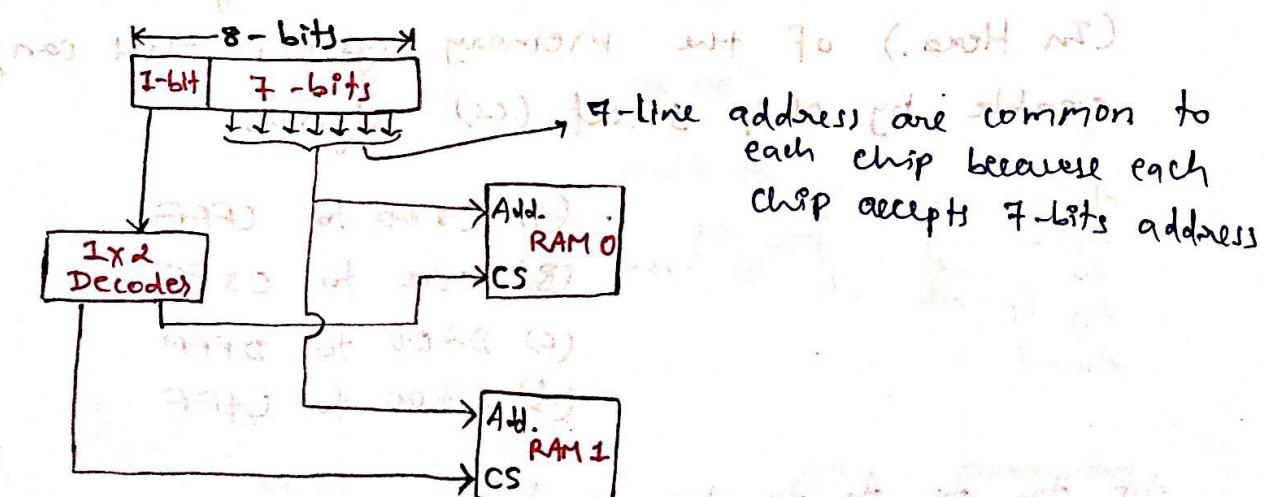
128×8 - bits

7-bit address \Rightarrow RAM 1

Total 256×8 - bits

8-bit address

		8-bit address							
		0	0	0	0	0	0	0	0
		127	0	1	1	1	1	1	1
		128	1	0	0	0	0	0	0
		256	1	1	1	1	1	1	1



7-line address are common to each chip because each chip accepts 7-bits address

- Q- (i) How many 128×8 bit RAM chips are needed to provide a memory capacity of 2048 Bytes.
- (ii) How many lines of the address bus must be used to access 2048 Bytes of memory. How many of these lines will be common to all chips.
- (iii) How many lines must be decoded for the chip select. Specify the size of decoders.

$$\rightarrow \text{Capacity} = 2048 \times 8 \text{ bits}$$

(i) RAM size = 128×8 bit

$$\# \text{RAM} = \frac{2^1 \times 2^3 \text{ bits}}{2^{10} \text{ bits}} = 16$$

(ii) Address lines = $\log 2048$

$$= 11 \text{ lines}$$

Common lines to all chips = 7 bit

(iii) Lines for decoded = $\log 16$

$$= 4 \text{ lines}$$

Size of decoder is 4×16

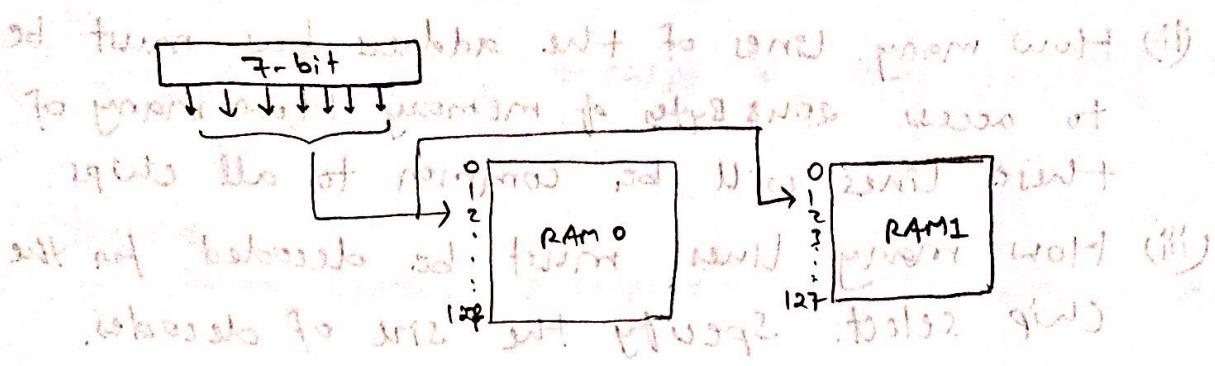
- Q- How many 128×8 RAM chips are needed to provide a memory capacity of 128×16 bits?

$$\rightarrow \frac{128 \times 16}{128 \times 8} = 2$$

Total capacity = 128×16 bits Add. Line = 7

1 RAM chip capacity = 128×8 bits Add. Line = 7

Horizontal Arrangement



NOTE: (i) If required no. of address are more, use vertical arrangement.

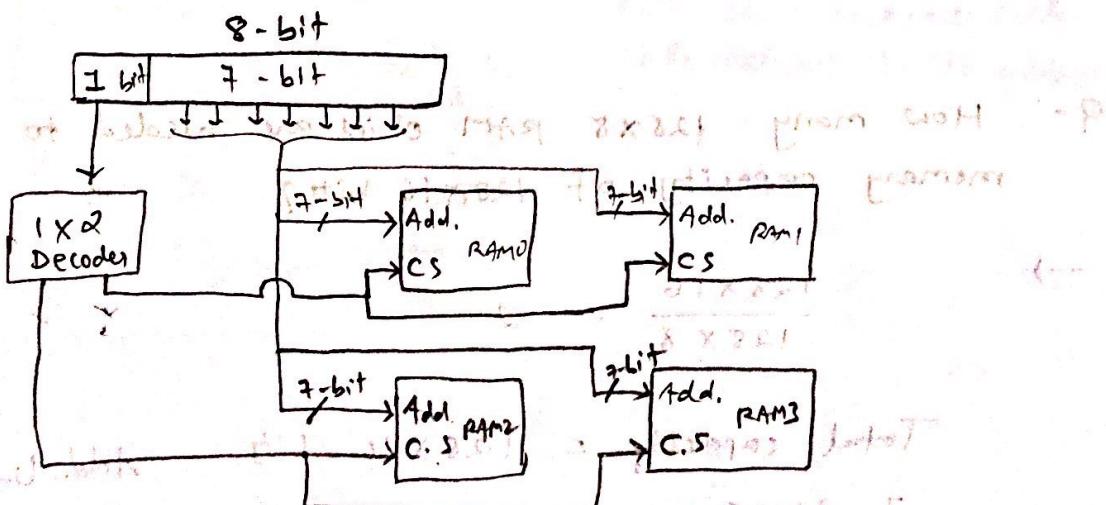
(ii) If required no. of data lines are more, use Horizontal arrangement.

(iii) If both are more, use hybrid arrangement (Vertical & Horizontal).

a. How many 128×8 RAM chips are needed to provide memory capacity of 256×16 bits?

Ans

$$\frac{256 \times 16}{128 \times 8} = 4$$



Q- Consider a M.M system that consists of 8 memory modules attached to the system bus, which is 1 word wide. When a write request is made, the bus is occupied for 100 ns. by the data, add and control signal. During the same 100 ns and for 500ns thereafter, addressed memory module executes 1 cycle accepting and storing the data. The (internal) operation of different memory modules may overlap in time, but only one request can be on the bus at any time. The max. no. of stores (of 1 word each) that can be initiated in 1 ms is ?

In 100 ns Z stored initiated after $\frac{1}{2}$ of dose.

$$1 \text{ ms} = \frac{1 \times 10^{-3}}{100 \times 10^{-9} \text{ s/cycle}} \text{ cycles} = 10^4 \text{ cycles}$$

= 10,000 stores initiated ± 100ns 3
3m

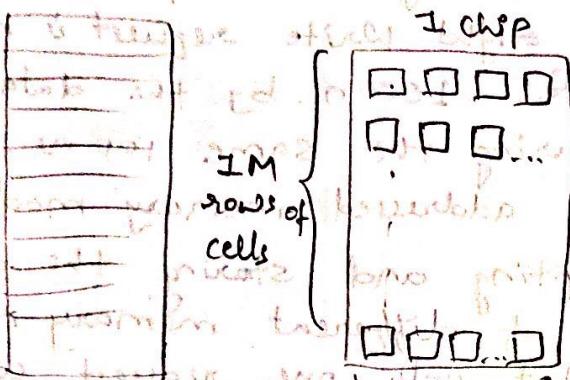
~~and due $10,000$ is now given as unit distance $\frac{7 \times 10^{-3}}{100 \times 10^{-9}}$~~ $10,000$

How many request can be successfully completed in 1 ms.

16-882

→ $\frac{10,000}{5} = 2,000$ miles from map A - P
814000 $\frac{5}{99.95}$ miles from map A - P
approximately 9600 miles from map A - P
The other 50 miles from map A - P are included in the 9600 miles from map A - P.
The distance between map A - P and map B - C is approximately 1000 miles.
The distance between map B - C and map C - D is approximately 1000 miles.
The distance between map C - D and map D - E is approximately 1000 miles.
The distance between map D - E and map E - F is approximately 1000 miles.
The distance between map E - F and map F - G is approximately 1000 miles.
The distance between map F - G and map G - H is approximately 1000 miles.
The distance between map G - H and map H - I is approximately 1000 miles.
The distance between map H - I and map I - J is approximately 1000 miles.
The distance between map I - J and map J - K is approximately 1000 miles.
The distance between map J - K and map K - L is approximately 1000 miles.
The distance between map K - L and map L - M is approximately 1000 miles.
The distance between map L - M and map M - N is approximately 1000 miles.
The distance between map M - N and map N - O is approximately 1000 miles.
The distance between map N - O and map O - P is approximately 1000 miles.
The distance between map O - P and map P - Q is approximately 1000 miles.
The distance between map P - Q and map Q - R is approximately 1000 miles.
The distance between map R - S and map S - T is approximately 1000 miles.
The distance between map T - U and map U - V is approximately 1000 miles.
The distance between map V - W and map W - X is approximately 1000 miles.
The distance between map X - Y and map Y - Z is approximately 1000 miles.

* Arrangement of cells in DRAM chip



In a DRAM chip, there are 1M rows of cells. Each row has 1K cells. So, total number of cells in one chip is $1M \times 1K$. This is represented by a 4x4 grid of cells. The grid has 4 rows and 4 columns. Each cell in the grid represents one bit of memory. The entire grid is called a 'chip'. The chip is then interconnected with other components like sense amplifiers, decoders, and drivers to form a complete memory system.

Assume DRAM chip of size 1GB, which has 1M rows of cells with 1K cells in each row.

→ In one refresh operation, 1 row of cells can be refreshed

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

if 'n' chips are there in one memory system.

$$\text{Memory system refresh time} = \text{1 chip refresh time}$$

Q - A main memory unit with capacity of 4MB is built using $1M \times 1$ bit DRAM chip. Each DRAM chip has 1K rows of cells with 1K cells in each row. The time taken for a single refresh operation is 100 ns. The time required to perform 1 refresh operation on all the cells in the memory unit is.

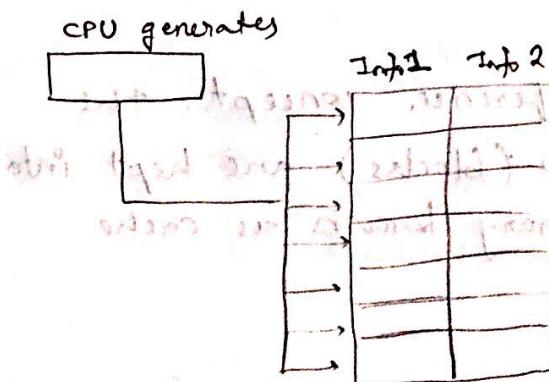
→

$$\frac{4 \times 2^{20} \times 2^3 \text{ bits}}{2^{20} \times 1} = 32 \text{ DRAM}$$

$$\begin{aligned} \text{1 chip refresh time} &= 1 \text{K} \times 100 \text{ ns} \\ &= 2^{10} \times 100 \text{ ns} \end{aligned}$$

$$\therefore \text{Memory refresh time} = \frac{\text{1 chip refresh time}}{\text{(1 chip)}} = \frac{100 \times 2^{10} \text{ ns}}{\text{(1 chip)}}$$

* Associative Memory



→ Cells do not have addresses

→ Every cell contains 2 info. 1. Info 1 = CPU knows
2. Info 2 = CPU requires

→ Searching performed based on content ⇒ Content addressable Memory

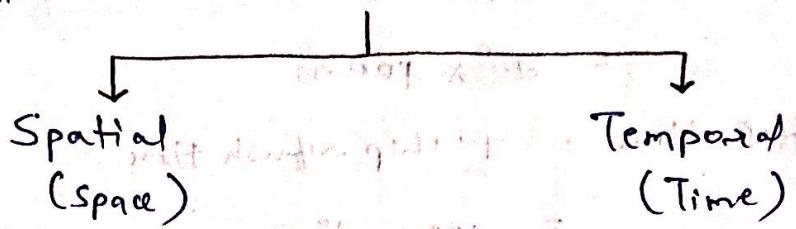
→ Content matching is done in parallel in all the cells, because each cell has one matching logic hardware.

→ Associative Memory is faster & expensive.

Locality of reference

If CPU refers the memory at specific address then the same address or the nearby address will be referred soon. This phenomenon is known as locality of reference.

Types of locality of Reference

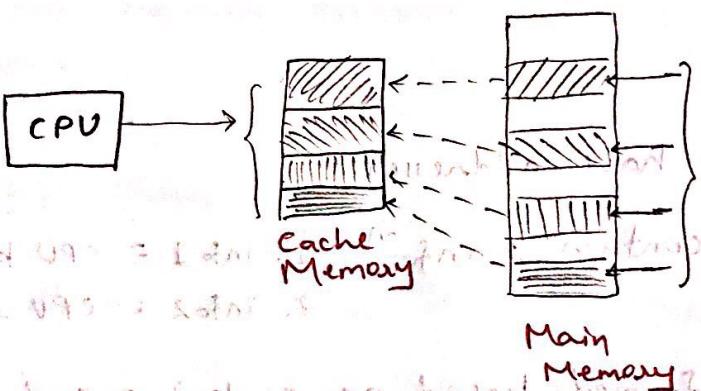


if nearby addresses
are referred soon

if same address is
referred soon.

* Cache Memory

→ Based on locality of reference concept, the current demanded localities (blocks) are kept into a smaller and faster memory known as cache memory.



Every time searching in MM for current demand block takes more time, so kept those blocks in Cache Memory.

→ Use of cache decreases the avg memory access time.

Working of Cache Memory

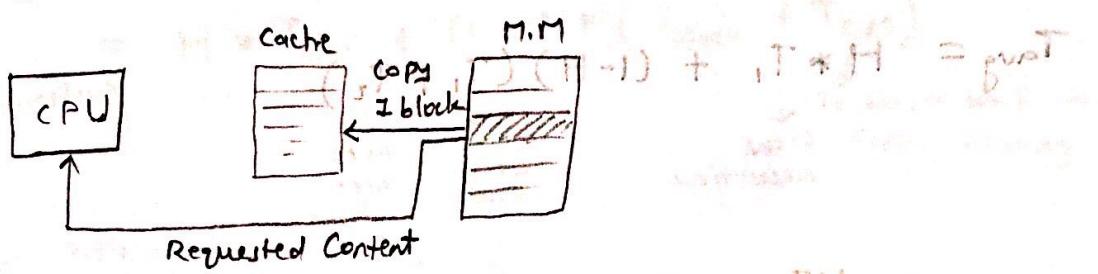
If requested content is present in cache memory then its cache hit.

If required content is not present in cache memory then its cache miss.

Performance of cache is given by cache hit ratio.

$$\text{Hit Ratio } H = \frac{\text{No. of hits}}{\text{Total references}}$$

In case of cache miss one block (which contains requested content) one block is copy from main memory to cache memory.

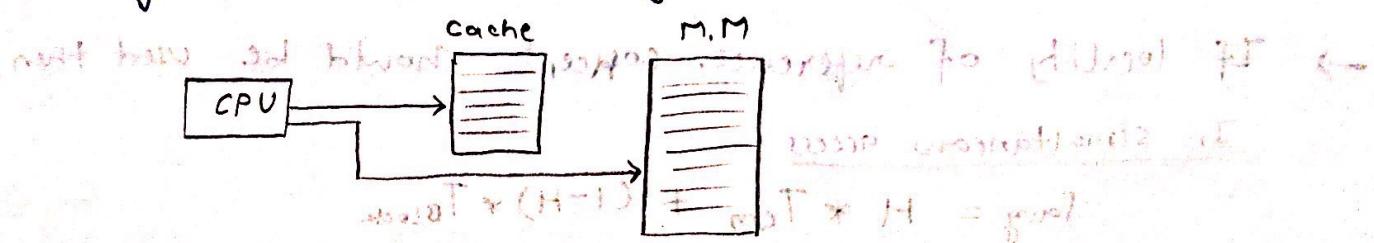


$$\text{Avg. Memory Access Time} = \text{Hit} * \text{Mem. access time when hit in cache} + (1 - \text{H}) * \text{Mem. access time when miss in cache}$$

* Types of Cache Access

1. Simultaneous Access

The requests for cache and main memory are generated simultaneously.



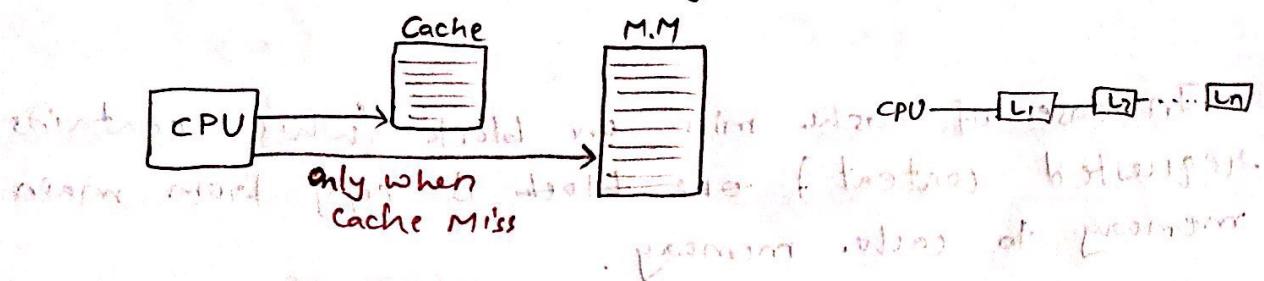
$$\text{Avg. Access Time} = H * T_1 + (1 - H) * T_2$$

↓ ↓
Cache access time M.M. access time

2. Hierarchical Access

Only the faster memory is accessed first.

If cache miss, then only then search in M.M.



$$\begin{aligned}
 T_{avg} &= H * T_1 + (1-H) (T_1 + T_2) \\
 &\quad \downarrow \quad \downarrow \quad \downarrow \\
 &\quad \text{Cache} \quad \text{Cache} \quad \text{M.M} \\
 &\quad \text{access time} \quad \text{access time} \quad \text{access time} \\
 \\
 &= HT_1 + T_1 - HT_1 + T_2 - HT_2 \\
 &= T_1 + (1-H) T_2
 \end{aligned}$$

→ If Cache memory & M.M access time not given
use generalized formula

- If cache memory & M.M access time are given
 - If any one of 2 words given → Hierarchical level
 - else, simultaneous access
- If locality of reference concept should be used then,

In simultaneous access

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

$$HT + (H-1) + (T+H)T_{block} = \frac{\text{Block size}}{\text{Block size}} * T_{m.m}$$

In Hierarchical Access

$$T_{avg} = H * T_{c.m} + (1-H) (T_{c.m} + T_{block})$$

When, block has to be transferred from M.M to cache
then,

In simultaneous access

$$T_{avg} = H * T_{c.m} + (1-H) * (T_{block} + T_{c.m})$$

To write back in
cache memory

In Hierarchical Access

$$T_{avg} = H * T_{c.m} + (1-H) (T_{c.m} + T_{block} + T_{c.m})$$

↓
 direct cache miss ↓
 Block access in M.M ↓
 write back in cache memory

Q- In two level mem. hier. if the top level has an access time of 8ns. & bottom level has an access time of 60ns. What is the hit rate on top level required to give an avg. access time of 10ns.

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

$$\rightarrow 10 \text{ ns} = H * 8 \text{ ns} + (1-H) (8 \text{ ns} + 60 \text{ ns})$$

$$10 \text{ ns} = 8H \text{ ns} + 68 \text{ ns} - 68H$$

$$\frac{-58}{-60} = H$$

$$\therefore \frac{58}{60} = H$$

$$H = 96.6\%$$

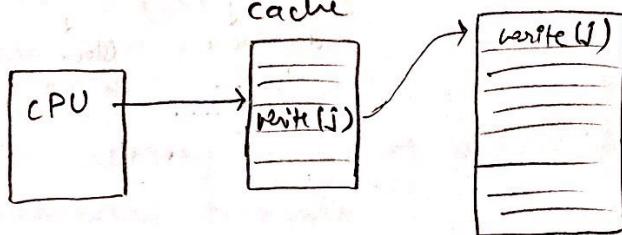
Q - Uncached mem. access takes 7 times longer than accessed to cache. If cache has hit ratio 0.9. The ratio of cached mem. acc. time to uncached mem. access time is?

\rightarrow at first need to access and already hold in cache
 Cache = x
 Uncached = $7x$

$$\frac{T_{Cache}}{T_{Uncache}} = \frac{0.9 \times x + 0.1(7x)}{7x} = \frac{0.9 + 0.7}{7} = 0.228$$

(assuming 1 miss per 10 access)

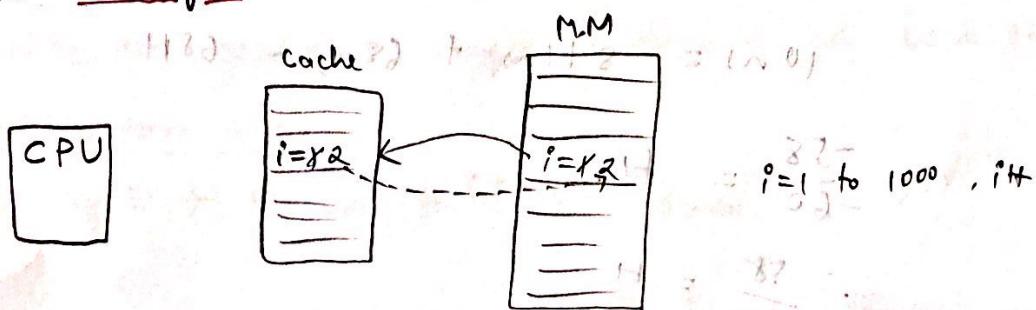
* Cache Write



last bus got set to main memory level out at -P
 no bus bottlenecked to 2.28 p worst access no

If C.P.U performs a write operation on cache contents (duplicate copy) then the respective main memory content (original copy) should also be updated.

I. Write Through



As soon as the content of cache is updated, the respective main memory content is also updated. Which means for the write operation CPU generates the request for both cache and main memory.

$$\text{Time writethrough} = \max(T_{\text{cm}}, T_{\text{m.m}})$$

\downarrow
Time for cache Time for m.m

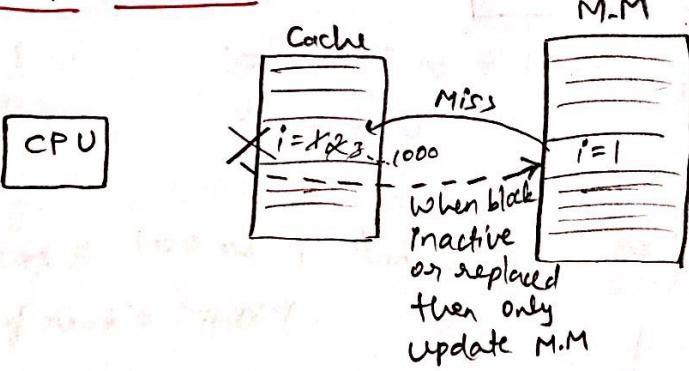
Adv.

Main memory always contains the updated content.

Dis adv.

Time consuming as every time when we update in cache we have to access M.M also to update.

II Write Back



Main memory is updated only when the cache block is inactive or to be replaced.

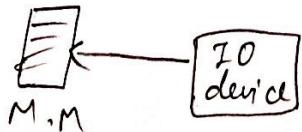
Adv.

Time saving

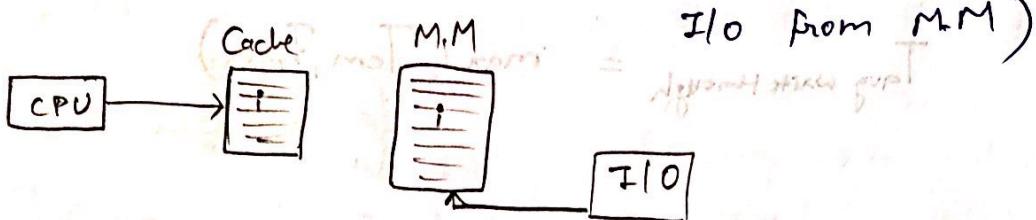
Dis Adv.

For the same content, cache & M.M may contain different value.

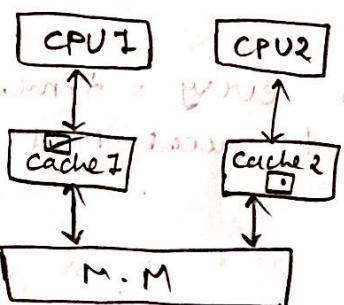
Ex. if I/O device wants to access memory content then it may read inconsistent data which is called **Memory coherence**.



Memory Coherence: If two different devices (CPU & I/O) trying to access the same content from two different memories (CPU from cache & I/O from M.M.)



Cache Coherence: Whether a specific program will



Multiprocessor System

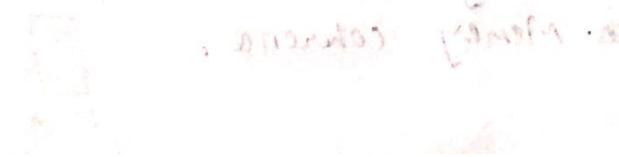
about section II



For TMR (Triple Modular Redundancy), the redundant modules must be synchronized to prevent race conditions. This can be achieved by using a distributed consensus algorithm like Paxos or Raft.

mission you can offload a task to another node with less network traffic.

With today's growth of data, clients are getting more and more data to process from the cloud.



Write through (Tavg overall): Total time taken for write operation
 Total time = time of read + time of write
 $T_{avg\ write\ through} = \max(T_{cm}, T_{m.m})$

Tavg read = Tavg formula using hit ratio etc.
 Tavg read = $\frac{t_{cm} + t_{m.m}}{hit\ ratio + miss\ ratio}$
 Tavg read = $t_{cm} * \text{hit ratio} + t_{m.m} * \text{miss ratio}$

$T_{avg\ overall\ read} = f_r * T_{avg\ read} + f_w * T_{avg\ write}$ where
 (read/write)

$$Ex. \quad t_{avg\ read} = 20 \text{ ns}$$

$$t_{avg\ write} = 50 \text{ ns}$$

Total 100 mem. reference

$$\text{read} = 80$$

$$\text{write} = 20$$

$$\begin{aligned} \therefore T_{avg\ overall} &= 0.8(20 \text{ ns}) + 0.2(50 \text{ ns}) \\ &= 16 \text{ ns} + 10 \text{ ns} \\ &= 26 \text{ ns} \end{aligned}$$

$$Q- \quad t_{cm} = 100 \text{ ns}, \quad t_{m.m} = 1000 \text{ ns}, \quad H_{read} = 80\%$$

$$\% \text{ of read} = 70\%$$

Write through policy employed

$$T_{avg\ read} = ?, \quad T_{avg\ write} = ?, \quad T_{avg\ overall} = ?$$

$$\rightarrow T_{avg\ write} = 1000 \text{ ns}$$

$$\begin{aligned} T_{avg\ read} &= 0.80(100 \text{ ns}) + 0.20(1000 \text{ ns}) \\ &= 80 + 200 \\ &= 280 \text{ ns} \end{aligned}$$

$$\begin{aligned} T_{avg\ overall} &= 0.7(280) + 0.3(1000) \\ &= 196 + 3000 \\ &= 496 \text{ ns} \end{aligned}$$

Q The memory access time is 1ns for a read opⁿ with hit in cache, 5ns for read opⁿ with miss in cache, 2 ns for write opⁿ with a hit in cache & 10 ns for write opⁿ with a miss in cache. Execution of sequence of instⁿ involves 100 IF opⁿ, 60 mem. operand read opⁿ and 40 memory operand write opⁿ. The cache hit ratio is 0.9. The average memory access time (in ns) in executing the sequence of instⁿ is?

→

$$\text{Cache Hit} = 0.9 \\ \text{Cache Miss} = 0.1$$

$$\text{Read Cache Hit time} = 1\text{ns}$$

$$\text{Read Cache Miss} = 5\text{ns}$$

$$\text{Write Cache Hit} = 2\text{ns}$$

$$\text{Write Cache Miss} = 10\text{ns}$$

$$\text{Read Op}^n = 4(60)(60 + 100) 2.0 = 11200 \text{ ns}$$

$$\text{Write Op}^n = 40$$

$$\text{Total} = 200 \text{ ns}$$

$$\text{Tavg. Read} = 0.9(1\text{ns}) + 0.1(5\text{ns})$$

$$= 0.9 \text{ ns} + 0.5 \text{ ns} = 1.4 \text{ ns}$$

$$= 1.4 \text{ ns}$$

$$\text{Tavg. Write} = 0.9(2\text{ns}) + 0.1(10\text{ns})$$

$$= 1.8 \text{ ns} + 1 \text{ ns}$$

$$= 2.8 \text{ ns}$$

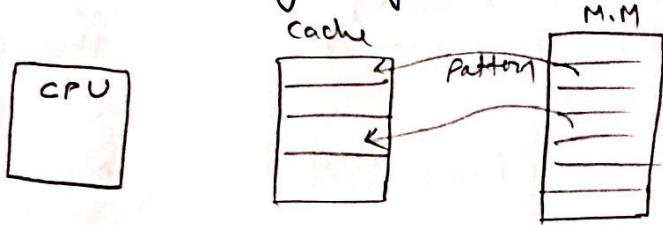
$$\text{Tavg.} = \frac{160}{200}(1.4) + \frac{40}{200}(2.8)$$

$$= 1.12 + 0.56$$

$$= 1.68 \text{ ns}$$

* Cache Mapping

CPU always generates Main memory address.



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

With the help of mapping, cache memory is searched using M.M access.

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

Types

→ Direct Mapping

→ Set-Associative Mapping

→ Fully Associative Mapping

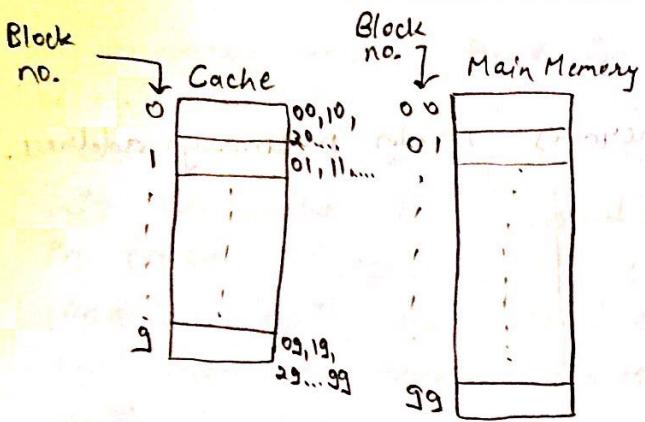
Direct Mapping

Assume,

Cache memory contains 10 blocks (0-9)

Main memory contains 100 blocks (00-99)

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



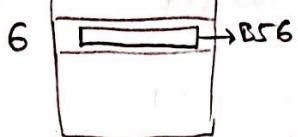
Cache Memory block no. = $k \bmod n = i$

for mapping with Main memory block no.

Ans. $k \bmod n = i$ \Rightarrow $i = \text{no. of cache line}$

$i = \text{cache memory line no.}$

CPU Reg. (M.M block no.)	C.M block no. for a ref.	Ans.
56	$56 \% 10 = 6$	Goto block no. 6 in Cache, & check. It's miss in cache. Bring Block no. 56 Content in cache at block no. 6
92	$92 \% 10 = 2$	Its miss in cache, Bring B92 content from M.M to Cache at 2.
26	$26 \% 10 = 6$	Goto block no. 6 in Cache & check. Block no. 56 is available in Cache. Hence, B26 is miss, Replace B56 by B26

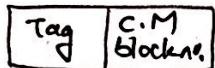


→ block Cache

What if memory
block no. 6 is present?

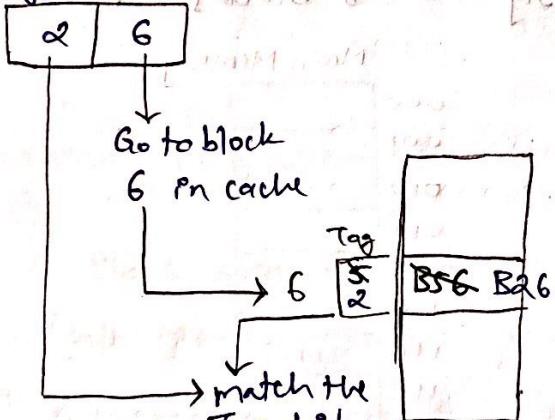
Not solution is to add
tag bits to all blocks.

M.M block no.



For request of M.M block 26

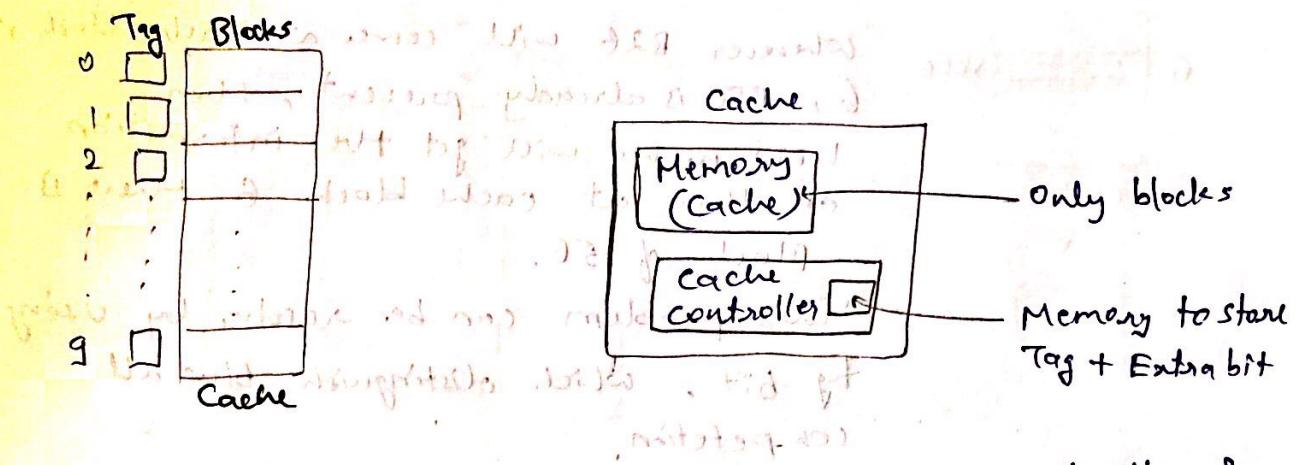
(a) If M.M block no. 26 is present



If Not matched

Cache Miss

Bring block 26 from M.M to cache in block 6 in cache and update tag bit by 2



A tag information is stored in cache controller for each block in cache.

Assume, cache memory = 14 blocks (000-111).

Main Memory = 8 blocks (000-111)

	Cache	Mapping
00		0, 4
01	I	B5
10		1, 5
11		2, 6
		3, 7

	Main Memory
000	
001	
010	
011	
100	
101	
110	
111	

CPU Request block no. 5 (101) of Main Memory

→ Mapping $5 \% . 4 = 1$

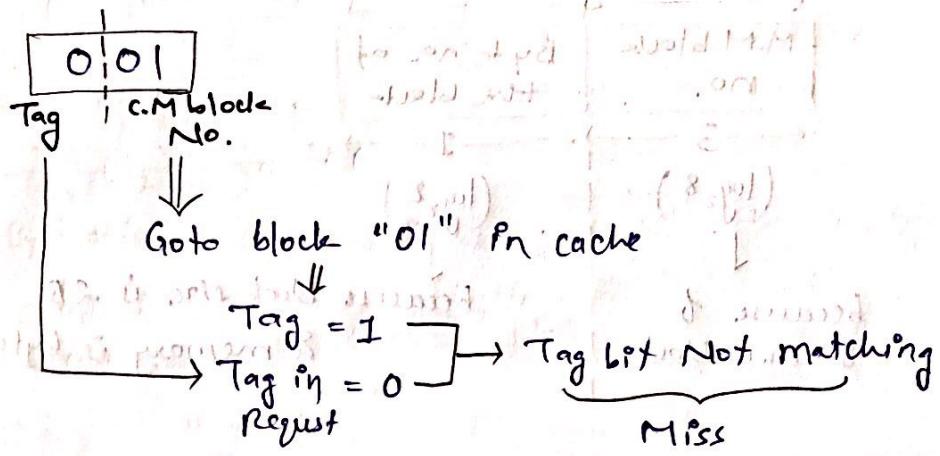
M.M. block No.

101

Tag of C.M. block No.

Goto Block "01" in cache & its a miss,
Bring B5 (101) in cache at block "01"
with tag "I".

CPU request for block no. 1 (001)₂ of Main Memory



Bring block "001" from main memory to cache at block no. "01" with tag "0" (by replacing block "101")

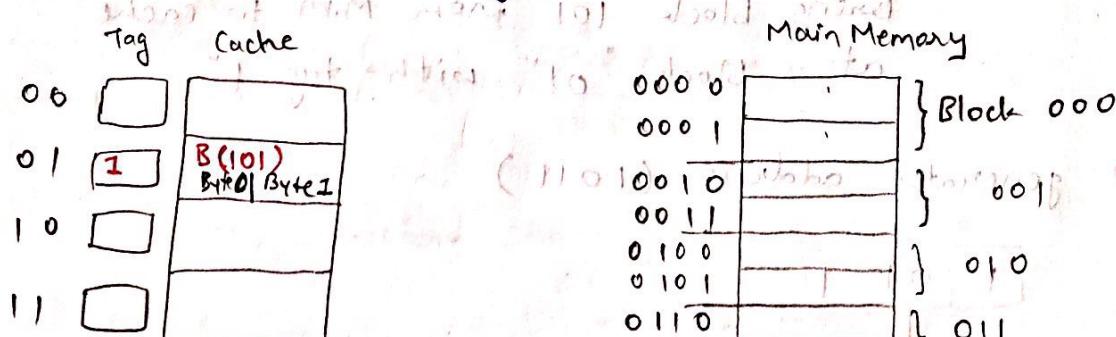
→ Comparator is used to compare tag bit.

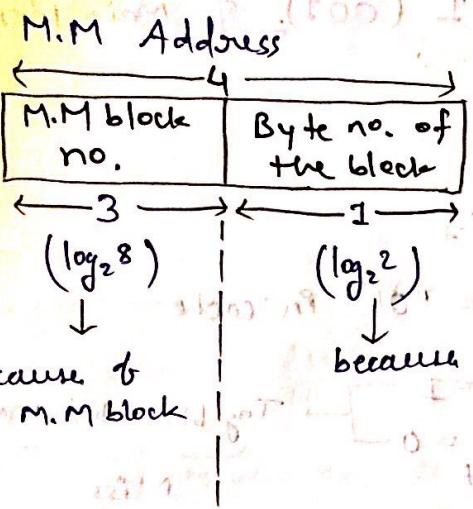
Assume, Block size = 2B

$$\text{Main Memory Size} = 8 * 2B = 16B$$

Main Memory is Byte Addressable

Main Memory address = 4-bit





Tag	C.M block No.	Byte no. of block
-----	---------------	-------------------

1st part of address = 1

→ CPU generates address (1010)

M.M Block no.
1 0 1 | 0

1 0 1 | 0

↓
Goto block "01" in cache

↓
Miss

↓
Bring block "01" from M.M to cache
at block "01" with tag "1"

→ CPU generates address (1011)

1 0 1 | 1

↓
Goto block no."01" in cache

↓
Tag true = 1] → Matching

Hit

↓
Send byte 1 of block
for access to CPU.

M.M address	
M.M block no.	Byte offset

Index bit		Byte offset
Tag	C.M Block no.	Byte offset
		$\log_2(\text{Cache size})$

$$\text{No. of bits for Byte offset} = \log_2(B.S) \cdot \text{bits}$$

$$\text{No. of bits for C.M block no.} = \log_2(\#C.B) \cdot \text{bits}$$

$$\text{Tag Memory size} = \text{no. of cache block} * (\text{Tag-size} + \text{Extra bits})$$

or

size of Tag Directory

or

size of meta-data

$$\text{Ex: Cache size} = 64 \text{ kB}$$

$$\text{Block size} = 16 \text{ B}$$

$$\text{M.M add.} = 32\text{-bit} \quad (\text{Byte add.})$$

Direct Mapping

$$\text{Tag size} = ? \quad \text{Tag directory size} = ?$$

Let us make the presented

→	32
using	16 12 4

$$\# \text{Cache line} = \frac{64 \text{ kB}}{16 \text{ B}} = 412$$

$$\text{Tag directory size} = 412 * 16 \text{ bit}$$

$$16 \text{ bit} \times 412 \text{ lines} = 4 \text{ kB} * 16 \text{ B}$$

$$= 8 \text{ kB}$$

16 bit = 2⁴ bits (4 bytes) in each of block, i.e.

each block contains 2⁴ = 16 blocks of 16 bit each

so 16 blocks = 16 * 16 = 256

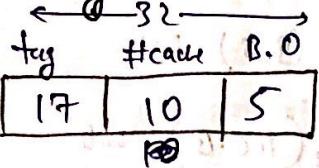
16 blocks = 16 * 16 = 256

16 blocks = 16 * 16 = 256

2905

- Q- Consider a Direct Mapped cache of size 32 kB with block size of 32 Bytes. The CPU generates 32 bit address. The no. of bits needed for cache indexing and no. of tag bits are.

→

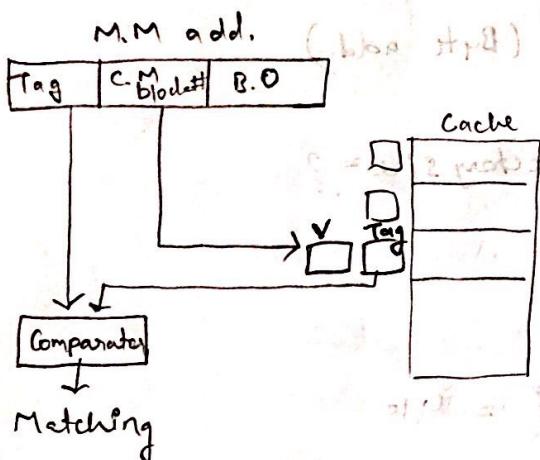


$$\# \text{ cache line} = \frac{32 \text{ kB}}{32 \text{ B}} = 1 \text{ K}$$

$$\text{Index bits} = 20$$

$$\text{tag bits} = 17$$

Cache Initialization

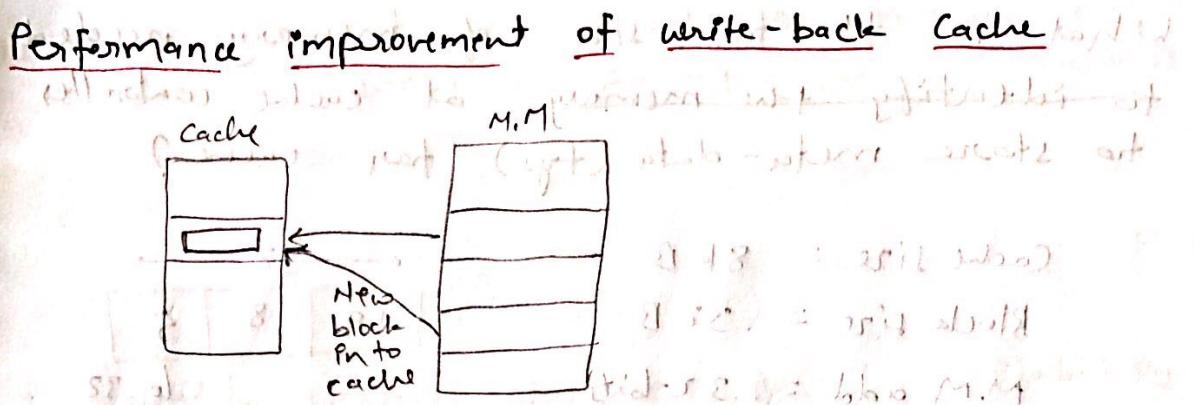


When system turn off, memory get flushed & some garbage value is available.

Whenever we turn on the system the CPU generate first request & compare the tag bit in cache if somehow they match then CPU will access unwanted data from cache.

- To solve this problem there is one valid bit in cache which gives the information whether tag field & block in cache is valid or not. So initial Valid bit is "0" when CPU generates request 1st time.

Valid bit 'V' → 0 invalid tag & cache block
→ 1 valid tag & cache block

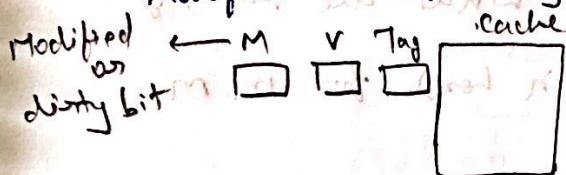


During read operation in cache we don't required any updation in M.M. So whenever conflict occurs in cache to replace blocks, we can just replace without updating in M.M.

During writer operation in cache, if any conflict occurs in cache to replace blocks, then first we have to update in M.M. before and then replace the block in cache.

So, for read operation we can save some amount of time & hence performance can be improved.

To improve performance we need to keep the information for read or write. So we required Modified bit or dirty bit in cache.



Q-6 An 8KB direct mapped write-back cache is organized as multiple block, each of size 32-bytes. The processor generates 32-bit add. The cache controller maintains the tag information for each cache block comprising of following:

- 1 valid bit
- 1 Modified bit
- As many bits as the minimum needed to identify the memory block mapped in cache.

What is the total size of memory needed to identify the memory at cache controller to store meta-data (tags) for cache?

$$\rightarrow \text{Cache size} = 8 \text{ kB}$$

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

$$\text{M.M add} = 32\text{-bit}$$

19	8	5
----	---	---

$$\log_2 32$$

$$\text{Number of cache lines} = \frac{8 \text{ kB}}{32 \text{ B}} = 2^8$$

Valid bit = 1 bit
Mod bit = 1 bit

$$\text{Size of Meta data} = 2^8 \times (19 + 1 + 1)$$

Size of Meta data = $2^8 \times 21$ bits

Total size of meta data = 5376 bits

Q-19 Consider machine with byte add. main memory of 2^{20} Bytes, B.size = 16 B & Direct Mapped cache of having 2^{12} cache lines. Let add. of two consecutive bytes in M.M be $(E201F)_{16}$ & $(E2020)_{16}$. What are the tag & cache line address (in hex) for M.M address $(E201F)_{16}$

$$\rightarrow \quad \begin{array}{c} \text{---} \\ \text{20} \\ \hline \boxed{4} \quad 12 \quad 4 \end{array} \quad \checkmark \text{(a) } E, 201$$

$$\text{(b) } F, 201$$

$$\text{(c) } E, E20$$

$$\text{(d) } 2, 01F$$

$E201F$

4	12	4
E	201	F

Cache B.Offset

Tag
Line add.

Q-15 A 64 word cache & 256 word M.M are partitioned into 16 word blocks. The tag information is

Block	Tag
0	10
1	10
2	00
3	01

Identify correct statement w.r.t availability of M.M words in cache

- (a) words 50 & 132 are available in cache
- (b) words 150 & 132 are available
- (c) words 178 & 150 are available
- (d) words 35 & 50 are available

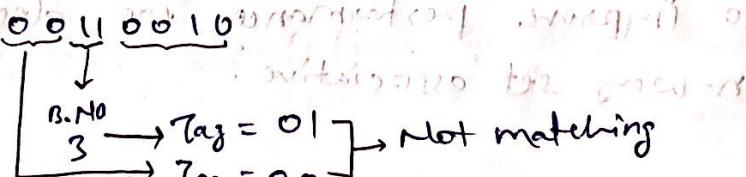


2	2	4
Tag	C.M block no.	word offset

$$\# \text{cache block} = \frac{64}{16} = 4$$

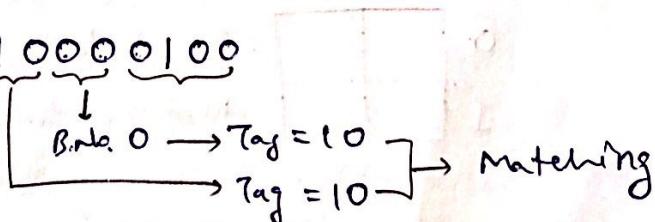
Word address : $\log_2 4 + \log_2 16$ → 8 bits
Word offset : 4 bits
C.M block no. : 2 bits
Tag : 2 bits

Word no. 50 : 00110010 (Binary)



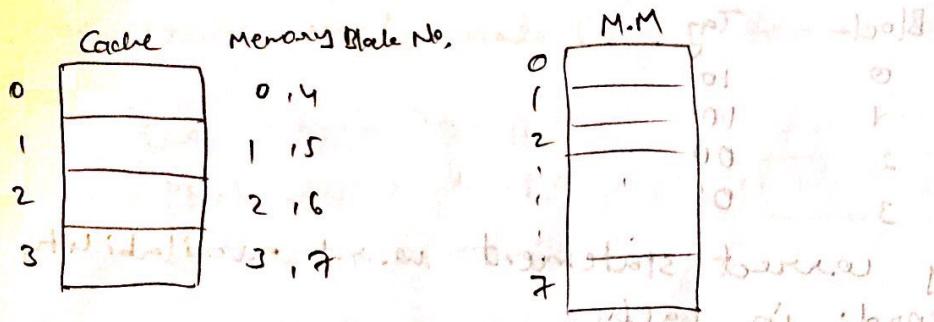
So, 50 is not in cache

word No. 132 : 10000100



So, 132 is present

* Set - associative Mapping



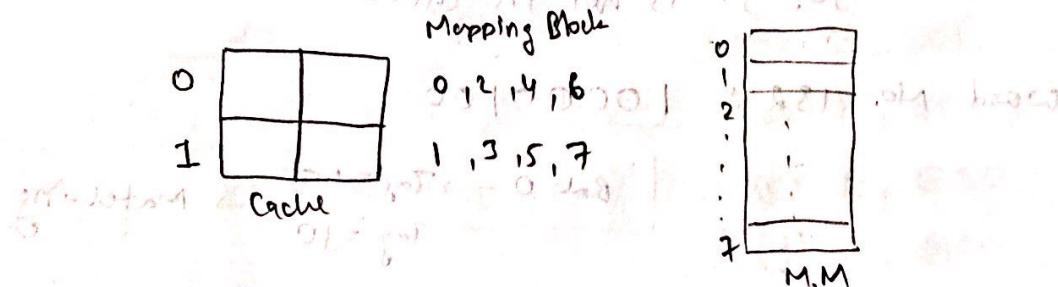
CPU Req. M.M block No. Maps Hit or Miss Replacement

CPU Req. M.M block No.	Maps	Hit or Miss	Replacement
1	→ I	Miss	Brkly B2
5	→ I	Miss	Brkly B5
1	→ I	Miss	Brkly B1
5	→ I	Miss	B5
1	→ I	Miss	B2

Here in direct mapping if request is generated by CPU as above then every time it leads to miss which consumes more time in replacement.

To improve performance we design cache in n-way set associative.

In 2-way set associative Mapping



M.M block req : 1, 5, 1, 5, 1, 5
 ↓ Miss Miss All Hit

In set associative mapping each cache line can be hold multiple blocks at a time.

$$\frac{\text{No. of index}}{\text{No. of sets}} = \frac{\# \text{cache block}}{k-\text{way}}$$

k = no. of main memory block are possible in each cache

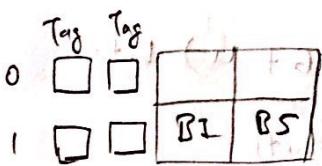
(so 2 per cache line, not others)

Main Memory Address

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

Mapping function = M.M block no. % no. of set in cache

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



For every block we required tag.

Tag directory size = #cache block * (Tag + Extra bits)

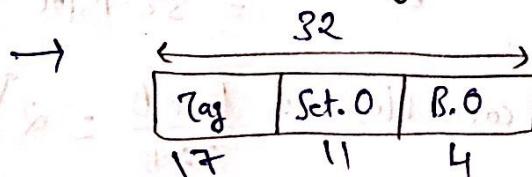
Ex: If M.M address = 32-bits

Cache size = 128 kB

block size = 16 Bytes

4-way set associative cache

Size of tag directory.



\log_2^{17} \log_2^{11} \log_2^{16}

$$\# \text{block in cache} = \frac{128 \text{ kB}}{16 \text{ B}} = 2^{13}$$

$$\text{No. of set} = \frac{2^{13}}{2^2} = 2^{11}$$

: Size of tag = $2^{13} * 17$

directory

$$= 17 \times 8 \text{ kB} \\ = 136 \text{ kB}$$

Q-7 Cache size = 256 kB, 4-way set associative.
 Block size = 32 Byte. What is block offset?
 M.M. add = 32-bit
 Tag directory entry contains in addition to tag,
 2 valid, 1 Mod., 1 replacement bit.
 #bits in tag field?
 Cache tag memory size?

\rightarrow

32		
16	11	5
\log_{32}		

#cache line = $\frac{2^2 \cdot 2^8 \cdot 10}{2^5} B$
 $= 2^{13}$
 No. of sets = $\frac{2^{13}}{2^{2+1}} = 2^9$

tag field = 16 bits
 tag memory size = $2^{13} * (16 + 4)$ bits
 $= 20 * 8 \text{ kbit}$
 $= 160 \text{ kbit}$

Q-8 A 4-way set-associative cache memory unit with capacity of 16 kB is built using block size of 8 words. The word length is 32-bits. The size of physical address space is 4 GB. The no. of bits for the TAG field.

\rightarrow

block size = $8 \times 32 \text{ bits} = \frac{2^3 \text{ bits}}{2^3} = 2^5 \text{ Byte}$
 $= 32 \text{ Byte}$

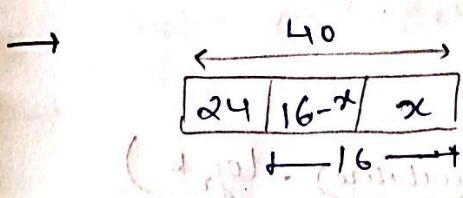
$\log 4 \text{ GB} = 32$

20	7	5
tag	set offset	block offset

#cache block = $\frac{16 \text{ kB}}{32 \text{ B}} = 2^9$
~~#set~~ #set = $\frac{2^9}{2^2} = 2^7$

TAG field = 20 bits

Q-22 The width of physical address on machine is 40 bits. The width of tag field in 512 kB 8-way set associative cache is — bits.



$$\# \text{cache} = 512 \text{ kB}$$

$$\# \text{set} = \frac{512 \text{ kB}}{8} = 2^{16}$$

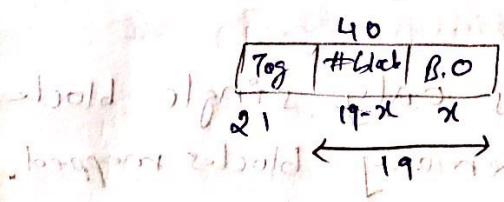
$$\frac{2^{19}}{2^3}$$

Let B. size = $2^{12} x$

$$\# \text{cache line} = 2^{19-x}$$

$$\# \text{set} = \frac{2^{19-x}}{2^3} = 2^{16-x}$$

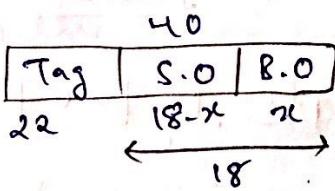
Direct Mapping



$$\text{Block size} = 2^x$$

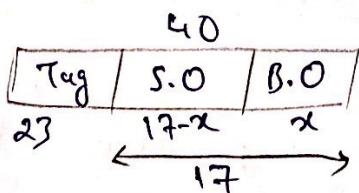
$$\# \text{blocks} = \frac{2^{19}}{2^x} = 2^{19-x}$$

2-way set associative



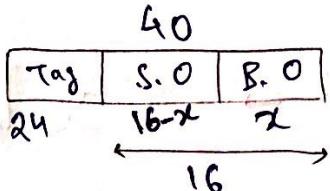
$$\# \text{sets} = \frac{2^{19-x}}{2^1} = 2^{18-x}$$

4-way set associative



$$\# \text{sets} = \frac{2^{19-x}}{2^2} = 2^{17-x}$$

8-way set associative



$$\# \text{sets} = \frac{2^{19-x}}{2^3} = 2^{16-x}$$

Hot 0

For k-Way Set Associative

Tag	Set off.	Block off.
-----	----------	------------

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

$$\therefore \text{Tag} = \text{M.M add size} - (\log_2(\text{cache size}) + \log_2 k)$$

$$= \text{M.M add size} - \log_2(\text{cache size}) + \log_2 k$$

* Fully Associative Memory

In this mapping there is only single block on which all the main memory blocks mapped.

Direct

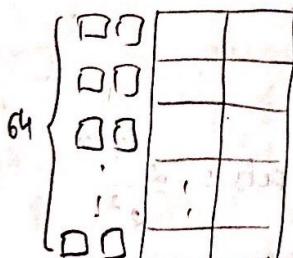
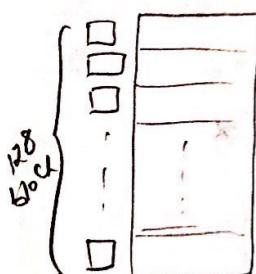
2-way

Fully Associative

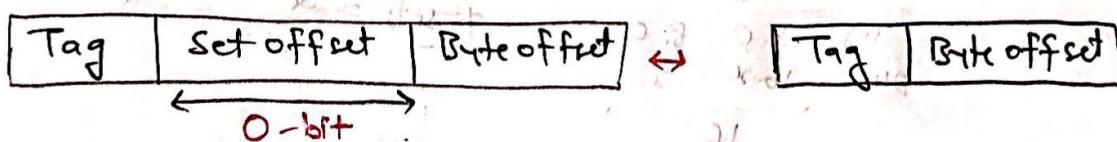
Tag	C.M block	B.O
7		

Tag	S.O	B.O
6		

Tag	S.O	B.O
0		



→ There is no index bit (set offset), in fully associative



No. of set = 1

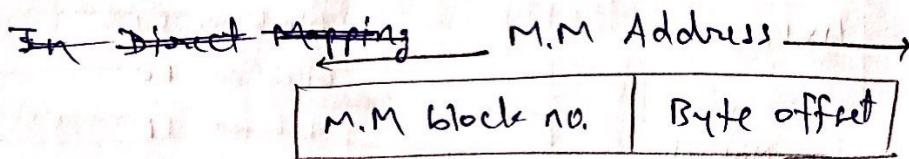
∴ Set offset is 0, $2^0 = 1$

Q. A fully associative cache contains 2^{10} blocks. The main memory has 2^{16} blocks of the same size as the cache. The size of memory required at cache controller to store the cache directory?



$$2^{10} \text{ cache block} = 2^x = \frac{\text{Cache size}}{2^{10}}$$

$$2^{16} \text{ M.M block} = 2^x = \frac{\text{M.M size}}{2^{16}}$$



In Fully Associative

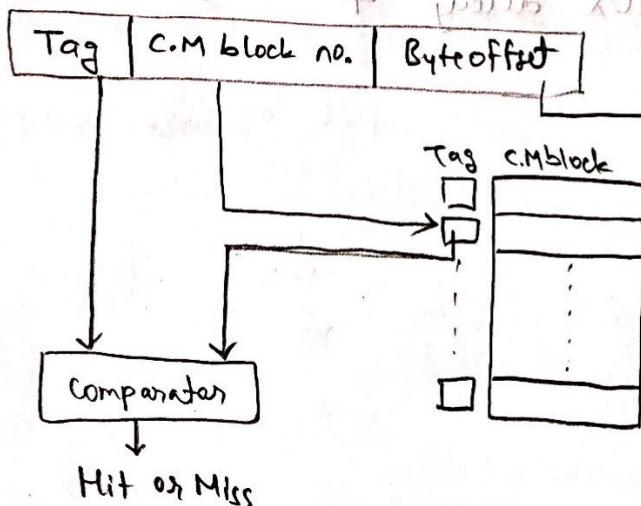
Tag	Byte offset
16	

$$\therefore \text{Size of tag directory} = 2^{10} * 16 = 16 \text{ k bits}$$

→ In fully associative memory, Tag field is nothing but main memory block no.

* Hardware Implementation of Mapping

I. Direct Mapping



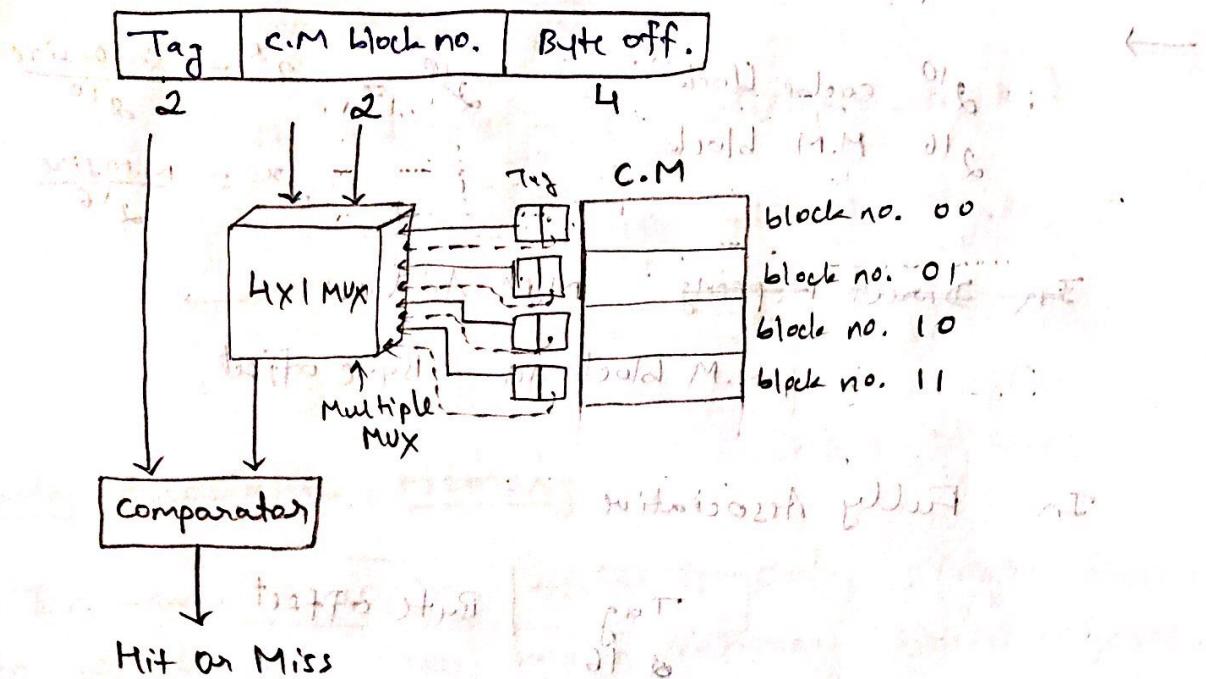
→ Which byte no. CPU will access from cache block

Ex: No. of blocks in Cache Mem. = 4

Main. Mem. address = 8 - bits

Block size = 16 Byte

Syndrome along with write of 16 locations into



No. of Mux = $31 + 0$ → 31×2 → 62 bits

MUX can forward only 2-bit of tag, so we

required multiple MUX with 2 bits output

No. of Mux = No. of tag bits

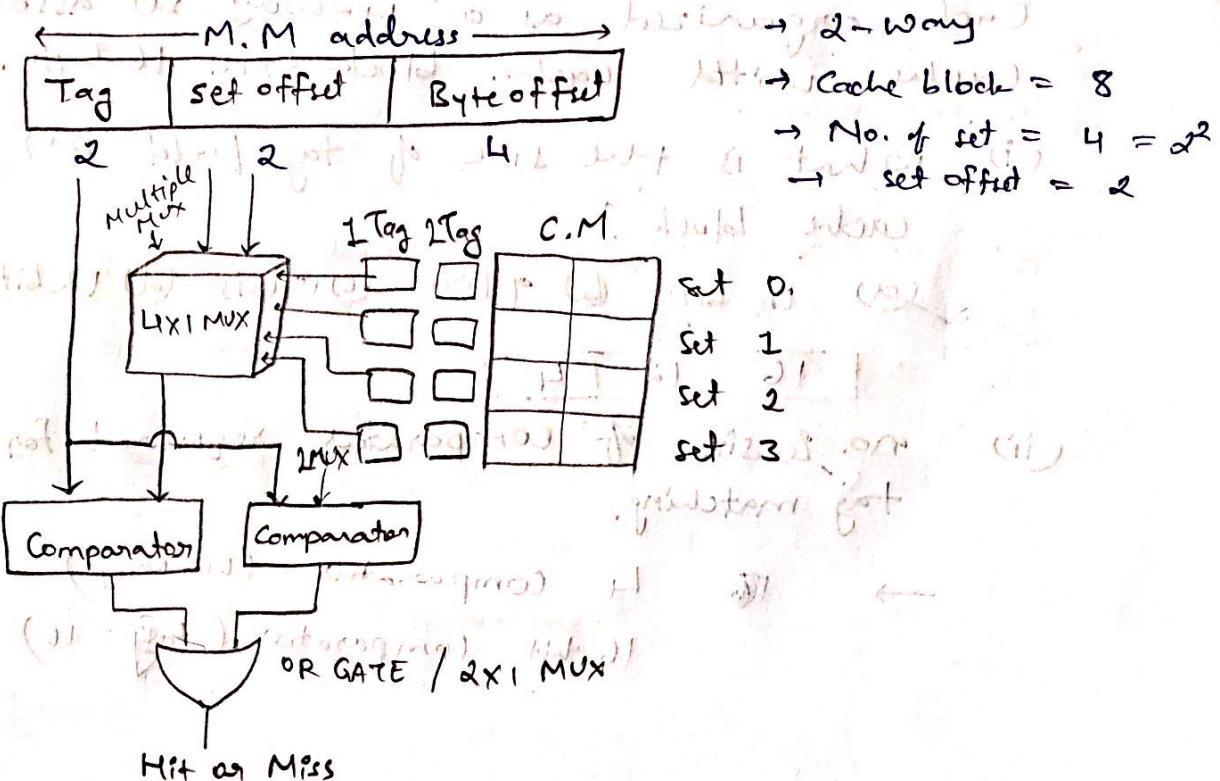
Size of Mux = No. of cache blocks \times 1

Size of comparator = No. of tag bits

Time required to declare hit or miss = MUX delay + Comparator delay

Hit latency (Hit time)

II. 2-Way Set Associative



Q-1 CPU has 32-bit memory add. and 256 kB Cache organized as a 4-way set associative cache with cache block size 16 byte.

(i) What is the size of tag field per cache block?

- (a) 16 bits (b) 9 bits (c) 19 bits (d) 12 bits

16	12	4
----	----	---

(ii) no. & size of comparators required for tag matching.

- 4 comparators (4-way)
16-bit comparator (tag = 16)

Q-11 1st cache → 32 kB, 2-way set ass.

Block size = 32 B \Rightarrow 5 bits

2nd cache → same size about direct mapped

Address = 32 bits.

2×1 MUX latency = 0.6 ns

K-bit comparator latency = $K/10$ ns

Hit latency of set associative is h_1

Hit latency of direct mapped is h_2

$h_1 = ?$, $h_2 = ?$

→ 2-way set

18	9	5
----	---	---

$$\# \text{MUX} = \frac{2^9}{2 \times 1}$$

$$= 256$$

~~18~~
~~9~~
~~5~~

~~# MUX~~

Direct Mapped

17	10	5
----	----	---

$$h_1 = \text{Comparator delay} + 2 \times \text{1 MUX delay}$$

$$\text{Time} = \frac{18}{10} + 0.6 \quad (\text{R}=0) \quad \text{Total time} = 2.4 \text{ ns}$$

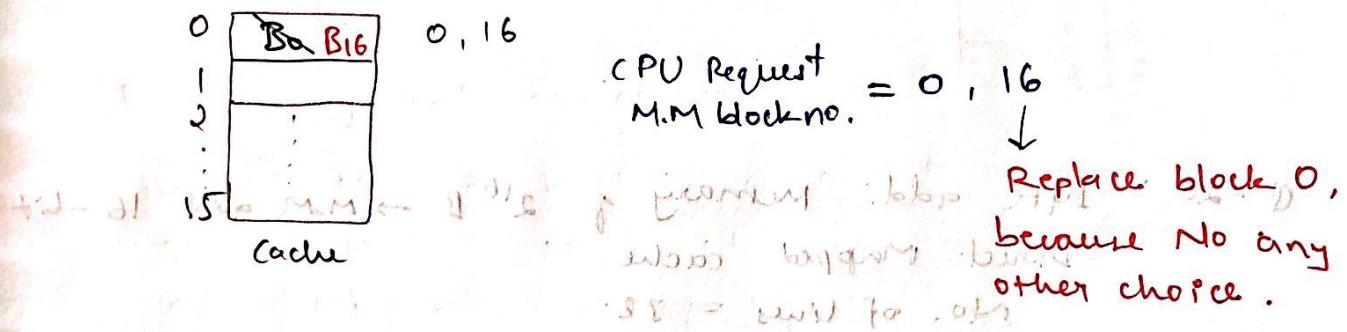
$$h_2 = \text{Comparator delay} + 2 \times \text{2 MUX delay}$$

$$\text{Time} = \frac{17}{10} \text{ ns} \quad \text{for } h_2, \text{ add to } 2.4 \text{ ns}$$

$$= 1.7 \text{ ns}$$

* Cache Block Replacement

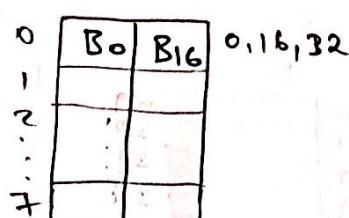
Direct



→ In Direct mapping, No replacement policy needed.

Set Associative

2-way



H 0011 = invalid address

start invalid process

M.M block = 0, 16, 32
request ↓

which block
is replaced ↓

Selected by
block replacement
Policy

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

Q - Consider direct mapped cache with 8 cache blocks (0-7) if M.M. block required are in order,

3, 5, 2, 8, 0, 63, 1, 16, 20, 17, 25, 18, 30,

24, 2, 63, 5, 82, 17 (in 24)

which of the following block will not be in cache at the end of sequence?

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

0	8	16	3524
1	X	X	317
2	X	18	282
3	3		
4	20		
5	5		
6	30		
7	63		

translating block address to Cache

block

Q-20

Q-20 Byte add. memory of 2^{16} B \rightarrow M.M. add 16-bit

Direct Mapped cache

No. of lines = 32

block size = 64 Bytes (ignoring bank offset)

50x50 2D array of bytes \Rightarrow 2500 elements

\Rightarrow 2500 Bytes

Base address = 1100 H

Array access twice



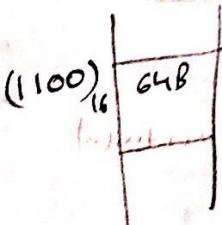
M.M. add. = 16

Tag	C.M. Block no.	B.O.
	5	6

double address

0001000100000000

block No. - 4



3 111111
2 111111
1 111111
0 111111

29	
30	
31	
32	
33	
34	X
35	X
36	X
37	X
38	X
39	X
40	X
41	
42	
43	
44	
45	
46	
47	
48	
49	
50	
51	
52	
53	
54	
55	
56	
57	
58	
59	
60	
61	
62	
63	

$$\text{No. of blocks required} = \frac{2500B}{64B} \\ \approx 40$$

After block 32, from 33 to 40 will be miss & replace at block no 4 to 11.

In 2nd access, block 1 to 8 will be miss & will replace block no. 33 to 40. = 8 miss

From block no. 9 to 32 will be hit.

From block 33 to 40 will be miss & replace the block no. 1 to 8. = 8 miss

$$\begin{aligned}\text{Total miss} &= 40 + 8 + 8 \\ &= 56\end{aligned}$$

$$\text{No. of blocks required} = \frac{2500B}{64B} \\ \text{for array} \\ \approx 40$$

After block 32, from 33 to 40 will be miss & replace at block no. 4 to 11.

In 2nd access, blocks 1 to 8 will be miss & will replace block no. 33 to 40. = 8 miss

From block no. 9 to 32 will be hit.

From block 33 to 40 will be miss & replace the block no. 11 to 8. = 8 miss

$$\text{Total miss} = 40 + 8 + 8 \\ = 56$$

Q5 A CPU has 32 kB direct mapped cache with 128B block size. Suppose 'A' is a two D array of size 512x512 with elements that occupies 8 B each.

Consider code segments P_1 & P_2 .

P_1 : for ($i=0$; $i < 512$; $i++$)

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

P_2 : for ($i=0$; $i < 512$; $i++$)

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

Suppose P_1 & P_2 are executed independently with same

initial state, namely the array A is not in the cache & P_1 & P_2 in registers. Let no. of cache misses experienced by P_1 be m_1 & for P_2 be m_2 .

$$(i) m_1 = ?$$

$$(ii) \frac{m_1}{m_2} = ?$$

Cache size = 32 kB

Block size = 128 B

Array $A = 512 \times 512 = 262144 \text{ elements}$

Element size = 8 B

$$\# \text{ Cache Line} = \frac{32 \text{ kB}}{128 \text{ B}} = 2^8 = 256 \text{ lines}$$

$$\text{No. elements in 1 cache line} = \frac{128 \text{ B}}{8 \text{ B}} = 16 \text{ elements}$$

Default storage is now major in memory.

(i) In P_1 , array is accessed in row major.

For 1st element whole block will bring in to cache i.e., (16 elements) (1 miss)

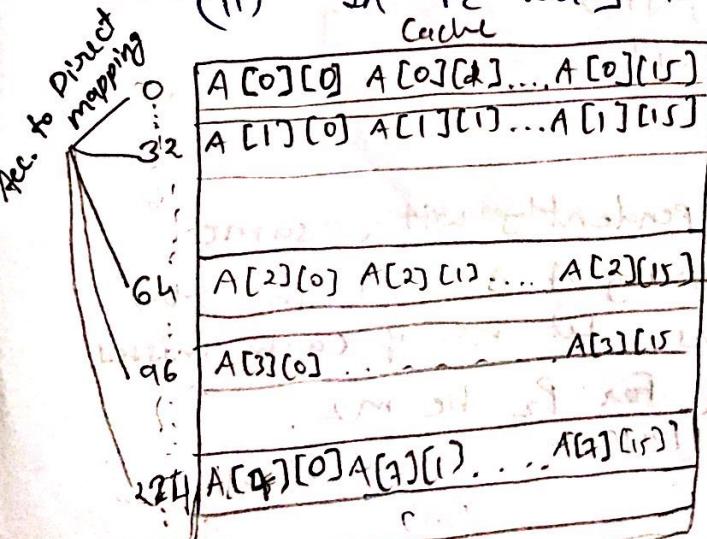
For next 15 element it will be ~~row~~ wt.

$$\therefore \text{Total miss} = \frac{512 \times 512}{16} = 16384 \text{ miss}$$

$$16384 \times 1 = 16384 \text{ miss}$$

$$\therefore m_1 = 16384$$

(ii) In P_2 array is accessed in column major



For 1st iteration we required $A[0][0]$ which is miss

& 16 elements will be bring from memory to cache

For 2nd iteration we required $A[1][0]$ which is miss

& 16 elements will be bring from memory to cache

\therefore In column major access add in every iteration there will be miss always.

$$\therefore \text{Total miss} = 512 \times 512 \\ = 262144 \text{ misses}$$

$$\therefore M_2 = 262144 \text{ Misses per column access}$$
$$\therefore \frac{M_1}{M_2} = \frac{16384}{262144} = 0.0625$$

Division of 16384 by 262144 gives 0.0625

- Q- A sequence of cache block addresses is of length n/N and contains 'n' unique block addresses. The no. of unique block addresses between two consecutive accesses to the same block addresses is bounded above by k . What is the miss ratio if the access is passed through a cache of associativity $A \geq k$.
Exercising LRU policy.

(a) n/N (b) $1/N$ (c) $1/A$ (d) k/n

Ans: n/N (b) $1/N$ (c) $1/A$ (d) k/n

block address = N

'n' unique block add. after k add.

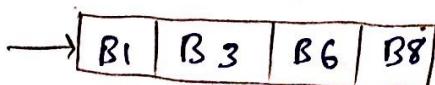


1 3 6 3 6 1 3 8 9 3
↑ ↑ ↑ ↑ ↑ ↑ ↑ ↑ ↑ ↑
 $\therefore K=4$ = no. of unique block access between 2 consecutive access of a block, should be less than or equal to k .

atleast 4-way set associative cache

min 'n' no. of misses will occur definitely

Assume all blocks map to same set



$\therefore n/N$

there will be not be any cache miss because of replacement

* Cache Miss Penalty

The amount of time CPU spends to service the cache miss.

Time required to bring the missed block from main memory to cache memory.

Assume,

No. of cycles needed to send add. to mm = 1

No. of cycles needed to access m.m = 10

No. of cycles needed to transfer 1 cell = 1

Cycle \cong C.P.U cycle

Cache block size	M.M cell size	cache miss penalty
4 words	1 word	$1 + 4 \times 10 + 4 \times 1 = 45$ cycle
4 words	2 word	$1 + 2 \times 10 + 2 \times 1 = 23$ cycle
4 words	4 word	$1 + 10 + 1 = 12$ cycle

If CPU clock rate is given,

$$1 \text{ cycle time} = \frac{1}{\text{clock rate}}$$

∴ time = no. of cycles * 1 cycle time

→ CPU sends address to memory only once

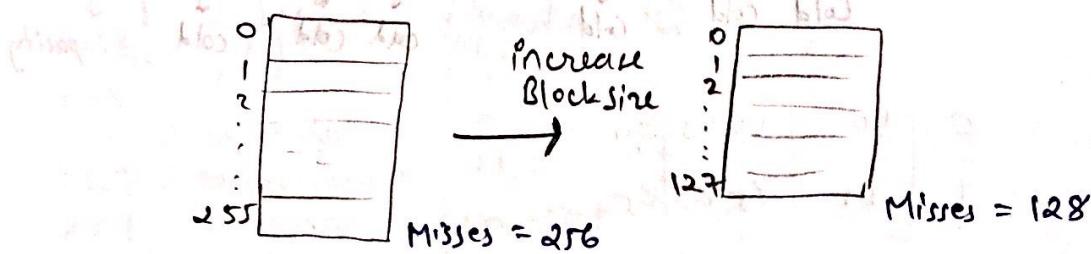
→ In one access from main memory 1 cell data can be accessed.

number of address bits = number of address bits

number of address bits

* Types of cache Misses:

1. Cold or Compulsory Miss: Cache miss occurs due to first access of any block, then it will be compulsory Miss.



→ To reduce the compulsory Misses increase the block size, which leads to less no. of cache lines. Therefore no of misses will be less.

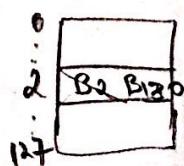
2. Capacity Miss: Cache miss occurs when the cache memory is full, then that is known as capacity miss.
[It should not be a compulsory miss]

→ To reduce miss then increase capacity of cache.

3. Conflict Miss: The cache miss occurs due to mapping constraints. Or the cache miss occurs due to tag mismatch.

→ It should not be compulsory miss and capacity miss.

→ Conflict miss occurs when cache is not full, but 2 blocks mapping to same cache line.



Example : 4 sets LRU replaces to what?

4 blocks
LRU replacement, Initially cache empty.

M.M. block req. = 0 4 0 8 0 4 1 5 1 3 1 5
↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓
Cold Cold hit Cold hit conflict cold cold hit cold hit capacity

0	B0	B4 B8 B4	B4
1	B1	B5 B3 B5	B5

not present in M. M. program so will have conflict

Q-23 Consider 2-way set associative cache with 256 blocks and uses LRU replacement. Initially the cache is empty. Conflict misses are those misses which occur due to contention of multiple blocks for same cache set. Compulsory misses occur due to first time access to block.

misses = (0 128 256 128 0 128 256 128 1 129
257 129 1 129 257 129)

This sequence is repeated 100 times.

Find no. of conflict misses experienced by Cache).

→ No. of blocks = 256

2-way set associative

$$\# \text{set} = \frac{256}{2} = 128$$

Using this it can be seen that there are 128 sets of 2 blocks each and 128 blocks of 2 sets each.

0	0 256 0 128	128
1	128 X 256 128 2	257
2	257 128	129
3		
...		
126		
127		

	1st access	2nd access	3rd access
0	Compulsory	conflict	Same
128	Compulsory	Hit	as previous
256	Compulsory	conflict	remaining 8 access
128	Hit	Hit	
0	Conflict	conflict	
128	Hit	Hit	
256	conflict	conflict	
128	Hit	Hit	
1	Compulsory	conflict	
129	compulsory	Hit	
257	compulsory	Conflict	
129	Hit	Hit	
1	Conflict	Conflict	
129	Hit	Hit	
257	conflict	Conflict	
129	Hit	Hit	

Compulsory = 6

Conflict = 4

Conflict = 8

$$\text{Total conflict} = 8 * 9 + 4$$

$$\text{misses} = 72 + 4$$

$$= 76$$

$$\text{Total compulsory miss} = 6$$

NOTE: There can not be any conflict miss in fully associative cache. Because there is only one set in fully associative cache, so block will be replaced only when set is full.

$$[(read \cdot (hit+1) + write \cdot hit)] \cdot (hit+1) + write \cdot hit = just$$

* Goal of Using Cache

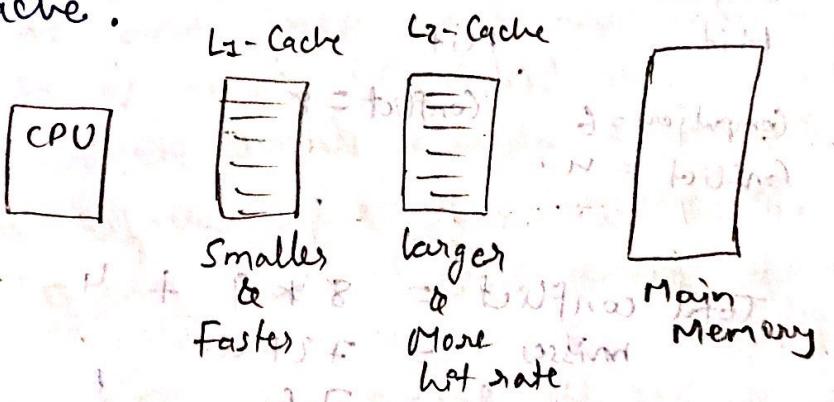
1. Minimum average memory access time
2. Maximum cache hit rate
3. Minimize cache miss penalty

→ To minimize average memory access time

we require small cache memory size.

And to maximize cache hit rate we require larger cache memory.

So to achieve both we implement multi-level cache.



→ To minimize cache miss penalty we use technique like Non-blocking cache, multipart or multibank cache.

Multilevel Cache

Simultaneous Access:

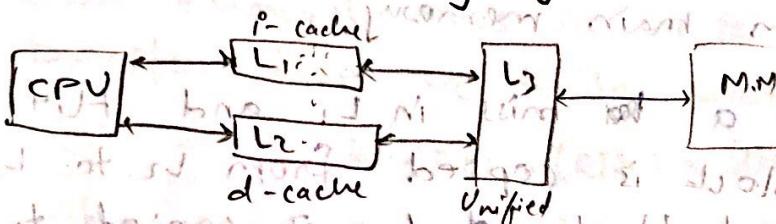
$$t_{avg} = H_1 \cdot t_1 + (1 - H_1) [H_2 \cdot t_2 + (1 - H_2) t_{m.m}]$$

Hierarchical Access:

$$t_{avg} = H_1 \cdot t_1 + (1-H_1) [H_2 \cdot (t_1+t_2) + (1-H_2) (t_1+t_2+t_{m,m})]$$

$$= t_1 + (1-H_1) (t_2 + (1-H_2)t_{m,m})$$

Q - Consider a multilevel memory hierarchy. The hit ratio of level L₁, L₂ & L₃ and MM are 0.8, 0.85, 0.9, I-O respectively. The access time of L₁, L₂, L₃ & MM are 10 ns, 10 ns, 50 ns & 500 ns. Among total memory references 60% are data access. What is the avg access time of the memory system.



$$t_{avg} = 0.8 * 10 + 0.2 (0.9 * (10+50) + 0.1 (10+50+500))$$

$$\text{data} = 0.85 * 10 + 0.15 (0.9 * (10+50) + 0.1 (10+50+500))$$

$$= 25 \text{ ns}$$

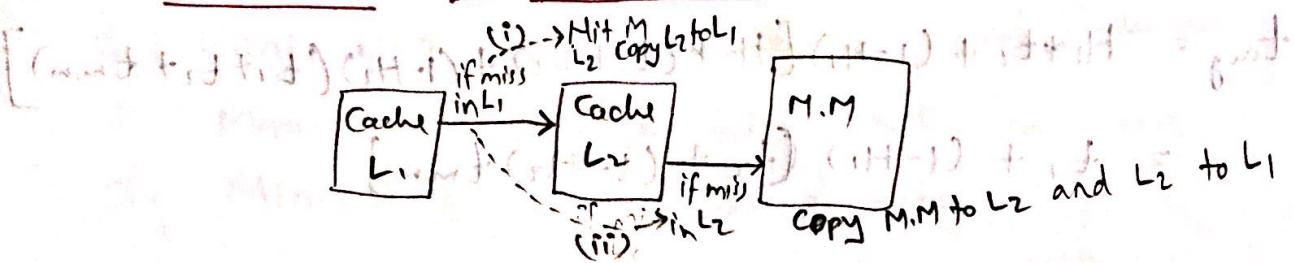
$$t_{avg} = 0.6 * t_{avg_data} + 0.4 * t_{avg_i/o}$$

$$= 0.6 * 25 + 0.4 * 30$$

$$= 15 + 12$$

$$= 27 \text{ ns}$$

* Inclusion Vs. Exclusion



Inclusion: Whatever content present in cache L₁ is also present in L₂. So content will be copied from L₂ to L₁.

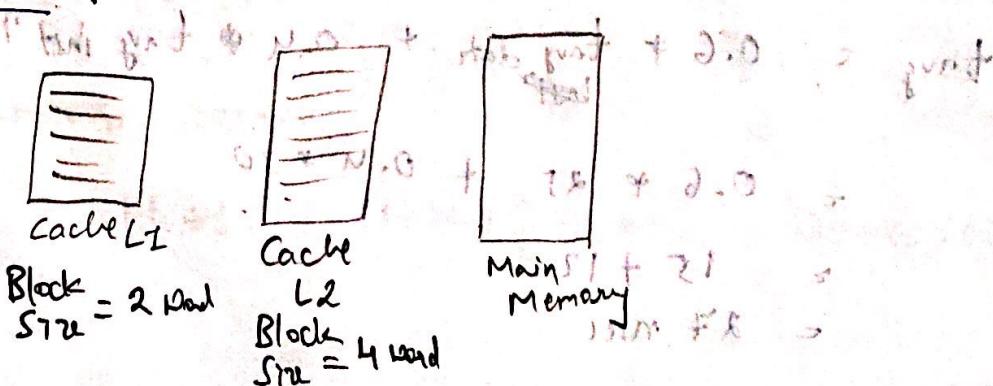
Exclusion: Content present in cache L₁ is not necessarily present in cache L₂. So content in L₁ will be bring from main memory.

→ If there is a miss in L₁ and hit in L₂ then a block is copied from L₂ to L₁ and replaced block of L₁ is copied to L₂.

→ If there is a miss in L₁ and L₂ both, then a block is copied from main memory to L₁ and replaced block of L₁ is copied to L₂.

→ L₂ cache is also known as victim block.

Inclusion:-



if a block copied from main memory \Rightarrow Size of block in L₂ (4 words)

if block copied from L₂ to L₁ \Rightarrow size of block in L₁ (2 words)

Q

Computer system has an L₁ & L₂ cache & M.M unit. Connected as shown below. Block size in L₁ is 4 word, block size in L₂ is 16 words & the memory access time are 2 ns, 20 ns, & 200 ns

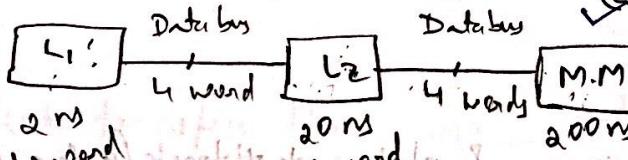
(i) When there is miss in L₁ & hit in L₂ a block is transferred from L₂ to L₁ cache. Time taken?

(a) 2 ns (b) 20 ns (c) 22 ns (d) 88 ns

(ii) When there is miss in both L₁ & L₂ cache.

First a block is transferred from main memory to L₂ cache & then L₂ to L₁. Time taken?

(a) 222 ns (b) 888 ns (c) 902 ns (d) 968 ns



(i) We have to take only transfer time & checking if miss in L₁ & hit in L₂ is already done.

\therefore Transfer time = L₂ \rightarrow L₁ access time.

$$\begin{aligned} \text{Opacity of bus} &= 20 + 2 \\ \text{Time for data access} &= 22 \text{ ns} \end{aligned}$$

$$(ii) \text{From M.M to L}_2 = 4 * (200 + 20) = 880$$

$$\text{From L}_2 \text{ to L}_1 = 22 \text{ ns}$$

$$\begin{aligned} \therefore \text{Transfer time} &= 880 + 22 \\ &= 902 \text{ ns} \end{aligned}$$

Disk Access Time:

$$\text{1 Disk Access Time} = \text{Seek time} + \text{Rotational latency} + \text{1 sector transfer time}$$

of disk controller + delay time

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

Rotational latency: Time required to rotate the desired sector under read/write heads.

1 sector Transfer time: Time required to read / write sector

$$\text{Avg. Rotational latency} = \frac{\text{Rotational latency}}{2}$$

NOTE : In one rotation entire track can be transferred

$$1 \text{ sector transfer time} = \frac{1 \text{ rotation time}}{\text{No. of sectors per track}}$$

- Q. Consider a disk with 16 platters, 2 surfaces/platter, 1K track/surface, 2K sectors/tracks & 2048 Byte/sector. Disk rotates with 3000 rpm. Seek time = 10 ms.

- (i) Capacity of Disk in GB?
- (ii) #bits for disk addressing?
- (iii) Find 1 Disk access time.
- (iv) Disk data transfer rate.

$$\text{Data transfer rate} = \frac{2048 \times 8 \times 10^3}{3000}$$

$$21.8 \text{ MB/s} = 2.18 \text{ GB/s}$$

platter = 16, 2 surface/platter

1K track / surface

2K sector / track

2048 Byte / sectors

3000 rpm.

seek time = 10 ms

at max seek time after bringing head to target disk

(i) Capacity = $2^4 \times 2 \times 2^{10} \times 2 \times 2^{10} \times 2$

blocks per track $\times 2^{37}$ of bytes in each block (total bytes)

result = 128 GB values taken

(ii) bits per byte $\log_2 2^{26} = 26 \approx 6$ bits

(iii) rotational latency = $\frac{6k}{3000} = 0.02 \text{ sec}$

rotational latency = $\frac{3000}{20000} = 0.15 \text{ msec}$

Avg. rotational latency = 10 msec

time for reading 1 track = 20 msec

1 sector ?

time for reading 1 sector = $\frac{1 \text{ sec} \times 20 \text{ msec}}{20000 \text{ bytes/sector}} = 0.00001 \text{ sec} = 0.01 \text{ msec}$

total work so, total time = 21 msec to read 1 sector + 10 msec to read 1 track + 10 msec to read 1 sector = 20.01 msec

(iv)

20 msec

1 sec

$\frac{2048 \times 2048}{0.02 \text{ sec}} = 209.71 \text{ MB/s}$

$2^1 \times 2^{10} \times 100 = 200 \text{ MB/s}$

If 'n' no. of sectors are to be transferred from the disk.

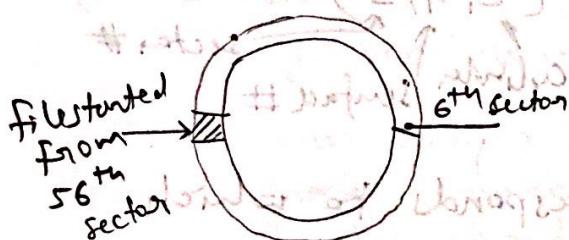
→ If sectors are stored sequentially,

$$\text{Disk Access time} = \text{Seek time} + \frac{\text{Avg. Rotational latency}}{\text{Time}} + n * \frac{\text{1 sector transfer time}}{\text{Time}}$$

→ If sectors are stored randomly,

$$\text{Disk Access Time} = n * \left[\text{Seek time} + \frac{\text{Avg. Rotational latency}}{\text{Time}} + \frac{\text{1 sector transfer time}}{\text{Time}} \right]$$

Ex - P



500 sectors in 1 track

10 msec is 2 rotations time

current head position = 6th sector

sector to transfer = 56th sector

∴ Rotational latency = time req. to rotate 50 sectors only

$$(21 \text{ at } 0) \cdot 3000 = \text{rotations} \cdot 30 = \frac{10 \text{ msec}}{500} = 1 \text{ msec}$$

Cylinder! To save seek time.

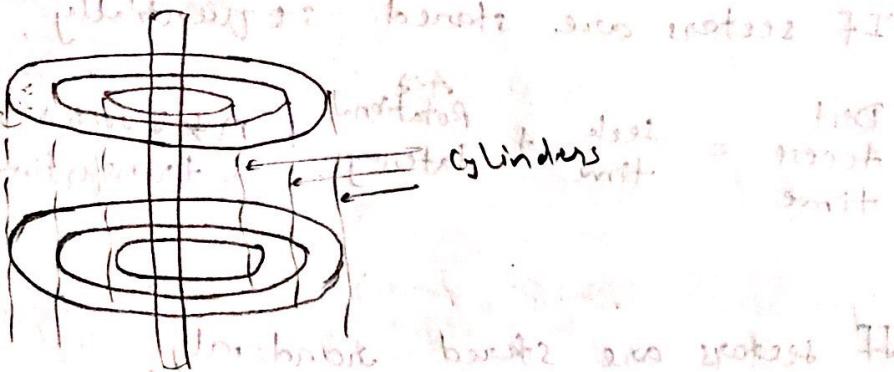
→ Respective tracks of every surface makes cylinder.

#cylinder on disk = #track per surface

$$M.D = \frac{P.C}{S.S} \rightarrow \text{higher surface if } S.S \text{ is constant}$$

constant ... n = max.

Collection of the tracks of the same radius from all the surfaces form a cylinder.



Q-3

63 sectors/track
10 platters with 2 surface
1000 cylinders.

Address of sector 13 (c, h, s)

sector # = 13
cylinder # = 13
surface # = 1
unit (sector no) = 002
unit (cylinder no) = 01
unit (surface no) = 00

(i) (400, 16, 29) corresponds to which sector no.
sectors in 400 cylinders (0 to 399)
= $400 \times 20 \times 63$

Ans: 02 sectors at 16 surface = 504000
Ans: 16 sectors M: 16 surface = 1008 (0 to 15)

Ans: add. of (400, 16, 29) = $504000 + 1008 + 29$
= 505037

Now we have to find for which surface no. 505037

(ii) Add. of 103.9^{th} sector = 505037
surfaces required = $\frac{103.9}{63} = 16.49$

We require 16 surface
 \therefore # sectors in 16 surface = 16×63
= 1008

Now remaining sectors are = $1039 - 1008$

$$= 31$$

∴ address will be $(0, 16, 31)$

* Parallel Processing

Technique for simultaneous data processing.

→ Array processing

→ Vector processing

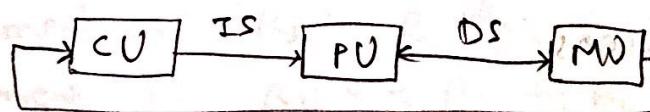
→ Pipeline Processing

Flynn's Classification

(i) SISD

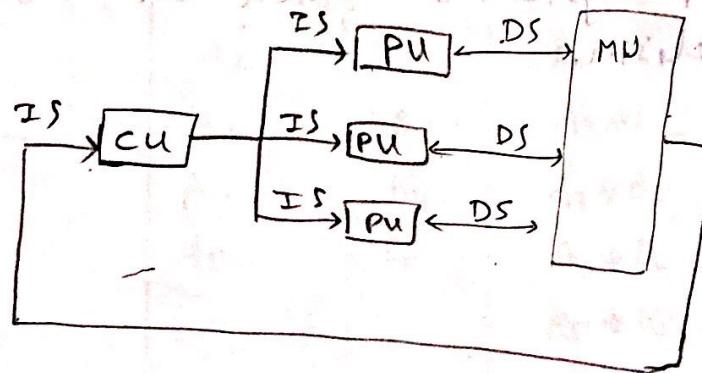
e.g. Von-Neumann

- Instruction & data both are stored in same memory after the operation the result is also stored in same memory.
- Sequential execution



- Single Instruction & Fetch & Single Instruction Execution.

(ii) SIMD :

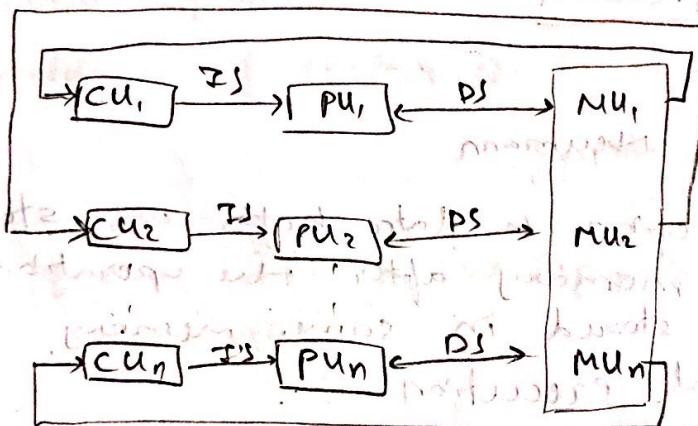


- Single instruction fetch in one cycle & multiple instruction execution.
 Eg. Pipeline processing

MISP:

- Multiple instruction fetch in one cycle & single execution.
 ∵ It is not implemented.

MIMD:



Eg. Multiprocessors

- Multiple instruction is fetched in single cycle & multiple instruction execution.
 → Multiple pipeline
 → Such type of system is based on Instruction level parallelism.

* Pipeline Processing

It is technique to divide a sequential operation into sub operation.

Pipeline is applicable only when same processing is performed over multiple input.

Each sub operation is performed in a separate unit known as segments. All the segments can perform their respective subop? parallelly on different different inputs.

An operation performed in all the segments is known as a task.

Example: $A_i * B_i + C_i$

$i=1 \text{ to } 5$; A, B, C are memory operand.

Segment 1: $R_1 \leftarrow A_i$; $R_2 \leftarrow B_i$

Segment 2: $R_3 \leftarrow R_1 * R_2$, $R_4 \leftarrow C_i$

Segment 3: $R_5 \leftarrow R_3 + R_4$

q) * 1 : dual 3 stage of pipeline unit

Segments cycle	seg 1	seg 2	seg 3	
	R_1	R_2	R_3	
1.	A_1	B_1		
2.	A_2	B_2	$A_1 * B_1$	C_1
3.	A_3	B_3	$A_2 * B_2$	C_2
4.	A_4	B_4	$A_3 * B_3$	C_3
5.	A_5	B_5	$A_4 * B_4$	C_4
6.			$A_5 * B_5$	C_5
7.				$A_5 * B_5 + C_5$

#task	#cycle in pipeline	#cycle in sequential system
$n=5$	7	$3+5 = 15$
$n=6$	8	$3+6 = 18$

→ Pipeline accept the new input before completion of the old input that means the new input are executed along with the old input.

Therefore, pipelining allows the overlapping execution of two tasks to fill the pipeline.

* General Consideration about Pipeline

Cycle time in Pipeline: Amount of time in which all the segments can perform their respective suboperation.

Consider a k-segment pipeline with cycle time t_p to perform n -task.

Time required to perform 1st task : $k \cdot t_p$

Time required to perform remaining tasks : $(n-1) \cdot t_p$

∴ Execution time of pipeline = $(k + (n-1)) t_p$

Consider non-pipeline system, which takes t_n time to perform one task.

∴ Execution time of non-pipeline system = $(n * t_n)$

Performance of pipeline

$$\text{Speed up} = \frac{\text{Performance of Pipeline}}{\text{Performance of non-pipeline}}$$

$$= \frac{1}{\frac{\text{E.T pipeline}}{\text{E.T non-pipeline}}}$$

$$= \frac{\text{E.T non-pipeline}}{\text{E.T pipeline}}$$

$$\text{Speedup (S)} = \frac{n t_n}{(k + (n-1)) t_p}$$

When $n \gg k$

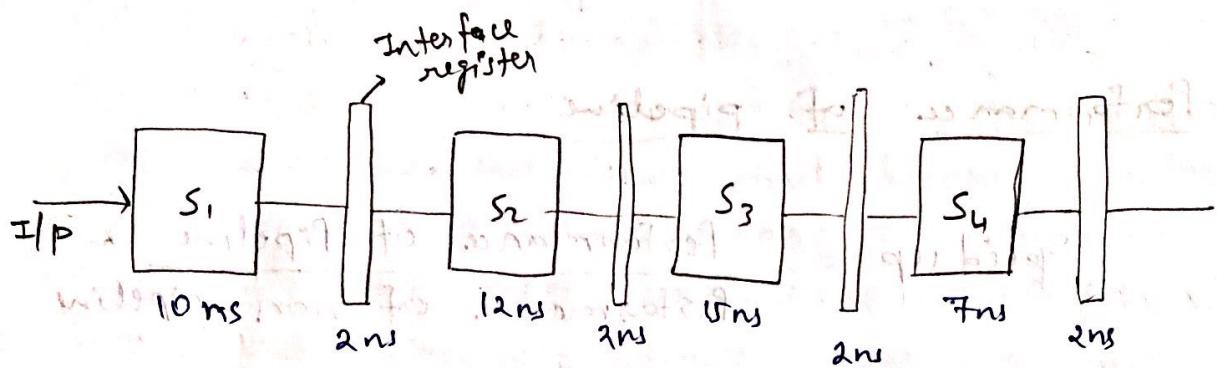
$$S_{\text{ideal}} = \frac{n t_n}{n t_p} = \frac{t_n}{t_p}$$

→ If times required to perform one task in pipeline and non-pipeline system are same,

$$t_n = k * t_p$$

$$\therefore S_{\text{ideal}} = \frac{k t_p}{t_p} \quad \text{assuming no pipeline overhead}$$

(adding) $S_{\text{ideal}} = k$



→ When buffer delay is included

$$t_p = \max(\text{segment delay}) + \text{intermediate buffer delay}$$

$$\begin{aligned} \therefore t_p &= \max(10, 12, 15, 7) + 2 \text{ ns} \\ &= 15 + 2 \\ &= 17 \text{ ns} \end{aligned}$$

→ If setup time overhead is included

$$t_p = t_p + \text{setup time}/\text{skew time}$$

- Q - A non-pipeline system takes 50 ns to process a task. The same task can be processed in 6 segment pipeline with cycle time of 10 ns. Speedup? for pipeline for 100 task.

$$\rightarrow \text{Speedup} = \frac{\frac{100 * 50 \text{ ns}}{105 * 10 \text{ ns}}}{105 * 10 \text{ ns}} = 4.76$$

$$\text{Speedup (max)} = \frac{t_n}{t_p} = \frac{50 \text{ ns}}{10 \text{ ns}} = 5$$

After applying 3 clock cycles, the output will be ready.

Q-19 A 4-stage pipeline has stage delays of 150, 120, 160, 140 ns resp. Interstage buffer delay is 5 ns. Registers used constant clock rate. The total time taken to process 1000 data items on this pipeline is —

$$\rightarrow \text{Total time req.} = (k + m) t_p$$

$$= 4 \cdot (4 + 1000 - 1) (160 + 5)$$

$$= 1003 \cdot 165$$

$$= 165495 \text{ ns}$$

$$= 165.5 \mu\text{s}$$

Q-20 CPU takes 12 cycle to complete an inst. The corresponding pipelined CPU uses 6 stage with execution times of 3, 2, 1, 5, 4, 6 and 2 clocks. What is the speed up assuming that a very large no. of instructions to be executed?

$$\rightarrow \text{Ideal Speedup} = \frac{t_n}{t_p}$$

$$\text{and total number of cycles} = 12$$

$$\text{so Speedup} = \frac{12}{6} = 2$$

$$(t_n/t_p) \cdot (1/t_p) = \text{Benefit factor} = 2$$

So Speedup = 2

Q-19 We have to design D_1 & D_2 for synchronous pipelined CPU. D_1 has 5 stages with execution time 3ns, 2ns, 4ns, 2ns, 3ns, while D_2 has 8 pipeline stage each with 2ns execution time. How much time can be saved using D_2 over D_1 for executing 100 inst?

$$\text{Execution Time, } D_1 = (5+99) * (4) \\ = 416 \text{ ns}$$

$$\text{Execution Time, } D_2 = (8+99) * (2) \\ = 214 \text{ ns}$$

$$\text{Saved Time} = 416 - 214$$

∴ Total no. of cycles of D_1 = 202 ns (slot 0 to slot 199)

→ Total no. of cycles of D_2 = 103 ns (slot 0 to slot 102)

Consider a 4-segment pipeline with segment delays 50ns, 30ns, 45ns & 45ns. The intermediate register delay is 5ns.

(i) How long would it take to process 100 m/p. in this pipeline?

(ii) How can we reduce total time about $\frac{1}{2}$ of the time calculated in part 1.

$$\rightarrow (i) \text{ reg. time} = (4+99) * (95+5) \\ = 985 \text{ ns} \quad 103 * 100 \\ = 10300 \text{ ns}$$

$$(ii) \frac{1}{2} \times 9785 = (k + n - 1) t_p$$

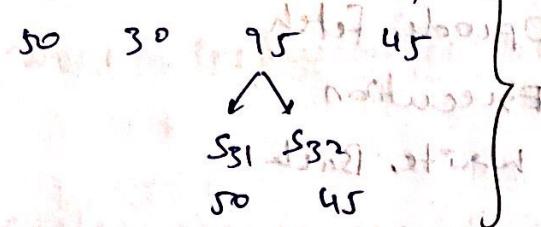
$$4892.5 = (4 + 99) * (t_p + 5)$$

$$\frac{4892.5}{108} = (t_p + 5)$$

$$45 = t_p + 5$$

$$\text{target address } 1112.5 = t_p + 5 = 45$$

not increasing number of stages → ACT



We can increase the stage while maintaining equal stage delay to increase the performance.

Q3

Consider a non-pipelined processor with a clock rate of 2.5 GHz and average cycle per instr of four. The same processor is upgraded to pipelined processor with five stages; but due to the internal pipeline delay, the clock speed is reduced to 2 GHz. Assume that there are no stalls in the pipeline. The speedup achieved in this pipelined processor is _____

→

$$k = 5$$

$$t_{Rn} = \frac{1}{2.5 \text{ GHz}} = 0.4 \text{ ns} = 0.4 \times 4 = 1.6 \text{ ns}$$

$$t_p = \frac{1}{2 \text{ GHz}} = 0.5 \text{ ns}$$

$$\text{Speedup} = \frac{t_n}{t_p} = \frac{1.6}{0.5} = 3.2$$

* Instruction Pipeline

- Pipeline processing is implemented for instruction cycle.
- Assume 5 segment pipeline as follows:

IF - Instructions fetch

DA - Decode & Address calculation

OF - Opcodes fetch

EX - Execution

WB - Write Back

clock cycle 1 2 3 4 5 6 7 8 9 10 11 12

Instⁿ

I₁ IF DA OF EX WB

I₂ IF DA OF EX WB

I₃ IF DA OF EX WB

stall cycle

I₄ IF DA OF EX WB

I₅ IF DA OF EX WB

I₆ IF DA OF EX WB

I₇ IF DA OF EX WB

I₈ IF DA OF EX WB

In cycle 4 CPU decodes that I₃ is branch instruction.

I₁ and I₂ will be completely executed.

remove all instruction (I_n) from pipeline, fetched after I₃.

→ After EX phase of I₃, branch result (target add.) will be available [in cycle no. 6] → In cycle no.

(6 + 2) = 8 = target add. will be available for next instⁿ is fetched.

NOTE: After execution phase of the branch instⁿ, target Address is available.

Actually no. of instructions executed = 6 (n)

$$k = 5 \text{ (No. of instructions)}$$

$$\text{No. of cycles req. in usual} = k + n - 1$$

$$= 5 + 6 - 1$$

After 5 executed normal inst : ~~target add~~ + reading of stall cycle because of branch instⁿ = 3

$$\text{Total no. of cycles} = 10 + 3$$

$$= 13$$

NOTE: If ~~first~~ after ith stage of the branch instⁿ executed then, no. of stall cycle is (i-1)

Stall Cycle: Extra cycle in the pipeline

because of any hazard (problem).

Q-

FI - 5 ns

Total instⁿ = 12

DZ - 7 ns

In 13 branch instⁿ & target D Zg.

FO - 10 ns

FX - 8 ns

WB - 6 ns

Storage buffer = 17 ns

$$\text{Total cycle} = (k+n) + \text{stall cycle}$$

* Total no. of total instruction will be execute = $(I_1 - I_{n_1}) + (I_2 - I_{n_2})$
 $= 8$

$\therefore k+n = 5+8-1 = 12$ (Total no. of instruction)

$\therefore \text{Total cycle} = 12 + 3$
 $= 15$

(iii) Total time req. = $15 * 10 \text{ ns}$
 $= 150 \text{ ns}$
 $I_1 - I_{n_1} = 5 \text{ ns}$ (Time taken by first instruction)
 $I_2 - I_{n_2} = 8 \text{ ns}$ (Time taken by second instruction)

* Pipeline Hazard: The reason because of which the smooth conduct of pipeline is disturbed.

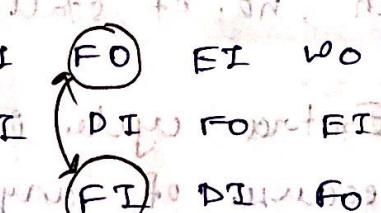
Types

(i) Resource Conflict: If the same resource is required in two stages simultaneously.

Ex 1

I_1 FI DI FO EI WO

I_2 FI DI FO EI WO



Same resource is required

Required for both

Ex 2 1 cycle in Execution \rightarrow ADD

4 cycle in Execution \rightarrow MUL

MUL FI DI FO EI EI EI WO

ADD FI DI WFO EI EI WO

Execution for ADD
will be available after MUL

(ii) Data Dependency : If the result of an instruction is used as input in the next instruction.

Ex:

$$i : R_1 \leftarrow R_2 + R_3$$
$$i+1 : R_4 \leftarrow R_1 + R_6$$

i : IF ID OF EX WB

i+1 : IF ID $\xrightarrow{=}$ OF

Can not perform OF
as R_1 is available after WB

→ Regular pipeline can not detect the data dependency.

→ To minimize the no. of stall in pipeline due to data dependency

S/W Solⁿ : Delayed load

H/W Solⁿ : H/W interlock or By pausing

Delay load! Compiler inserts no operation instruction are independent instⁿ before the dependent instruction.

i : $R_1 \leftarrow R_2 * R_3$ IF ID OF EX WB

i+1: NOP

i+2: NOP

i+3: $R_4 \leftarrow R_1 * R_6$ IF ID OF EX WB

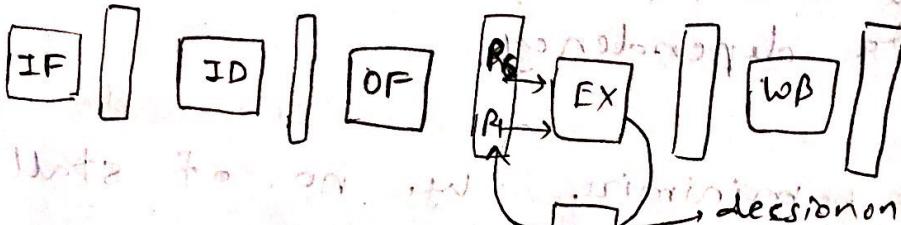
H/w interlock: Extra stall cycle created in execution by H/w.

IF ID OF EX WB
IF ID \leftarrow OF EX WB
stall cycle added by H/w.

operand forward!

$$R_1 \leftarrow R_2 + R_3$$
$$R_4 \leftarrow R_1 + R_6$$

IF ID OF EX WB
IF ID OF EX WB



Condition (because every instⁿ b/w happens): We don't require operand forwarding

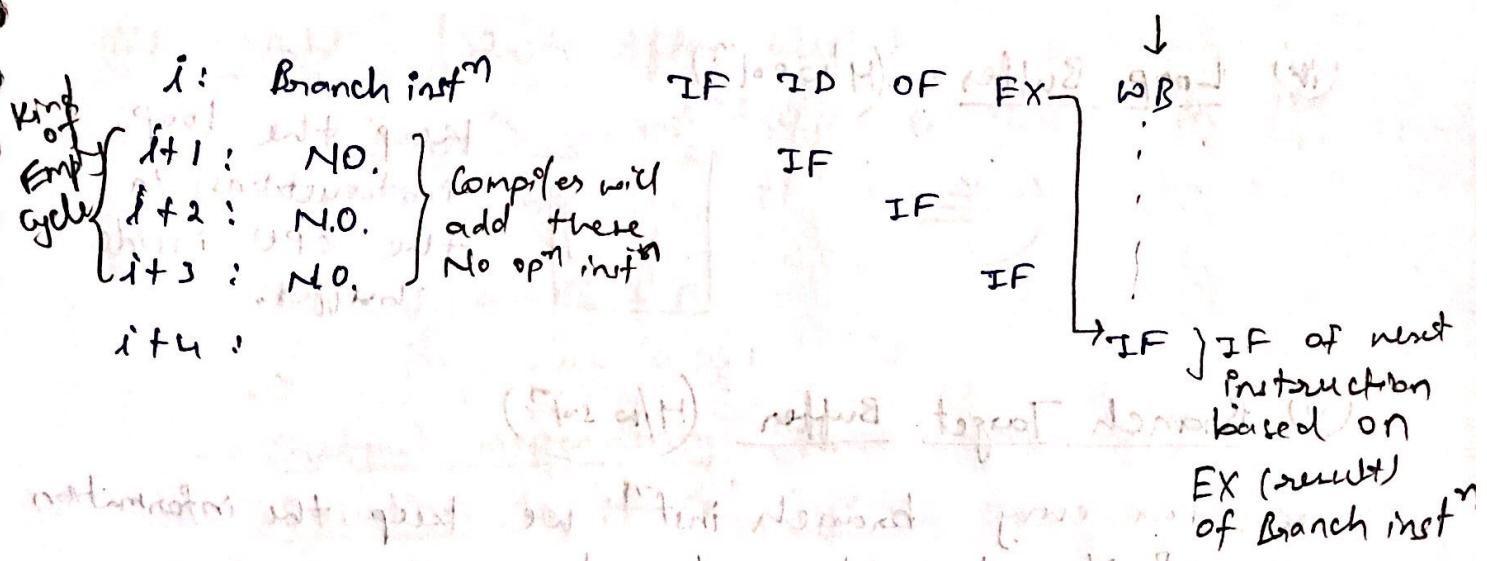
If operand Forwarding is used then no any stall cycles for the data dependency.

(iii) Branch Difficulty / Control Dependency

If occurs because of branch instⁿ

(i) Delayed branch (slow solⁿ)

→ Provided by compiler



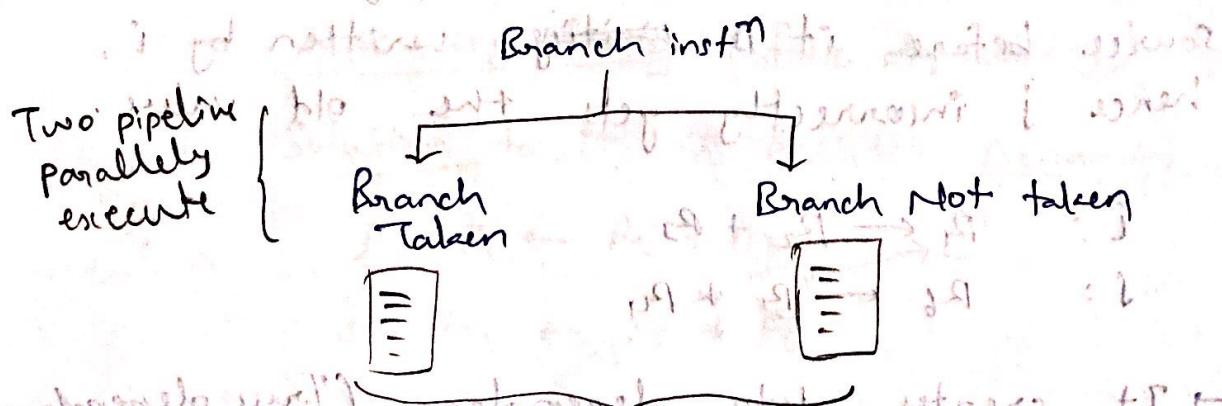
(ii) Branch Prediction: (H/W solⁿ)

get from this stored part of pipeline of Branch predictor, predict whether branch is taken or not.

If branch is predicted (taken), condition will be verified, if true then there is no extra stall cycle required.

If it is false rollback to corresponding instruction & execute based on target address available which required extra stall cycle.

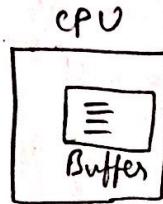
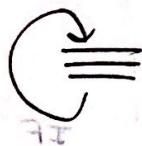
(iii) Prefetched Target Instruction: (H/W solⁿ)



(whichever will be true to go to next stage)

After execution of Branch instⁿ we come to know condⁿ is true or false.

(iv) Loop Buffer: (H/w sol^m)



keep the loop instructions in the CPU mode buffer.

(v) Branch Target Buffer (H/w sol^m)

For every branch inst^m we keep the information of target address based on its previous execution. ~~(for all H) & consider instruction (i)~~

If previously taken branch will not be the target address then rollback.

not taken, (not Branch inst^m) Target Add

if went wrong, If it is taking If it is taking wrong data write on

* Data Hazard Classification

Consider two inst^m i, j; if i execute before j.

(i) RAW (Read After Write)

Instruction j is trying to read a source before it is writing written by i, hence j incorrectly gets the old value.

$$i: R_1 \leftarrow R_2 + R_3$$

$$j: R_6 \leftarrow R_2 * R_4$$

→ It creates data dependency. (True dependency) which can be solved by delayed load, operand forwarding.

(ii) W.A.W : (Write After Write)

Instⁿ j is trying to write in the destination before it is written by i.

$$i: R_4 \leftarrow R_2 + R_3$$

$$j: R_4 \leftarrow R_5 * R_6$$

→ Output dependency (false dependency)

Solution to it is Register Renaming (H/w so/m)

$$i: R_4 \leftarrow R_2 + R_3$$

$$j: R_{4'} \leftarrow R_5 * R_6$$

$$R_4 \rightarrow R_{4'}$$

(iii) W.A.R (Write After Read)

Instⁿ j is trying to write a destination before it is read by i, hence i incorrectly gets new value.

$$i: R_4 \leftarrow R_2 + R_3$$

$$j: R_4 \leftarrow R_5 * R_6$$

→ Anti-dependency (false dependency)

Solution to it is Register Renaming - (H/w so/m)

$$i: R_4 \leftarrow R_2 + R_3$$

$$j: R_{4'} \leftarrow R_5 * R_6$$

1st iteration → $\Sigma \leftarrow \text{sum of } R_2 + R_3$

2nd iteration → $\Sigma \leftarrow \text{Vid of } R_4$

3rd iteration → $\Sigma \leftarrow \text{Vid of } R_4'$

F

$$\text{Total time} = 5 + 8 = 13 \text{ clock cycles}$$

Q-6 IF = 1
 ID = 2
 OP = 2n
 PO = 2
 WO = 1

For ADD PO \rightarrow 2 ns
 For SUB PO \rightarrow 1 ns
 For MUL PO \rightarrow 3 ns
 For DIV PO \rightarrow 6 ns

operand forwarding is used in pipeline.

$I_0 : MUL R_2, R_0, R_4 \quad R_2 \leftarrow R_0 \times R_4$

$I_1 : DIV R_5, R_3, R_4 \quad R_5 \leftarrow R_3 / R_4$

$I_2 : ADD R_2, R_5, R_2 \quad R_2 \leftarrow R_5 + R_2$

$I_3 : SUB R_5, R_2, R_6 \quad R_5 \leftarrow R_2 - R_6$

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15

$I_0 : IF ID OF PO PO PO (W0)$

$I_1 : IF ID - OF PO PO PO PO PO PO W0$

$I_2 : IF ID - OF PO PO W0$

$I_3 : IF ID - OF PO W0$

15 cycle required.

#hazⁿ = 4

K = 5

Usual no. of cycles without hazard = K + n - 1

$$= 5 + 4 - 1 = 8$$

All hazⁿ takes 2 cycle for PO, ~~for~~.

for MUL \rightarrow 3 \rightarrow extra 2

for DIV \rightarrow 6 \rightarrow extra 5

for ADD/SUB \rightarrow 2 \rightarrow extra 0

$\frac{7}{7}$

$$\therefore \text{Total w/cycles} = 8 + 7 = 15$$

For 'n' instruction,

$$\# \text{cycles} = k + n - 1$$

Without Hazard $CPI_{Avg.} = \frac{k+n-1}{n}$

With Hazard $CPI_{Avg.} = \frac{k+n-1+x}{n}$ [x = extra stall cycle]

In Ideal Condition,

for 'n' instruction, #cycles = n [Ignore k-1]

$$CPI_{Avg.} = \frac{n}{n} = 1$$

Without Hazard, $CPI_{Avg.} = 1$

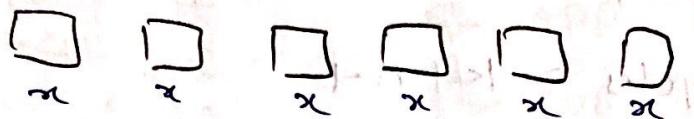
With hazard, $CPI_{Avg.} = \frac{n+x}{n}$

Average instruction execution time = $CPI_{Avg.} * t_p$

Q-16 Consider 6-stage "inst" pipeline, where all stages are perfectly balanced. Assume that there is no cycle-times overhead of pipelining. When an app. is executing on this 6-stage pipeline, the speedup achieved with respect to non-pipelined execution, if 25% of inst incur 2 pipeline stall cycles.

Even if first hazard occurs, there is no additional stall.

$$k = 6$$



$$t_p = x$$

$$t_n = 6x$$

freq. of stall instⁿ = 25%

stall cycle = 2

$$S = \frac{t_n}{(1 - 2 \text{ stn}) \cdot CPI_{avg} + t_p}$$

$$CPI_{avg} = 0.75 * I^{197} + 0.25 (I+2)$$

$$I = 0.75 + 0.75 \\ = 1.5$$

$$\frac{t_n}{S} = \frac{6x}{1.5 * x}$$

$$6x = 4 \\ x = 0.6666666666666666$$

Q-9: Instⁿ pipeline has five stages where each stage takes 2 ns. & all instⁿ use all five stages. Branch instⁿ are not overlapped, i.e., the instⁿ after the branch is not fetched till the branch instⁿ is completed. Calculate the avg. instⁿ execution time assuming that 20% of all instⁿ executed are branch instⁿ. Ignore the fact that some branch instⁿ may be conditional.

$$CP_{\text{Lang}} = 0.80 \times 2 \text{ m} + 0.20 (1.5 \times 2 \text{ m})$$

$$= 1.6 \text{ m} + 0.20 \times 2 \text{ m}$$

$$= 3.6 \text{ m}$$

Q-10. From above question, if a branch instⁿ is a conditional branch instⁿ, the branch need not be taken. If branch is not taken, the following instⁿ can be overlapped. When 80% of all branch instⁿ are conditional branch instⁿ, & 50% of the conditional branch instⁿ are such that the branch is taken, what is the average instⁿ execution time?

$$= 0.20 \left[0.80 \left(0.50 \times 20 + 0.50 \times 5 \right) + 0.20 \times 1 \right]$$

$$= 0.8 + 0.2(0.8(3) + 1)$$

= ~~A~~ 1.48 cycle

$$t_p = CPJ_{avg} * t_{p_{\text{initial}}}$$

$$k_{\text{trap}} t = 1.48 \times 2 \text{ ns}$$

$$= 2.96 \text{ ms} \text{ (out of balance)}$$

* Efficiency & Throughput

if efficiency is 100% which means the speedup of pipeline is maximum speedup.

Throughput is $\frac{1}{t_p}$ because 1 cycle $\rightarrow 1$ instruction and been delayed after pipeline registers and total time is t_p so throughput is $\frac{1}{t_p}$.

* Dynamic Pipeline

S_1	X			X	
S_2			X		
S_3		X			X

Execution Sequence : $S_1 S_3 S_2 S_1 S_3 + 8.0$

Initiation: Providing a new input to the pipeline is known as initiation.

Latency: Number of time units between two initiation.

or After how many cycles next input provided to the pipeline.

Collision: Resource conflict is known as collision.

Permissible Latency: The latency which does not cause collision.

Forbidden Latency: The latency which causes collision.

Example

	1	2	3	4	5	6	7	8	9	10	11	12
S_1	X					X		X				
S_2		X		X								
S_3			X	X		X						

Find the number of collisions in 8 lines not →

After 8 lines there is no collision yet. If →

I_1 $S_1 \ S_2 \ S_3 \ S_2 \ S_3 \ S_1 \ S_3 \ S_1$ 0101010101010101

I_2 $S_1 \ S_2 \ S_3 \ S_2 \ S_3 \ S_1 \ S_3 \ S_1$ 0101010101010101

I_3 S_1 0101010101010101

I_4 S_1 0101010101010101

I_5 S_1 0101010101010101

I_6 S_1 0101010101010101
because of resource conflicts → Collision (S_1 is being used by I_1)

i.e., S_2 is in use

i.e., after 5 cycle we can not insert new instruction

∴ forbidden latency = 5 cycle

Forbidden Latency : (6-1), (8-1), (8-6), (4-2), (5-3), (7-3),

$$= 5, 7, 2, 2, 2, 4, 2$$

$$= 7, 5, 4, 2$$

Permissible Latency : 8, 6, 3, 1

Collision Vector:

with 8 bits, position 0 is 1st bit from left
 for collision (for forbidden) = 1
 for no collision (for permissible) = 0

$c_8 \ c_7 \ | c_6 \ c_5 \ | c_4 \ c_3 \ | c_2 \ c_1$

0 1 0 1 | 1 0 1 0

State Transition Diagram

01011010

→ For every 0 in collision vector, right shift it by position of 0' and take bit-wise OR with original collision vector.

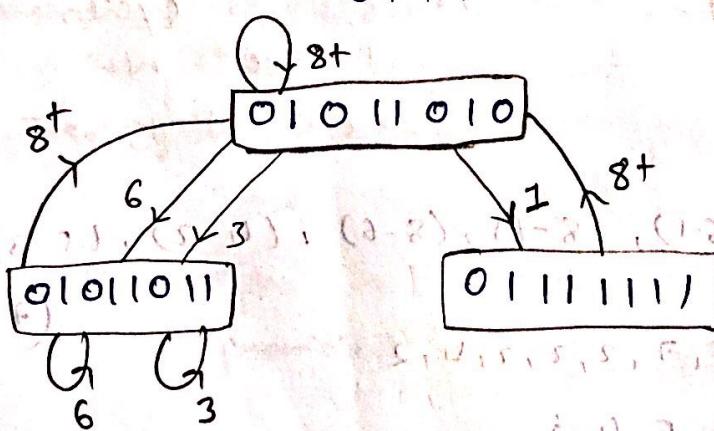
For state 01011010

(i) For '0' at position 1:

00101101 → Right shifted by position

OR 01011010 → original state

01111111



(ii) For '0' at position 3:

$$\begin{array}{r} 0000101 \\ \text{OR } 01011010 \\ \hline 01011011 \end{array}$$

(iii) For '0' at position 6:

$$\begin{array}{r} 00000001 \\ \text{OR } 01011010 \\ \hline 01011011 \end{array}$$

(iv) For '0' at position 8:

$$\begin{array}{r} 00000000 \\ \text{OR } 01011010 \\ \hline 01011010 \end{array}$$

For state 011111

(i) For '0' at position 8:

$$\begin{array}{r} 00000000 \\ \text{OR } 01011010 \\ \hline 01011010 \end{array}$$

for state 01011011

(i) For '0' at position 3

$$\begin{array}{r} 00001011 \\ \text{OR } 01011010 \\ \hline 01011011 \end{array}$$

(ii) For '0' at position 6:

$$\begin{array}{r} 00000001 \\ \text{OR } 01011010 \\ \hline 01011011 \end{array}$$

Simple Cycle: A cycle in which every state appears only once. (except the starting state)

(8), (1,8), (3,8), (6,8), (3), (6)

Greedy Cycle: The simple cycle in which the minimum latency from every state is considered.

(1,8) (3)

$$\text{Average latency} = \frac{1+8}{2} = 4.5 \quad \left\{ \begin{array}{l} \text{Latency to '8' not '0'} \\ \text{Min Avg. latency} = 3 \end{array} \right.$$

$$\text{Avg. latency (3)} = \frac{3+8}{2} = 5.5$$

* Floating Point Representation

- A larger range of numbers can be represented using limited number of bits as compared to the fixed point representation.
- In this representation every number is represented as follows:

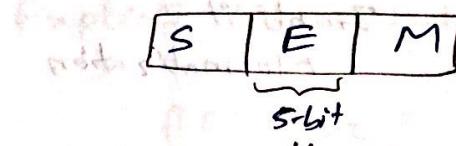
S	E	M
Sign bit	Exponent	Mantissa

Sign bit always 1-bit

0 → -ve no.

1 → +ve no.

- Exponent is represent in biased form
- Mantissa is a normalized (explicit or implicit) fraction number.



2) Biased exponent

Range of sign no. = -2^4 to $+2^4$ / -16 to $+15$

$1010.1 = 14$ $10101.0 = 31$

$1010 = 14$ $10101.0 = 31$

$s = 1$ $s = 0$

biased {represent in unsigned} not i.e., the
0 to 31

Original exponent from number (e)

stored Exponent (E)

$$S + M \cdot 2^{e-127}$$

-16 0

-15 1

Excess-16

∴ 0 transposed sign, 1 is bias code

Bias = 16

$$E = e + \text{bias}$$

$$\text{Bias} = 2^{k-1}$$

if k-bits are used to represent Exponent

Normalization

(Result as follows) Normalization is needed to → 101.01 estimate answer

↓
Explicit
Normalization

→ Number should be
in form $0.\underline{\hspace{2em}}$

i.e., 0.10101×2^3

$M = 10101$ [No. after
the point]

$e = 3$

$E = e + \text{bias}$

↓
Implicit
Normalization

→ Number should be
in form $1.\underline{\hspace{2em}}$

i.e., 1.0101×2^2

$M = 0101$

$e = 2$

$E = e + \text{bias}$

$\text{Value} = (-1)^s \cdot M \times 2^{E-\text{bias}}$

$\text{Value} = (-1)^s \cdot M \times 2^{E-\text{bias}}$

→ Conventional floating point can not represent '0'.

Q. Consider a 16-bit register used to represent floating point numbers. Mantissa is normalized sign magnitude fraction number. Exponent is in Excess - 32 form. What is the 16-bit pattern representing $(-10.25)_{10}$ in this register.

→ Exponent is Excess - 32

bias = $32 - k = 2^{k+1}$

$\therefore k = 6$ [k represent Exponent]

value = 9 = 9

value = 9 = 9

S	E	M
1	6	9

$$(10.25)_{10} = (1010.01)_2$$

Explicit Normalization has (e) = 3.0000

$$(1010.01)_2 = \frac{0.101001}{\text{normalize}} \times 2^4$$

$$M = 101001$$

$$e = 4$$

$$E = e + \text{bias}$$

$$= 4 + 32$$

$$E = 36$$

$$= (100100)_2$$

$$\therefore S E M = 100100$$

$$1 \quad 100100 \quad 110100 \quad \underbrace{1}_{\text{padding bit}}$$

Q- What is the maximum value which can be represented in above register.

→ For max. no. - $S=0$

$$E = E_{\max} = 11..1111$$

$$M = M_{\max} = 11..11$$

$$\therefore \text{Value} = (-1)^0 * 0.1111111 * 2^{63-32} = 0.1111111 * 2^{31}$$

$$= 0.1111111 * 2^{31} = 511$$

$$(1011001) = 1.1111111 * 2^9 * 2^3$$

$$= (2^9 - 1) * 2^{22}$$

$$= 511 * 2^{22} = 1023 * 4194304$$

Q-4 Mantissa is a pure fraction in signed magnitude form. The decimal no. 0.239×2^{13} has the following hexadecimal representation without normalization and rounding off

- (a) 0024 (b) 0D4D (c) 4D0D (d) 4D3D

Sign	Exponent	Mantissa
1	14	870

Excess-64



$$S = 1\text{-bit} + 0 = 0$$

$$E = 7\text{-bit} + 13 = 20$$

$$M = 8\text{-bit} = 0.239$$

$$0.239 \times 2^{13} = (001001)_2$$

$$0.239 \times 2 = 0.478$$

$$0.478 \times 2 = 0.956$$

$$0.956 \times 2 = 1.912$$

$$0.912 \times 2 = 1.824$$

$$0.824 \times 2 = 1.648$$

$$0.648 \times 2 = 1.296$$

$$0.296 \times 2 = 0.592$$

$$0.592 \times 2 = 1.184$$

$$0.184 \times 2 = 0.368$$

$$0.368 \times 2 = 0.736$$

$$0.736 \times 2 = 1.472$$

$$0.239 \times 2^{13} = 0.0011101 \times 2^{13}$$

$$M = 0011101$$

$$E = 13 + 64 = 77$$

$$E = 13 + 64 = 77 = (1001101)_2$$

$$\therefore \underline{01001101} \underline{0011101} \times (1001101)_2$$

$$\times 4 \quad 03D$$

Q-5 The normalized representation for the above format is specified as follows. The mantissa has an implicit preceding the binary (radix) point. Assume that only 0's are padded in while shifting a field. The normalized representation of the above

$$(0.239 \times 2^{13})_{10}$$

$$\rightarrow 0.0011101 \times 2^{13} = 1.11101000 \times 2^{10}$$

$$M = 11101000$$

$$e = 10$$

$$E = 10 + 64 = 74$$

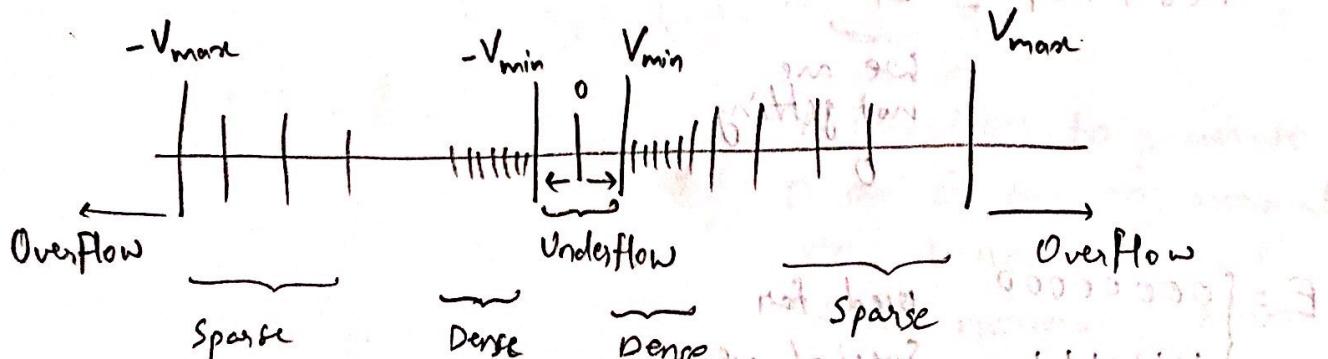
$$= 74 = \cancel{0100000}$$

$$= (1001010)_2$$

$$\therefore \underline{01001010} \underline{11101000}$$

$$0x4AE8$$

Distribution of Nos.



More bits in exponent \rightarrow Large range

More bits in Mantissa \rightarrow Better precision

* IEEE-754 Floating Point Standard

Single Precision

Double Precision

S	E	M
1	8	23

S	E	M
1	11	52

$$\text{Bias} = 127$$

$$\text{Bias} = 1023$$

$$E = 8 - \text{bit}$$

$$\therefore \text{Range} = -128 \text{ to } 127$$

$$127 + 127 = \underbrace{254}_{\text{e + bias}}$$

we are not getting 255

$$-128 + 127 = \underbrace{-1}_{\text{0}}$$

we are not getting 0

$$E = \begin{cases} 00000000 \\ 11111111 \end{cases}$$

used for
Special no.
representation

Special cases

S	E	M	
0	00...0	00...0	+0
1	00...0	00...0	-0
0/1	11...1	00...0	$\pm\infty$
0/1	111...1	$M \neq 0$	Not A Number
0/1	00...0	$M \neq 0$	Fraction or denormalized no.
0/1	$E \neq 00..0$ and $E \neq 111..1$	$M = XX.X$	Implicit normalized no.

Denormalized No. : A very very small number, which can not be implicitly normalized. To format numbers correctly, add 8, might not add 1.

For any number biffles $e = -127$

biased exponent $\Rightarrow E = -127 + 127$

$$E = 0 \quad \underline{S(0.011)} = 0 \quad \underline{E(0.011)}$$

Restriction to generate 0 as E for any normal no., hence restrict the normalization upto only -126 position.

$$\text{Eg. } 0.00\dots1 / S + 10011.1 =$$

$$1.1 * 2^{-127} \rightarrow \text{Not allowed}$$

$$(01)_{26.00} 0.11 * 2^{-126} \rightarrow \text{Store as denormalized}$$

→ If after normalization upto ± 126 position, we get implicit normal form then that no. is implicit normalized no., otherwise it will be denormalized no. w/ fractional part.

$$\text{Value} = (-1)^S \times 0.M \times 2^{-126} \text{ or } 2^{-1023} \text{ (double)}$$

denormalized

$$\text{Value} = (-1)^S \times 1.M \times 2^{E_{bias}} \text{ (Implicit)}$$

Q-7 Using single-precision 32-bit floating point format of IEEE-754 standard that uses 1-bit for sign, 8 bits for biased exponent & 23-bit for mantissa. A float type variable X is assigned the decimal value of -14.25 . X in hexa is ____.

$$(14.25)_{10} = (1110.01)_2$$

Showing of distribution

S	E	M
1	8.000	23

bits = 127

midas@laptop: ~%

Writing as $(1.110.01)_2 \rightarrow$ Implicit representation

$$= 1.11001 \times 2^3$$

$$M = +1.11001 \times 1.11111111111111111111111$$

$$E = 3 + 127 = 130$$

$$\text{bit sequence is } 0 \text{ to } 31 \rightarrow = (10000010)_2$$

$$\therefore \underline{1.1000010} \underline{11001000...0}$$

0x C164 0000

Q- Consider the following 32-bit number that denotes single precision floating point no. What is the decimal value represented by it.

$$0.10000011 \ 11000...0$$

→

$$S = 0$$

$$E = 1000\ 0011$$

$$= \del{100} 131$$

$$131 - 127 = 4 \Rightarrow e = 4$$

$$M = 11000$$

$$\text{Value} = (-1)^0 * 1.11 * 2^4$$

$$\text{Value} = (1110)_2$$

$$= 28$$

Q- 0 0000 0000 1100...0

$$S = 0$$

$$E = 0$$

$$\therefore e = 0 - 127 = -127 \rightarrow \text{denormalized}$$

so $e = -128$

~~please~~

$$\text{Value} = (-1)^0 * 0.11 * 2^{-126}$$

$$= 11.0 \times 2^{-128}$$