

Cache Mapping

Q1 Explain direct mapping, associative mapping & set associative mapping with brief.

Ans Direct Mapping

The simplest technique, which maps each block of main memory into only one possible cache line.
or In Direct mapping, assign each memory block to a specific line in the cache.

If a new line is previously taken up by a memory block when a new block needs to be loaded, the old block is trashed.

The cache is used to store the tag field whereas the rest is stored in the main memory.
Direct mapping's performance is directly proportional to the Hit ratio.

Formula used for mapping is :-

~~Random~~

~~Direct~~ (100)

Where: i = cache line number

j = main memory

$$i = j \% m$$

~~Cache replacement algorithm~~ ~~with block number~~
 m = number of lines in the cache.



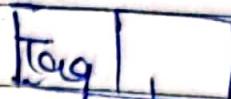
Associative / Fully Associative mapping

→ Associative mapping overcomes the disadvantage of direct mapping by permitting each main memory block to be loaded into any line of the cache. In this case the control logic interprets a memory address simply as a tag & a word field.

Adv: → fast & easy to implement.

→ we need replacement algo to maximize hit ratio.

- Disadv :- \rightarrow Complex circuitry req. to examine the tags of all cache blocks parallel.
- \rightarrow No tag comparison increases.



Set Associative mapping :-

- \rightarrow It is a mixture of direct mapping & fully associative mapping by arrangement lines of a cache into sets.
- \rightarrow The sets are persistent using a direct mapping scheme. However the lines within each set are treated as a small fully associative cache where any block that can scan to the set can be stored to any line inside the set.

Tag	# of sets	Block size
-----	-----------	------------

Fig: memory box.

Q.2 For a direct mapped cache design with a 32 bit address & byte addressable memory the following bits of the address are used to access the cache : Tag (22 bits), Index (5 bits), Offset (5 bits). Ans the following:-

Q. what is the cache line size (in words)?

Ans Given :- Offset = 5 bits & {block offset} = cache line size

Q. Size of cache line = 2^5 words

Since 1 word = 32 bytes \therefore 32 words

Q. How many entries (cache lines) does the cache have? Cache lines are also called blocks.

Ans Given Index Bits = 5 bits

Index bits = bits required to prevent the number of cache lines,

Q. Num of cache, when $2^5 = 32$ lines

- Q3) Consider a direct mapped Cache of size 16 KB with 16 blocks of size 256 bytes. The size of main memory is 128 KB. Find the number of bits in tag.

Ans) Given: Cache size = 16 KB

Bits required to represent cache size = 2^{10}

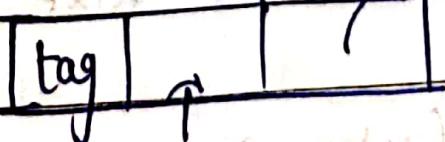
$$= 10$$

14 bits = tag

Block Size = 256 bytes
 Main memory size = 128 KB

Tag bits = ?

On direct mapping



of bits required for block size = 8 bits

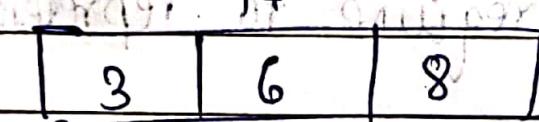
PA size (Ans bits) = 17 bits

of cache lines = Cache size \Rightarrow 16 KB

Ans bits in program from block sizes = 256 B

$$= \frac{2^{14}}{2^8} = 6 \text{ bits}$$

$$\therefore \text{Ans } \frac{2^8}{2^8} = 96 \text{ words}$$



\downarrow Ans bits = 17
Tag

∴ Tag bits = 3 bits

Q) Consider a direct mapped cache of size 512 KB with block size 1 KB. Then, what bits are the tag. And the size of main memory?

Ans

Given :- Cache size = 512 KB

Block size = 1 KB

Tag bits = 9 bits.

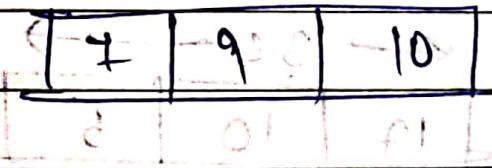
bits for block size = 10 bits
 # bits

$$\# \text{ of Cache Lines} = \frac{\text{Cache size}}{\text{Block size}} = \frac{512 \text{ KB}}{1 \text{ KB}} = 9 \text{ bits}$$

PA is denoted by :- of cache line



Tag bit .



$$\therefore \# \text{ of bits in PA (main mem)} = 7 + 9 + 10 = 26 \text{ bits}$$

$$\therefore \text{main mem size} = 2^{26} = 2^6 \cdot 2^{10} = 64 \text{ MB}$$

Q5) Consider a direct mapped cache of size 32 KB with block size 32 bytes. The CPU generates 32 bit address. Cal. the number of bits needed

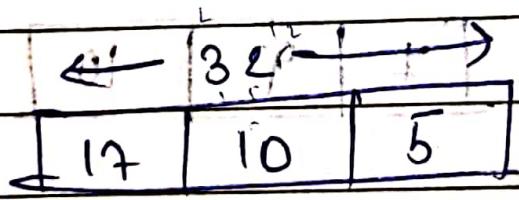
for cache indexing & number of tag bits?

Ans) Given Cache size = 82 KB
Block size = 32 bytes

$$\# \text{ of Cache lines} = \frac{82 \text{ KB}}{32 \text{ KB}} = \underline{\underline{1000}}$$

$$= 2^{10} = \underline{\underline{10 \text{ bits}}}$$

Given: Main memory is represented by 32 bits.



∴ Tag bits = 17 bits

Cache indexing bits = 10 bits

Q6) Consider a 2 way set associative mapped cache of size 16 KB with block size 256 bytes. The size of main memory is

128 KB. Find number of tag bits?

Ans

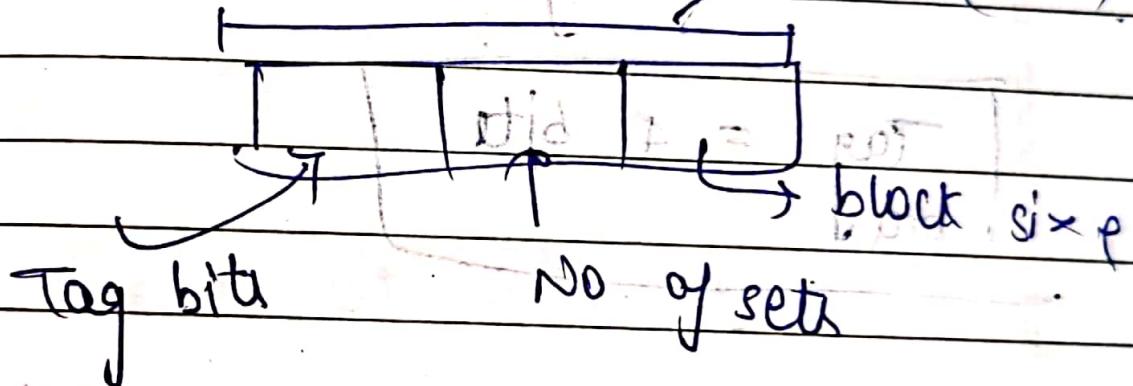
Given: Cache size = 16 KB

Block size = 256 bytes

Main Mem. size = 198 KB.

In set associative mapping:-

PA is denoted as :- $\text{PA} \rightarrow \text{PA (bits)}$



$$\text{PA in bits} = 2^7 \cdot 2^{10} = 17 \text{ bits}$$

~~$$\text{Number of cache lines} = \frac{\text{Cache size}}{\text{Block size}} = \frac{16 \text{ KB}}{256 \text{ B}} = 64$$~~

$$\begin{aligned} \# \text{ of cache lines} &= \frac{16 \text{ KB}}{256 \text{ B}} = \frac{2^{14}}{2^8} = 2^6 \\ &= 64 \text{ lines} \end{aligned}$$

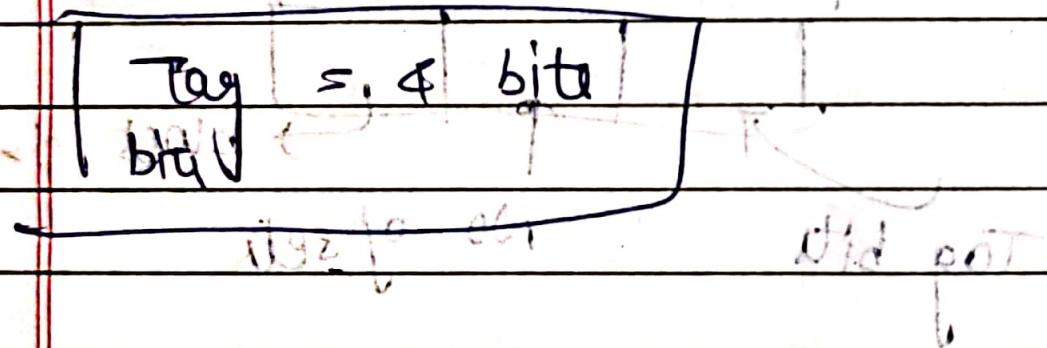
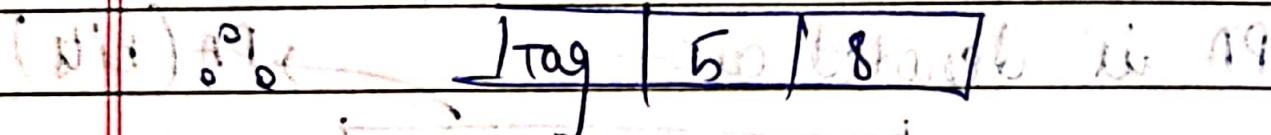
$$\# \text{ of sets} = \frac{64}{2} \quad \left\{ \begin{array}{l} \# \text{ of cache lines} \\ \text{set size} \end{array} \right\}$$

$$64 \times 16 = 9 \times 12 \quad \text{and} \quad 16 \times 16$$

$$64 \times 16 = 9 \times 12 \quad \text{and} \quad 16 \times 16$$

$$64 \times 16 = 9 \times 12 \quad \text{and} \quad 16 \times 16$$

of reg. to represent # of sets = 5 bit



Q7) Consider a 8 way set associative mapped Cache of size 512 KB with block size 1 KB. There are 7 bits in the tag. Find size of main memory.

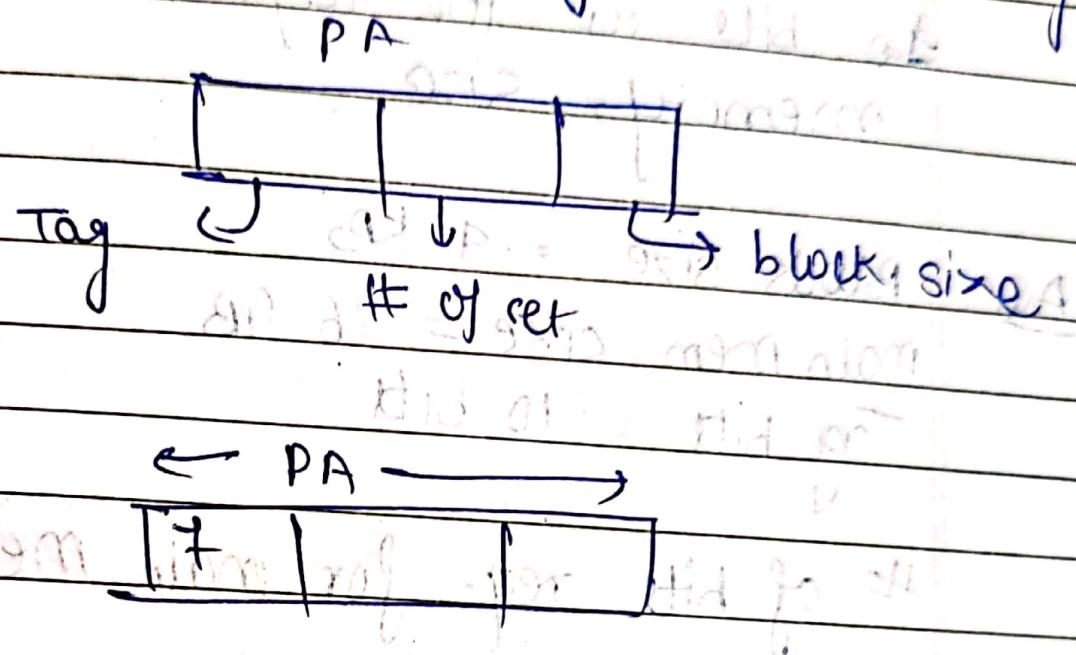
Ans: Given: Cache size = 512 KB

Block size = 1 KB

Tag bits = 7 bits

PA size = ?

In set associative mapping PA is denoted by:-

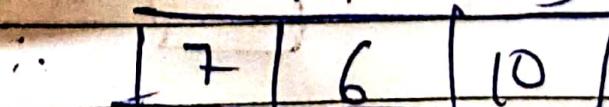


Block size in bits = ~~10~~ \otimes 1 KB = 10 bits

of cache lines = $\frac{512 \text{ KB}}{1 \text{ KB}} = \underline{\underline{512}}$

of sets = $\frac{\text{Cache Line size}}{\text{Set size}} = \frac{512}{8} = \underline{\underline{64}}$

No. bits reqd. for sets = $\alpha \sim 6$ bits



$$PA = 7 + 6 + 10 = 23 \text{ bits}$$

$$PA \text{ size} = 2^{23} = 8 \text{ MB} /$$

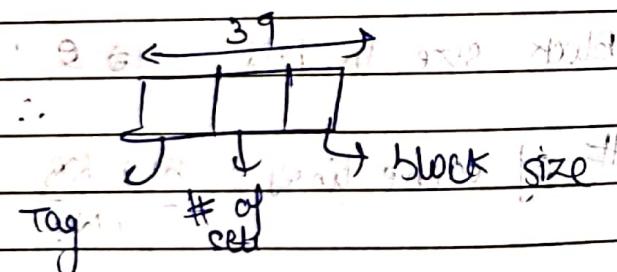
Q8) Consider a 4 way set associative mapped cache with block size 4 KB. The size of main memory is 16 GB & there are 10 bits in the tag. Find the cache memory size.

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

$$\text{main mem size} = 16 \text{ GB}$$

$$\text{Tag bits} = 10 \text{ bits}$$

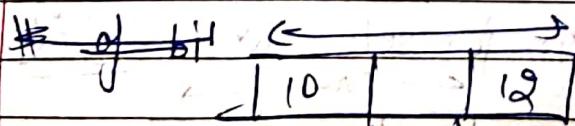
$$\# \text{ of bits req. for main mem} = 2^4 \cdot 2^{30} \frac{\text{bit}}{\text{bit}}$$



$$\# \text{ of bits for block size} = 2^2 \cdot 2^{10} = 2^{12}$$

$$= 2^{12} \rightarrow 12 \text{ bits}$$

$$39$$



$$\therefore \# \text{ of bits req. for } \# \text{ of sets} = 34 - 22$$

$$= 12 \text{ bits}$$

$$\therefore \# \text{ of sets} = 2^{12}$$

$$\# \text{ of sets} = \frac{\# \text{ of cache lines}}{\text{Set size}}$$

$$2^{12} = \frac{\# \text{ of cache lines}}{4}$$

$$2^{12} \cdot 2^2 = \# \text{ of cache lines}$$

$$2^{14} = \# \text{ of cache lines.}$$

$$\therefore \text{Cache mem. size} = \# \text{ of cache lines} + \text{Block size}$$

$$= 2^{14} \times 2^{12}$$

$$= 2^{26} = 2^6 \cdot 2^{20} = 64 \text{ MB}$$

Q9) Consider a 4 way set associative mapped cache with block size 4 KB. The size of main memory is 16 GB & there are 10 bits in the tag. Find the size of cache memory?

Given :- Block size = 4KB

main mem. Size = 16 GB

Tag bits = 10

To find. Size of cache memory = ?

Given: 4-way set associative map.

∴ the set associative mapping

main mem will be represented as :-

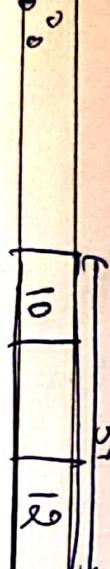


PA

$$\text{No. of bits req. to rep. PA} = 2^4 + 2^{10} = 34 \text{ bits}$$

$$\text{block size (in bits)} = 2^9 \cdot 2^{10} = 12 \text{ bits}$$

of cache blocks = .



$$= 34 - 10 = 24 \text{ bits}$$

solved Previously

Ques Consider a system with 2 level cache. Access time of level 1 cache, level 2 cache and main memory are 1 ns, 10 ns and 50 ns respectively. The hit rates of level 1 & level 2 caches are 0.8 & 0.9 respectively.

what is the avg access time of sys. ignoring the search time within the cache?

$$\text{Ans. Level 1 Access Time (H)} = 1 \text{ ns.}$$

$$\text{Level 2 Access Time (L)} = 0.8 \text{ ns.}$$

level 2 access time (L) = 10 ns

hit ratio of L = 0.91

Mem. access time = 50 ns

EAR = ~~(1 - H1)(1 - H2) + H1(1 - H2) + H1H2~~

~~(1 - H1)(1 - H2) + H1(1 - H2) + H1H2~~

$$\text{EAR} = (1 - H_1)(1 - H_2) + H_1(1 - H_2) + H_1H_2$$

$$= (0.9)(0.1) + (0.2)(0.9) + (0.2)(0.1)$$

Pipelining

- Q) What is arithmetic pipeline & instruction pipeline?
- Ans) Arithmetic Pipeline

- Arithmetic pipeline divides an arithmetic problem into various sub-problems for execution in various pipeline segments.
- Its used for floating point operation, multiplication & various other operations.

Instruction Pipeline

- It is a sequence of instructions which can be executed by overlapping fetch, decode & execute phase of an instruction cycle.
- It is divided into four stages:
- Fetch, Decode, Execute & Write back.

This type of technique is used to increase throughput of computer system.

An instruction pipeline needs data from the memory while previous part are being executed in other segments of the pipeline. Thus we can execute multiple instructions simultaneously.

- Q) Explain floating point addition stages using pipeline.
- Ans) Floating point addition is performed by arithmetic pipeline.

The add operation is being divided to various suboperation as follows:-

1) Compare the exponents by subtraction
Align the mantissas
Add the mantissas
Normalize the result.

Let's take an example :-

Step 3: Add mantissas :-

$$\text{let } X = A \times 10^a = 0.9509 \times 10^3 \\ Y = B \times 10^b = 0.8200 \times 10^2$$

$$X = 0.9509 + 0.8200 \times 10^3 \\ X = 1.0329 \times 10^3$$

where $a > b$ = mantissa

$a - b$ = exponents

Step 4: Compensating exponents by subtraction :-

The diff of exponents: $3 - 2 = 1$
determines how many times the mantissa associated with the smaller exponent must be shifted to right.

Step 5: Align the mantissas :-

The mantissa associated with smaller exponent is shifted according to diff of exponents determined in Step 4.

$$\therefore X = 0.9509 \times 10^3$$

$$Y = 0.8200 \times 10^2$$

Step 6: Normalize the result:-

$$X = 0.9509 \times 10^4$$

Explains the following terms: latency, throughput

throughput in ideal or max condn, speedup ratio, speedup time in ideal or max condn, efficiency.

Latency :- After how much time the new input

is given to the system.

Throughput :- Number of operation performed unit time.

$$\text{Throughput} = CK + n - 1 \cdot HP$$

In ideal case :- $K \ll n$ in other phone ($K \neq 1$)

$$\text{Throughput (ideal)} = \frac{1}{t_p^2}$$

Speed up ratio :- Performance of pipeline
is given by speed up ratio.

$S = \frac{\text{Non Pipeline Processing time}}{\text{Pipeline Processing time}}$

$$S = \frac{t_n}{t_p} = \frac{\text{Stage width}}{(C_k + n - 1)t_p}$$

Efficiency :- It's the %age of time when processor are busy over total time taken.

→ In ideal cond'

$$= n \cdot \frac{1}{t_p} \cdot K \ll n$$

$$= [K + n - 1] \cdot \frac{1}{t_p} \cdot K \ll n$$

Ideal

Q4 Explain Instruction Hazards & their solutions with neat & clean diagrams

i) Resource Conflict Structural Hazard :-

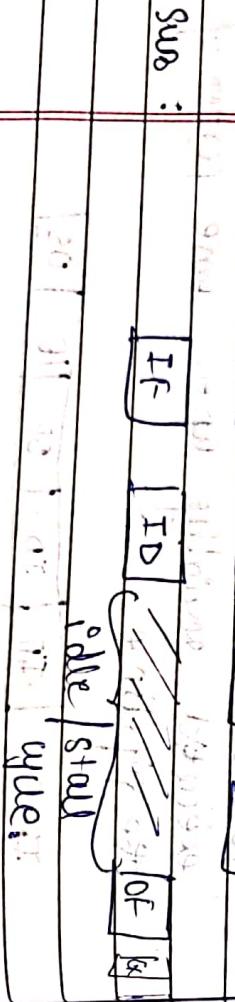
When 2 pipeline instructions or even more want to access the same resource it results in resource hazards. It is also known as structural hazard. A solution to this hazard is that there must be an overlap serially upto some portion of the pipeline.



Data hazard :- It's a cond' when accessing an operand location creates conflict. Consider that you have 2 insts :-
 ADD R₁, R₂, R₃
 SUB A₁, R₁.

obs. the sub. the result of add is stored to R₁ after execution.

R₁ act as an operand for subd. inst.



The sub. inst. can fetch the operand at t₃ but others. It depends on R₁. & R₁ will be updated by ADD at t₄. Sub has to wait for idle until it is executed.

Control hazard / Branch difficulty :- occurs when pipelining fails to predict branches. i.e. the first resulting instruction entering the pipeline that must be discarded. A control hazard is often referred to as a branch hazard.

Q5) Show Prediction = 0011 0111 1000
 Q6) Dynamic Branch Prediction
 Q7) Recording Instructions
 Q8) Pipeline stalls due to time = 10000 items / 1000 ns = 10 ms

A 4 stage pipeline has the stage delays as 150 ns, 120 ns, 160 ns and 100 ns respectively. Registers that are used between the stages have a delay of 5 ns each. Assuming constant clocking rate, the total time taken to process 10000 data items in this pipeline will be _____ ns?

Ques

Given: A pipeline with four stages and three registers.



P₁: Four stage with latencies 1ns, 2ns, 2ns, 1ns
P₂: 1ns, 1.5ns, 1ns, 1ns

Total time taken to process 10000 data

$$\text{Steps} = 8$$

$$10000 \times 8 \times 1 \text{ ns} = 80000 \text{ ns}$$

$$P_4: 5 \text{ stage pipeline with latencies } 0.6 \text{ ns}, 0.6 \text{ ns}, 0.6 \text{ ns}, 0.6 \text{ ns}, 0.6 \text{ ns}$$

Total cycle time = $(K+n-1)tp$

where $K = \text{no. of stages}$, $n = \text{no. of data}$

$$n = \text{no. of data} = 1000$$

$$tp = (\text{sum of all stage delay}) +$$

initial setup time. (Req. Delay) $\approx 2 \text{ ns}$

Any lets assume there are n task for each process.

$$tp = (160 + 5) \text{ ns}$$

$$\therefore \text{Total cycle time} = (4 + 1000 - 1) + 165$$

$$= 103 \times 165 \text{ ns}$$

$$= 16995 \text{ ns}$$

$$\therefore \text{Cycle time for P1: } tp = 2 \text{ ns}$$

$$\text{Given: } K = 4$$

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

$$\therefore \text{Cycle time} = (K+n-1)tp$$

$$= (4+1) \times 2 \text{ ns} / 1$$

Cycle time for P₂ :-

Given : K = 4

$$t_p = 1.1 \text{ ns}$$

We get that the cycle time of P₂ is at least

$$\text{Cycle time} = (K + n - 1) t_p$$

= 11 ns

Cycle time for P₃ :-

$$t_p = 1.1 \text{ ns}$$

Given : K = 5

$$\text{Cycle time} = (K + n - 1) t_p$$

$$= (5 + 4 - 1) \cdot 1 \text{ ns}$$

Q3

Explain all the reasons of stall cycle.

Cycle time for P₄ :-

$$t_p = 1.1 \text{ ns}$$

Given : K = 4

What is data dependency & possible solution for data dependency?

Cycle time = $(K + n - 1) t_p$

Ans
When an instruction is modified in different stages of a pipeline with the

help of instructions that data dependency in this case, the data hazard will occur.

→ When instructions are read/write, the registers that are used by some other instructions in this module, the hazard will occur resulting in the delay in the pipeline.

The WAR & WAW hazard will not cause the delay if a processor uses the same pipeline for all the instructions & executes these instructions in the same order in which they appear in a program.

Consider a pipelined processor with the following 4 stages, IF, ID, EX, WB. The IF ID WB stage take one clock cycle to complete the operation. The number of cycle for the EX stage depends on the instruction. ADD & SUB need 1 clock cycle and MUL need 2 clock cycle. Pipeline do not use any operand forwarding. What is the number of cycle take to complete the following sequence of instruction.

→ RAW (Read After Write)
WAR (Write After Read) or读后写
WAW (Write after Write) 或写后写

The WAR is very uncommon or impossible.

To avoid RAW & WAW hazard processor should use different register to generate the output of each instruction. A number of registers

Date _____
Page _____

Date — / — /
Page —

APPENDIX 2. *W. B. H. G. 16. 2. 9. 8.*

MUL 1F ID EX EX GND.WB

1996-1997
Year 1996-1997

Total clock cycles required is 8.

Q19 What (PPI & Consid'') = Stress

Q19 What CPU? Consider a 5-stage pipeline which is executing a program of 100 instructions. Among all the instructions cause 3 stall cycles each.

Calculate CPI of Pipeline

- If pipeline cycle time is say then what is avg. exec. time - Cal. CPI and pipeline in ideal cond'.
- Offered load? If pipeline cycle time is one then what is avg. first-execution time in ideal cond'.

CPI stands for cycle per Instruction

Given: $K = 5$, $n = 100$, $\tau = 30(10 + 3)$

$$CPI = \frac{K + n - 1}{n} = \frac{5 + 100 - 1}{100} = \frac{104}{100} = 1.04$$

Given: $tp = 3\text{ns}$

$CPI_{avg} = \frac{k + n - 1 + \alpha}{n} = \frac{5 + 100 - 1 + 30}{100} = \frac{139}{100} = 1.39$

Avg exec. time = $CPI_{avg} * tp = 1.39 * 3 = 4.17\text{ns}$

Q11 CPI during 1st 100 ns = 130
 Ideal latency would be 100 ns
 Avg Enq time in ideal cond' = CPI * t_p

$$= 1.03 \times 3 \text{ ns}$$

$$= \underline{\underline{3.09 \text{ ns}}}$$

Non linear pipeline allows read forward
 read back connections in add to
 stream line connections.

Reservation Table:

It displays the time space flow of data
 through the pipeline for one func'
 evaluation.

- Q11 What is non-linear pipeline reservation
 table, permissible latency forbidden
 latency, collision vector & trace diagram
 Example: cycle & greater cycle & cycle time

Any Non linear Pipeline

Permissible latency :-
 Latency that will not cause collision

- Non linear pipelines are dynamic pipeline
 because they can be very hard to
 perform valuable function diff
 times.

Collision vectors :- It is the representation of forbidden latency & permissible latency together.

Simple cycle :- A simple cycle in a latency cycle in which each state appears only once.

Greedy cycle :- A simple cycle whose edges are all made with max latency from the respective states touching states.

Q21 Consider the following execution table:

Stages	1	2	3	4	5	6	7	8
S1	X				X			X
S20		X		X				
S33			X		X			X

Now calculate permissible latency, forbidden latency, collision vector, state diagram, simple cycle & greedy cycle.

Collision vector $C_1 = [0 \ 1 \ 0 \ 1 \ 1 \ 0 \ 1 \ 0]$
 $L_{\text{forbidden}} = [1 \ 3 \ 6 \ 8 \ 10 \ 12 \ 14 \ 17]$ Permissible = 0

State diagram :- (01011010)
 $\begin{matrix} C_1 \\ C_2 \\ C_3 \\ C_4 \\ C_5 \\ C_6 \\ C_7 \\ C_8 \end{matrix}$

Shift right $0 \ 0 \ 1 \ 0 \ 1 \ 1 \ 0 \ / \ 1 \ 0 \ 1 \ 0 \ 1 \ 1 \ 1 \ 0 \rightarrow$ diff. bits as tell

C_2 as subscript is 3.
 C_3 here we need to shift right 3 times

$\begin{matrix} S_1 \\ S_2 \\ S_3 \end{matrix} \rightarrow 0 \ 0 \ 1 \ 0 \ 1 \ 1 \ 0$

$\begin{matrix} S_1 \\ S_2 \\ S_3 \end{matrix} \rightarrow 0 \ 0 \ 0 \ 0 \ 1 \ 1 \ 0$

state

we need to shift these

16 times 6 minutes 1966 Aug 10

卷之三

0 0 0 0 0 0 0

卷之三

(continued) → same as G

卷之三

○○○○○○○○

卷之三

1894-1895

卷之三

8
7
6
5
4
3
2
1

ANSWER

W
or

Date — / — /
Page —

Digitized by
Date — / — /

Simple cycle $\rightarrow \{1, 93, 83, 3, 83, 6, 83, 3\}$

Greedy cycle \rightarrow {33}, {1, 89}