

INTRODUCTION

Computer Architecture

- Programmer's Perspective
- Abstract View
- Programmer Knows that there's CPU, Memory, I/O & ways to interact with it say, using instructions.

Computer Organization

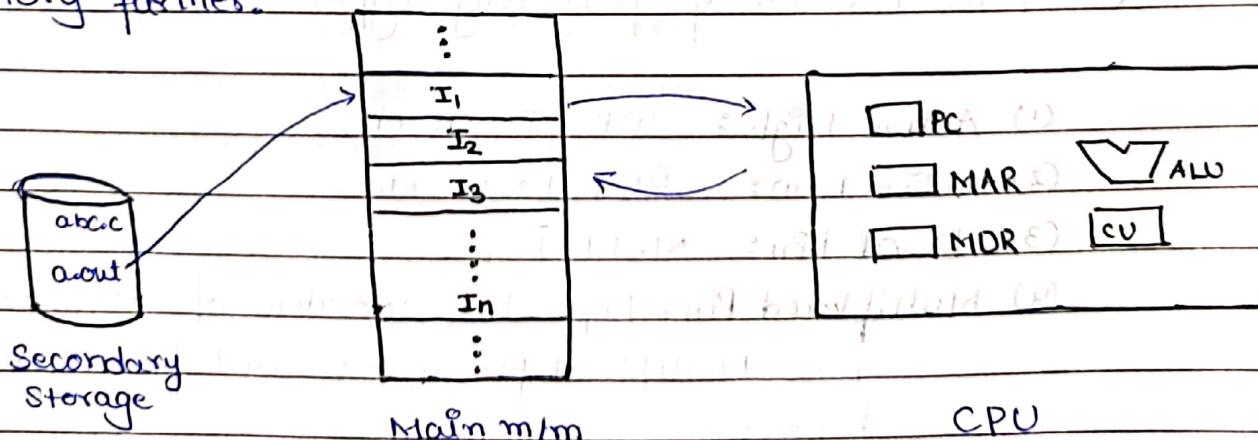
- Designer's perspective
- Hardware View
- Talks about, how to design low latency m/m, may use cache or not, etc, what instructions to provide, etc.

* Computer Organization

Story started as follows :

- (1) C Program, written by programmer
- (2) Compiler & linker converted to binary executable.
- (3) Loader loaded the program into memory.

Now CPU has access to instructions. COA continues the story further.



UNIT OUTLINE

* Program: Sequence of instructions + Data

* Memory is divided into equal sized cells called partition called cells & each cell has unique address.

• Each cell allows only two operations: Read or Write.

• CPU generates Memory Request to access data or instruction from memory.

Memory Read: <Address>, \overline{RD}

Memory Write: <Address>, WR, <value>

• We see two types of signals here: \overline{RD} & WR. \overline{RD} is active low & WR is active high.

• Pins like these are how two components can communicate with each other.

* Pins can be of following types:

(1) Active High: WR, INTR, etc.

(2) Active Low: \overline{RD} , INTR, etc.

(3) Dual Pin: MEM/ \overline{IO} ,

(4) Multiplexed Pin: Depending on value of some other pins Multiplexed pins can be used for different functions.

Ex: AD₀ ... AD₁₅

Can carry out 2 or more funcs.

* n -bits Processor: Tells that:

- (1) No. of Data lines: ~~n -bits~~ n
- (2) Data bus size: n -bits
- (3) Word Size: n -bits
- (4) ALU: n -bits
- (5) Accumulator: n -bits
- (6) Memory Data Register: n -bits
- (7) General Purpose Registers: n -bits.

Basically anything to do with data by the processor
is n -bits.

* Addressability of Memory:

(1) Byte Addressable: Each cell is of 1 Byte. OR address is given to each byte.

(2) Word Addressable: Each cell is of 1 word. OR address is given to each word.

Ex. MIM is 8 GiB. No. of address bits for

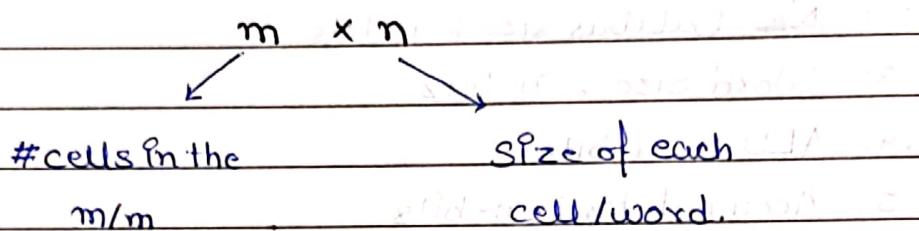
(1) Byte Addressable: 33 bits.

(2) Word Addressable ($w = 4B$)

$$\therefore \# \text{cells} = \frac{8 \text{GiB}}{4B}$$

$\therefore 31$ bits.

* Memory Chip Configuration



Ex. $64K \times 16$

$$\therefore \# \text{cells} = 64K$$

$$\Rightarrow \text{Size of address bus} = \log_2 64K \\ = 16$$

width of address bus = 16 bits

- If unit of cell size isn't mentioned, take as bit.

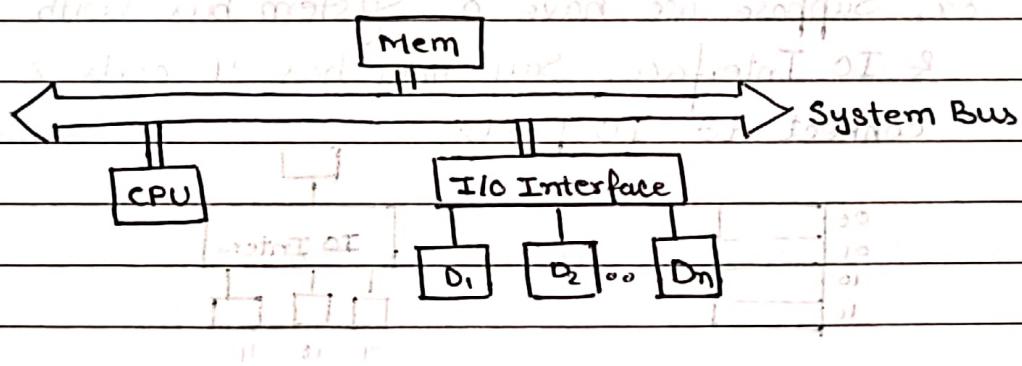
* Components of Computer System

Whenever a program is run, it would use / require three components:

- (1) CPU: Master component. Would execute the program.
- (2) Memory: Slave component. Slow. Stores program & data
- (3) I/O: Slave component. Slowest.

When we say fast / slow this is based on clock cycle time.

Faster component has lower cycle time.



* System Bus

Consists of three set of lines:

- (1) Address Lines (AL): Used to carry address from CPU to mem or I/O. Unidirectional.

Ex. 8085 has 16 bit Address Bus:

So, max addressable m/m = 2^{16} Bytes Cells

Ex. Mem is $4^{10} \times 8$

So, no. address lines = 10.

(2) Data Lines [DL]: Used to carry data from
CPU to M/M, M/M to CPU,

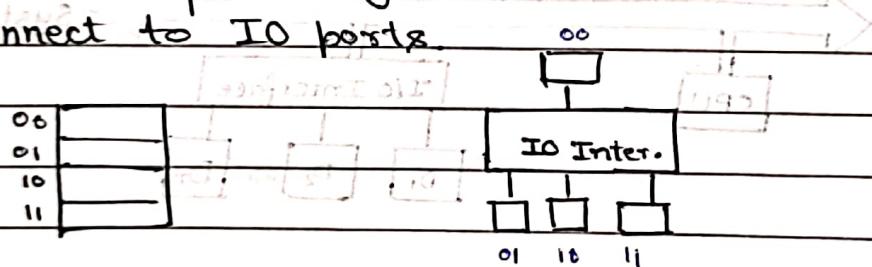
CPU to I/O & I/O/CPU.

(3) Control Lines [CL]: Used to carry control signals.

Control signals tell what type of operation
is being performed.

Memory never issues control Signal.

Ex. Suppose we have a system bus with CPU, M/M
& IO Interface. Say M/M has 4 cells & 4 IO devices
connect to IO ports.



Now if CPU issues, say a Memory read req:

01, RD

This would go on System bus & hence both to
M/M & IO Interface.

How is one to know if it is meant for IO read or
M/M read?

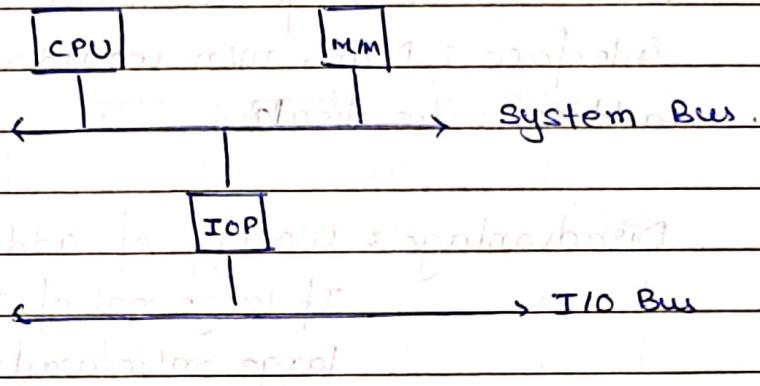
Since all three components share the same AL, DL &
CL, this problem is happening.

~~IO T hold data you can't access~~

- * There are three solutions to the problem:

(1) Separate I/O Processor

It's a separate processor which only performs I/O operation. DMA is an example of IOP.



∴ A separate bus for IO. All IO operations go through IOP.

We can use same addresses & control signal, but different buses.

(2) Isolated I/O (IO Mapped I/O)

We have separate read & write operation for IO & M/M.

IO: IORD & IOWR

M/M: MRD & MWR

So we use same ~~bus~~ system bus & same addresses but different control signals.

This is generally achieved using an additional signal

	M/I/O	RD	WR	Operation	Assembly Inst Ex.
	0	0	0	IO RD	IN R ₁ , 01
	0	1	1	IO WR	OUT 011, R ₂
	1	0	0	MM RD	LOAD R ₁ , 01
	1	1	1	MM WR	STORE 11, R ₂

(3) Memory Mapped IO:

We keep the same system bus, use same signal & use some addresses specifically for IO device addressing, and disable those addresses in M/m.

When req. meant for IO device goes to both IO Interface & M/m, m/m would reject it since the addresses are disabled.

Disadvantage: Wastage of address space in m/m.

If large no. of IO devices, then large m/m wastage.

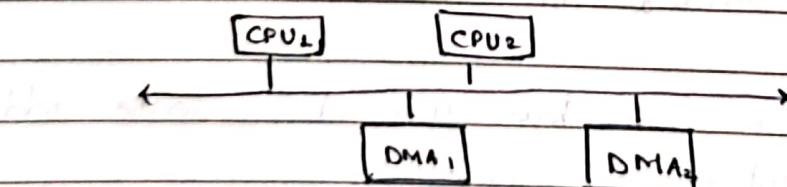
Ex. M/M has 8 bit addresses,

then m/m size = 2^8 cells = 256 cells.

Say we decide that any address of form $xyz1ab0c$ is reserved for IO port addressing.

(We disabled) $2^6 = 64$ cells in m/m.

Ex. Say we have such a system bus

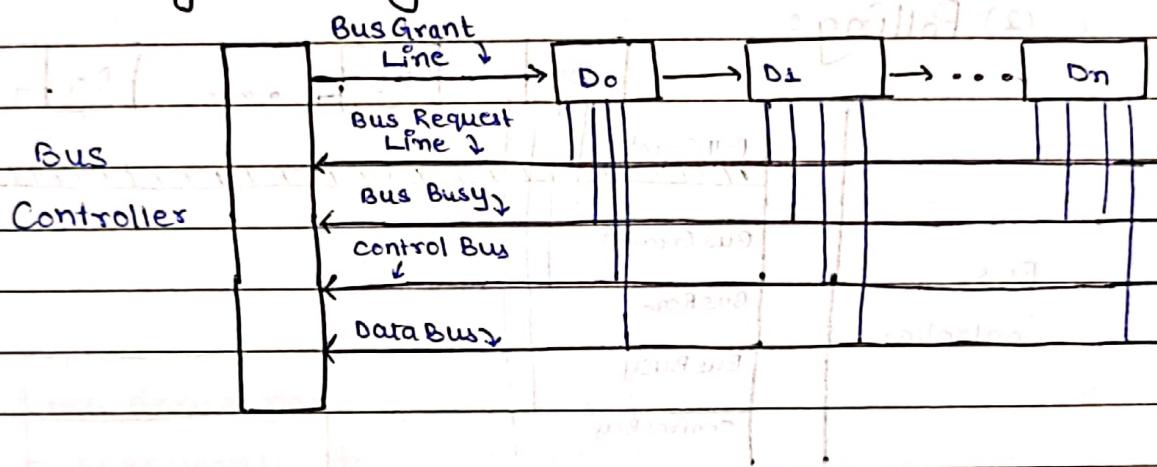


If multiple masters are trying to access system Bus for various operations at same time, then we need a method to choose the next master of system bus.

* Bus Arbitration: The process of deciding on the next master of System Bus is called as Bus Arbitration.

There are three methods of bus arbitration

(1) Daisy Chaining:



When at least one device activates BRL, the Bus Controller activates bus grant line of D₀. If D₀ made the request, it activates Bus Busy signal so that Bus controller & other devices know that it is accessing the Control & Data buses. Bus Controller grants access to D₀. If D₀ didn't make request it just activates Bus Grant Line of D₁ & the process repeats.

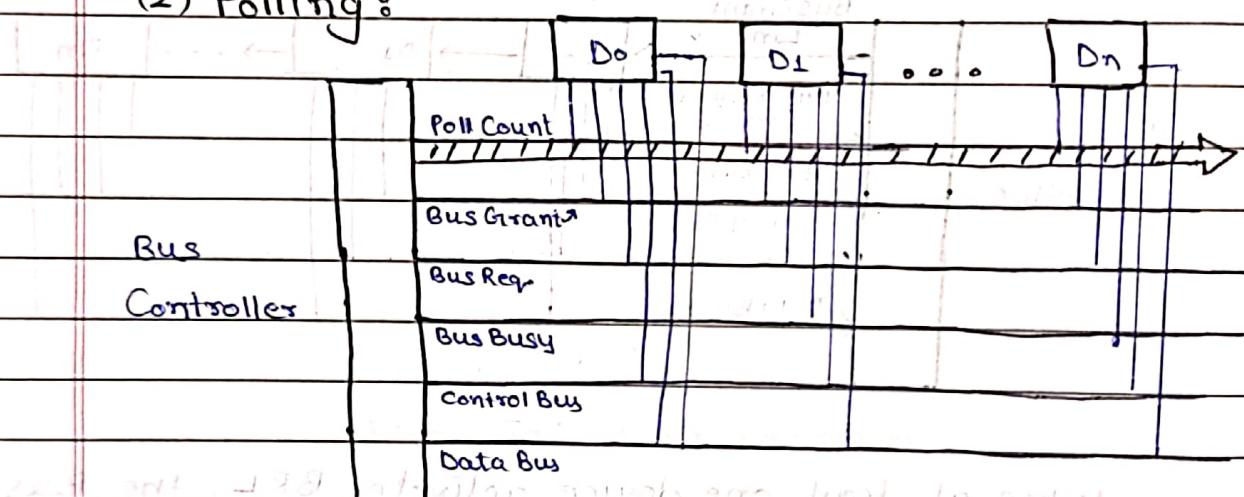
Note: Devices closer to Bus Controller get higher priority.

Disadvantages:

- (1) If one node fails whole system fails. Single point of failure
- (2) Higher unnecessary delay for devices further from Bus Controller.
- (3) Device priority is static & higher for closer devices to Bus Controller.

The task of Bus Controller also keeps shifting from one device to other. Sometime CPU can act as BC, or sometime CPU5, sometime DMA10. The control unit of the current BC decides which Master device would become BC next.

(2) Polling:



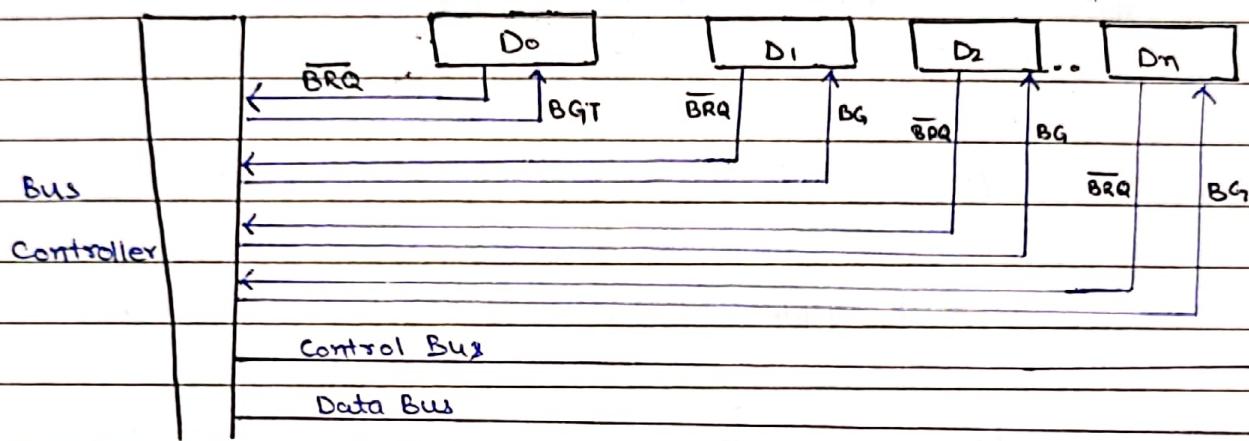
Poll Count is a K-lines/signal. $K \gg [log n]$ for n devices.

when one or more device activates Bus Request Line, the Bus controller generates Poll Count (or device index) in random order & queries if device generated the request. The devices are indexed based on dynamic priority which can change over time.

The device which did activate request line activates Bus Busy signal. Bus controller activates Bus Grant signal & grant access to Control & Data bus to that device.

So all ~~st~~ devices have fair chance & no single point of failure.

(3) Independent Bus Req. & Grant:



Each device has separate Bus Request & Bus Grant lines.

Independently they request access to system buses.

BC has a static priority for the devices based on which it chooses the next device to grant access.

BC activates the selected device's Bus Grant Line & allows access to Control Bus & Data Bus.

Instructions & Addressing Mode

* Instructions : Inside a machine, Instruction is nothing but sequence of bits which is processed by CPU.

Each Instruction is made up of some operation & operands.

Say we want to add values in register R₁ & R₃.

Human Readable form	Assembly Code	Corresponding Bits
$R_1 + R_3 \equiv ADD R_1, R_3$		$1010\ 0111$

The mandatory operational code or opcode.

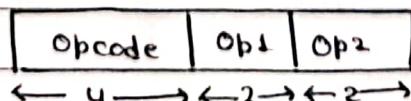
Then there are 0 or more operands.

Opcode	Op ₁	Op ₂	...	Op _n
--------	-----------------	-----------------	-----	-----------------

Ex. Say opcode field in a particular CPU is 3 bits, then total $2^3 = 8$ operations can be supported & we can encode whichever operation to whichever opcode. One possible assignment is

ADD: 000	LOAD: 010	MUL: 100	INCR: 110
SUB: 001	STORE: 011	DIV: 101	DECR: 111

Ex.



So, $2^4 = 16$ operations

$2^{4+2+2} = 256$ instructions possible

With each operation, $2^{2+2} = 16$ instr. are possible.

about parallel & concurrent

Ex. Instruction size = 20 bits M/M size = 256 Bytes

2 addr instructions. # operations?

(Assuming 256x8 m/m)

$$\text{Addr size} = \lceil \log_2 256 \rceil = 8 \text{ bits}$$

$$\therefore \# \text{opcodes} = 20 - 8 - 8$$

bits \Rightarrow 4 bits available for opcodes

$$\therefore \# \text{opcodes} = 2^4 = 16$$

and 3 bits for address bits left. Then, what?

Ex. Inst Size = 32 bits #operations = 256? Only 3 addr

inst. Size of m/m?

28 bits

8 bits

avail

$$\text{Rem. bits} = 32 - \lceil \log_2 256 \rceil = 24$$

$$\text{per addr} = 24/3 = 8$$

$$\therefore \text{m/m size} = 2^8 = 256 \text{ B (or cells)}$$

Inst. [28 bits] [8 bits] [8 bits]

at least 3 256 cells required or 8 bits address plus 8

8 bits for address + 8 bits for 256 cells makes

24 bits for address + 8 bits for 256 cells makes

32 bits for instruction + 8 bits for 256 cells

[28 | 8 | 8]

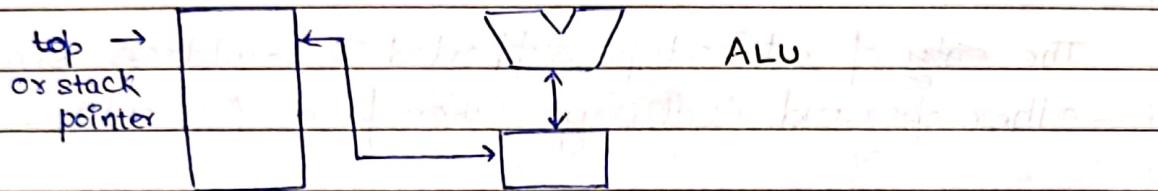
8 bits for address + 8 bits for 256 cells

and address + 8 bits for 256 cells

and address + 8 bits for 256 cells

* Instruction Format (Based on M/C Type)

• Stack Organization:



At each cycle, the top of stack is moved to accumulator, then operation is performed by ALU & moved to accum. Then pushed to stack top again.

(1) PUSH Acc → Stack Top

(2) POP Stack Top → Acc.

(3) Operations such as ADD, INCR, etc.

All stack org. instr's are 0-address.

Ex. POP

pops two elements from Top of Stack

POP adds them & pushes result back to stack.

PUSH

• Accumulator Organization:

Only one register is available: Accumulator.

One operand is available in Accum. always & the other

operand is in memory.

After the operation completes the result is stored in the accumulator.

This is also called LOAD & STORE machine.

The ~~only~~ operation type supported is 1-address, since other operand is always taken from Accumulator.

Ex. LOAD 3000 Acc \leftarrow [3000]

ADD 2000 Acc \leftarrow Acc + [2000]

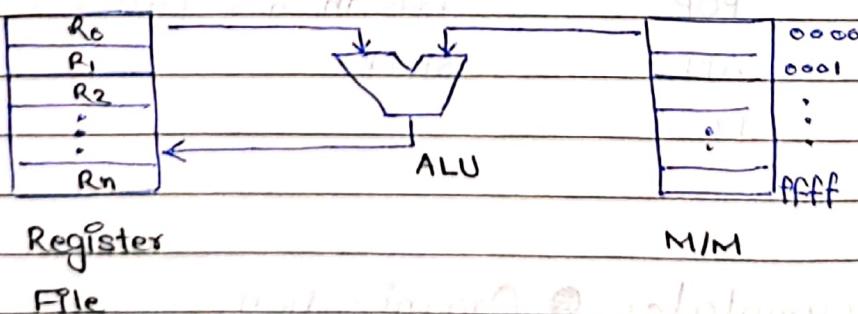
STORE 3002 3002 \leftarrow Acc

(Zero addr. operations are also possible).

General Register Organization:

(1) Register to M/M organization

Small set of registers. One operand is from register and another from the memory. Result would be stored in the register.



This type of machine supports upto 2-addr. instruction.

Ex. LOAD R₀, [2000] R₀ ← [2000]

STORE [3000], R₁ 3000 ← R₁

ADD R₂, [2000] R₂ ← R₂ + [2000]

(2) Register To Register Orgn:

- Both operands are register locations and the result is also stored in register.

- Processor supports more no. of registers.

- Since #registers is high, operands can be prefetched as well.

- It supports upto 3-address instructions.

Ex. LOAD R₁, 3000 R₁ ← [3000]

LOAD R₂, 2000 R₂ ← [2000]

ADD R₃, R₁, R₂ R₃ ← R₁ + R₂

STORE 3002, R₃ 3002 ← R₃

- Fastest of all orgn since operands are in registers.

* Instruction Format (Based on number of operands)

- Zero Addr:

opcode

- One Addr:

opcode	Op ₁
--------	-----------------

- Two Addr:

opcode	Op ₁	Op ₂
--------	-----------------	-----------------

- Three Addr:

opcode	Op ₁	Op ₂	Op ₃
--------	-----------------	-----------------	-----------------

Operand can be mem addr., register addr., or constants.

Ex. Instⁿ size : 20 bits. 64 reg. MM = 256 B.

If one addr is reg. addr, other is m/m addr. No. of
2addr opcodes supported?

$$\text{m/m addr size} = \log_2 256 = 8$$

$$\text{reg. addr size} = \log_2 64 = 6$$

$$\therefore \text{Opcode bits} = 20 - 8 - 6 = 6$$

$$\therefore \# \text{opcodes (2addr)} = 2^6 = 64$$

* There are two types of machines:

(1) Same Length Inst

The size of instruction is fixed.

Ex. Say there's a machine with upto 3-addr instructions,
with Mem. addr of 32 bits, and 256 operations.
then addr would be like

Opcode	Addr ₁	Addr ₂ A ₃
←8	→32	→32→A ₃

and size of instⁿ is 104 bits.

Even if we have a 0 addr instⁿ out of the 256 operations, say CMA or INCRA etc, we'd still have to use 104 bits to store it in the memory.

Disadvantage: Memory wastage.

(2) Variable Length Instructions :

We may have some 1-addr, some 2-addr, some 0-addr instructions that we want our machine to perform.

Using variable length instⁿ would create tradeoff b/w opcode bits and address bits.

Expanding Opcode Technique: This basically means, we are making the opcode field length variable & not the instⁿ length.

Ex: Say we have 16 bit instⁿ & 4 bit addr, then we can vary opcode field size as per no. of addresses:

Opcode	Addr ₁	Addr ₂	Addr ₃	2 ⁴ 3-addr ops possible
↔ 4 ↔				

Opcode	Addr ₁	Addr ₂	2 ⁸ 2-addr ops possible.
↔ 8 ↔			

Opcode	Addr ₁	Addr ₂	2 ¹² 1-addr ops
↔ 12 ↔			

Opcode	Addr ₁	Addr ₂	2 ¹⁶ 0-addr ops
↔ 16 ↔			

Ex. Possible Design.

0000	xxxx	yyyy	zzzz	15 3-addr inst ⁿ
0001	xxxx	yyyy	zzzz	
:				
1110	xxxx	yyyy	zzzz	

1111	0000	xxxx	yyyy	13 2-addr inst ⁿ
:				

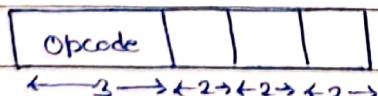
1111	1100	xxxx	yyyy	(3x16)=48 1-addr inst ⁿ
1111	1101	0000	xxxx	
:				
1111	1111	1111	0000	

1111	1111	1111	0000	10 0-addr inst ⁿ
1111	1111	1111	1010	

The CPU would be designed to recognize the following:

- (1) If first 4 bit are in range (0000 to 1110), then interpret as 3-addr Instn.
- (2) If first 4 bit in range (1111 to 1111) & the next 4 bit in range 0000 to 1100, then interpret as 2-addr Instn.
- (3) If first 4 bit in range (1111 to 1111), next four bit in range (1101 to 1111), and next 4 bits in range then 1-addr Instn.
- (4) Else 0-addr Instn.

Ex. Instn Size = 9 bits. Addr Size = 2 bits. 2 3-addr Instn, 20 2-addr Instn, and 8 1-addr Instn are required to be designed. How many 0-addr Instn are possible?



$$\# \text{possible 3-addr Instn} = 2^3 = 8$$

$$\# \text{Req.} - 11 = 2$$

$$\text{Remaining} = 6.$$

$$\# \text{possible 2-addr Instn} = 8 \times 2^2 = 24$$

$$\# \text{Req.} - 11 = 20$$

$$\text{Remaining} = 4$$

$$\# \text{possible 1 addr Pinst}^n = 4 \times 2^2 = 16$$

$$\text{Req. } 16 - 8 = 8$$

$$\# \text{possible 0 addr Pinst}^n = 8 \times 2^2 = 32$$

$$\therefore 32$$

Ex. Is following config possible:

Used = 14 2-addr Pinst

Used = 127 1-addr Pinst

Used = 60 0-addr Pinst

Opcode	Op1	Op2
← U →	← 6 →	← 6 →

$$\# 2 \text{addr possible} = 2^4 = 16$$

$$\text{Used} = 14$$

$$\text{Remain} = 2$$

$$\# 1 \text{addr Pinst possible} = 2 \times 2^6 = 128$$

$$\text{Used} = 127$$

$$\text{Remain} = 1$$

$$\# 0 \text{addr possible} = 1 \times 2^6 = 64$$

$$\text{Remain} = 64 - 60$$

$$\Rightarrow 4$$

∴ Yes possible.

* Addressing Modes

Instruction Cycle has two phases :-

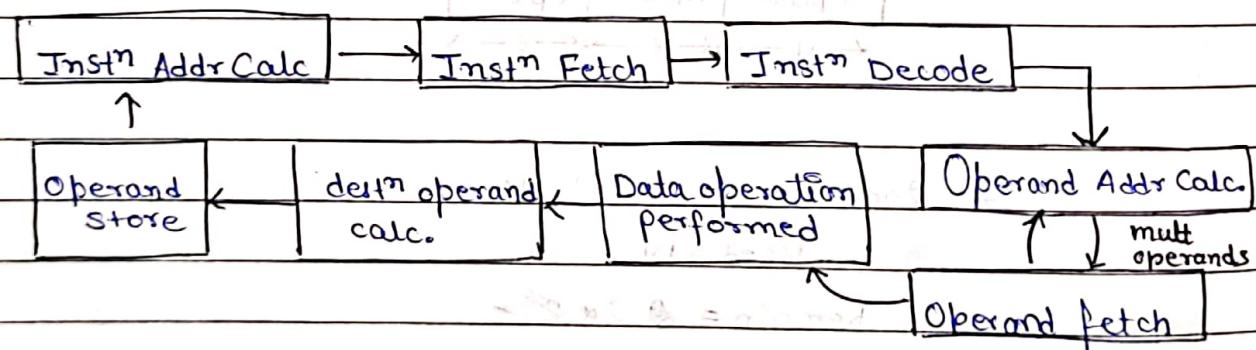
- Fetch Cycle
- Execution Cycle.

The first step during execution would be to calculate the address of $inst^n$, then we can go ahead & fetch $inst^n$ from m/m.

Once inst^n is in Inst^n register, we can decode it.

After decoding the operation, we require operands.

This requires an iterative cycle of operand addr. calculation & operand fetch for each operand.



So the addressing mode dictates where the operand would be found.

An operand can be found in four places:

- (1) Opcode: INCA, (Increment Accumulator), CLRC (Clear Carry), etc.
 - (2) Instruction: ADDA #0FFF (Add accum.) , etc.
 - (3) Registers: ADD R₁, R₂
 - (4) Memory: STOREA [3000]

- Any inst^n may use combination of multiple add. modes.
 - Some bits in the inst^n format are reserved for specifying addressing mode, per operand.

MOV: Used to move data b/w registers. It copies.

LOAD: Used to read from m/m.

STORE: Used to store data to m/m.

CLASSMATE

Date _____

Page _____

• Implied / Implicit Addressing Mode.

Operand is specified within the opcode. No effective address calculation. 0-addr instructions.

Ex. INCRA, CLRC, CLRA, etc.

• Immediate Addressing Mode:

Operand is available inside the instruction.

Sometimes specified using I at end of opcode.

Ex. ADD + I \Rightarrow ADDI

Ex. ADI 1000 adds value 1000 to Accum. Value.

Sometimes immediate values are prefixed with #.

Ex. ADA #1000 adds value 1000 to Accum. value.

Ex. MVI R1, 1000

• Register Addressing Mode.

Operand is available inside some register and the address of that register is specified in the instn.

Eff. Address = Reg. Name / Reg. Address

Ex. MOV R1, R2 $R_1 \leftarrow [R_2]$

• Direct (Absolute) Addressing Mode

Operand is stored in some memory location & the address of the operand is provided inside the Instⁿ.

$$\text{Eff. Addr.} = \text{Mem. Addr in Inst}^n$$

Ex. ADD 2000, 3000 EA = 2000, EA = 3000

$$[2000] \leftarrow [2000] + [3000]$$

∴ 2 m/m access during operand fetch

& 1 m/m access during write back.

• Indirect Addressing Mode

Used in case of pointers in High Level language.

Operand is stored in some m/m loc x & the address x is stored in another m/m location y.
The Instⁿ contains the value /x.

$$\text{Eff. Addr} = [x] = y.$$

Requires two memory accesses.

Ex. ADA @ 1000 say $[1000] = 5000$.

$$\text{then } EA = [1000] = 5000$$

$$\text{Acc} \leftarrow [\text{Acc}] + [[1000]]$$

$$\text{Poe. Acc} \leftarrow [\text{Acc}] + [5000]$$

• Register Indirect Addressing Mode

The operand is stored in some m/m loc. x , the addr. x is stored in some register & the instⁿ contains the register name/address.

$$\text{Eff. Address} = [R_n]$$

Ex. ~~AD~~ LOAD R₁, @R₁ & $[R_1] = 1001$

$$EA = [R_1] = 1000$$

$$R_1 \leftarrow [R_1]$$

$$R_1 \leftarrow [1001]$$

One Register Access & one m/m access.

• Indexed Addressing Mode

Used for Array access. Can't use array construct without this AM.

The Base Address of the array is stored in a special register: Index Register. In the instⁿ, the offset is provided.

$$EA = [IR] + \text{offset}$$

Ex. LOAD 16(R_n)

$$EA = [R_n] + 16$$

$$Acc \leftarrow [R_n] + 16$$

• Program Counter Addressing Mode

The base addr is taken from Program Counter.

In the instⁿ, the offset only is provided.

(aka. PC Relative)

(aka. Relative)

At the time the $inst^n$ is decoded, the Program Counter is already moved to the next $inst^n$ location.

Ex. LOAD -20

would result in:

$$EA = [PC] - 20$$

$$= 1002 - 20$$

$$\Rightarrow 982$$

Load Imm. offset	LOAD -20	1002

$$Eff. Address = [Program Counter] + offset$$

- Base Index Addressing Mode

$$Eff. Address = [Base Reg.] + [Index Register]$$

Ex. LOAD [Bx][Ri]

$$EA = [Bx] + [Ri]$$

This requires two register access to calculate the effective address, and one m/m access to get operand.

- Base-Index w/ Displacement AM

$$Eff. Address = [Base Reg] + [Index Reg] + Displ.$$

Ex. LOAD 20[Bx][Ri]

$$EA = [Bx] + [Ri] + 20$$

Also req. 2 register access & one m/m access

- Register Auto Incr./Decr. (AIM)

(1) POST

$$EA = [\text{Reg. Name} / \text{Addr}]$$

$$\text{Reg} \leftarrow [\text{Reg}] \pm \text{Inst Size.}$$

(2) PRE

$$EA = [\text{Reg}] \pm \text{Inst. Size}$$

$$\text{Reg} \leftarrow [\text{Reg}] \pm \text{Inst. Size.}$$

Ex. LOAD + (R₂) Inst size = 2.

$$EA = [R_2] + 2$$

$$R_2 \leftarrow [R_2] + 2$$

Ex. STORE (R₁) - Inst size = 4

$$EA = [R_1]$$

$$R_1 \leftarrow [R_1] - 4$$

So, one register access, one arithmetic operation & one m/m access.

- Base Addressing Mode.

$$EA = [BA] + \text{offset.}$$

^{Imp} This & PC relative modes are used to write Position Independent Code.

Position Ind. Code \rightarrow think PC relative, then Base Addr. Mode.

* Instruction Cycle

There are four cycles, two of which are optional:

(1) Fetch Cycle: The instruction is fetched from m/m based on PC value.

(2) Indirect Cycle [optional]: Required when using some form of indirection & the effective address has to be brought into the Instⁿ, such as in Indirect AM, Indexed AM, PC Relative AM, etc.

(3) Execution Cycle: Instruction is decoded, operands are fetched, operation is executed & result is written back.

(4) Interrupt Cycle [optional]: Executed when there's an interrupt to the CPU. There's an INTR line of CPU, which is activated by IO & m/m. After each instruction, CPU checks if this line is high, it runs this cycle.

(1) Acknowledges the interrupt. INTA

(2) Saves state

(3) Jump to ISR.

(4) Run ISR

(5) Jump back & resume state.

Fetch Cycle

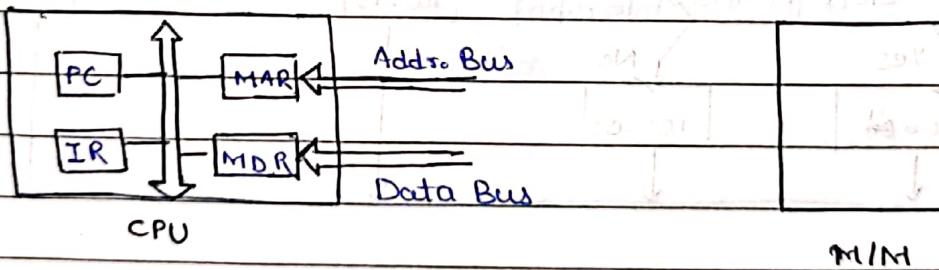
Instⁿ are stored in the m/m. CPU can't act of m/m directly, so it needs to be fetched to CPU.

Program Counter: Stores the address of the next ~~reg~~ instⁿ to be executed.

Instruction Register: Stores the instⁿ being executed.

Memory Address Register: Used to communicate address b/w system bus & internal bus.

Memory Data Register (Mem. Buffer Reg): Used to communicate data b/w system bus & internal bus.



(1) Transfers addr. of next Instⁿ from PC to MAR.

MAR floats this addr on system bus & m/m read signal is issued. $PC \rightarrow MAR, RD$

(2) The value from the m/m location is loaded into the MBR. At the same time (in the same CPU cycle) the value of PC can be incremented.

$$MBR \leftarrow Mem[MAR], PC \leftarrow PC + 1$$

(3) The value from MBR is transferred to IR.

$$IR \leftarrow MBR$$

- How does CPU know which cycle is executing & which control signals to generate?

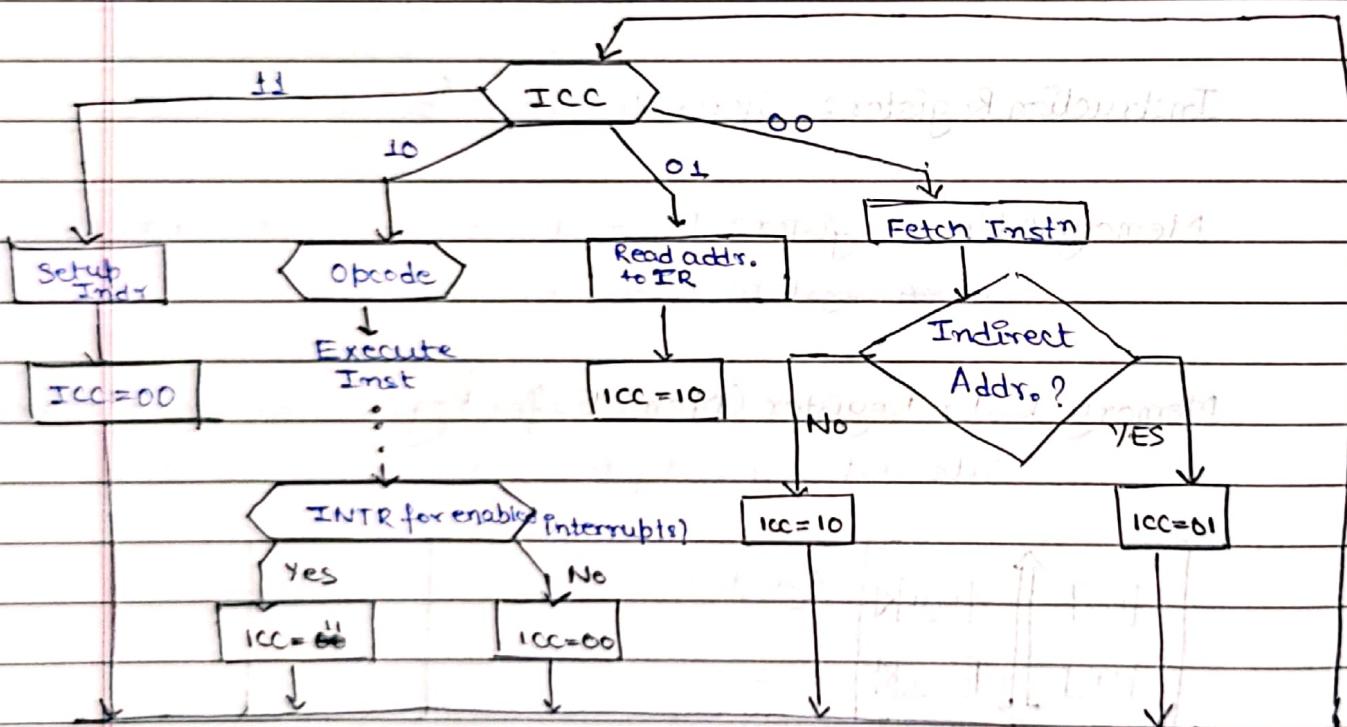
CPU has Instruction Cycle Code which has four values:

00 → fetch

10 → execute

01 → indirect

11 → interrupt.



- Indirect Cycle:** If one or more operand uses indirect addressing, then this cycle is executed. It checks the addr. mode of operands incrementally & if indirect, brings the effective addr. in the instr.

MAR ← Addr. part of IR

MBR ← Mem [MAR]

Addr. part of IR ← MBR.

Ex. LOAD @ 1000

MAR \leftarrow 1000

MBR \leftarrow [1000]

IR[Mem Addr] \leftarrow MBR

LOAD 2052.

2052	1000
20	2052

Note: The m/m cycles to get eff. address are included in Execution Cycle only unless stated otherwise.

Execution Cycle: The following steps can take place during this cycle.

(1) Instruction Decode

(2) Operand Fetch

(3) ALU operation

(4) Write Back.

Ex. ADD R₁, X

t₁: MAR \leftarrow X, RD

t₂: MBR \leftarrow Mem[MAR] No write back.

t₃: R₁ \leftarrow [R₁] + [MBR]

Instⁿ Decode happens by decoding the opcode using a n: 2^m decoder, based on which sequence of microinstⁿ is generated.

Ex. JMP d0(R₁)

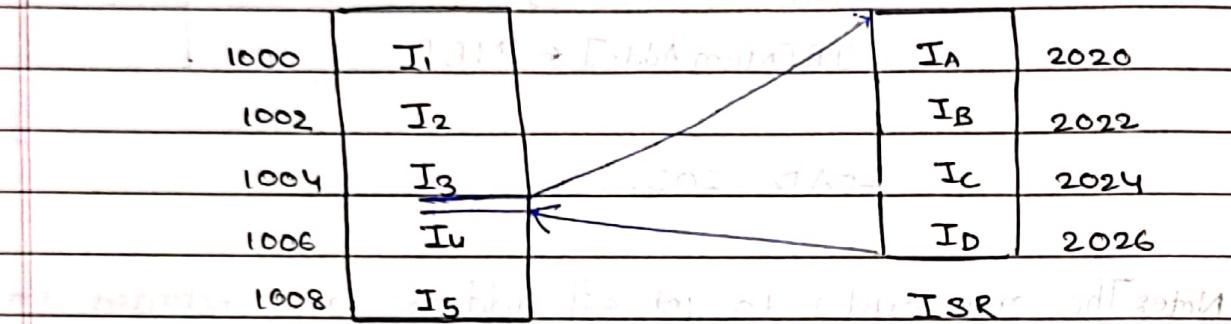
t₁: MAR \leftarrow [R₁] + 20, RD No ALU opⁿ or WB.

t₂: MBR \leftarrow Mem[MAR]

t₃: PC \leftarrow MBR

• Interrupt Cycle

Ex. Say the prog. $\langle I_1, I_2, I_3, \dots \rangle$ is executing.



Say I_2 is executing when an enabled interrupt comes, e.g., CPU's INTR line is activated.

At this point $PC = 1006$.

Once CPU completes executing I_2 it checks if an enabled interrupt has been activated.

- First the program state is stored/saved in m/m.
- Then PC is loaded with base addr. of ISR.
- Then ISR is run.
- Once finished the PC value is restored.

$t_1: MBR \leftarrow PC$

$t_2: MAR \leftarrow \text{Save_addr}, WR$

$t_3: PC \leftarrow \text{ISR BA}$

* Interrupts

Interrupts can be classified on various basis.

• Vectored & Non-Vectored Interrupt

The interrupts whose ISR is located (stored) in a fixed memory location are called Vectored Interrupts.

Non-Vectored interrupt comes or supplies the base address of the ~~ISR~~ ISR, i.e. the device supplies ISR addr. along with INTR signal while raising interrupt.

NV interrupt allows for the device to supply different ISR addresses.

• Software & Hardware Interrupt

Hardware interrupts come through hardware lines connected to the CPU.

Software interrupts are raised by the program, such as Divide by 0, or Segmentation fault, etc.

• Maskable & Unmaskable Interrupt

Maskable interrupts mean that they can be disabled and enabled.

Non-maskable interrupts can't be disabled.

- Interrupt Service Routines are part of OS & stored in Kernel space of m/m.
- Non-maskable interrupts have higher priority compared to maskable interrupts.
- Interrupt Handling Procedure:
 - (1) Save the PC value, flags & register values on the m/m stack.
 - (2) Load PC with beginning address of ISR.
 - (3) Finish ISR when RET inst. is executed.
 - (4) Return to the point of execution by restoring PC, flag & register state.

* Instruction Set

made

Set of all the instructions[↑] available to the programmer by the designer & performable by the CPU.

• Data Transfer Instⁿ

Used to transfer data (copy)

- (1) from Reg to Reg
- (2) from Reg to M/m
- (3) from Mem to Reg.
- (4) from Mem to Mem.
- (5) ~~from~~ Immediate to Mem
- (6) Immediate to Reg.

(a) MOV Rd, R_s

Move

$$Rd \leftarrow [Rs]$$

(b) MVI Rd, Imm

Move Immediate.

$$Rd \leftarrow Imm$$

(c) LDA Addr

Load to Accumulator

$$Acc \leftarrow [Addr]$$

(d) STA Addr

Store Acc. value to m/m addr

$$Addr \leftarrow [Acc]$$

(e) XCHG R₁, R₂

Exchange values stored in reg.

• Arithmetic Instⁿ:

Affects all flags: Sign, Carry, Aux Carry, zero, Parity, Overflow

(a) Add R, M

$$R \leftarrow [R] + [Mem]$$

(b) Add R

$$Acc \leftarrow [Acc] + [R]$$

(c) Add M

$$Acc \leftarrow [Acc] + [Mem]$$

(d) ADC R, M Add with Carry

$$[R] \leftarrow [R] + [M] + [C]$$

(e) ADI Imm.

$$Acc \leftarrow [Acc] + Imm$$

(f) ACI Imm.

$$Acc \leftarrow [Acc] + Imm + [C]$$

(g) SUB R, M

SUB R

SUB M

(h) SBB R, M Subtract w/ Borrow

(i) SUI Imm

(j) SBT Imm.

(k) INR R

Increment

INR M

(l) DCR R

Decrement

DCR M

• Logical Inst: (AND, OR, XOR, ROTATE, COMPARE, COMPLEMENT, SHIFT)

Affects Sign, Zero & Parity flags.

(a) ANA R/M AND accumulator

(b) ANI Imm. AND acc. with imm. value.

(c) ORA R/M

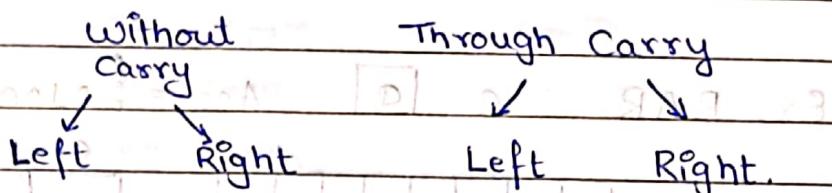
(e) XRA R/M

(d) ORI Imm.

(f) XRI Imm.

(g)

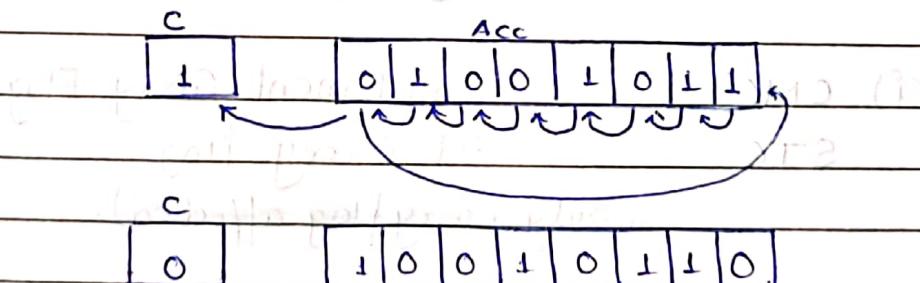
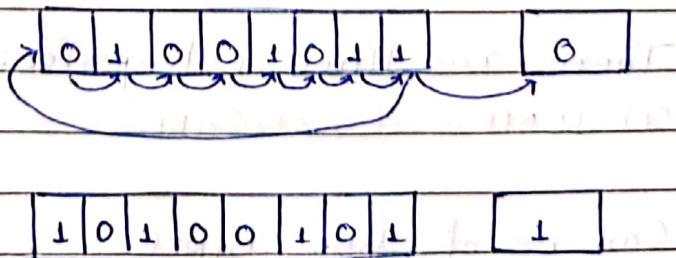
Rotate



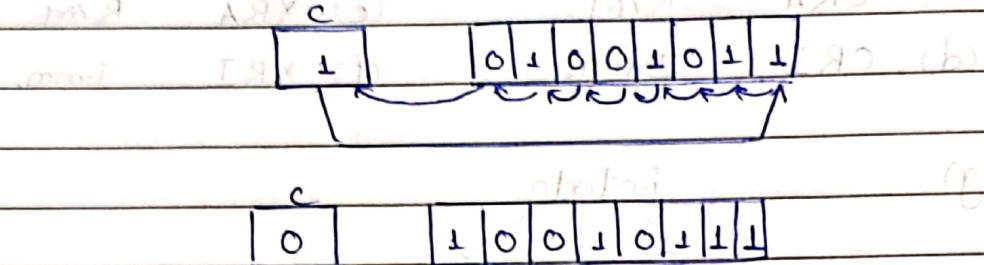
RLC: Rotate Acc. Left w/o Carry

RRC: Rotate Acc. Right w/o Carry.

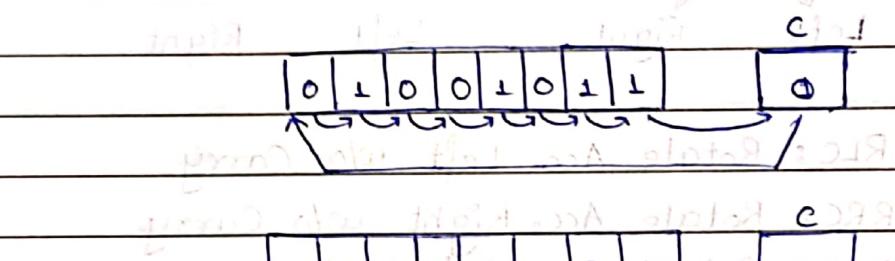
RAL: Rotate Acc. Left with Carry

~~RAR~~ RAR: Rotate Acc. Right with Carry.Ex. RLC C = 1 Acc = 01001011.Ex. RRC C = 0 Acc = 01001011.

Ex. RAL C = 1 Acc = 0100 1011.



Ex. RAR C = 0 Acc = \$01001011.



(h) CMA Complement Accumulator
(No flags affected).

(i) CMC Complement Carry Flag.
STC Set Carry flag
(Only carry flag affected).

• Branching Insts:

There are three categories:

- (1) JUMP
- (2) CALL
- (3) RET.

Can be of two types:

(1) Conditional

(2) Unconditional.

(A) Jump

(a) JMP addr (Unconditional Jump)

(b) J X addr

X here can be flag checks, like

JC: Jump if carry

JNC: Jump if not carry

JZ, JNZ

(9-1) JP, JNP;

(9-1)

(121)

JS, JNS

(Sometimes BUN: Branch Uncond.)

(B) Call

(a) CALL addr

(unconditional)

When this instn executes, it first stores the "return address" on the program stack & when RET instn executes, it pops the "return addr" & sets PC to it.

(b) CX addr

X here can be flag checks.

Prog. Stack = Call Stack = Runtime Stack.

(C) Return

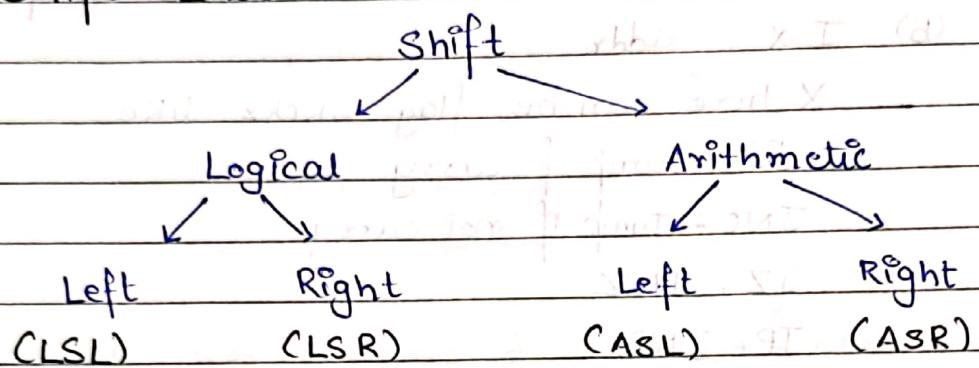
(a) RET

Pops the "return addr" from stack & sets PC to this value.

(b) RX

~~classmate~~

Shift Instⁿ



In Arth. Shift the sign bit mustn't change.

Ex. $A = \boxed{10000100}$ (d)

(a) ASI

Not Allowed. Sign bit changes.

(b) ASR

$\boxed{11000010}$

Arth. Shift the sign bit mustn't change.

Control Unit

- Every Program is a sequence of instructions. Each instⁿ has two cycles (in general), fetch & execute. Both the cycle require some microinstⁿ to run.

Say, fetch cycle:

- $t_1: M_{AR} \leftarrow PC, RD$, ~~constitutes fetch stage~~
- $t_2: M_{BR} \leftarrow Mem[M_{AR}], PC \leftarrow PC + 1$
- $t_3: IR \leftarrow M_{BR}$

For execution cycle, the microoperations depends on the opⁿ.

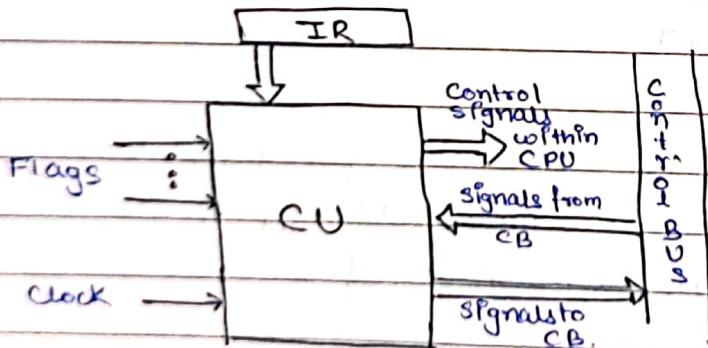
- CPU majorly consists of these components:
 - ALU
 - Register
 - Internal Data Path
 - External Data Path
 - Control Unit.

- Control Unit is responsible for generation of control signals which make the microinstⁿ happen. It is also responsible to maintain the sequence of execution.

Eg. For fetch cycle, equivalent control signals would generate in sequence:

- $t_1: PC_{out}, M_{ARin}, MEMRD$
- $t_2: M_{BRin}, PC-Incr$
- $t_3: IRin, M_{BRout}$

- Control Unit keeps track of current cycle of Instⁿ based on Instruction Cycle Code register. Based on the value of ICC, the microoperations are run, by CU generating control signals in sequence.
- Microoperations fall into one of following categories:
 - Transfer data from one reg to another.
 - Transfer b/w internal reg. to external interface.
 - Perform ALU operation using registers for inp & output.
- CU performs two basic tasks:
 - Sequencing:** The CU causes the processor to step through a series of micro-operations in the proper sequence, based on program being executed.
 - Execution:** The CU causes each popⁿ to be performed.



From the IR opcode field, the decoder decodes which operation is to be performed & respectively CU generates signal.

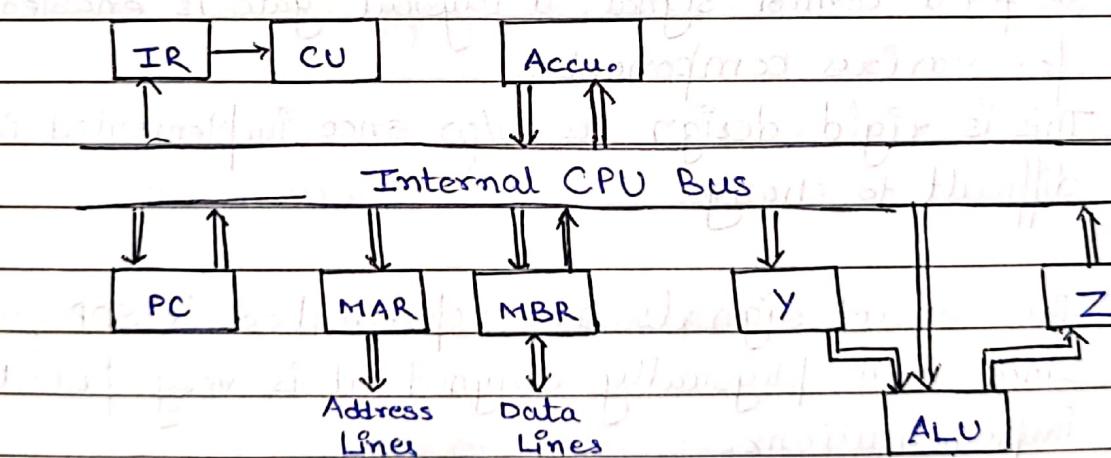
Flags allow for check values for Instⁿ like JNZ, CZ, etc.

Control bus supplies signals such as INTR, HOLD, etc.
A receives control signals for mem & IOs such as RD, IN, OUT, etc.

Control Signals within CPU allows for data transfer & ALU functions.

Clock simply generates pulses based on which microinstⁿs are run.

* Processor Organization



- A buffer (temp) register Z is required to store result of ALU, otherwise if ALU result is directly allowed to float on Internal bus, then ALU would again take it as input.

Ex: ADA 2002H

t₁: PCout, MARin, Yin, RD, INCR, Zin

t₂: MBR_(ext)in, ADD, Zout

t₃: MBR_(ext)out, TRin

t₄: Zout, Accin.

t₅: MBR_(ext)in, RD

t₆: MBR_(ext)in, Accout, Yin

* Control Unit is generally implemented in one of the two ways:

(a) Hardwired

(b) Microprogrammed.

* Hardwired Implementation

Was used in earlier systems. Suitable for RISC archs.

Control Signals are directly implemented by h/w circuits. So for a control signal, a physical gate is enabled disabled for various components.

This is rigid design, as h/w once implemented is very difficult to change.

The control signals are represented in SOP form. Since it is physically designed, it is ~~very~~ fastest of all implementations.

Ex. Designing a hypothetical CPU with 3 instructions & 4 control signals using Hardwired CU design.

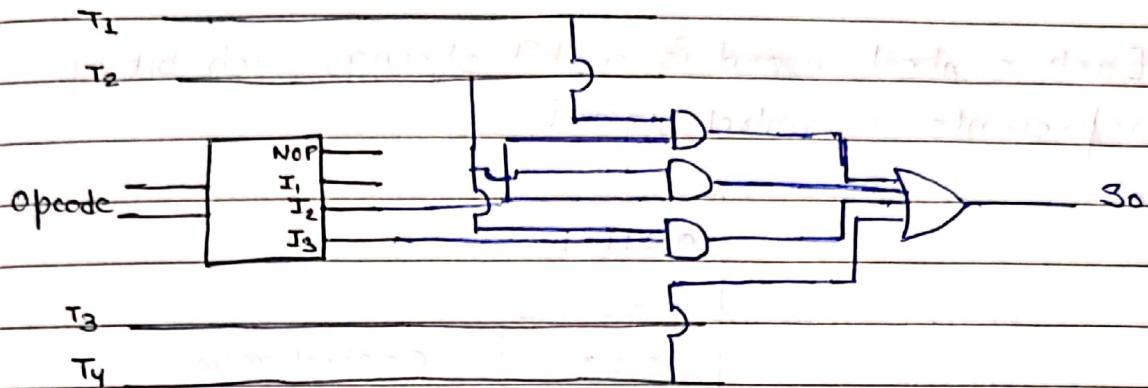
	T ₁	T ₂	T ₃	T ₄
I ₁	S ₁ , S ₃	S ₁ , S ₂	S ₁ , S ₃	S ₀ , S ₁
I ₂	S ₀ , S ₃	S ₀ , S ₂	S ₁ , S ₃	S ₀ , S ₂
I ₃	S ₀ , S ₁	S ₀ , S ₁	S ₁ , S ₃	S ₀ , S ₃

$$S_0 = T_1 I_1 + T_2 (I_2 + I_3) + T_4$$

$$S_1 = T_1 I_1 + T_2 (I_1 + I_3) + T_3 + T_4 I_1$$

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

$$S_3 = T_1 + T_3 + T_4 I_3$$



Generally a counter is used to generate control steps: T_1, T_2, T_3, \dots . Here, since there are four states, a mod-4 counter can be used, then output of the counter & output of opcode decoder can be used to realise control Signal SOP_2 .

* Microprogrammed CU

In Microprogrammed CU, the logic of the control unit is specified by a microprogram.

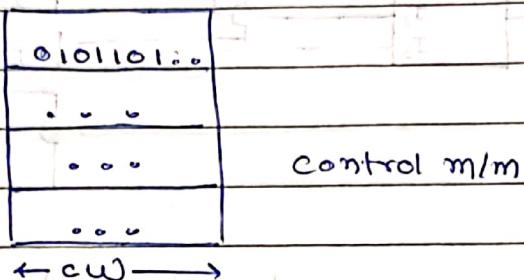
A microprogram consists of a sequence of instns in a microprogramming language. These are simple instns that specify microoperations.

Suitable for CISC architectures.

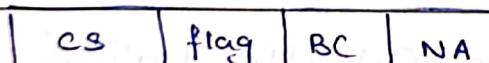
Microprogram resides somewhere b/w Hardware & S/W.
Firmwares are written in microprogramming language.

There is a control memory which stores Control Words.
We have Control Address Reg & Control Buffer Reg. &
Program Counter.

Each control word is a bit string, each bit of which represents a control signal.



- Control Word: Control Word consists of various fields:
 - Control Signals
 - Flags
 - Branch Condition
 - Next Address



- Horizontal pprg. CU: In this the control signals in CW are in decoded form, i.e. for each control signal there is one bit indicating active/inactive in the CW.

Ex. 64 Control Signals.

64 bit for control signals
per in CW.

This uses a lot of space for each CW since 1bit/control signal.

- Vertical μprog. CU: In this the control signals are stored in the ControlWord in encoded form.

Ex. 64 Control Signals

$$\therefore \lceil \log_2 64 \rceil = 6 \text{ bits only}$$

(for control signals in the CU).

Since the Control Signals are stored in encoded form, a decoder is also required to decode all the control signals.

cs	Flags	BC	Next Ad.	
		bcadr		
DCDR		return		DATA

This saves space but requires time.

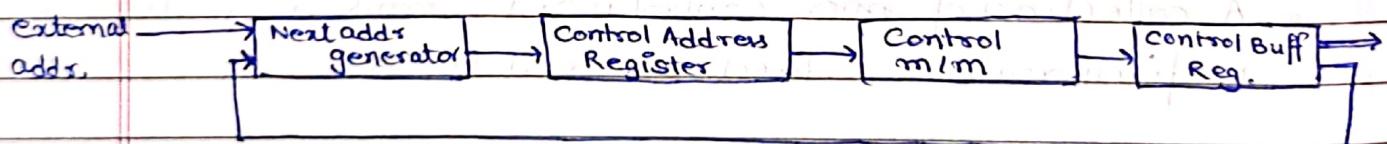
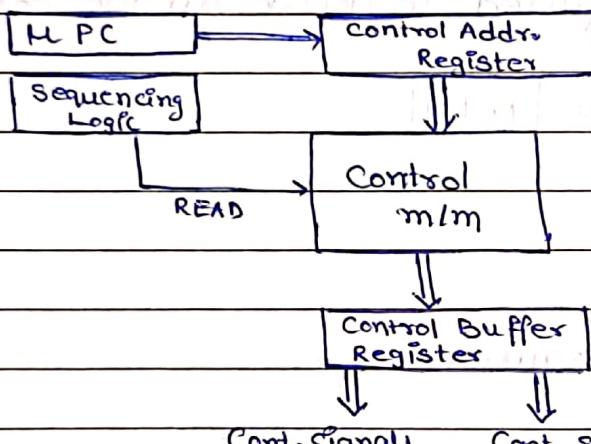
- If some branch condition is true, then the next address field is used to fetch the next control word from the control m/m. It is transferred to μPC.
- A collection of ~~the~~ control words forms a routine, such as fetch routine, indirect cycle routine, etc.
- Each microinstⁿ / cw is associated with a single CPU clock cycle, say a Instⁿ takes 4 clock cycles. It means there are total 4 control words associated with the Instⁿ.
- Since Horizontal μprog CU doesn't require decoding it is faster than Vertical μprog. CU.

Control Memory

The set of microinstⁿs is stored in Control memory.

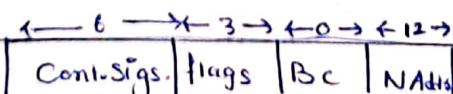
The control address reg. contains the address of next microinstⁿ to be read.

When microinstⁿ is read from CM, it is transferred to CBR. CBR is then connected with control lines (or decoders). Reading CW from CM is same as executing microinstⁿ.



Ex. A vert. pip. CU supports 256 Inst. System uses 8 flags & contains 48 Control Signal. Each Instⁿ req. 16 microoperations. Size of CM?

$$\# \text{Microops} = 256 \times 8 = 2^{12}$$



$$2^6 \times 2^3 \times 2^0 \times 2^{12} = 2^1 \times 2^9 \text{ B} = 107374 \text{ B.}$$

Cache Memory Design

- Cache m/m resides b/w CPU & main m/m. It is made up of static RAM.
- Main m/m can be accessed using physical ~~cache~~ address. Cache m/m can be on-chip (with CPU) or off chip. On chip Cache can be accessed by virtual addr. generated by the CPU. For off chip cache first a Memory Management Unit converts virtual address to physical address.

* Why Cache Memory

When a process is executed, it performs only two operations: fetching from main m/m & execute on CPU. Majority of the time is spent on fetching from m/m. If we could improve the memory access time, then we can improve overall execution time.

It is observed that not all parts (locations) associated with a process are accessed uniformly. Some are accessed more often than others. This is called Locality of Reference.

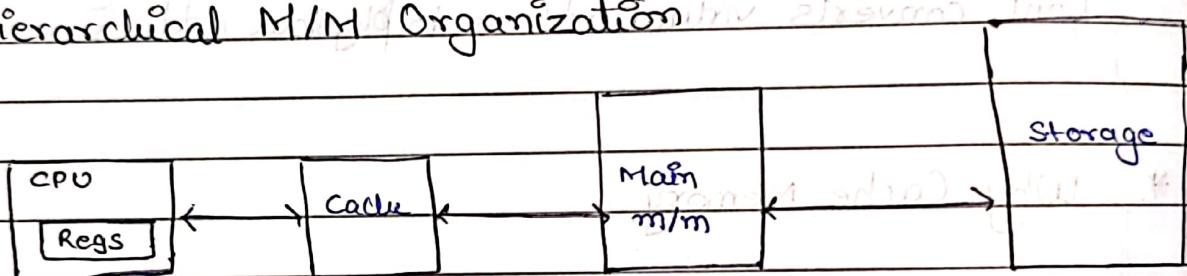
- Locality of Space (Spatial Locality): If CPU access location X at time t_k , then words near to X, such as $X-2, X+4$ would be accessed in near future, with high probability.
- Locality of Time (Temporal Locality): If CPU accesses loc. X at time t_k , then the same loc. would be accessed in near future with high probability.

Ex. for $(i=0; i<10; ++i)$
 $a[i]++;$

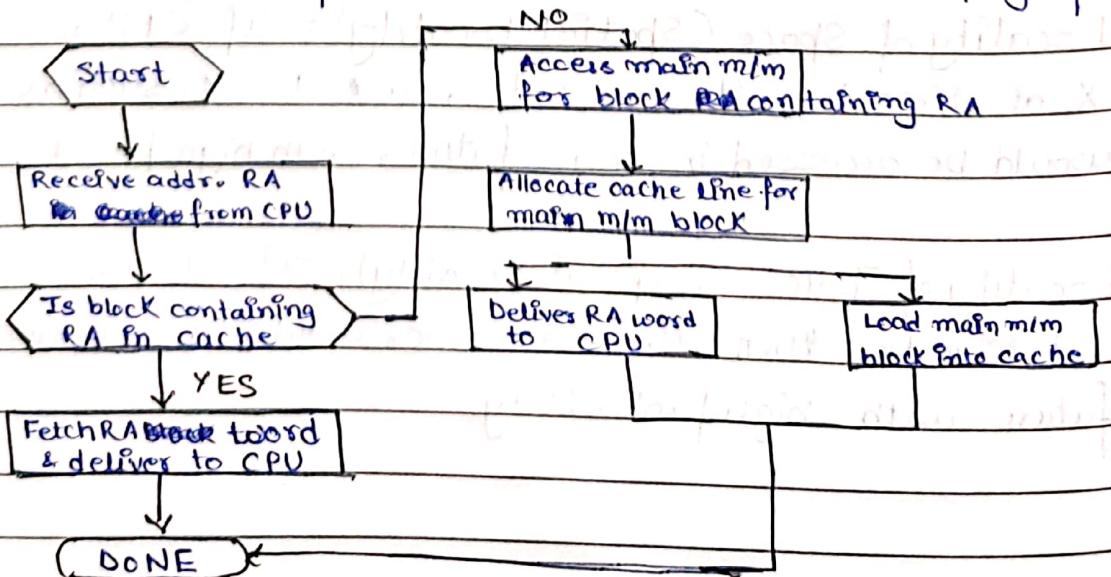
Here 'i' shows temporal locality.

' $a[i]$ ' shows spatial locality.

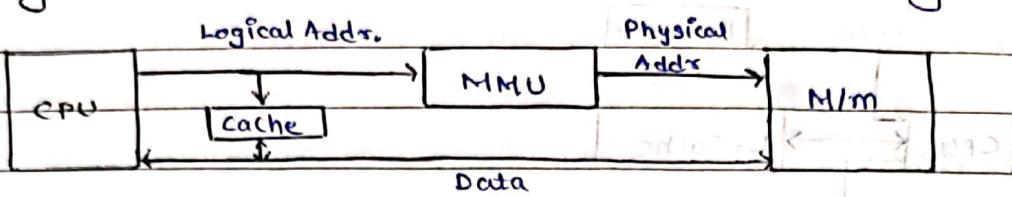
* Hierarchical M/M Organization



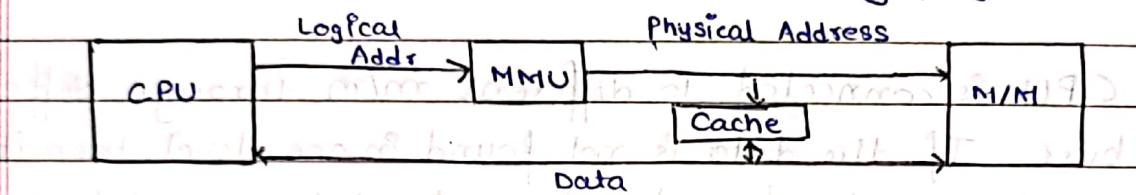
- Generally memory in a system is organized in hierarchical fashion.
- If data/instr isn't found in cache, its called Cache miss, if not found in main m/m its called page fault.



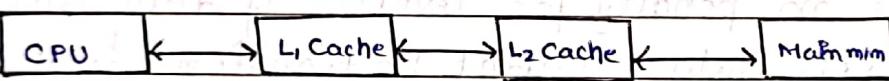
- Logical Cache: Can access cache m/m using Virtual address



- Physical Cache: Can access cache using physical address.



- Hierarchical Cache Orgn



Block/word travels only to its next neighbour for each miss.

(Hit Ratio: Prob. of finding a ~~hit~~ addr in m/m X)
(of m/m X)

(T_x : Access time of mem. X)

Case 1: L₁ Hit:

$$\text{Access Time} = H_{L_1} \times T_{L_1}$$

Case 2: L₁ Miss, L₂ Hit

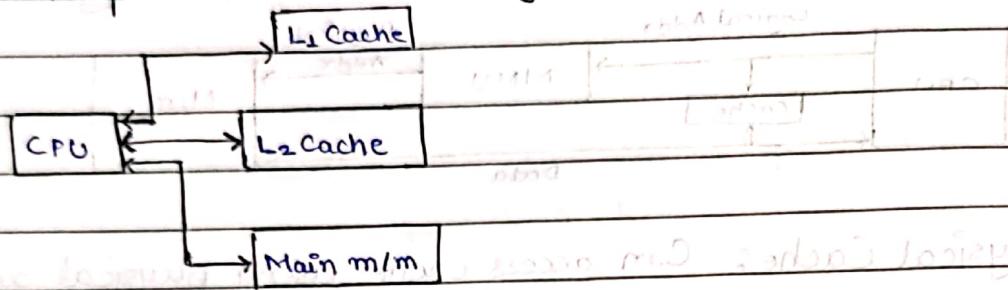
$$\text{Access Time} = (1 - H_{L_1}) \times H_{L_2} \times (T_{L_1} + T_{L_2})$$

Case 3: L₁ Miss, L₂ Miss

$$\text{Access Time} = (1 - H_{L_1})(1 - H_{L_2}) \times H_{M} (T_{m/m} + T_{L_1} + T_{L_2})$$

$$\text{Avg Access Time} = H_{L_1} \times T_{L_1} + (1 - H_{L_1}) \times H_{L_2} \times (T_{L_2} + T_{L_1}) \\ + (1 - H_{L_1})(1 - H_{L_2}) (T_{m/m} + T_{L_2} + T_{L_1})$$

Independent Cache Orgⁿ:



CPU is connected to different m/m through different buses. If the data is not found in one level, then it would never be found in that level, because data is never copied to that level from higher level, directly transferred to CPU.

There is no penalty for level(i) if data isn't found in that level. Data wouldn't be transferred from level (i+1) to (i), only to CPU.

C1: L₁ Hit

$$\text{Access Time} = H_{L1} \times T_{L1}$$

C2: L₁ Miss, L₂ Hit

$$\text{Access Time} = (1 - H_{L1}) \times H_{L2} \times T_{L2}$$

C3: L₁ Miss, L₂ Miss, M/m hit

$$\text{Access Time} = (1 - H_{L1})(1 - H_{L2}) T_{m/m}$$

$$AT_{avg} = H_{L1} \times T_{L1} + (1 - H_{L1}) H_{L2} \times T_{L2} + (1 - H_{L1})(1 - H_{L2}) T_{m/m}$$

* Cache Structure

Cache m/m is divided into equal sized units called Blocks / Lines. At a time b/w cache & ~~m/m~~^{m/m} the least size of transfer is one block. (one word at a time)

Ex. 16B cache. 4B cache line.

00	00 01 10 11	# Lines = $\frac{16B}{4B} = 4$
01	00 01 10 11	
10	00 01 10 11	∴ 4 lines, each of 4B.
11	00 01 10 11	

* Cache Mapping

Initially the cache is empty (or lines are invalid) when CPU requests a word from m/m, it is brought into the cache as well.

Cache Mapping techniques answer the question, which block of m/m will be mapped to which cache line.

Fully Associative

CPU generates an address to access main m/m. Fully Associative Cache will divide the address into two parts: Tag & Offset. Offset is the location (index) of word within a cache block.

In fully associative cache, any m/m block can be put in any available (empty) cache line. Each cache line additionally

stores tag bits. Each cache line has a comparator attached. When CPU sends address to Cache to check for hit/miss, Cache takes the tag part of the address and compares it with tag part of all the cache lines. If any comparator returns true, then cache hit.

TAG	OFFSET	DATA	VALID	WE
-----	--------	------	-------	----

Address fmt for
Fully-Assoc.

OFFSET size: No. of bits req. to index \rightarrow words in block.
 $\Rightarrow \lceil \log_2 (\text{block size in words}) \rceil$

TAG size: All bits other than offset. / No. of memory blocks
 $\Rightarrow \text{Addr len} - \text{offset size.}$

All tags are compared in parallel.

Tag Storage Size = #tag bits \times #cache lines.

If a tag matches, then \uparrow all the words from that block go to Multiplexer, with the offset bits of addr. acting as select line.

comparators = #cache lines

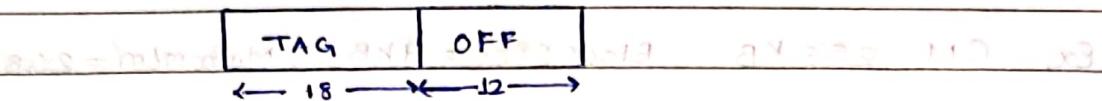
bits in comparator = size of tag field

multiplexer = 8.

Ex. CM = 256 KB 4 KB Blocksize 1 GB m/m.

$$\# \text{ cache lines} = \frac{256 \text{ KB}}{4 \text{ KB}} = 2^6 = 64 \quad \# \text{ words per block} = 4 \text{ K} \Rightarrow 2^{12}$$

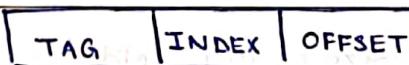
$$\# \text{ blocks} = \frac{2^{30} \text{ B}}{2^{12} \text{ B}} = 2^{18}$$



• Direct Mapping

Rule: Cache line no. for m/m block = m/m block no. % no. of cache lines

• Mem. address is divided into three parts:



INDEX: Cache line index.

TAG: Index of m/m blocks within the cache line

OFFSET: Word offset.

The Direct mapped cache stores m/m blocks based on modulo of m/m block no. by no. of cache lines.

When addr. from CPU comes, cache directly jumps to the specified cache line using INDEX. TAG represents, out of all the m/m blocks which can possibly exist in the particular cache line, which one is currently residing there. OFFSET is word offset within the cache block.

Ex. CM = 256 KB Block Size = 4 KB Main m/m = 2 GB.

$$\# \text{blocks} = \frac{2 \text{ GB}}{4 \text{ KB}} = \frac{2^{31} \text{ B}}{2^{12} \text{ B}} = 2^{19}$$

$$\# \text{cache lines} = \frac{256 \text{ KB}}{4 \text{ KB}} = 2^6$$

$$\text{tag dir. size} = 2^6 \times 13 = 104 \text{ B}$$

TAG	INDEX	OFF
← 13 →	← 6 →	← 12 →

One multiplexer required to fetch the tag bits of particular cache line with INDEX bits as select line. OR 13 bits
One comparator req. to compare TAG bits.

$$\# \text{multiplexers} = 1$$

$$\# \text{comparator} = 1$$

$$\# \text{size of comparator} = \text{size of tag field.}$$

$$\text{Tag dir. size} = \# \text{cache lines} \times \text{size of tag field.}$$

After tag matcher of the specified index in add. all words of that index go to multiplexer with offset bits acting as select lines.

Ex. CM = 4 MB Block Size = 64 B

Mem Addr = 0xA B12CF. Find Cache line no. of this word address?

$$\# \text{Offset bits} = 6$$

$$\# \text{cache lines} = \frac{4 \text{ MB}}{64 \text{ B}} = \frac{2^{22} \text{ B}}{2^6 \text{ B}} = 2^{16}$$

TAG	IDX	OFF
← 2 →	← 16 →	← 6 →

$$\therefore \text{Block No. of cache} = 10 1011 0001 0010 11$$

Set Associative Mapping

A combination of Direct & Fully Assoc. mapping.

We take all the cache lines & divide them into sets of K-lines.

K-way Set Associate means that each set is made up of K lines.

For a m/m block the Set Associative Cache finds the set index by : m/m block no. % no. of sets in cache.

Within the designated set the block can reside in any of the "K" cache lines. Each cache line consists of a tag bits indicating that out of all possible m/m blocks which could be mapped to the particular cache set, which ones are residing in the set currently, specifically in that block.

When CPU generates & sends address to the cache, it is interpreted in three parts:

TAG	SET INDEX	OFFSET
-----	-----------	--------

A set is selected from the cache based on set index. All the blocks within the set (K) send their tag to K different comparators to compare with Tag field. If matched with any of the K blocks, its a hit & offset is used to select a particular word from the block.

comparator = K

size of comparator = size of tag field.

Tag DPr. Size = # cache lines x size of tag field.

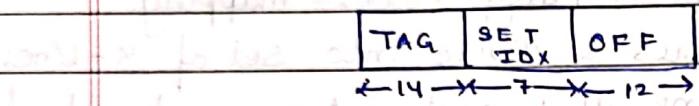
mux = 2 (one to select block, another to select word).

Ex. CM = 4 MB

Block Size = 1 KB

Main m/m = 8 GB

8-way Set Assoc.



$$\# \text{cache lines} = 2^{10}$$

$$\# \text{sets} = 2^{10} = 2^7$$

$$\# \text{m/m blocks} = 2^{21}$$

$$\# \text{comps} = 18 \quad \text{word size comp} = 14 \text{ bits}$$

$$\text{Tag dr size} = 2^{10} \times 14 = 1792 \text{ bits}$$

$$\# \text{muz} = 2^3 : 1 \quad (\text{Block Selector})$$

$$\text{Muz1} : 2^3 : 1 \quad (\text{Word Selector})$$

- Along with storing tag info, cache lines also need to store a few extra bits, such as Valid / Invalid, Dirty, etc.

* Cache Miss

When CPU demands a particular word from cache & the containing block is not in the cache, it's called a Cache Miss & the cache has to incur a miss penalty.

There are three types of cache misses:

- **Compulsory Misses:** The very first access to a block can't be in the cache, so the block must be brought into the cache. Compulsory misses are those that happen even if you have infinite cache.
- **Capacity Miss:** If a block is being discarded from cache due to the cache being full and later retrieved, it is known as capacity miss.

Fully Associative cache would have high capacity miss.

- **Conflict Miss:** If the mapping strategy isn't fully associative, then conflict misses would happen (when it's not compulsory & cache has space) when a block is discarded due to conflict with another block & later retrieved.
Ex. In direct map, even if cache has empty blocks lines a block maps to a particular line, and if it's occupied then the occupying block would be discarded, causing conflict miss.

- * **Cache Latency:** The time required to access the first word of the block.

* Cache Read.

When CPU tries to read a word from the cache:

- (1) Block is available in cache: $T_{C \times H_C}$
- (2) Block is not available in cache: The block would be transferred from m/m:
 - (a) Read Allocate: The block is parallelly transferred to both CPU & cache. $[C_1 - H_C] (T_{mm} + T_C)$
 - (b) No Read Allocate: The block would directly go to CPU only. $[C_1 - H_C] T_{mm}$

* Cache Write

When CPU is trying to write to a particular address:

- (1) Block is available in cache: $[T_{C \times H_C}]$
- (2) Block isn't available in cache:
 - (a) Write Allocate: The block is transferred from m/m to cache. $[C_1 - H_C] (T_C + T_{mm})$

Ex. Say CPU wants to add two nos.

LOAD $r_0, 1000$

LOAD $r_1, 2000$

ADD r_2, r_0, r_1

STORE $1000, r_2$

LOAD $r_3, 4000$

Let's say 4 block Direct Cache & 8 block main m/m.

DM Cache		Main m/m	
b ₀	1000	b ₀	20 1000
b ₁	2000	b ₁	40 4000
b ₂		b ₂	30 2000
b ₃		b ₃	

LOAD $r_0, 1000$: Read Allocate.

$r_0 \leftarrow b_0, mm \rightarrow b_0, cache \rightarrow r_0, CPU$

LOAD $r_1, 2000$: Read Allocate.

$r_1 \leftarrow b_2, mm \rightarrow b_2, cache \rightarrow r_1, CPU$

ADD r_2, r_1, r_0 : $r_2 \leftarrow [r_1] + [r_0]$

$$r_2 \leftarrow [r_1] + [r_0]$$

STORE $1000, r_2$: $r_2 \rightarrow b_0, cache$

$r_2, CPU \rightarrow b_0, cache$

At this point block in Cache contains value 50 for addr 1000 & main m/m contains value 20.

This is Cache Coherence problem.

- Cache Coherence: Same address in two different memories contains different values.

Key takeaways:

LOAD $r_3, 4000$:

The block b_0 gets overwritten with b_4 & the value at loc. 1000 is lost due to cache coherence.

- Write Allocate: A block to be written, on a cache miss is brought from m/m to cache, followed by a write hit. Write Alloc. acts as Read Allocate.
- No Write Allocate: Write miss don't affect the cache. Block is directly modified in main m/m.
- Write Through: The information is written both to the block in cache & in main m/m.
- Write Back: The information is written only to the cache. The modified cache block is written to main m/m only when it is replaced.
- Dirty Bit: A bit is associated with each cache block, which indicates, while in cache if the block has been modified or not. During replacement only the blocks which have been modified are written back.

- With Write Through by default No write allocate is used.

$$T_{avg(r)} = H_c T_c + (1-H_c) (T_{mm} + T_c)$$

$$T_{avg(r)} = H_c T_c + T_{mm} + (1-H_c) T_{mm}$$

$$(w) \Rightarrow H_c (\max(T_c, T_{mm})) + (1-H_c) (\max(T_c, T_{mm}))$$

$$\Rightarrow H_c (\max(T_c, T_{mm})) + (1-H_c) T_{mm}$$

(Hierarchical)

$$T_{avg(r)} = H_c \times T_c + (1-H_c) T_{mm} \quad (\text{parallel})$$

$$T_{avg(w)} = T_{mm} \quad (\text{parallel})$$

- With Write Back by default Write Allocate is used

$$T_{read} = T_{write} = H_c \times T_c + (1-H_c) [\% \text{ of dirty bits} \times (T_{mm} + T_{mm} + T_c) + \% \text{ of not dirty} \times (T_{mm} + T_c)]$$

For the dirty blocks, first the block is written to main m/m

then new block is accessed from m/m & then written to cache

Ex. Consider a two level m/m hierarchy. L₁ (cache) has

AT = 10 ns & main m/m has a AT = 20 ns. Writing & updating into L₁ & m/m takes 10 ns & 30 ns resp. H_{L1} = 20%.

Avg ~~AT~~ in case of Write Through?

$$T_{avg(w)} = H_c \times (\max(T_{mm}, T_c)) + (1-H_c) \times T_{mm}$$

$$= 0.8 \times 30 + 0.2 \times 30$$

$$\Rightarrow 30$$

read

- In Hierarchical access in case of a miss, the block is first read from m/m & put into cache (read allocate) & then "inst" is restarted & this time its a read hit & word is read from cache.
 $\therefore T_{mm} + T_c$

If no read allocate, then T_{mm} ,

(not in cache)

- For write if its a hit both cache & m/m are updated simultaneously in Write Through, & on miss, only m/m is updated.

Ex. A cache line is 64B. The main m/m has latency of 32ns & bandwidth of 1GB/s. The time req. to fetch entire cache line from the main m/m is?

$$1\text{GB/s} = 1 \times 10^9 \text{ B/s} = 1 \times 10^9 \text{ B/s}$$

$$\begin{aligned}\text{Time req.} &= \text{Time to access first word} \\ &\quad + \text{Time to transfer all subsequent words}\end{aligned}$$

$$\Rightarrow 32 \times 10^{-9} \text{ sec} + \frac{64}{10^9 \text{ B/s}}$$

$$\Rightarrow 32 \text{ ns} + 64 \text{ ns}$$

$$\Rightarrow 96 \text{ ns.}$$

* Virtual Cache

- Till now what we've seen is cache being accessed through Physical m/m address: Physical Tag Physical Index.

CPU sends addr. to MMU to convert virtual addr. to Physical addr. either using TLB or page table, & then the addr is sent to cache.

- Virtual cache allows us to simultaneously search the cache & translate Vaddr to Paddr. in MMU.
- Virtual addr. can use: Virtual Tag, Virtual Index.

CPU sends V.A. to cache, which interprets Tag & Index fields from it, and to MMU to translate to P.A.

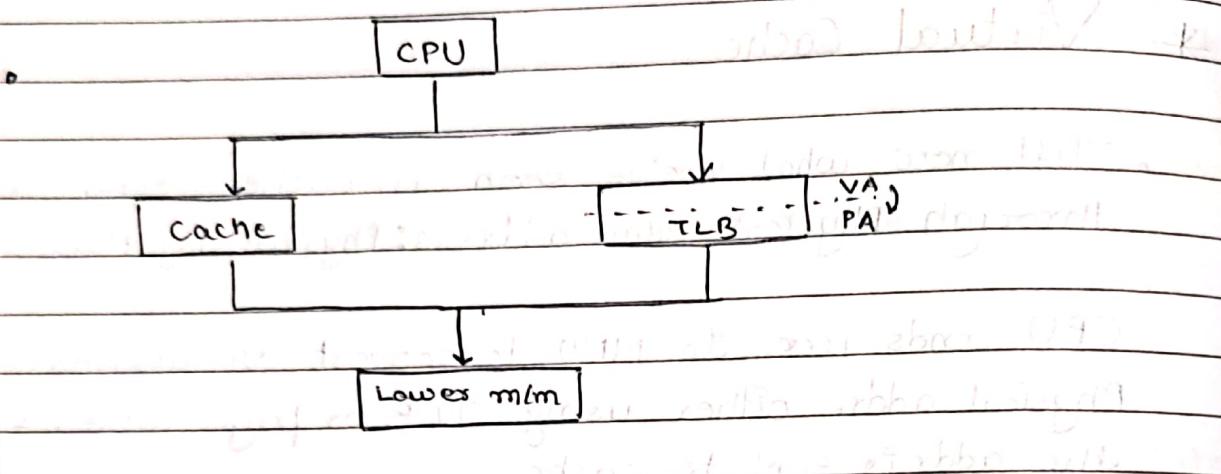
- Tag can be Physical & Virtual.
Index can be Physical & Virtual.

Based on this two types of Cache designs are used

- (1) Virtual Index Physically Tagged
- (2) Physically Index Virtual Tagged [NOT used]

- Problem with VA: VTVI

- (1) In case of context switch, cache has to be flushed
- (2) Diff. VA referring to same PA : Synonym problem.
- (3) Diff. PA referring to same VA : Homonym problem.



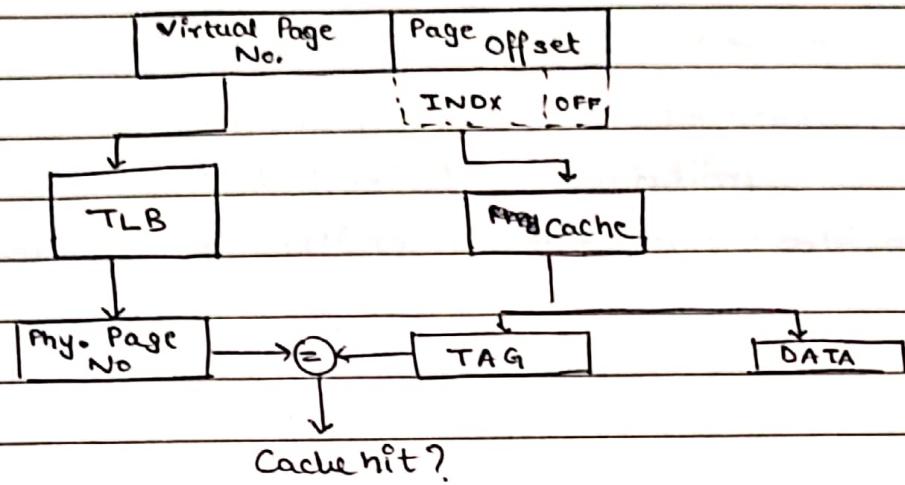
- **Synonym Problem:** Two virtual ~~page~~ addr. refer to same physical location. If both are in cache, then writing to one line is invisible to other. So if other line is read, it will read wrong info.

- **Homonym Problem:** A single VA, for two different process can refer to different Physical loc. If CPU generates such an address, then cache won't know which location's data to bring from m/m.
 Soln: (A) Store PID identifiers along with VA.
 (B) Flush on Context Switch.
 (C) Physical Tags

- **Virtual Index Physical Tag**

Cache size can't exceed page size, so as to avoid synonym problem.

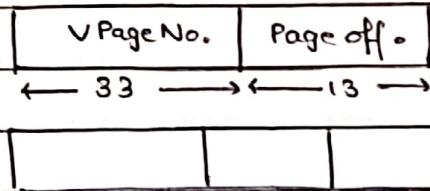
$$\text{Page size} \times \text{Assoc.} \geq \text{Cache Block size}$$



In case, $C > P \times A$, then bits from Virtual Page No. field would be used & would cause Synonym problem.

Ex. $VA = 46$ bits Page Size = 2^{13} B. 1 MB cache with 16 way

Assoc.



$$C = 1 \text{ MB} = 2^{20} \text{ B} \quad A = 16 = 2^4 \quad P = 2^{13}$$

$$\therefore 2^{20} > 2^4 \times 2^{13}$$

$$20 > 17.$$

\therefore 3 bits from the Virtual Page No. field would be used to for set index. Hence 2^3 virtual pages would map to same set.

Hence 8 page colors req. to avoid Synonym.

For each color cache would return 1 physical addr.

Pipelining

Pipelining is an implementation technique where multiple instⁿs are overlapped in their execution. It takes advantage of parallelism of actions needed to execute an instⁿ.

Ex.

	IF	ID	OF	EXE	WB
	IF	ID	OF	EXE	WB

- * Cycle Time: The time period b/w two rising edges of a clock.
- * Frequency of Clock: $\frac{1}{\text{Cycle Time}}$

- * There are various stages to execute an instⁿ. Each stage requires / interacts with different component independent component.

Ex: IF : Mem

ID : Instⁿ Decoder / CU

OF : Mem / Register

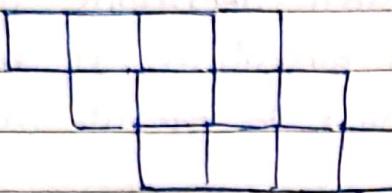
Op. : ALU

WB : Mem. / Register.

The cycle time of CPU is set large enough to make sure each stage can be executed in almost 1 clock cycle.

Pipelining takes benefit of independence of these components, to start executing further instⁿ before finish of the first without hampering the order.

Ex. Say K-step pipeline for n instⁿ.



a K

1

1

∴ for first instⁿ \Rightarrow K cycles.

for remaining (n-1) instⁿ = 1

$$\therefore \# \text{cycles} = (K + (n-1))$$

Say cycle time is t_p

$$\therefore t_p \times (K+n-1)$$

Say K = 5

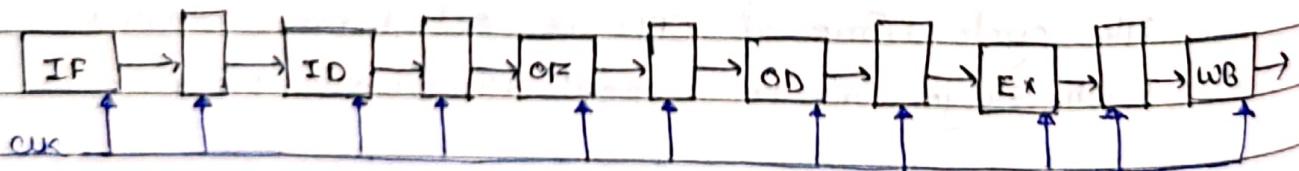
$$n=10 : CPI = (10+5-1)/10 = 1.4$$

$$n=20 : CPI = (20+5-1)/20 = 1.2$$

$$n=100 : CPI = (100+5-1)/100 = 1.04$$

∴ If we keep increasing instⁿ CPI ≈ 1 .

* Pipeline Design



All the stages are followed by buffers to store the output of the stage. Each buffer has a common clock & hence a sync.

∴ The clock period would depend on the slowest stage delay + buffer delay.

$$\therefore t_p = \max(\text{staged delays}) + \text{buffer delay}$$

$$\text{Ex. } N = 1000, K = 5, t_p = 12 \text{ ns}, t_{NP} = 2 \text{ ns}$$

$$\# \text{cycles}_P = 1000 - 1 + 5 = 1004 \quad \therefore \text{Exec Time}_P = 2008 \text{ ns}$$

$$\# \text{cycles}_{NP} = N \times K = 1000 \times 5 = 5000. \quad \therefore \text{Exec Time}_{NP} = 10000 \text{ ns}$$

• Speed Up = $\frac{\text{Performance of Pipeline}}{\text{Performance of Non PL}}$

∴ For each stage, $\frac{1}{t_p}$ part of total work is completed.

∴ $\frac{1}{t_p} \Rightarrow$ $\frac{1}{ET_P}$ part of total work is completed

$\frac{1}{t_{NP}}$ part of total work is completed

∴ Total work is completed

$$\Rightarrow \frac{ET_{NP}}{ET_P}$$

$$\Rightarrow n \times t_{NP} = n \cdot t_{NP}$$

$$(n+K-1) \cdot t_p = (n+K-1) \times t_p$$

• Ideal Condition : $n \gg K-1$.

$$\therefore \text{Speedup}_{ideal} = \frac{n \cdot t_{NP}}{n \cdot t_p} = \frac{t_{NP}}{t_p}$$

• Theoretically, $t_{NP} = K * t_p$

$$\text{Speedup}_{\text{max}} = \frac{\text{Time}_{\text{seq}}}{\text{Time}_{\text{par}}} = \frac{k \times t_p}{t_p} = k$$

Note: If all stages have same time then t_p

- Efficiency (η) = $\frac{\text{Speedup}}{\text{No. stages}} = \frac{k}{n}$

- Throughput = $\frac{\# \text{ tasks}}{\text{time taken}} = \frac{n}{(n+k-1)t_p}$

- Throughput_{ideal} = $\frac{n \cdot m}{m \cdot t_p} = \frac{n}{t_p}$

* Types of Pipelining

- Uniform Delay Pipeline: Every stage of the pipeline executes in same cycle time (t_p), i.e. each stage of pipeline takes exactly one cycle & no less.

$$T_p = t_p + \text{buffer delay}$$

- Non-Uniform Delay Pipeline: Stages have different cycle times, and hence, $t_p = \max(\text{stage delays})$.

$$T_p = \max(\text{stage delay}) + \text{buffer delay}$$

Ex. Consider an n stage pipeline with speedup factor of 10 & operating with 80% efficiency.

$$\begin{aligned} S &= \frac{\eta \cdot k \cdot t_p}{(k-1) \cdot t_p} \\ &\Rightarrow \frac{10 \cdot 80}{(10-1)} = \frac{800}{90} = \frac{80}{9} \end{aligned}$$

$$\eta = \frac{S}{K} = \frac{10}{K} = 0.8$$

$$K = \frac{10}{0.8} = \frac{100}{8} = 12.5 \approx 13$$

Ex.

Proc. A

$$K = 8$$

uniform delay

$$N = 100$$

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

Proc. B.

$$K = 5$$

Non-uniform delay.

$$N = 100$$

$$\text{Stage Delays} \rightarrow 2 \text{ ns}, 3 \text{ ns}, 1 \text{ ns}, 2 \text{ ns}, 2 \text{ ns}$$

Find Speedup if A proc. is used instead of B?

Proc. A

$$ET = (107) \times 2 \text{ ns}$$

Proc. B.

$$ET = (104) \times 3 \text{ ns}$$

$$\text{Speedup} = \frac{1}{ET_A}$$

$$\frac{1}{ET_B}$$

$$\frac{104 \times 3}{107 \times 2}$$

$$\Rightarrow 1.45$$

* RISC vs CISC

RISC

(1) Fewer instⁿ, fewer addressing modes.

(2) Simple instⁿ.

(3) Fixed len instⁿ.

(4) 1 cycle per instⁿ.

(5) Easy to implement using hardwired CU.

(6) Register to Register ALU opⁿs only.

(7) Suitable for pipelining

CISC

(1) More instⁿ & more addressing modes.

(2) Complex instⁿs.

(3) Variable length instⁿ.

(4) Multiple cycles per instⁿ

(5) Difficult with hardwired CU.

(6) Reg to Mem & Reg to Reg ALU opⁿ.

(7) Unsuitable.

- In RISC systems, the opⁿs can only be performed b/w registers, so if an opⁿ needs some data from some m/m loc., it first needs to be loaded into some register, then opⁿ is performed, then the result is stored back into register. The result may be moved from reg. to m/m.

Ex. LOAD R₁, 2010

ADD R₁, R₂

STORE R₁, 2010.

* RISC Pipeline

The RISC instr set consists of the following instr's.

- (1) LOAD & STORE
- (2) ALU Inst's
- (3) Branch Inst's

RISC Pipeline Stages

- (1) Instr fetch cycle : Fetch instr & update PC.
- (2) Instr decode & register fetch cycle : decode the instr & fetch from registers parallelly. All registers can be read from in parallel.
- (3) Execution / Effective Addr Calc. : ~~ALU operation~~
ALU operates on the operands from previous cycle depending on the instr type:
 - (a) Memory Reference : The ALU adds base & offset register to form effective address.
 - (b) Register-Register ALU instr : The ALU performs the op specified by the opcode in the instr on values read from registers.
 - (c) Register-Immediate ALU Inst : The ALU performs operation specified by the ALU opcode, on register & immediate value.
- (4) Memory Access : If the instr is a load, the mem. does a read using the effective address computed in the previous cycle. If it is store, then the m/m writes the data from the register read in previous cycle using the eff. addr.
- (5) Write Back :

- (a) Register to Register ALU instr or LOAD instr :

Write the result into the register file.

* Basic Performance Issues

- For fewer instⁿ, pipelining ~~reduces~~ increases the execution time.
- Imbalance of stage delays results in reduced performance since the clock can run no faster than the slowest stage.

* Pipeline Hazards:

- Anything causing stall cycle is called Hazard.
- Stall Cycles:** Waiting cycle in pipeline in which no work is done, i.e. cycle at end of which no instⁿ completes or cycles in which no new instⁿ enters the pipeline.
- Structural Hazards:** Hazards arising due to structural or resource conflicts.
- Data Hazards:** When an instⁿ depends on the result of previous instⁿ in a way that it affects pipelined execution.
- Control Hazards:** Arise from the pipelining of branch & other instⁿ, which change PC content.
- The CPI gets hampered due to stall cycles:

$$\text{CPI} = \text{CPI}_{\text{ideal}} + \# \text{stall cycles.}$$

$$\therefore E_{Tp} = (n-1+K+\text{stalls}) \times t_p.$$

Ex. [Structural Hazard]

LOAD:	IF	ID	EX	M	WB	
I ₁ :	IF	ID	EX	M	WB	1st cycle
I ₂ :		IF	ID	EX	M	2nd cycle
I ₃ :			IF	ID	EX	3rd cycle
I ₄ :				IF	ID	4th cycle

In cycle 4, both LOAD & I₃ access memory simultaneously
hence this is Structural Hazard.

∴ We'd have to stall.

LOAD:	F	D	E	M	W	
I ₁ :	F	D	E	M	W	ULL
I ₂ :	F	D	E	M	W	ULL
I ₃ :		S	F	D	E	M
I ₄ :	S	F	D	E	M	ULL

∴ instⁿ takes 10 cycles.

Solⁿ to this structural hazard is to keep separate instⁿ & data memory.

Ex. [Data Hazard]

ADD R₁, R₂, R₃

SUB R₄, R₁, R₅

AND R₆, R₁, R₂

OR R₇, R₁, R₉

XOR R₁₀, R₁, R₁₁

Types of Data Dependence

- (1) Read After Write: Instⁿ j tries to read some value written by Instⁿ i, before Instⁿ i is able to write it. ($i < j$)
- (2) Write after Write: Instⁿ j tries to write to some reg. before ith which Instⁿ i also writes. ($i < j$)
- (3) Write after Read: Instⁿ j tries to write to some reg. before ith which Instⁿ i has read from. ($i < j$)

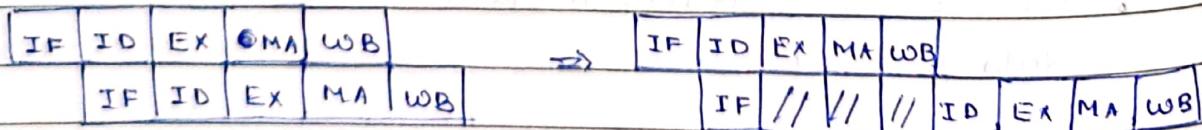
WAU & WAR are called name dependence, and can be removed by using different registers.

RAW is known as true dependence.

Ex. [RAW]

ADD R₁, R₂, R₃

SUB R₅, R₁, R₄



Ex. [WAU]

ADD R₁, R₂, R₃

SUB R₁, R₄, R₅

IF	ID	EX	MEM	WB
IF	ID	EX	MEM	WB

Since both update the value of R_1 in WB, so order maintained.

Ex. [WAW]

Say m/m takes 3 cycles to read

LOAD R₁, 0x3FFE2
ADD R₁, R₂, R₃

IF	ID	EX	M	M	M	WB
IF	ID	EX	M	WB		

load	[id]	[ex]	[M]	[M]	[M]	WB
add	[id]	[ex]	[M]	[M]	[M]	WB

This creates problem, as newer value is overwritten.

Ex. [WAR]

ADD R₁, R₂, R₃
SUB R₂, R₅, R₆.

IF	ID	EX	M	WB
IF	ID	EX	M	WB

∴ No hazard.

- Only RAW dependency is sure to cause hazard.
- WAR is called anti-dependency as it rarely causes stall.
- Register renaming can solve both WAR & WAW dependencies.

• Operand Forwarding / Data Forwarding / Bypassing

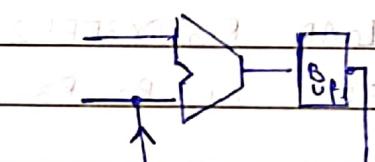
Result of previous operations can directly be transferred from ALU buffer to ALU directly for next operation.

Ex ADD R₁, R₂, R₃

SUB R₄, R₁, R₅

this looks like:

IF	ID	EX	M	WB
IF	ID	EX	M	WB



This helps to prevent stalls.

Ex. LDAD R₁, 10(R₂)

SUB R₄, R₁, R₅

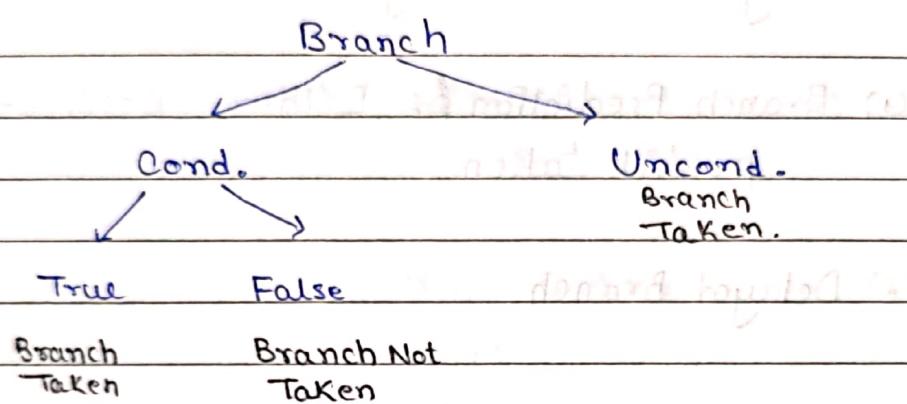
AND R₆, R₁, R₇

OR R₈, R₁, R₉

IF	ID	EX	MM	WB	
IF	ID	//	EX	MM	WB
IF	ID	//	EX	MM	WB
IF	ID	//	EX	MM	WB

Bypassing can't handle all RAW hazards. There can still be stalls.

- Control Hazards: Cause huge performance loss. When a branch is taken it may/maynot take PC to current value +1.



Ex. 1000: ADD R₁, R₂

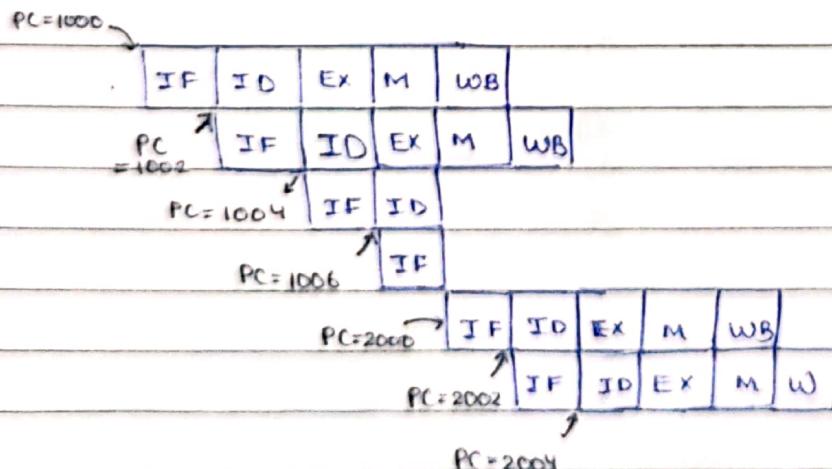
1002: JMP 2000

1004: SUB R₂, R₄

1006: MUL R₃, R₄

2000: ADD R₂, R₄

2002: SUB R₃, R₄



Instⁿ at loc. 1002 calculates the target of JMP in EX cycle & then takes a jump, till this happen two Instⁿ at 1004 & 1006 have already been fetched. Once branch is taken the results of these 2 Instⁿs need to be flushed from buffers before the next Instⁿ from 2000 can be executed.

- Branch Penalty : The penalty incurred due to fetching wrong instrn following branch.
- Reducing Branch Penalty :

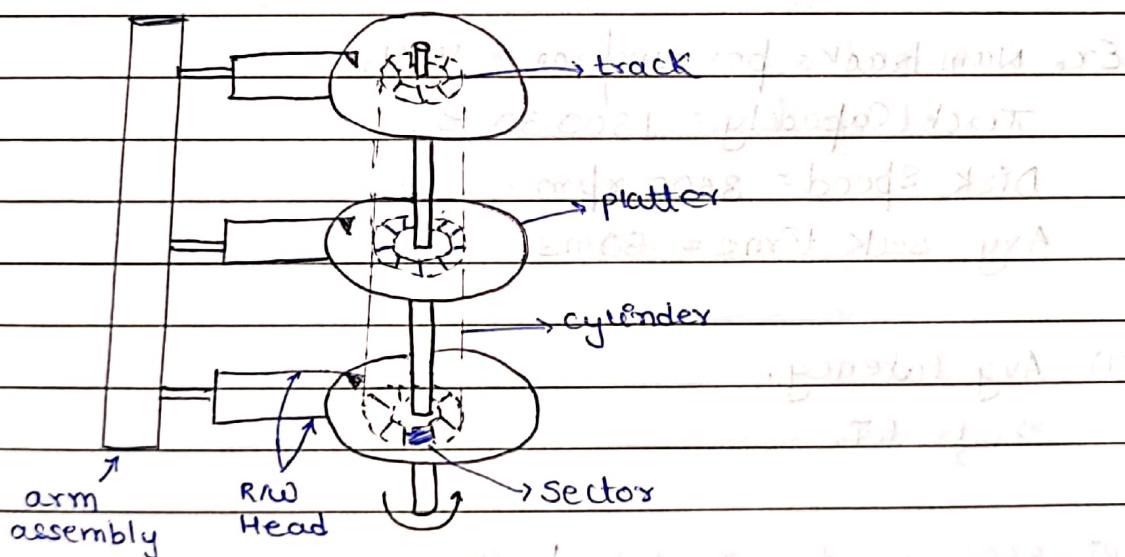
(1) Branch Prediction : Either predict-not taken or predict taken

(2) Delayed Branch

Disk Structure

Main m/m is small & volatile & hence we need Disk Storage to sustain our data.

- **Sector:** The smallest unit of data that can be read or written to on a disk.
- **Track:** Made up of multiple sectors.
- **Platter:** Collection of tracks.
- **Cylinder:** Collection of single track across platters.



- Platters can be single sided or double sided.

- **Rotation Time:** Time to one complete rotation.

Rotational Latency: Half of Rotation Time.
(Avg Latency)

$$\frac{RT}{2}$$

- Seek Time: Time to align R/W head over the desired track.
- Sector Transfer Time: $\frac{1}{\text{Rotation Time}} \times \text{No. of sectors on the track}$
- Disk Access Time = Rotational Latency + Seek Time + 1 sector access time + any additional delay (if given)

Ex. Num tracks per surface = 404.

Track Capacity = 130030 B.

Disk Speed = 3600 rpm.

Avg seek time = 80 ms.

(1) Avg latency.

$$\Rightarrow \frac{1}{2} RT$$

$$RT = \frac{3600}{60} = \frac{1}{x} \Rightarrow x = \frac{1}{60} \text{ sec.}$$

$$\therefore AL = \frac{1}{120} \text{ sec.} = \frac{25}{3} \text{ ms}$$

$$(2) \text{ Storage Capacity} = \# \text{ surfaces} \times 404 \times 130030 \text{ B} \\ \approx 50 \text{ MB} \times \# \text{ surface}$$

(3) Data Transfer Rate.

1 Rot \rightarrow 1 track

$\Rightarrow 50 \text{ ms} \Rightarrow 1300.30 \text{ B}$

3

$\Rightarrow 1 \text{ sec} \Rightarrow 1300.30 \times 60 \text{ B}$

$\therefore \text{Data transfer rate} = 7618.94 \text{ KB/s}$

one sector (or block)

• Disk Transfer Time: Time to transfer \uparrow
~~of a sector or block~~

• Avg disk access time = Avg Seek Time + Avg Rot Latency
 $+ \text{Disk Transfer time} + \text{Controller Delay}$
 $= 10.7 \text{ ms} + \text{Queuing delay}$

(Section 13) HAB

* Disk Scheduling

Goal of Disk Scheduling is to minimize head movement, since it affects seek time.

So given a queue of track ids we want to schedule the seeks in a way that it minimizes head movement.

• FCFS

Serve the track req. as per the arrival time.

• Shortest Seek Time First

From the current head position serve the track request which is nearest.

• SCAN (Elevator)

From its current pos. the R/W head goes to the one end of surface servicing all the req. on its path.

Then it reverses its direction & it goes to the other end, again servicing the requests in its path.

• ~~C-SCAN~~ C-SCAN,

From its current pos. the R/W head goes to the one end of the surface servicing all the req. on its path. Then it reverses its direction & goes to the other end

without servicing any req. From the other end it reverses its direction & starts servicing.

- LOOK

Same as SCAN but in one direction It only goes to last request & not the surface end.

- CLOOK

Same as CSCAN but in a direction It only goes to last request & not the surface end.

FLOATING POINT REPR

- This is used to store very large or very small numbers.

Ex. Say we have the following hex no.

$$(897\ 000\ 000\ 000)_H$$

If we were to store this, it would take $12 \times 4 = 48$ bits.
We can come up with a different representation of the same number.

$$+ 8.97 \times 16^{11}$$

So we can just store the values +, 8.97, 11 into a smaller sized register & later interpret it back.

This representation is known as Mantissa - Base - Exponent Representation.

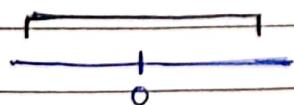
$$\pm M \cdot B^{\pm E}$$

This representation stores the following info:

Sign: 1 bit for +/ve (0 for +ve, 1 for -ve).

Mantissa: Value after decimal point

Biased Exponent: Since our exponent can also be +ve or neg., we shift our exponent value by a bias towards pos. side.



Original range
of E



New range of E

$$\therefore E' = E + \text{Bias}$$

The bias value depends on the length of bias field in our representation.

Ex.	<table border="1"> <tr> <td>S</td><td>BE</td><td>M</td></tr> </table>	S	BE	M
S	BE	M		
	← 1 ← S → 10 →			

In the above repr. BE is of 5 bits meaning it can store values in range:

$$0 \text{ to } 2^5 - 1$$

So, since E' is of range 0 to $2^5 - 1$,

$$\begin{aligned} E \text{ would be in range } & -(2^4 - 1) \text{ to } (2^4 - 1) \\ & \Rightarrow -15 \text{ to } 15 \end{aligned}$$

∴ For n bit Biased Exponent field, range

of B Exponents is

~~$$-(2^{n-1} - 1) \text{ to } +2^{n-1}$$~~

and the bias value is $(2^{n-1} - 1)$.

Normalized Form: For a binary number, xyz.ABC, the normalized form is x.yzABC, i.e. only 1 bit to the left of decimal.

Ex. 1001.0101

$$\Rightarrow 1.0010101 \times 2^3$$

Ex. 0.0010101

$$\Rightarrow 1.0101 \times 2^{-3}$$

Ex. Format:

S	BE	M
1	5	10

Num: 10010.11101 $\times 2^4$

$$\Rightarrow 1.001011101 \times 2^8 = 59 \quad 10000001000000000000000000000000 = 59 \times 2^8$$

$$\therefore M = 0010111010 \quad (10 \text{ digit bits})$$

$$S = 0 \quad (\text{+ve}) \quad (1 \text{ bit})$$

$$E = 8 \quad E' = 8 + (2^{n-1} - 1) \quad n=8 \quad \Rightarrow 8 + (16 - 1)$$

$$\Rightarrow 101110_2 = 27 \quad (5 \text{ bit})$$

$$0|10111|0010111010 = 27BA$$

* IEEE - 754 Single Precision Format (32-bits)

S	BE	Mantissa
← 1 →	8	→ 23 ←

$$\therefore \text{Bias} = 2^7 - 1 = 127$$

$$\& \text{Exponent Range} = -(2^7 - 1) \text{ to } +(2^7 - 1) = -127 \text{ to } 128$$

* IEEE-754 Double Precision Format (64-bits)

S	BE	Mantissa
← 1 →	11	→ 52 ←

$$\therefore \text{Bias} = 2^{11-1} - 1 = 2^{10} - 1 = 1023$$

$$\& \text{Exponent Range} = -(2^{10}-1) \text{ to } +(2^{10}-1) = -1023 \text{ to } 1024$$

* Conversion To Value: Given S, BE & M.

$$B = BE - \text{Bias} \quad \text{Value} = (-1)^B \times 1.M \times 2^B$$

$$\text{Ex. } R_1 = 0x42200000 \quad R_2 = 0xC1200000$$

$$R_3 = \frac{R_1}{R_2} \quad \text{in single Prec. format.}$$

$$R_1 \Rightarrow 0|00\ 0010\ 0010\ 0000\ 0000\dots \quad S = 0$$

$$\therefore S = 0 \quad BE = 132 \quad M = 0100\dots$$

$$\begin{aligned} E &= 132 - 127 \\ &\Rightarrow 5 \end{aligned}$$

$$\therefore \text{Value} = +1.01000\dots \times 2^5$$

$$\Rightarrow +(101000)_2$$

$$\Rightarrow (+40)_{10}$$

$$R_2 \Rightarrow 1|100\ 0001\ 0010\ 0000\ 0000\ 0000\dots$$

$$\therefore S = 1 \quad BE = 130 \quad M = 01000\dots \quad E = 3$$

$$\therefore \text{Value} = -1.01000\dots \times 2^3$$

$$\Rightarrow -(1010)_2$$

$$\Rightarrow (-10)_{10}$$

$$\therefore R_3 = \frac{R_1}{R_2} = \frac{+40}{-10} = -4.$$

$$R_3 \Rightarrow \underline{\underline{-}} \quad - (100.0)$$

$$\Rightarrow - (1.00) \times 2^2$$

$$\therefore S = 1 \quad E = 2 \quad BE = 129 \quad M = 000\dots$$

$$1|100\ 0000\ 1|0000\dots$$

$$\Rightarrow (C0800000)_H$$

Ex. $(3 \text{ E } 6 \text{ D } 0000)_H$. Single Prec. Format.

$0b11\ 1110\ 1101\ 0000\ 0000..$

$S=0 \quad BE=124 \quad M=11011010000..$

$$E = 124 - 127 = -3.$$

$$+ (1.11011010) \times 2^{-3}$$

$$\Rightarrow +(0.0011101101)_2$$

\Rightarrow

* Special Values

Although the BE field in (say) single prec. format allows for values in range: -127 to 128 by biasing them (with 127) to 0 to 255.

0 is 0000 0000 & 255 is 11111111

These two values are used for special purposes, so the range usable is (-126 to 127).

S	BE	M	Remark
+/-	All 0s	All 0s	± 0
+/-	All 0s	Not All 0s	Denormalized Nums
+/-	All 1s	All 0s	$\pm \text{Infinity}$
+/-	All 1s	Not all 0s	NAN

* Implicit Normalization: 1. Mantissa format. 1 is implicit.

Explicit Normalization: 0.1 Mantissa. 1 is explicitly included in the mantissa.