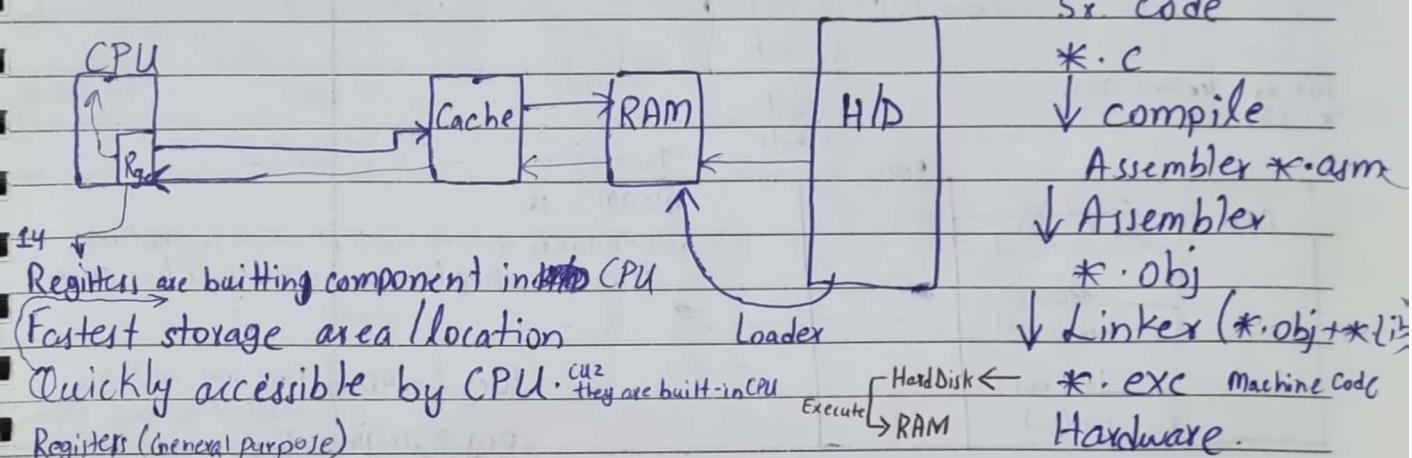
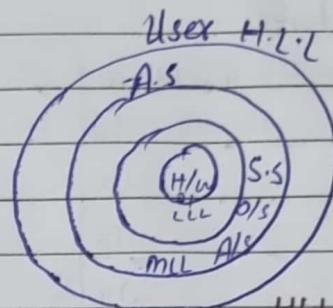
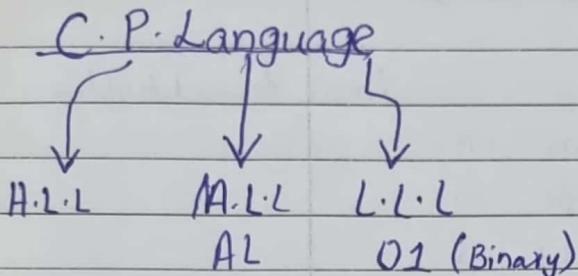


$A \oplus 2$ → instructions.
 $2 + 3$ → opcode

DATE ___ / ___ / 20___

System Software → embedded
 Programming
 ↓
 Programming to
 run hardware

COAL



① accumulator (a)	code segment	instruction
② accumulator extended (ax)	data "	Pointer ip
③ base (bx)	GP Segment	stack pointer Sp
④ counter / code (cx)	extended segment	f] Flag Register
⑤ data (dx)	source index (si)	b] Base Register.
		Destination Index (di)

Efficiency of CPU → Accumulator: Input / Output, Operations.

Base: Holds address of Data.

Counter: Used in loop

Data: Holds data for output.

Source Index: Point the source Operand.

Destination " : Point the destination Operand.

Instruction Pointer: Holds the next instruction.

Code Segment: Holds addresses of code segments. Stack Pointer: Point the current top of the stack.

Data Segment: Holds address of data segments.

Stack Segment: Holds " of stack management.

Extra segment: Holds address of data segment.

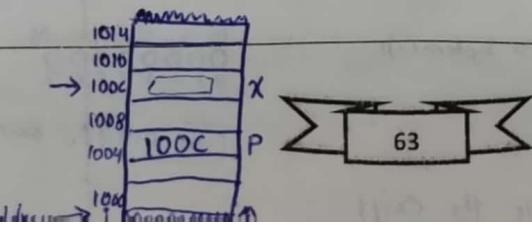
Provides additional memory addressing capability beyond the limitation of the Data Segment in 16 bit real mode.

Flag Register: Holds current status of program.

Base Pointer: Base the top of stack

```

int x;
int *p;
x = 25;
p = &x;
  
```



$$\begin{array}{r}
 & \text{borrow} \\
 & \text{1} \\
 & \text{0} \quad | \\
 & \text{1} \quad | \\
 & \text{1} \quad | \\
 & \text{0} \\
 \hline
 \end{array}$$

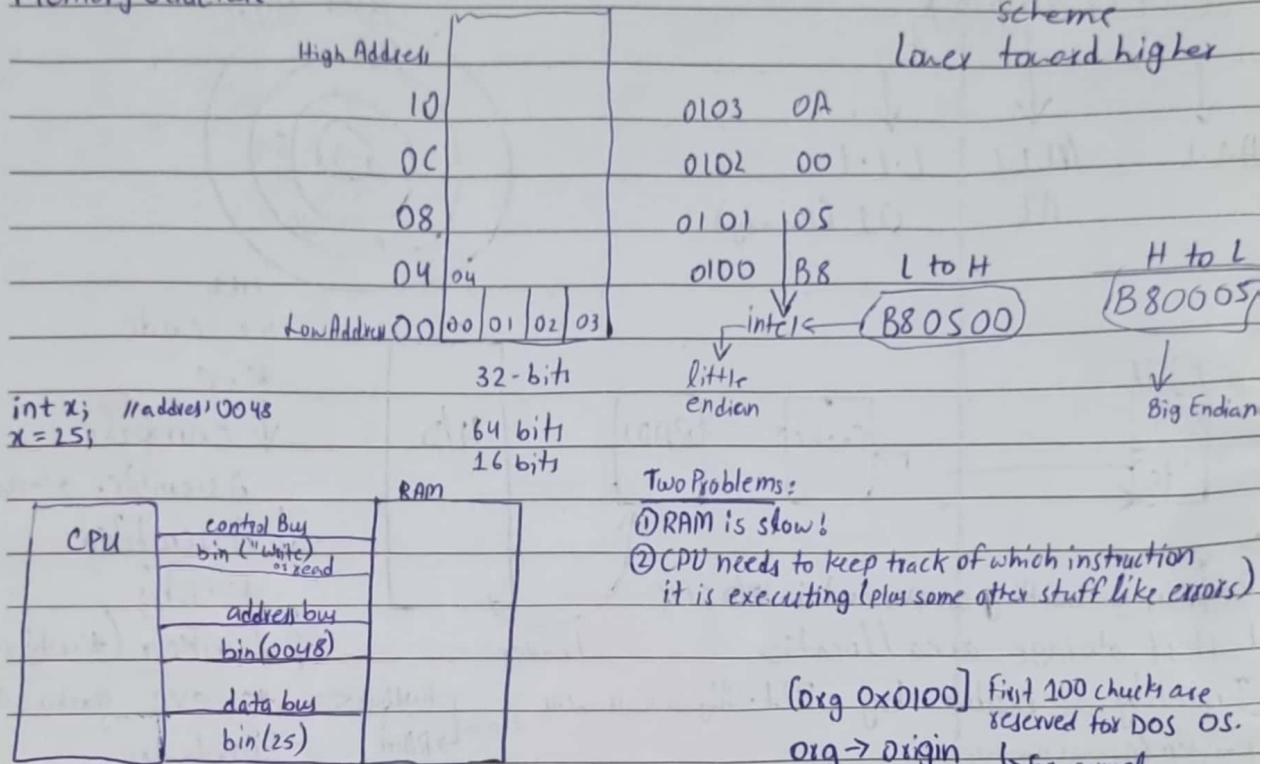


GO GETTERS

storage high speed slow.

Pause
Pose
DATE 1/120

Memory Structure



In general, all are based on:

Von-Neumann Architecture:

- > Code & data both stored in memory
- > General Purpose
- > Fetch, decode, execute.

mount c c:/MP

c:

edit file.asm

First Code → ax=5

register long 0x0100]
instruction mov ax, 5 → operand.

mov bx, 10

add ax, bx ∵ ax = ax + bx

→ nasm file.asm -f file.com

[CPU keeps track of]

afd file.com

mov bx, 13

add ax, bx

Press F1 for next step

mov ax, 0x4C00 ; exit...

nasm file.asm -f file.list -o file.com
edit file.list

int 0x21 ; ... is what the OS should do for me.

↳ interrupt

Request Pause the Program.

Request Signal → System Call

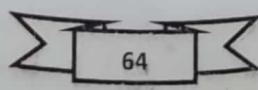
Application ↓
(User Space) ↳ OS
(Kernel Space)

↳ to handle the exit.

B8 04 00 B8 03 00 01 D8
00 00 00 00 00
0 1 2 3 4 5 6 7

GO GETTERS

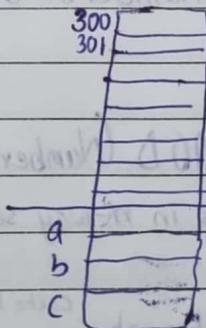
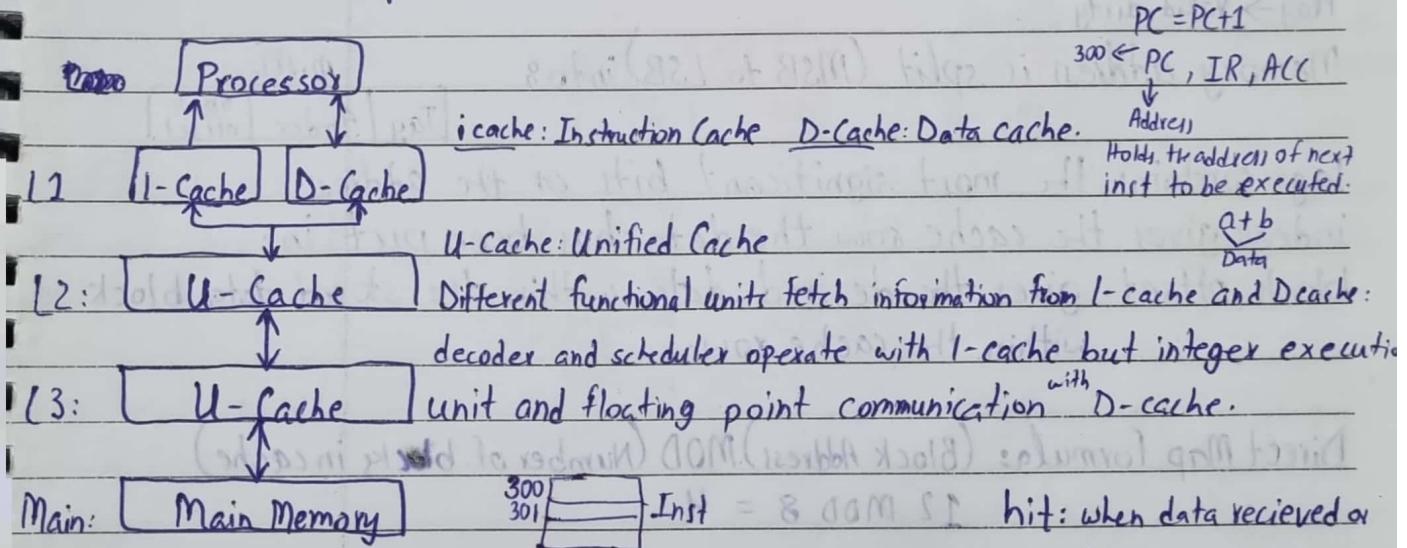
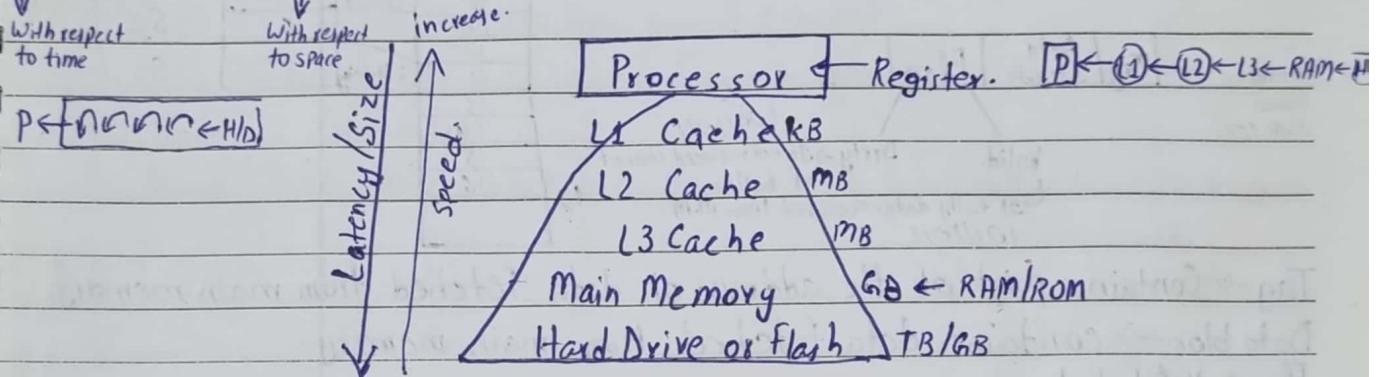
mi 0100 To BoundIndex.



Memory Hierarchy Design

Sharing of memory
by layering.

- Programmer want unlimited amount of memory with low latency.
- Fast memory technology is more expensive per bit than slower memory
- Solution: Organize memory system into a hierarchy
 - Entire addressable memory space available in largest, slowest memory
 - Incrementally smaller & faster memories, each containing a subset of the memory below it, proceed in steps up towards the processor.
- Temporal and spatial locality insures that nearly all references can be found in smaller memories
 - With respect to time
 - With respect to space
 - increase.
- Givei the illusion of a large, fast memory being presented to the Processor.



When a word is not found in the cache, a miss occurs.
 ↳ Fetch word from low level in hierarchy requiring high latency refer.
 ↳ Low level maybe cache or the main memory.

Miss Rate:

↳ Fraction of cache access that result is a miss

↳ Causes of miss (3C's)

↳ Compulsory ↳ Capacity ↳ Conflict.

GO GETTERS

Compulsory: First ref to a block - would happen even for infinite Cache.
Capacity: Block discarded & later retrieved.

Conflict: Program make repeated references to multiple addresses from different blocks that map to the same location in the cache.

- Placement Policies:
- o Main memory has a much larger capacity than cache
 - ★ o Mapping between main and cache memories.
 - ★ o Where to put a block in cache.

$$2^8 = 256$$

$$2^3 = 8$$

DATE ___ / ___ / 20___

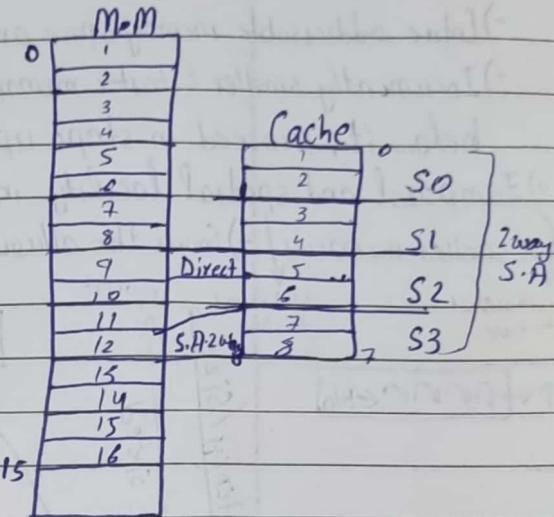
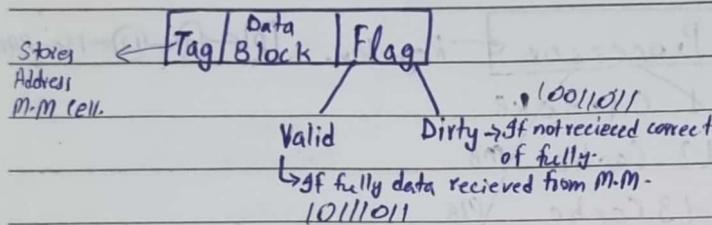
How cache can communicate with Main Memory

3- Approaches to match cache memory with main

- ① Direct Map
- ② Set Associative
- ③ Full Associative

→ Whole cache will go to one chunk of M.M.
→ A block can be placed in any location in cache.

Cache Row:

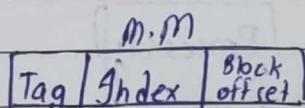


Tag → Contains part of the address of data fetched from main memory.

Data block → contains data fetched from main memory.

Flag → Valid, dirty

Memory address is split (MSB to LSB) into 8



Tag: contains the most significant bits of the address.

index: gives the cache row the data has been put in.

block offset: gives the desired data within the stored data block within the cache row.

Direct Map Formula: (Block Address) MOD (Number of blocks in cache)

$$12 \text{ MOD } 8 = 4$$

Set Associative Formula: (Block Address) MOD (Number of sets in cache)

→ A block can be placed in one of n locations in n-way set associative cache.

Writing to cache: two strategies

o Write-back

→ Only update lower levels of hierarchy when an updated block is replaced.

o Both strategies use write buffer to make writes asynchronous.

o Write-through

↳ Immediately update lower level of hierarchy

cache & memory (update)



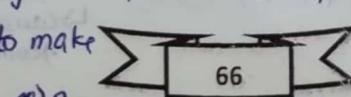
o Why write buffer?

so CPU doesn't stall.

o Why buffer, why not just one register.

Burst of writes are common.

GO GETTERS



o Are read after write (RAW) hazard an issue for write buffers?

Yes! Drain buffer before performing any read if stall.

(Q) You have 32 blocks of M.M and 16 blocks cache determine the cache block to map with 10 block of MM using

(2 way) SA -

$$10 \% \cdot 8 = 2$$

Hit: data appears in some block in the upper level (example: Block X)

Hit Rate: the fraction of memory access found in the upper level.

Hit Time: Time to access the upper level which consist of RAM access time + Time to determine hit/mis

• Hit Time << Miss Penalty (500 instruction on 212646)

Miss: data needs to be retrieved from a block in the lower level (Block Y)

• Miss Rate = 1 - (Hit Rate)

• Miss Penalty = Time to replace a block in the Upperlevel + Time to deliver the block to processor.

A [] = 5
↓
memory

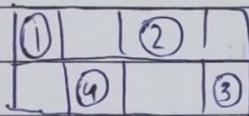
Address
Age = 20;
Variable

dw : define word
↳ size is 16 bits/2 bytes.

db: define bytes.

↳ to define space, specific space as per needed

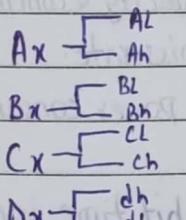
[arg 0x0100]



dw → 0005
db → 0A05

mov ax, [num1]
mov bx, [num1+1]
add ax, bx

mov bx, [num1+2]
add ax, bx



Byte stor.
variable.
00000000, 00000000,
AH AR

first ~~pro~~ loader

finds the global

mov [num1+3], ax

xor ax, ax

Mov [num1], [num2], illegal

we can't assign memory to memory

Atk wat mac nik hi address carry k.

mov ax, 0x4c00

int 0x21

(ZF) Zero Register: tells either it's 0 or 1

↳ tell the status of program

→ code executed if ZF=0 (result is not zero)

→ code executed if ZF=1 (result is zero).

num1: dw 5, 10, 8, 0

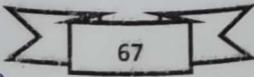
↓

Global variable: In AL, global variable are typically accessed using memory addressing.

The memory address of global variable is known at compile time and the variable is accessed using its address.



GO GETTERS



67

How global variable seen by assembly language.

Listing:

PreCodeList

nasm work.asm -l work.lst -o work.com \Rightarrow to see real address of
edit work.lst memory.

Hit rate: fraction found in that level.

~~It~~ So high that usually talk about MissRate.

Average memory-access time = HitTime + Miss rate \times Miss Penalty (ns or clock)

Miss Penalty: time to replace a block from lower level, including time to replace in CPU

access time: time to lower level = f(latency to lower level)

transfer time: time to transfer block = f(BW between upper & lower levels, blocks)

Improve Performance by: 1) Reduce the miss rate.

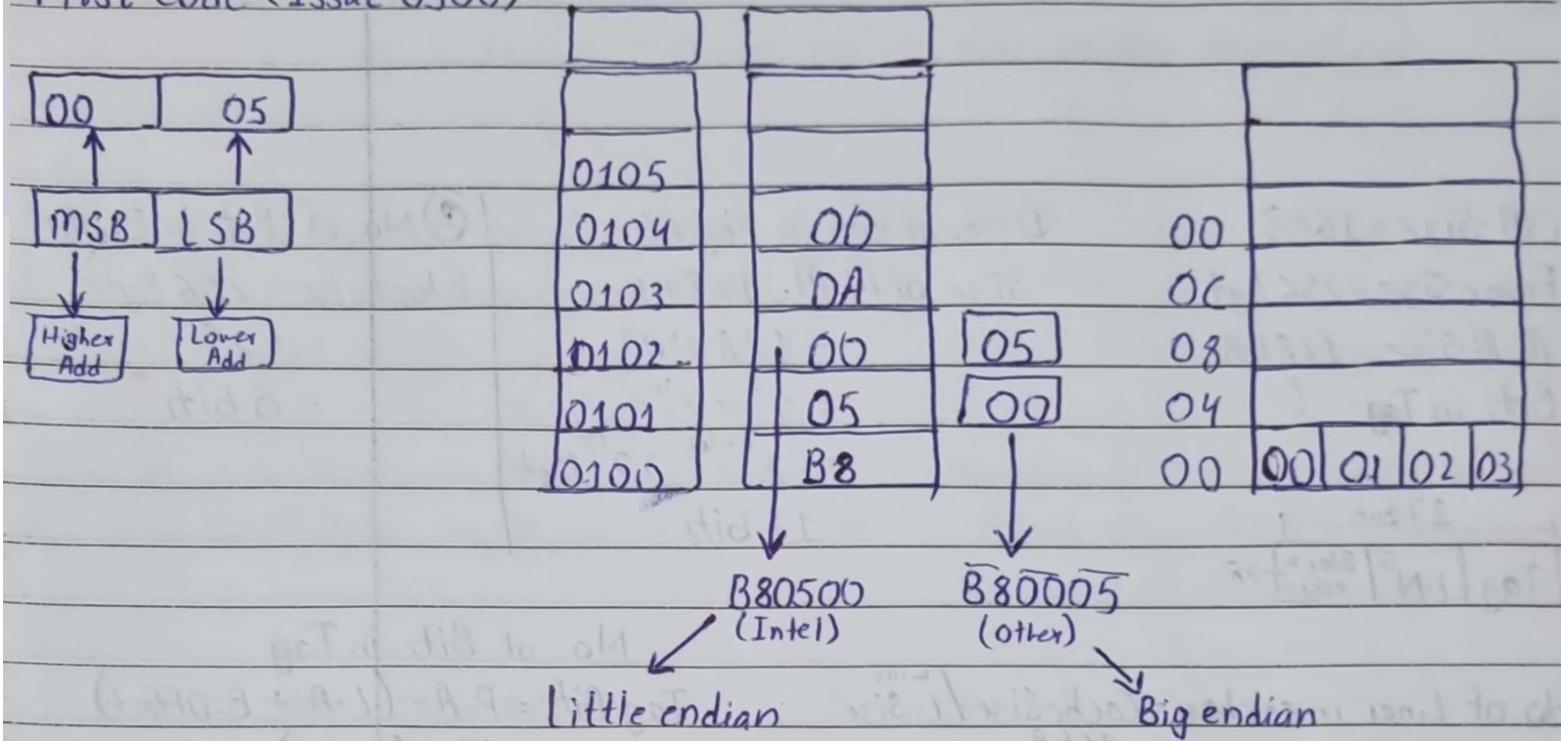
2) Reduce the miss penalty, or

3) Reduce the time to hit in the cache.

Six basic cache optimizations:

- Larger block size
 - \hookrightarrow Reduce compulsory misses, \Rightarrow Increases capacity & conflict misses, increases miss penalty.
- Larger total cache capacity to reduce miss rate.
 - \hookrightarrow Increases capacity hit time, increases power consumption
- Higher associativity
 - \hookrightarrow Reduce conflict misses \Rightarrow Increases hit time, increase power consumption.
- Higher number of cache levels.
 - \hookrightarrow Reduces overall memory access time
- Giving priority to read misses over writes.
 - \hookrightarrow Reduce miss penalty.
- Avoiding address translation in cache indexing
 - \hookrightarrow Reduce hit time

First Code (Issue 0500)



#Assemble and run

```
nasasm c.asm -o c.com
afld c.com
```

#Assemble with listing

```
nasasm c.asm -l c.lst -o c.com
```

Numerical :-

(Q.1) Consider a direct map

$$CM \cdot Size = 16KB$$

$$Frame Size = 256 Byte$$

$$M \cdot M \cdot Size = 128KB$$

$$Bit in Tag = ?$$

① No. of Bits in Phy Add

$$\text{Size of } M \cdot M = 128KB$$

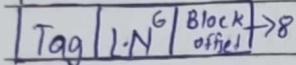
$$= 128 \times 100$$

$$= 2^7 \times 2^{10}$$

$$= 2^{7+10} \Rightarrow 2^{17} Byte$$

$$= 17 \text{ bit}$$

17 bit



② No. of Bits in B-offset

$$\text{Block size} = 256 Byte$$

$$= 2^8$$

$$= 8 \text{ bit}$$

$$\text{No. of Lines in cache} = \text{Cache Size} / L \cdot \text{Size}$$

$$= \frac{16KB}{256 Bytes}$$

$$= 2^4 \times 2^{10} \Rightarrow 2^6 Bytes$$

$$= 64$$

$$= [64 \text{ bit}]$$

No. of Bits in Tag

$$\text{Tag Bits} = P.A - (L.A + B.Offset)$$

$$= 17 - (8 + 6)$$

$$= 3 \text{ Bits}$$

(Q.2)

Number of Bits in Block offset

$$\text{Block Size} = 1KB = 2^{10} bytes$$

$$= 10 \text{ bits}$$

Number of bits in Line Number

$$\text{Line in cache} = \text{Cache Size} / \text{Line Size}$$

$$= 512KB / 1KB$$

$$= 2^9 \text{ lines}$$

$$= 9 \text{ bits. } \cancel{[9]}$$

$$P.A = 7 + 9 + 10 = 26 \text{ bit}$$

$$\therefore KB \times KB = MB$$

$$10^3 \times 10^3 = 10^6$$

$$M \cdot M = 2^{26} \Rightarrow 2^{16} \times 2^{10}$$

$$165000 \times 1000$$

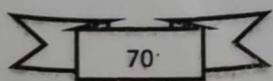
$$2^{10} \times 2^{10} \times 2^6$$

$$K \times K \times 2^6$$

$$64 KB \times KB \Rightarrow [64 MB]$$



GO GETTERS



Q3) Consider a direct mapped cache with block size 4KB - The size of main memory is 16GB & there are 10 bits in the tag. Find Size of cache memory.

Solution:

$$\text{Block size} = \text{frame size} = \text{line size} = 4\text{KB}$$

$$\text{Size of main memory} = 16\text{GB}$$

$$\text{Number of bits in tag} = 10 \text{ bit}$$

$$(10^9 = \text{Giga (GB)})$$

$$2^4 \times 2^{10} \times 2^{10} \times 2^{10}$$

$$\text{Size of main memory} = 16\text{GB} = 2^{34} \text{ bytes}$$

Number of Bits in Block offset

$$\text{Sum of bits in PhyAddress} = 34 \text{ bits}$$

$$\text{Block size} = 4\text{KB} = 2^{12} \text{ bytes}$$

$$\text{Number of bits in block offset} = 12 \text{ bits}$$

$$34 - (10 + 12)$$

Size of Cache Memory =

$$\text{Tot num of lines in cache} \times \text{Line Size}$$

$$2^{12} \times 4\text{KB}$$

$$2^{12} \times 2^{12} = 2^{10} \times 2^{10} \times 2^4 = 2^{24}$$

$$= [16\text{ MB}]$$

Problem # 06

$$\text{Cache Memory Size} = 8\text{KB}$$

Memory Size Needed At Cache Controller =

$$\text{BlockSize} = \text{FrameSize} = \text{LineSize} = 32 \text{ bytes}$$

Number of lines in Cache \times (1 Valid bit +

$$\text{Bits in Physical Address} = 32 \text{ bits}$$

1 modified bit +

19 bits to identify block

Bits in Blockoffset:

$$32 \text{ bytes} = 2^5$$

$$2^8 \times 2^1$$

5 bits

$$2^8 \times 2^1 \text{ bits}$$

5,376 bits

Option 4 is correct

Bit in Line Number:

$$\frac{\text{Cache memory size}}{\text{Line Size}} \Rightarrow \frac{8\text{KB}}{32 \text{ bytes}} \Rightarrow \frac{2^{13} \text{ bytes}}{2^5 \text{ bytes}}$$

$$2^8 \text{ lines} \Rightarrow 8 \text{ bits}$$

Bit in Tag:

$$32 - (8 + 5)$$

19 bits



Set Associative:

Problem-01: Consider a 2-way set associative mapped cache size of 16KB with block size 256 bytes - The size of main memory is 128KB - Find Number of bits in tag.

Sol: Given: SetSize = 2, Cache memory size = 16KB

$$\text{Block size} = \text{Frame Size} = \text{Line Size} = 256 \text{ bytes.} \quad P.\text{Ad} = 17 \quad \begin{array}{|c|c|c|c|} \hline \text{Tag} & \text{Line Set} & \text{B.off} \\ \hline \downarrow 5 & & \downarrow 8 \\ \end{array}$$

Bit of Phy Add :

$$\begin{aligned} \text{MM size} &= 128 \text{ KB} = 128 \times 2^{10} \\ &= 2^7 \times 2^{10} = 2^{17} \end{aligned}$$

Bit of Block offset:

$$\begin{aligned} \text{Block Size} &= 256 \text{ B} \\ &= 2^8 \text{ B} = 8 \text{ B} \end{aligned}$$

Bit of Phy Add = 17 bit

Total No. of sets in cache:

$$\begin{aligned} \text{CacheLine/setSize} &= 64/2 \Rightarrow 32 \text{ sets} \\ &= 2^5 \Rightarrow 5 \text{ bit} \end{aligned}$$

No. of line in cache:

$$\begin{aligned} \text{Cache/LineSize} &= 16 \text{ KB} / 256 \text{ B} = 16 \times 2^{10} / 2^8 \\ &= 2^4 \times 2^{10} / 2^8 = 2^6 \\ &= 64 \text{ lines.} \end{aligned}$$

$$\begin{aligned} \text{Tag Bits} &= \text{Phy Add} - (\text{Line S}) + (\text{B.off}) \\ &= 17 - (5+8) = 4 \text{ Bits} \end{aligned}$$

Problem-07: Consider a direct mapped cache with 8 cache block (0-7) - If the memory block request are in the order.

3, 5, 2, 8, 0, 6, 3, 9, 16, 20, 17, 25, 18, 30, 24, 2, 63, 5, 82, 17, 24.

Which of the following memory blocks will not be in the cache at the end of the sequence? (1) 3 (2) 18 (3) 20 (4) 30

Also calculate the hit ratio & miss ratio.

0	8, 0, 16, 24
1	9, 17, 26, 17
2	2, 18, 2, 82
3	3
4	20
5	5
6	6, 30
7	63

$$3/8 = 3, M$$

$$[H.R = 3/21] \quad [M.R = 18/21]$$

$$5/8 = 5, M$$

2/8 = 2, M 18 will not be in the cache.

$$8/8 = 0, M$$

$$0/8 = 0, M$$

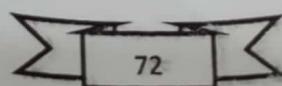
$$6/8 = 6, M$$

$$3/8 = 3, H$$

$$4/8 = 1, M$$

$$16/8 = 0, M$$

GO GETTERS



DATE ___ / ___ / 20___

Problem 08:

In 5 Block of cache → 7 of M Block

0	4,45
1	8,22
2	25
3	8
4	19,3
5	X,7
6	16
7	35

$$H \cdot R = 5/17 \quad M \cdot R = 12/17$$

Problem 09: $16/4 \Rightarrow 4$

0	0,48	$0 \cdot 4 = 0, M$
1	4,32	0 Set $255 \cdot 4 = 3, M$
2	8	$1 \cdot 4 = 1, M$
3	216,92	$4 \cdot 4 = 0, M$
4	1	$3 \cdot 4 = 3, M$
5	133	1 Set $8 \cdot 4 = 0, M$
6	129	$133 \cdot 4 = 1, M$
7	73	$159 \cdot 4 = 3, M$
8		$216 \cdot 4 = 0, M$
9		2 Set $129 \cdot 4 = 1, M$
10		$63 \cdot 4 = 3, M$
11		$8 \cdot 4 = 0, H$
12	285,155	$48 \cdot 4 = 0, M$
13	3	3 Set $32 \cdot 4 = 0, M$
14	159	$73 \cdot 4 = 1, M$
15	63	$92 \cdot 4 = 0, M$
		$155 \cdot 4 = 3, M$

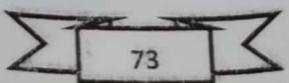
$$(H \cdot R = 1/17)$$

$$(M \cdot R = 16/17)$$

216 is not present
So D is correct



GO GETTER



Pipelining:

Sequentially it will take more time, while pipelined take less time / reduce time.

Throughput: No. of unit produced in a unit time.

I = instruction

stage $S_1 = IF$ (instruction fetch) Fetching instruction from memory.

$S_2 = ID$ (instruction decode) It decodes, read source reg & gen. generate control signals.

$S_3 = EX$ (Execution) Compute an R-type outcome.
Eg: Performing A+B

$S_4 = MFM$ (Memory) Read & Write the memory.

$S_5 = WB$ (Write Back) \rightarrow in HD Write back to hard drive
store a result in destination reg.

(1)	IF	ID	EX	WB	t_1	t_2	t_3	t_4	t_5	t_6	t_7	t_8	t_9	t_{10}	t_{11}	t_{12}	t_{13}	t_{14}	t_{15}
I_1	2	1	2	$3 = 8$	IF	I_1	I_1	I_2	I_3	I_4	I_4								
I_2	1	2	3	$2 = 8$	ID		I_1	I_2	I_2	I_3	I_3	I_4							
I_3	1	2	2	$1 = 6$	EX			I_1	I_1	I_2	I_2	I_3	I_3	I_4	I_4	I_4	I_4	I_4	
I_4	3	1	3	$2 = 9$	WB				I_1	I_1	I_1	I_2	I_3	-	-	I_4	I_4	I_4	

$$InPr = 31 \text{ O.C}$$

↳ sequential

$$P_T = 15$$

↳ After Pipeline.

Criteria for checking the performance of pipeline.

Speed-up: how much faster the pipeline execution is as compared to non-pipeline execution.

$$\text{Speedup} = \frac{\text{Non-pipeline execution time}}{\text{Pipeline exec time}} \quad S \cdot P = \frac{1}{P_T}$$

efficiency: Accurate result in less time (kam time mae kitna acha output data hai)

$$\text{efficiency}(n) = \frac{\text{Speed up}}{\text{no. of storage in pipeline Archtire}} \Rightarrow \frac{\text{Speed up}}{\text{Total no. of Phase}}$$

Throughput: how many task you achieve in a unit time.

$$\text{Throughput} = \frac{\text{No. of inst executed}}{\text{Total time taken}}$$

Non-Pipeline

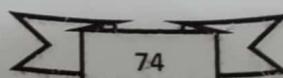
$$T = n * t_n$$

Pipeline (k stages)

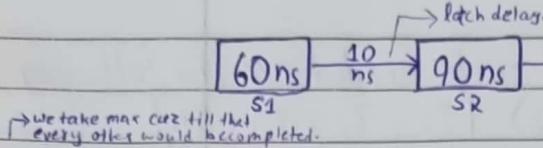
$$T_k = (k + n - 1) * t_p$$

t_p : Clock cycle (time to complete each suboperation)

$$\text{Frequency of the clock (f)} = \frac{1}{\text{Cycle Time}}$$



- (Q) Consider a pipeline having 4 phases with duration 60, 50, 90, 80 ns - Give latch delay is 10 ns - Calculate (1) Pipeline cycle time (2) Non-Pipeline exec time (3) Speedup ratio (4) Pipeline time for 1000 tasks (5) Sequential time 1000 task (6) throughput.



$$(1) P_t = \text{Max delay due to any stage} + \text{delay due to } i5 \text{ regx}$$

$$\max[60, 50, 90, 80] + 10$$

$$P_t = 100 \text{ ns}$$

$$(2) n P_t = \Sigma (\text{all stage})$$

$$= \Sigma (60w + 50w + 90w + 80w)$$

$$[n P_t = 280 \text{ ns}]$$

$$(3) \text{Speed-up ratio} = S.P. = \frac{n P_t}{P_t} \Rightarrow \frac{280}{100}$$

$$[S.P. = 2.8 \text{ ns}]$$

$$(4) \text{Pipeline time 1000 task}$$

P_t of 1 task + Time taken for 999 task

$1 \times 4 \text{ clock cycle} + 999 \times 1 \text{ C.C}$

$$4 + 999 \Rightarrow 4 \times 100 + 999 \times 100$$

$$400 + 99900 \Rightarrow [100300 \text{ ns}]$$

- (5) Sequential Time (Non-Pipeline)

$1000 \times n P_t$ for one task

$1000 \times \text{Time taken for one task}$

$$1000 \times 280 \text{ ns}$$

$$[280000 \text{ ns}]$$

- (6) Throughput : $T.p. = \frac{\text{No. of inst per unit time.}}{\text{executed.}}$

Total time

$$\frac{\text{Pipeline}}{100300} = 1000$$

$$\frac{\text{Non Pipeline}}{280000} = 1000$$

$$280000$$

- (Q) Consider an inst pipeline with four stages (S_1, S_2, S_3 & S_4) each with combinational circ only. The pipeline registers are required b/w each stages & at the end of the last stage. Delay for the stage & for the pipeline registers are as given in the figure.

$$n P_t = 30 \text{ ns}$$

$$P_t = \text{Max}[5, 6, 11, 8] + 1$$

$$= 11 + 1$$

$$P_t = 12 \text{ ns}$$

$$\text{Speed-up} = \frac{n P_t}{P_t} \Rightarrow \frac{30}{12}$$

$$= [2.5]$$

$$\text{Pipe exe time} = 1 \text{ clock cycle}$$

$$= 12 \text{ ns.}$$



GO GETTERS

$$\begin{aligned}
 & \frac{1.25 \times 10^3}{1.67 \times 10^3} - 100 \\
 & \frac{1.67 \times 10^3}{1.67 \times 10^3} - x \\
 & x = 133.6 \rightarrow \text{increase}
 \end{aligned}$$

DATE ___ / ___ / 20 ___

(Q.) Consider a 4 stages pipeline processor. The number of cycle needed by the four instructions I₁, I₂, I₃ & I₄ in stage S₁, S₂, S₃ & S₄ is shown below-

	IF	ID	EX	WB
I ₁	2	1	1	1
I ₂	1	3	2	1
I ₃	2	1	2	3
I ₄	1	2	2	2

What is the no. of cycle needed exc the following loop?
for(i=1 to 2) [I₁; I₂; I₃; I₄]

	t ₁	t ₂	t ₃	t ₄	t ₅	t ₆	t ₇	t ₈	t ₉	t ₁₀	t ₁₁	t ₁₂	t ₁₃	t ₁₄	t ₁₅	t ₁₆	t ₁₇	t ₁₈	t ₁₉	t ₂₀	t ₂₁	t ₂₂	t ₂₃
S ₁ I	I ₁	I ₁	I ₂	I ₃	I ₃	I ₄	I ₁	I ₁	I ₂	I ₃	I ₄												
S ₂ D		I ₁	I ₂	I ₂	I ₃	I ₄	I ₄	I ₁	I ₁	I ₂	I ₂	I ₃	I ₃	I ₄	I ₄								
S ₃ E			I ₁		I ₁	I ₂	I ₃	I ₃	I ₄	I ₄	I ₁	I ₁	I ₂	I ₂	I ₃	I ₃	I ₄	I ₄					
S ₄				I ₁		I ₂		I ₃	I ₃	I ₃	I ₄	I ₄	I ₁	I ₁	I ₂	I ₃	I ₃	I ₃	I ₄	I ₄	I ₄	I ₄	

Problem: 4

$$\begin{aligned}
 \text{1st Stage : } P_t &= \max\{800, 500, 400, 300\} + 0 \\
 &= 800 \text{ picoseconds}
 \end{aligned}$$

$$\text{throughput} = \frac{1}{800}$$

$$\begin{aligned}
 \text{2nd Stage : } P_t &= \max\{600, 300\} + 0 \\
 &= 600 \text{ Picosec}
 \end{aligned}$$

$$\text{throughput} = \frac{1}{600}$$

$$\begin{aligned}
 & \text{Throughput increase} \\
 & [(\text{Final throughput} - \text{initial throughput}) / \text{initial throughput}] \times 100 \\
 & [(\frac{1}{600} - \frac{1}{800}) / (\frac{1}{800})] \times 100 \\
 & \boxed{33.33\%}
 \end{aligned}$$

Problem: 9

IF, ID & WB stage take one clock cycle. EX depends on instruction, ADD & SUB need 1 clock cycle & MUL need 3 clock cycle in EX stage

	1	2	3	4	5	6	7	8
IF	A	M	S					
ID		A	M	S				
EX			A	M	M	M	S	
WB				A		M	M	S

8 clock cycle.



GO GETTE

Interrupts

→ Hardware → Maskable → Non-Priority based.
 → Non-maskable → Priority Based Interrupt
 → Software (exceptions)

- An unscheduled event that disrupts program execution is called exception.
- An exception that comes outside of computer is called interrupt.

In many computers exception are called interrupt but more precisely.

-) Exception is a signal from a program that requires OS to decide what to do next.
-) While interrupt is a signal from a device attached to computer.

Why Interrupts?

Device & program needs CPU but when? We can't predict.

Each device or program can give a signal to CPU in the form of interrupt.
 CPU can determine which program or device needs its attention.

Hardware Interrupts

↳ A signal to CPU through external device or hardware. Eg: key stroke, Mouse Movement.

Maskable: This type of interrupt can be delayed by the processor.

Non-Maskable: This type of interrupt can not be delayed by the processor.

Software Interrupt (Exception)

It is caused by a special instruction or exceptional condition.

E.g: Divide a number by zero, access to unavailable memory.

Interrupt Handler: A piece of code that is run as a result of an exception or interrupt.

Multiple Interrupts: An interrupt event when processor is already busy in handling other interrupt.



GO GETTER

(Q) Consider a main processor consist on 10 instructions & required 20ms mili sec for execution. Find out the size of processor, time required by for each instruction to executed, waiting time for 1 inst, & 3 inst & 6 inst.

MP

I ₁	Size of M.P = 10 inst
I ₂	T. MP = 20ms
I ₃	T/inst = $\frac{20}{10} = 2\text{ ms}$
I ₄	Waiting of I ₁ = 0ms
I ₅	" " I ₃ = I ₁ + I ₂ = 2 + 2 = 4ms
I ₆	" " I ₆ = 5 × 2 = 10ms
I ₇	
I ₈	
I ₉	
I ₁₀	

(Q) Using the above numerical, an interrupt occurs during the execution of third instruction of main process, the size of interrupt is 5 instruction. Each instruction of interrupt takes 1ms.

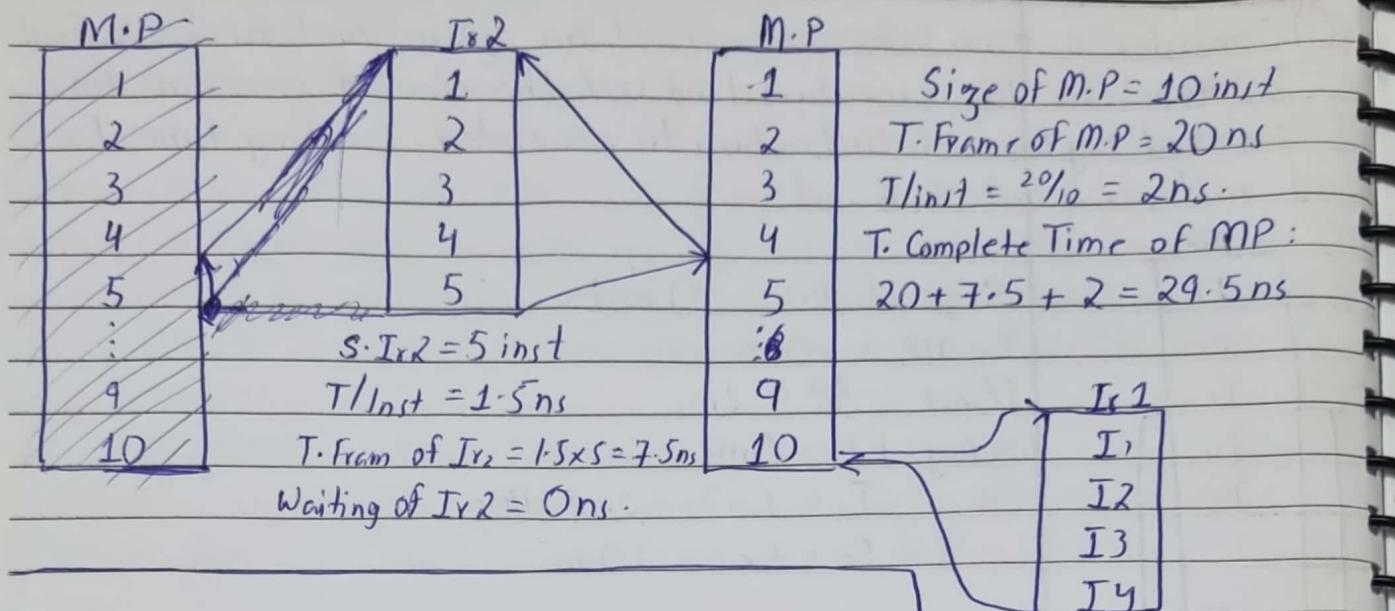
I.

I ₁	Size. Inst = 5 inst
I ₂	T/Inst = 1ms
I ₃	Total T = 5ms
I ₄	Waiting time I _{intup} = 0ms
I ₅	W.T of I ₆ = (5 × 2) + 5 = 15ms " of I ₄ = (2 + 2 + 2) + 5 = 11ms T. Comp. T. M.P = 20 + 5 = 25ms

(Q) Consider a process of 10 instruction in time frame require in 20ns. At the execution of 4th instruction 2 interrupt occurred - Interrupt number 2 is non-markable having 5 instruction, each inst takes 1.5ns for execution - Interrupt 1 consist of 4 insts with time frame 2ns - Do the necessary calculation.

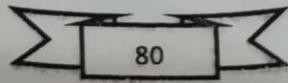
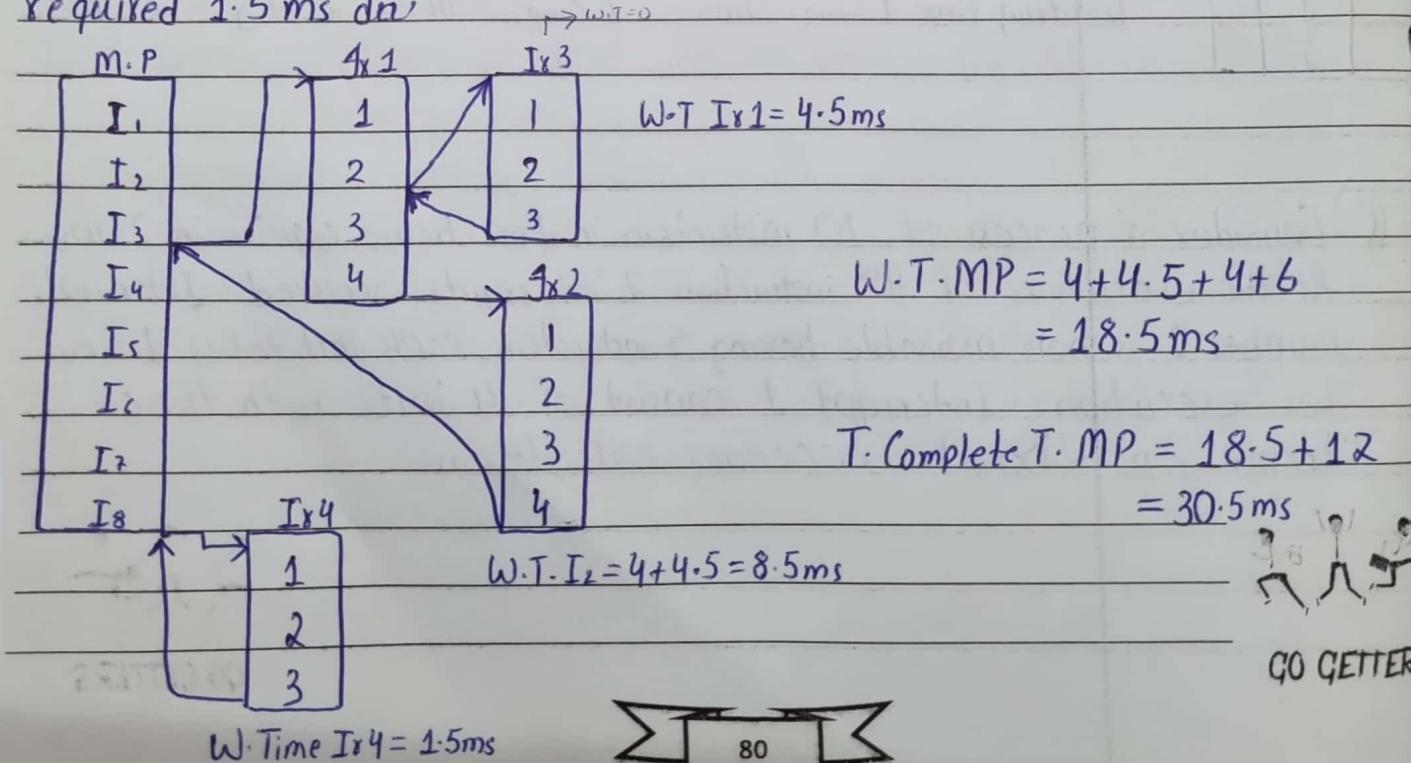


GO GETTERS



Q) Consider a M.P. of 8 inst each inst requires 1.5 ms during the execution of inst 3 two interrupt occur & both interrupt are non-maskable both are of same size having 4 inst with Timefram 4ms During the execution of 7th inst of M.P another interrupt occur with size 3 inst having T.F 6ms In the execution of interrupt 1, at second inst a mini interrupt ^{non maskable} occur with size 3 inst & each inst required 1.5 ms drv

$$\begin{aligned} S.I_{x1} &= 4 \text{ inst} \\ T.\text{Frame of } I_{x1} &= 2 \text{ ns} \\ T/\text{Inst} &= 0.5 \text{ ns} \\ W.T. I_{x1} &= 7.5 + 12 = \\ &= 19.5 \text{ ns.} \\ \text{Total Complete Time of } I_{x1} &= 19.5 \text{ ns} + 2 = 21.5 \text{ ns} \end{aligned}$$



GO GETTER

Pipeline Hazards: is a situation in which next instruction cannot execute in the following clock cycle.

There are 3 types of hazards.

① Structural Hazards: Cannot execute an instruction combination due to unavailability of a resource. If MIPS had only one memory - The same memory was used to: → fetch instruction → Load/store data.

1 IF ID Exe Mem WB

We have had a structural hazard with 4 insts

2 IF ID Exe Mem WB

→ First inst wants to get data from memory

3 IF ID Exe Mem WB

→ Fourth inst want to fetch inst from mem.

4 (IF) ID Exe Mem WB Without two memories, we will experience a structural Hazard.

Structural dependency occurs due to resource conflict, when different instruction collide while trying to access the same resource in the same segment of the pipeline, then it cause structural hazard.

S-Hazard are minimized using hardware technique is used called renaming, which splits the memory into two independent sub-modules to store instruction and data separation.

→ RAW

② Data Hazard: An occurrence in which a planned instruction cannot execute in proper clock cycle cuz data that is needed to execute the instruction is not yet available.

Data Hazard occurs due to data dependency, on instruction which is still in pipeline. Example: $Sum = a + b$

1 2 3 4 5

$Sub = Sum - c$] We cannot take sum as input until the sum operation is not completely performed.

IF ID EXF MEM [WB]

Sum insts writes value of sum in clock 5

IF [ID] EXE MEM WB

However sub insts reads it in clock 3.

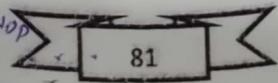
Solution: Insert nops, or bubbles , No inst fetched during nop

Three nop inserted, Not a good solution as waste of time

1 2 3 4 5

IF ID EXD MEM WB

Nop Nop Nop Nop Nop



81

IF ID EXF MEM [WB]

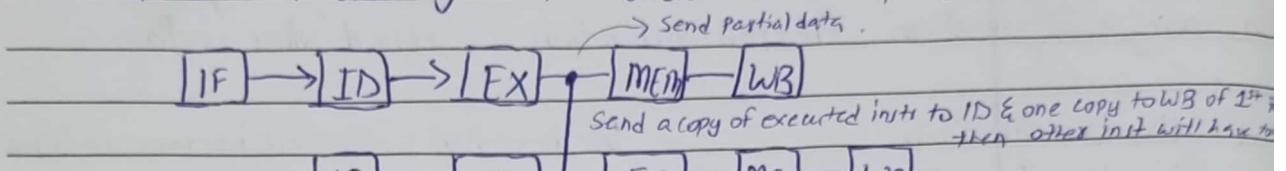


GO GETTERS

Data Hazard Sol

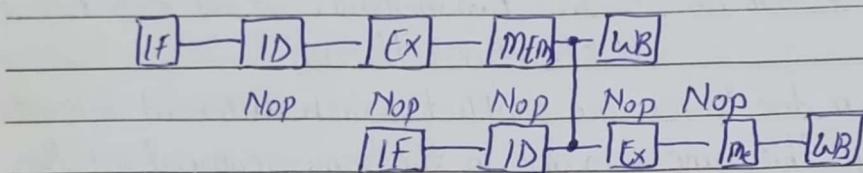
- ② Rely on optimizing compiler, can be used but not reliable
(as dependency can be data dependency as well)

Sol(3) Implement forwarding or bypassing - It is a method to resolve data hazard by retrieving the missing data element from internal buffer rather than waiting for it to arrive from programmer-visible register.



Forwarding can't solve all stalls, like this one

will take after memory so that a copy of data can be stored in this scenario we need to wait only one cycle



Lw = Load Word Sw = Store Word.

⇒ Reordering code to avoid Pipeline stalls.

Stall / Bubble required
As hardware must detect these cases and implement these stalls - We can't rely on software (comp) to resolve their issue

Control Hazard: Operation can't execute in proper clock cycle because the instruction fetched is not the instruction to be executed

Sol(1): Update ID phase in such manner that perform task in ID phase

⇒ Test register ⇒ Calculate branch address ⇒ Update PC

Even with this we need one stall.

Sol(2): Predict about the branch, when prediction is correct, continue as it when prediction is wrong, insert stall.

Sol(3) Delay branch decision.

Assembly lang: use mnemonics
to create instruction.

DATE ___ / ___ / 20___

Compiler: Convert code as whole batch in A.I

Word \rightarrow 16 bit size

Assembler: Convert that file in asm file

doubleword \rightarrow 32 bit

batch Processing / Packet switching: compiling as whole.

quad \rightarrow 4

Real time / circuit switching: On spot processing.

d-quad \rightarrow double of quad⁸

QuadWord \rightarrow 64

d-quadword \rightarrow 128

Hit & Miss

Occurrence of Hit - in a level when the address that an operation reference is found in that level of the memory hierarchy. Otherwise a Miss

Neither the hit rate nor the miss rate count the references that are handled by higher levels in the hierarchy.

For example: request that hit in the cache of our example memory hierarchy do not count in the hit rate or miss rate of the main memory.

Average Access time at a level

$$T_{avg} = (T_{hit} \times P_{hit}) + (T_{miss} \times P_{miss})$$

T = Time taken P = Probability.

Hit rate of the lowest level in the hierarchy = 100% (all requests that reach the bottom level are handled by the bottom level)

Starts at the bottom level & work up to compute the avg access time of each level in the hierarchy.

(Cache hierarchy level) \rightarrow 3 levels

(Q) Hit rate = 75%, $T_{hit} = 12\text{ns}$, $T_{miss} = 100\text{ns}$

$$\begin{aligned} T_{avg} &= (12 \times 0.75) + (100 \times 0.25) \\ &= 34\text{ns} \end{aligned}$$

Bottom to top

CPU

Cache

M.M → Virtual Memory

HD

1ms = 1000000ns

DATE ___ / ___ / 20___

(Q) If memory system contains the cache, main memory, & virtual memory.
the access time of the cache = 5ns , Cache hit rate = 80%.

the access time of the main memory = 100ns , Main Memory hit rate = 99.5%.

The access time of the virtual memory = 10ms

$$\text{Avg Access Time of M.M} = (100\text{ns} \times 0.995) + (1000000\text{ns} \times 0.05)$$
$$= 50,099.5 \text{ ns}$$

$$\text{Request to reach cache} = (5\text{ns} \times 0.8) + (50,099.5 \times 0.2)$$
$$= 10.024 \text{ ns.}$$

$$(A - A \div B \times B) \uparrow (A \bmod B)$$

DATE / /

Direct Map $\downarrow 17$

P1)	<table border="1"> <tr> <td>Tag</td><td>Line No</td><td>Block offset</td></tr> <tr> <td>3</td><td>2^8</td><td>8</td></tr> </table>	Tag	Line No	Block offset	3	2^8	8
Tag	Line No	Block offset					
3	2^8	8					

$$B.O = 256 \times 8 \Rightarrow 2^8$$

$$= 8 \text{ bits}$$

$$MM = 128 \Rightarrow 2^7 \times 2^{10}$$

$$= 18 \text{ bits}$$

$$\text{Line} = 2^4 \times 2^3 \Rightarrow 2^6$$

$$= 6 \text{ bits}$$

$$\begin{aligned} \text{Tag} &= 17 - (6+8) \\ &= [3 \text{ bits}] \end{aligned}$$

P4) S.A

(Q) 8-way set associative mapped cache, size of cache memory 512
10 bits in tag. Find size of main memory.

10	x	y
Tag	Line	B.off

$$\text{Size of M.M} = 10 + x + y$$

$$= 10 + 16$$

$$512 = x \times y \times 2^3$$

$$= 26$$

$$2^6 \times 2^9 = 2^n \times 2^8 \times 2^3$$

$$2^6 \times 2^9 \times 2^{10}$$

$$19 = x + y + 3$$

[QPRB]

$$x + y = 16$$

Q64MB

(P5)

(Q) 4-way SA cache. Size of M.M is 64MB

10 bits in tag, Find size of cache memory.

10	x	y
Tag	Line	B.off

$$2^{26} = 26 \text{ bits}$$

$$26 = 10 + x + y$$

$$\text{Cache Memory} = 2^x \times 2^y \times 4$$

$$x + y = 16$$

$$= x + y + 2$$

$$= 16 + 2 = 18$$

$$= 2^{18}$$

(Cache Mem = 256KB)

(P07)

0	8, 0, 16, 24	•
1	16, 8, 24, 17	3/21
2	2, 18, 8, 2	
3	3	•
4	20	18/21
5	5	•
6	30	
7	63	

(P8)

X45	1	H = 5/17
822	1	M = 12/17
25	11	
8	11	
193	1	
X7		
• 16	1	
35		

(P10)

0	808	SO	H = 1/5
1	12		M = 4/5
2			S1
3			

0
1
2
3
4

DATE ___ / ___ / 20___

SA

(P9)

0	048	
1	432	
2	8	
3	2492	
4	1	
5	133	S ₀
6	129	
7	73	
8		
9		S ₁
10		
11		
12	255155	
13	3	S ₂
14	159	
15	63	

(P11)

FIFO	542	11	542	11
	1261	11	1219	11
	13	11	13	11
	1719	11	1743	11

8443	54261	1111
12861	1219	1111
13	13	1111
1713	1743	1111

Hit FIFO: 5/15 Hit: LRU = 6/15

0	12412	110	12	10
1	15181718171861	111	12	11
2	2	111112	517124319	111111
3	174313	111113	1361	

Hit Direct Map: 1/15 Hit 2-Way SA = 5/15

(Q) Suppose you have a direct map cache with 8 blocks, if you receive following memory block request

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

Which of the following block will remain in cache at the end.

- (A) 18 (B) 30 (C) 3 (D) 20.

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

Q

$$3 \div 8 = 3,$$

$$5 \div 8 = 5$$

$$2 \div 8 = 2$$

$$8 \div 8 = 0$$

$$0 \div 8 = 0$$

$$6 \div 8 = 6$$

$$3 \div 8 = 3, H$$

$$4 \div 8 = 1$$

$$16 \div 8 = 0$$

$$20 \div 8 = 4$$

$$17 \div 8 = 1$$

A Option A will not remain in cache $25 \div 8 = 1$

Option B, C, D will remain in cache $18 \div 8 = 2$

$$30 \div 8 = 6$$

$$24 \div 8 = 0$$

$$2 \div 8 = 2$$

$$63 \div 8 = 7$$

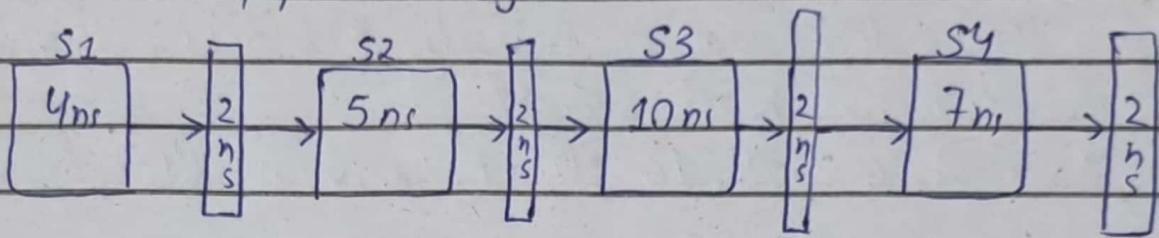
$$5 \div 8 = 5, H$$

$$82 \div 8 = 2$$

$$17 \div 8 = 1$$

$$24 \div 8 = 0, H$$

(Q.2) Consider a pipeline system



What will be the approx speed up.

Pipeline

$$\text{Max } \varepsilon(4+5+10+7) + \text{time}$$

$$10+2$$

$$12 \text{ ns}$$

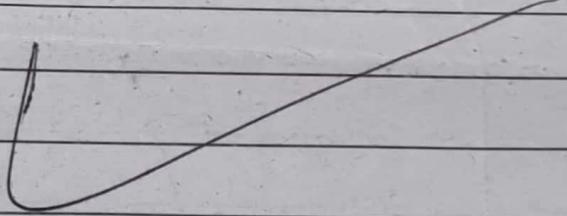
Non-Pipeline

~~Max~~

$$\varepsilon(\text{all stages})$$

$$\varepsilon(4+5+10+7)$$

$$26 \text{ ns}$$



$$\text{Speed up} = \frac{\text{Non-pipeline}}{\text{Pipeline}} = \frac{26}{12}$$

$$\begin{aligned} & \frac{13}{6} \\ & = [2.166 \text{ ns}] \end{aligned}$$