

Types of Instructions

① Data transfer Instructions

→ Use to move the data from

-) Register to m/m
-) m/m to register
-) register to register

② Data manipulation Instructions

-) Arithmetic opⁿ
-) Logical opⁿ
-) Shift opⁿ

③ Program Control Instructions

-) Use to make decisions
-) execution transfer statements

Data transfer instructions

We have 3 type of components

-) memory •) register •) I/O device

Generally I/O devices also consider as a registers.

① Load: Load instruction from m/m to register.

→ Generally use in loading accumulator.

② Store: register to m/m store

→ Use in case of Accumulator

③ Exchange: register to register or memory to registers.

→ Just like swap.

④ Input / output : Put input

data in register and from register take out as outputs

⑤ Push / Pop: for stack implementation.

Each instruction can generate many instruction on the basis of mode and operands.

Ex Load can be modified as -

LDI → Load Immediate

LD R, X → Direct addressing mode

Arithmetic Instruction

Basic operation → Addition, Subtraction
We can implement any operation using above two.

Increment ← Doesnt require CPU
Decrement ←

Multiply , Divide , Add with carry
Subtract with borrow, Negate.

Each instruction can have various execution on the basis of mode and addressing operands.

Negate gives 2's complements

Logical Instructions

Shift Instructions

Program Control Instructions

copied from book

Logical Instructions

→ Operations work on bits called logical operations.

Logical 'AND' or 'OR'. (A/V)

→ 'X-OR' operations.

→ 'Clear' done at register level. It means change all bits to '0'.

→ 'complement' means 1's complement.

XOR of any no. with all 1's gives 1's complement.

→ carry is a flag (1 bit data)

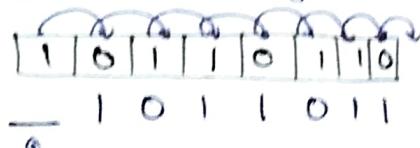
→ clear carry / set carry are operations works on flag bit.

→ Complement carry instruction is useful in implement subtraction.

→ Enable / Disable Interrupts opns. have interrupt flag separately.

Shift Instructions:

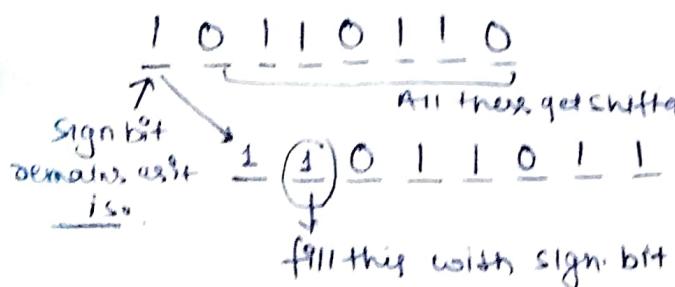
① Logical shift Right / Left



→ fill it with zero, and this zero is not due to LSB. It's just to fill the vacant bits.

Similarly works "Logical shift left."

② Arithmetic shift Right / Left



③ Ex:

01001101

After >> Shift

00100110

Arithmetic and Logical Left Shift both are same.

Shift instruction has

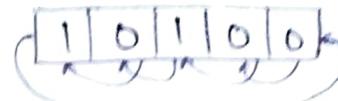
Opcode	Type of shift	Opand	new many shifts	Count mode

Rotate Right / Left



rotate right

01010 → o/p



rotate left

010110 → o/p

Rotate right / left through carry



rotate right through carry

→ we have an extra carry bit.

As right shift we can do left shift also.

Program Control Instructions

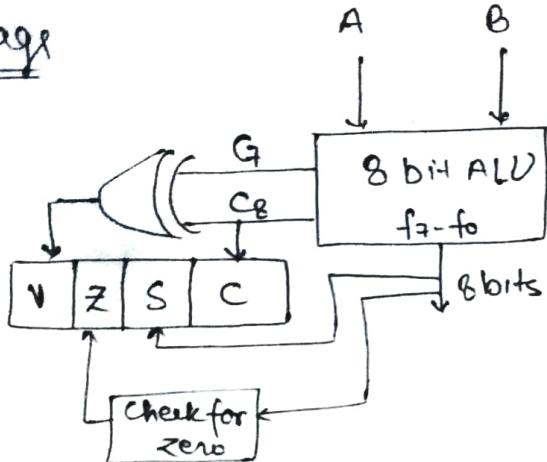
Branch, Jump, skips, (BR) (JZ, JNZ)

Call, return, compare, test

→ Skip blocks to ignore the next instruction.

→ Compare and test are similar kind of operations.

flags



V → Overflow flag

Z → Zero flag checks result
is zero or not.

$C_7 \oplus C_8 = 1$ means overflow.

→ We can also verify the value
of particular bit using zero
flag and test instruction.

S → Sign flag.

C → Carry flag.

On A - B	
<u>Signed</u>	<u>Unsigned</u>
$A = B, Z = 1$	$A = B, Z = 1$
$A > B, S = 1$ $Z = 0$ $V = 0$	$A > B, C = 1$ $Z = 0$
Otherwise $A < B$	

Flags can be used in comparing
two numbers.

Example

→ Generally we perform subtraction using adder.

$$A = 1100$$

$$B = 0010$$

$$A - B = ?$$

$$\begin{aligned} 2\text{'s comp. of } B &= 11.01 + 0001 \\ &= 1110 \end{aligned}$$

$$A - B = A + 2\text{'s comp. of } B$$

$$= \boxed{\quad} 1010$$

Carry \downarrow → MSB

Program status word

Z	V	S	C
0	0	1	1

$$V = C_{in} \oplus C_{out} = 1 \oplus 1 = 0$$

So, now we can compare A & B
using all the flags values.

$$A = 1100 \mid 12 \mid -4$$

$$B = 0010 \mid 2 \mid 2$$

\downarrow
Unsigned
 $(A > B)$

\downarrow
Signed
 $(A < B)$

In signed values if,

$$\textcircled{1} \quad V = Z = 0, S = 1 \quad (A < B)$$

In Unsigned number,

→ we don't care about overflow

$$\Rightarrow Z = 0, C = 1 \quad (A > B)$$

$$\text{Ex} \rightarrow A = 0010 \mid 2 \mid 2$$

$$B = 1100 \mid 12 \mid -4$$

$$\begin{aligned} A - B &= 0010 + 0100 \\ &= \underline{0110} \end{aligned}$$

PSW

0	0	0	0
Z	V	S	C

Signed: $V = Z = 0, S = 0 \quad (A > B)$

Unsigned: $Z = 0, C = 0 \quad (A < B)$

→ Carry is complement of borrow
It means if $C=0$ then $B=1$

Summary

Signed Number

V	Z	S	Remark
1	X	X	Invalid
0	1	X	$(A=B)$
0	0	0	$(A>B)$
0	0	1	$(A<B)$

Unsigned Number

Z	C	Remark
1	X	$A=B$
0	0	$(A < B)$ $\text{if } C=0 \Rightarrow B=1$ in 2's complement subtraction
0	1	$(A > B)$ (No borrow)

Conditional Branch Instruction

BZ : Branch if zero ($Z=1$)

BNZ : Branch if no zero ($Z=0$)

BC : Branch if carry ($C=1$)

BNC : Branch if no carry ($C=0$)

BP : Branch if positive ($S=0$)

BM : Branch if negative ($S=1$)

BV : Branch if overflow ($V=1$)

BNV : Branch if no overflow ($V=0$)

Unsigned

BHI → Branch if higher ($A>B$)

BNE → — higher or equal ($A \geq B$)

BLT → Branch if lower ($A < B$)

BLE → Branch if lower or equal ($A \leq B$)

BE → Branch if equal ($A=B$)

→ $(A \neq B)$

Signed

BGT → Branch if greater ($A > B$)

BGE → — greater or equal

BLT → — lower ($A < B$)

BLE → — lower or equal

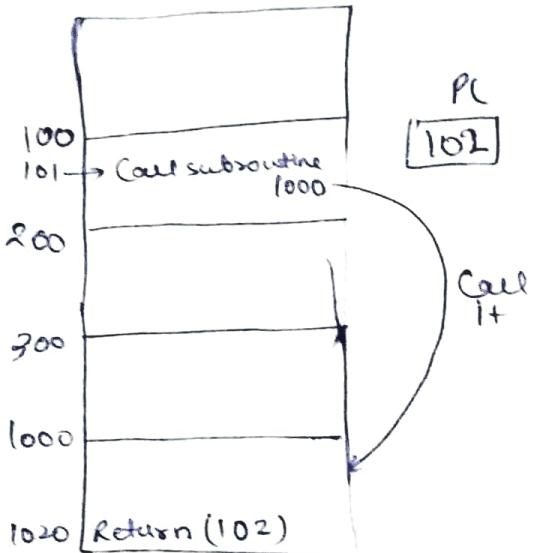
BE → $A=B$

BNE → $(A \neq B)$

Call and return instructions

→ Both are program control instructions.

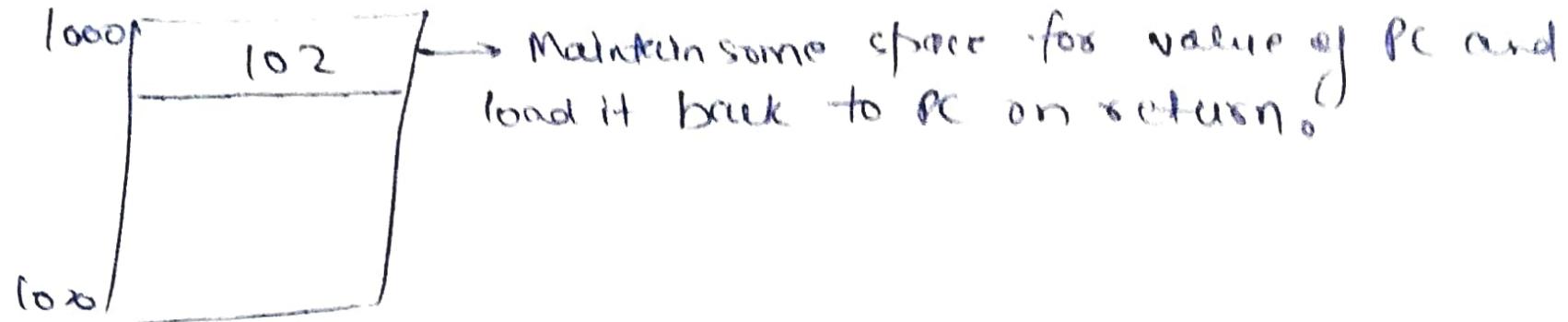
→ They manage the flow of programs



After completion we get return to 102 but if during execution of subroutine calling function get shifted to 200 then (102) add. creates inconsistency b/c 102 we need to return on address 202. for this we have soln like

- ① Store the value of 102 to temp. register and when return back load the value of temp to PC again but this was using long back and it was not much appropriate soln.

② Store in subroutine like



③ Store in any fixed memory location.

④ Store in ~~mem~~ stack and ^{it is the best method known till now}

Note → Read this topic from book too or from any other YT Video,
if we get issue while solving problem.

Pipelining

→ We can execute more than one instruction at a time using one CPU.

Operations perform on instruction -

- ① Instruction fetch
- ② Instruction Decode
- ③ Execute
- ④ Write Back.

→ Each instruction takes some time in terms of clock cycles

Suppose an instruction performs all four opn in 4 units then
for 3 instructions we require $3 \times 4 = 12$ units bcoz we
have only single CPU but we can do it in less time using
pipelining.

Phase time diagram

→ We are trying to manage CPU in such a way it get divided
into 4 parts and each will perform single operation so
that we can achieve parallelism.

Clock	1	2	3	4	5	6	7	8	9	10	11	12
IF	I ₁	I ₂	I ₃	I ₄	I ₅	I ₆	I ₇	I ₈	I ₉	-	-	-
ID	-	I ₁	I ₂	I ₃	I ₄	I ₅	I ₆	I ₇	I ₈	I ₉	-	-
EX	-	-	I ₁	I ₂	I ₃	I ₄	I ₅	I ₆	I ₇	I ₈	I ₉	-
BW	-	-	-	I ₁	I ₂	I ₃	I ₄	I ₅	I ₆	I ₇	I ₈	I ₉

Here, we complete even more than 4 instn in 12 clocks.

CPI → Clock per Instruction

Main aim of Pipelining is to keep, CPI ≤ 1 which is ideal.

1 Instruction = 1 Clock, 1 Instr. = $\frac{1}{f}$ sec, $\Rightarrow \frac{1}{f} \text{ sec} = f \text{ Instr.}$

Performance is measured using utilization or efficiency and speedup.

Load = # of instruction that can be perform in given time.
In above example Load = 9 I

$$\text{efficiency} = \frac{36}{12 \times 4} = \frac{3}{4} \frac{\text{(Occupant box)}}{\text{(Total box)}}$$

$$\text{Speedup (load is given)} = \frac{T_{\text{without pipeline}}}{T_{\text{with pipeline}}} = \frac{9 \times 4}{12} = 3$$

$$m \times \text{efficiency} = \text{speed up}$$

$$\underline{m \times e = s}$$

Where, 'm' is no. of operations required to perform an operation.

→ In practice we will never get 100% efficiency, because some blocks are empty in tables ($S < m$)

→ but in theory ($S \approx m$) nearly equal.

Suppose we execute 'n' instructions then time taken in

$$\begin{aligned}\text{Pipelining} &= 1 \times m + (n-1) \times 1 \\ &= \frac{m + (n-1)}{1}\end{aligned}$$

$$\text{without pipeline} = \underline{m \times n}$$

$$\text{Speed up} = \boxed{\frac{m \times n}{m + n - 1}}$$

$$= \frac{m}{1 + (m-1)/n}$$

Suppose, $n \rightarrow \infty$ very large

$$\text{speed up} \approx \frac{m}{1} \quad \text{In theory}$$

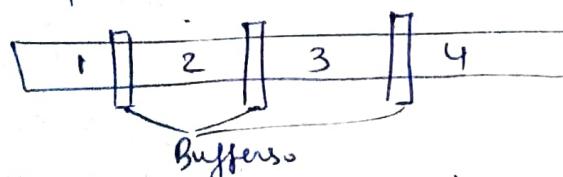
So, T_{wp} (without pipeline) $< m \times T_p$ → Just remember only,

Speed up → how fast we perform with pipelining?

So, second aim of pipelining is to keep $\boxed{S \approx m}$ ideal.

In practical value of 'm' can lie between 3 → 24 but various costs and delay add up in between. So generally we go for 5 (experimentally).

In non-pipeline mode, CPU acts as a single unit, but in pipeline mode we divide CPU to various stages to perform various operations as suppose each instruction is performing 4 opns then,



we store the result of previous opn in buffers

→ It is not always possible that all the opns takes equal time so it might provide some delays in loading to buffer or loading from buffer. So we even say that

$$\boxed{T_p \geq T_{wp}}$$

→ for a single instruction
not for various instructions.

Gates 99

	F	D	E	W
1	1	2	1	1
2	1	2	2	1
3	2	1	3	2
4	1	3	2	1
5	1	2	1	2

- Soln
- There are four operations and 5 instructions
 - Each instruction doesn't take equal no. of clock cycles
 - Don't start two instruction operation of same type at the same clock cycle

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
I ₁	F	D	D	E	W										
I ₂	F	X	D	D	E	E	W								
I ₃	F	F	X	D	X	E	E	E	E	W	W				
I ₄		F	X	D	D	D	D	X	E	E	W				
I ₅		F	X	X	X	D	D	D	D	X	E	W			

complete I₁

complete I₂

complete I₃

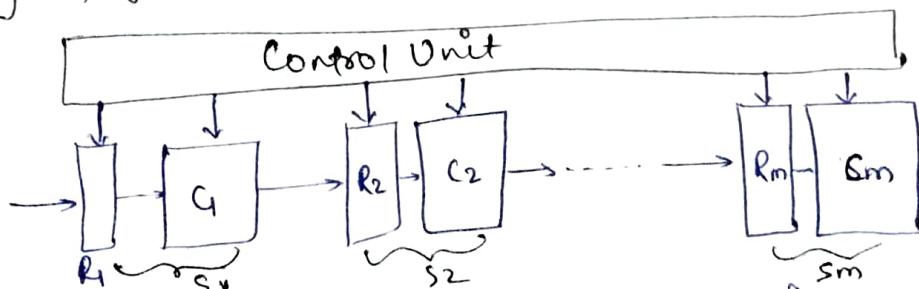
complete I₄

complete I₅

→ Instead of 15 we can take more than 9 but it will waste of space.

→ Choose max. value of cycle (15 here) as the max. value given in option.

Designing of pipeline (Hardware)



There are 'm' stages or operation that is managed by control unit with the help of clock signals.

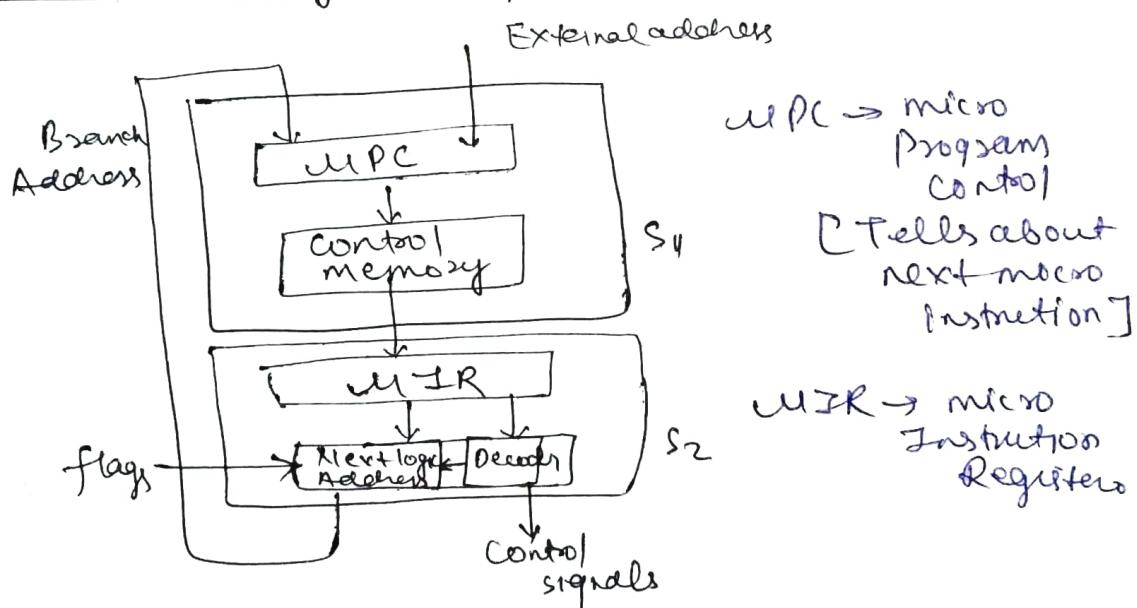
R_i → Registers used as buffer, C_i → Combinational ckt performing its tasks.

We designed it, to achieve CPI ≈ 1 (ideally)

- If clock is divisible then data will not move from one C_i to another C_j .
- If the stages in pipelining are having uneven or non-uniform delays then we will set the clock cycle according to largest delay, and in this case we will not achieve $(CPI = 1)$.

To achieve $CPI = 1$ we have to ensure that delay should be uniform and each opⁿ should take same time but this is not possible.

Here is internal str. of 2 stage pipeline



Q A 4-stage pipeline has the stage delays 150, 120, 165, 140 ns

Clock 2004

Sol Since we are using constant clocking rate (given),
So we will go for clock with max. time = 165 ns

for 1000 data items / instructions in 4 stages

$$\begin{aligned}
 & \frac{1 \times 4 \times 165 + 999 \times 165}{= 165.495 \mu s} \\
 & = \underline{\underline{165.5 \mu s}}
 \end{aligned}$$

Ques 2011: Consider an instrⁿ pipeline with 4-stages (S₁, S₂, S₃, S₄). What is the approximate speed up of the pipeline.

Soln Speedup = $\frac{T_{wp}}{T_p} = \frac{\text{Sum of time of all stages without buffer}}{\text{Max stage time with buffer}}$

$$= \frac{30}{12} = 2.50$$

Doubt: You might get confused the the question is talking about ideal condition so, S \approx m(4) here. But it is possible if delays are same but here at each stage delays are uneven.

But if they talk about max. possible speedup with giving any delay then speedup = 4 (here).

Note → High freq. means → more no. of instruction executed per clock or more no. of clocks generated.

Our aim is to get CPI ≤ 1 means clock per instruction = 1 but it is not possible due to hazards or dependencies.

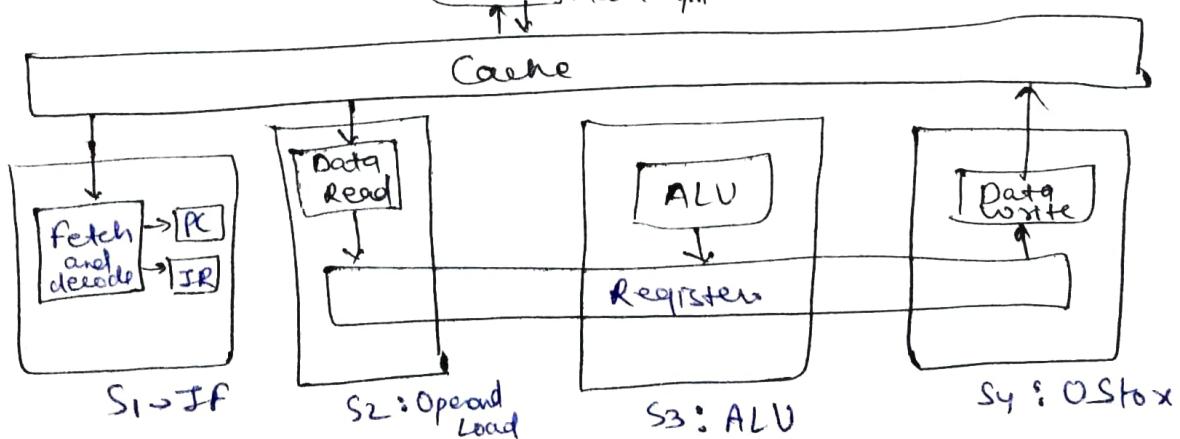
Hazards / dependencies:

- 1) Structural Hazards
- 2) Control Hazards
- 3) Data Hazards.

→ Structural dependencies occurs when two operations tries to access the same resource but only one can get it called as resource restriction.

We can avoid it using resource duplication or replication.

Ex



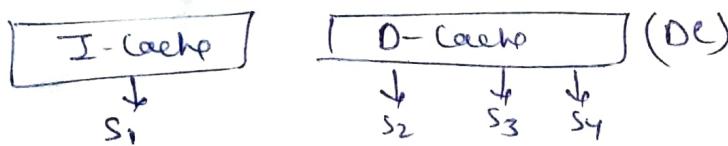
Pipeline

IF	OL	ALU	OS
IF	OL	ALU	OS

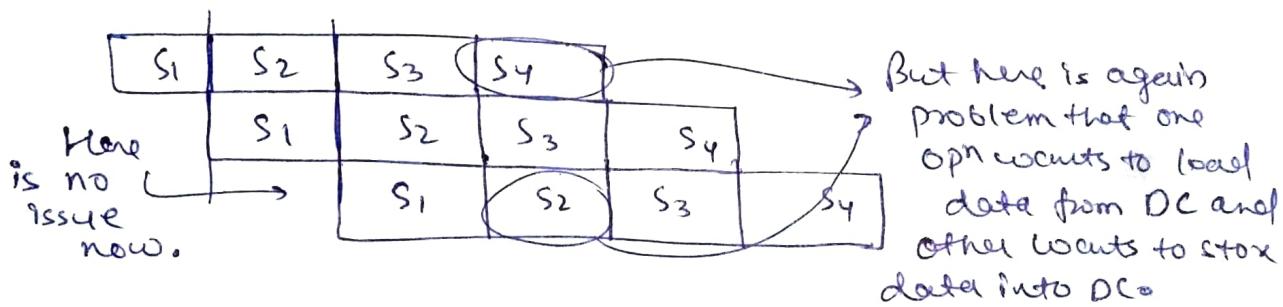
At the same time two operations are trying to access the same memory.

→ In resource duplication we will create multiple copy of Cache.

We divide it into 2 parts ① Instruction Cache (I-C) ② Data Cache (D-Cache)



But here, is again an issue



→ This way we need to devide cache into more parts, but we can divide it upto some extent only. So for this we have a solution that we should rearrange the stages in such a way that similar operations can work with synchronizations

Not in gate but reg

→ Now, we are discussing how (MP) handled this hazard. We are considering here (MP) MIPS R4000 or specifically MIPS R2000/R3000 as example, and these MP have the various pipelining instrⁿ given below.

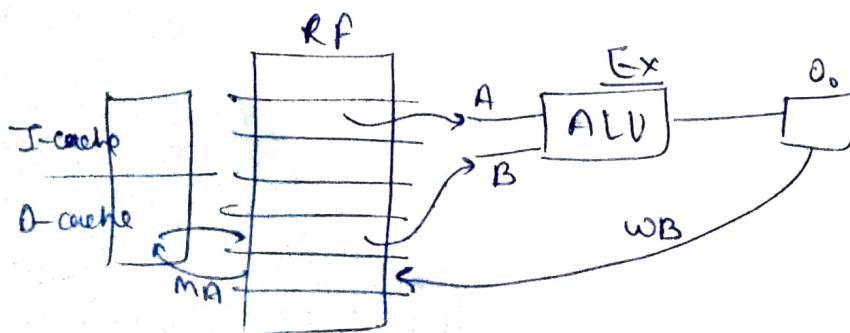
If → Instruction fetching using I-Cache

RD → Operand loading (reading) from the register file (RF) while decoding the fetched instruction.

Ex → Data processing using the ALU and RF as needed.

MA → Operand accessing (load or store) using the D-Cache.

WB → Operand storing (writing back) to RF.



Here we execute Data Read and Data Write as single instruction MA

Control dependency → When flow of control become unpredictable.

- ① conditional branches.
- ② Unconditional _____.
- ③ function calls/return.
- ④ Program Control Instructions that may effect PC.

Throughput → Instruction per sec.

Throughput from CPI: Given $CPI = 1$

$\Rightarrow 1 \text{ clock} \rightarrow 1 \text{ Instruction}$

Say 1 clock = 1 usec

$\Rightarrow 10^6 = 1 \text{ Instr}$

$\Rightarrow 1 \text{ usec} = 10^6 \text{ Instruction} = \underline{\text{Throughput}}$

~~Waste of time~~ → flushing of pipeline and branch penalty. Due to

~~Waste of branch penalty~~ we can't make $CPI = 1$.

~~Pipeline~~ Stalling the pipeline (As you identify that the instruction is branching then wait the pipeline till ~~Waste of~~ fast instruction get executed).

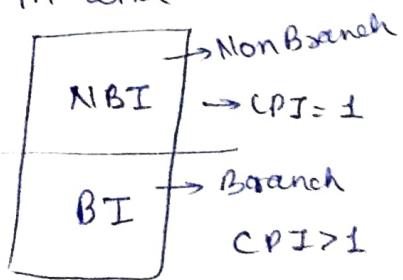
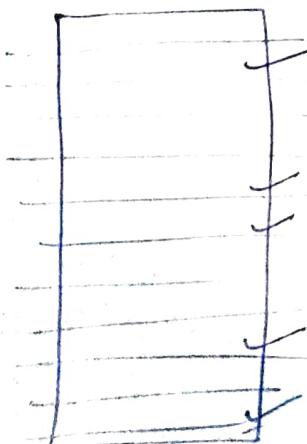
~~Pipeline~~ NOP (No operation instruction). → Our compiler will be smart enough.

~~Best method~~ Reschedule or Rearrange → Our compiler should be smart to fill NOP operations with the valid operations so that it doesn't effect the sequence of operation of other instruction.

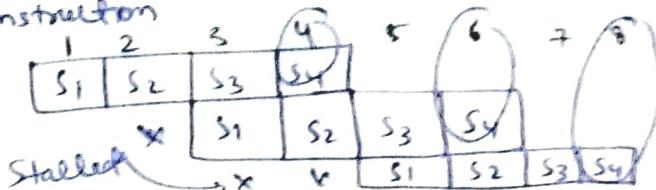
Analysis → Effect of branching on CPI.

~~Another~~ Problem: A program contains 1000 Instruction with 200 branching instruction. If result of branching Instⁿ known after the ~~the~~ second stage. What will be the effective CPI.

Assume branch instructions are distributed non-uniformly, but for our easiness we distribute in like



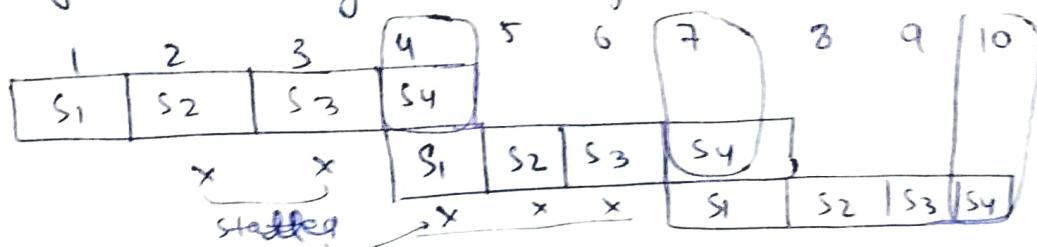
So, In branch instruction



→ CPI = 2, because every instruction came out after 2 clock cycles, after exit of 1st instruction.

$$\text{So, avg CPI} = \frac{80 \times 1 + 20 \times 2}{100} = 1.2$$

Suppose → on changing the question that result of branching get to know after 3rd stage then,



$$\underline{\text{CPI} = 3}$$

$$\text{So, avg CPI} = \frac{80 \times 1 + 20 \times 3}{100} = 1.4$$

→ Similarly we can extend it up.

Generic formula for CPI

Branch

Conditional Unconditional

→ If branching is unconditional then we will definitely move to another instruction Ex. goto 4. [100% chance]

→ In conditional branching its not certain becoz if condition fails we will not move anywhere [if-else] [50% chance]

→ Unconditional branching have less penalty because we can directly access the next instruction just after decoding the branch instructions.

→ But in conditional you have to execute the instruction after fetching it.

Ex

JNZ 200

as you decode
you will get the
next address
200.

if (A > b)

you have to
execute (A > b)
after decoding it
to verify the condition.

So, Conditional
branching provide
more penalty
than unconditional
branching.

So, CPI + 1 for Unconditional and conditional with true statements

If conditional is false then no penalty CPI = 1.

Case: Suppose there is 'n' instruction and 'p' is probability that it is conditional statement and 'q' is probability that condition is true. Pipeline have 'm' stages and we get to know about target statement after 'm' stages.

Soln Find avg. or effective CPI.

$$\text{No. of conditional branch instructions} = np$$

$$\text{Total instructions which branch} = npq$$

$$\text{Here, CPI} = m$$

$$\begin{aligned}\text{Total instruction which doesn't branch} &= (n - npq) \\ &= n(1 - pq)\end{aligned}$$

$$\text{So, effective CPI} = \frac{n(1 - pq) \times 1 + npq \times m}{n}$$

'm' → after which we get to know about branch o/p.

$$\begin{aligned}&= \frac{n - npq + mpq}{n} \\ &= \frac{1 - pq + pqm}{n} \\ &= \underline{\underline{1 + pq(m-1)}}\end{aligned}$$

Fate 2006: $f = 1 \text{ GHz} \Rightarrow T = \frac{1}{10^9} = \frac{1 \text{ ns}}{10^9}$

$$n = 10^9 \quad p = 0.2 \quad q = 0.8 \quad m = 3$$

$$\begin{aligned}\text{Time} &= \underbrace{\left(10^9 \times \frac{20}{100}\right)}_{\text{for conditional}} \underbrace{\left(3 \times 10^{-9}\right)}_{\text{for remaining}} + \underbrace{\left(10^9 \times \frac{80}{100}\right) \times \left(1 \times \frac{1}{10^9}\right)}_{\text{for remaining}} \\ &= \underline{\underline{1.4 \text{ sec}}}\end{aligned}$$

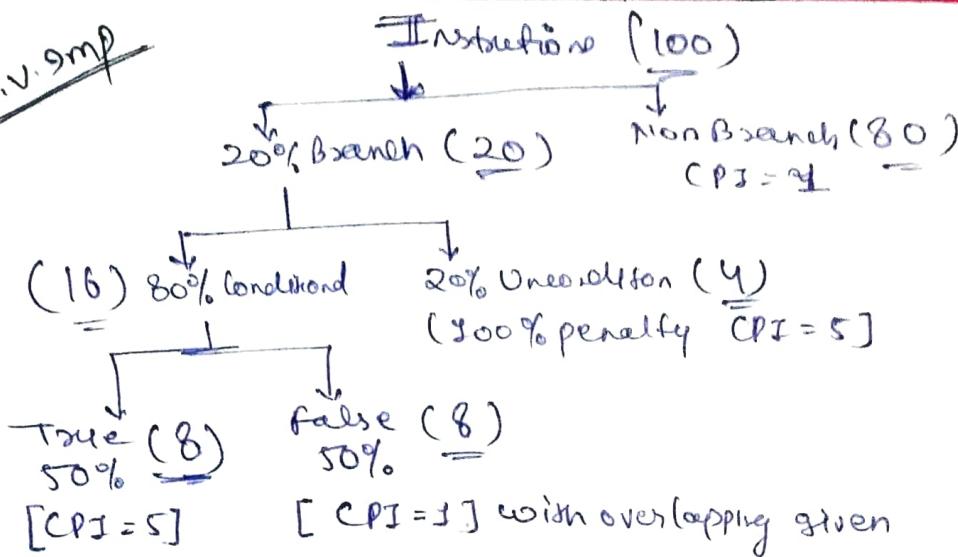
Cafe 2000: (a) $T_{avg} = \frac{n \times 0.2 \times 5 \times 2 \times 10^{-9} + n \times 0.8 \times 1 \times 10^{-9} \times 2}{n}$

$$= 2 \times 10^{-9} + \underline{\underline{0.6 \times 10^{-9}}}$$

$$= \underline{\underline{3.6 \text{ ns}}}$$

$$\text{Penalty} = 3.6 - 2 = \underline{\underline{1.6 \text{ ns}}}$$

(b)
V.V.9mp



Says 100% \rightarrow 100 Instructions.

$$\begin{array}{l} \xrightarrow{4 \times 5 \rightarrow 20} \\ \xrightarrow{80 \times 1 \rightarrow 80} \\ \xrightarrow{8 \times 5 \rightarrow 40} \\ \xrightarrow{8 \times 1 \rightarrow 8} \end{array} \quad \frac{148}{100} = \underline{\underline{1.48 \text{ clocks}}} \\ = \underline{\underline{2.96 \text{ ns}}}$$

Sometimes, time for Unconditional branches given separately.

→ Balanced stages means, all stages are taking same time.

Ways to avoid stall cycles

→ These ways works 80% of time not 100%.

① Branch Prediction & speculative execution.

→ flushing is involved

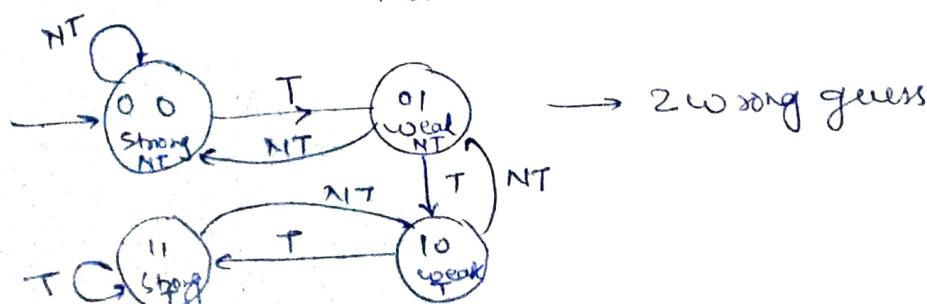
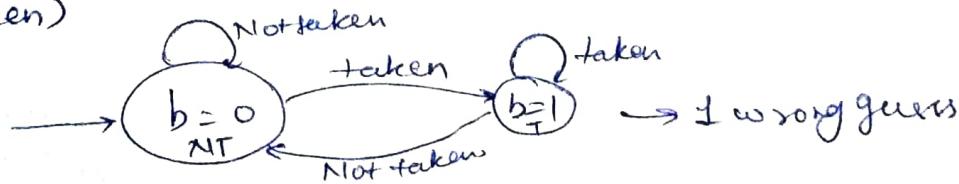
→ Two strategies : Branch always taken, Branch never taken

→ Dynamic Prediction

→ Bit implementation

→ Branch prediction Buffer implementation.

In bit implementation if $b=0$ (Branch Not Taken), $b=1$ (branch is taken)



Delayed branching

	I01	I02	I03	I04	I05	I06	I07	I08	I09
S ₁	I ₁	□	□	I ₂					
S ₂		I ₁	□	□	I ₂				
S ₃			I ₁	□	□	I ₂			
S ₄				I ₁	□	□	I ₂		

Suppose we can get the branched add. after S₃. So next instruction I₂ can't start before

I₁ completes the S₃ stage. This waste lots of space in pipeline but we can use it using smart compilers. Smart compilers can perform load and store kind of operations in vacant slots. So that whenever new Instⁿ come, it will get execute smoothly.

since I₂ is delayed by 2 clock cycles that is why it is called as delayed branching.

- Program meaning shouldn't change
- A smart compiler needed.

Data Hazards / Dependency : Due to read/write data.

- ① RAW : Read after Write (True Dependency)
- ② WAR : Write — Read (Anti Dependency)
- ③ WAW : — — Write (Write output dependency)

Ex-1 I_i: $R_1 \leftarrow R_2 + R_3$
I_j: $R_4 \leftarrow R_1 + R_5$

I _i	S ₁	S ₂	S ₃	S ₄
I _j	S ₁	S ₂	S ₃	S ₄

Here I_j is reading R₁ which is written by I_i. So in pipeline we need to wait for I_j until I_i gets executed because I_i will write back R₁ in S₄. So final pipeline will be look like

$S_1 \rightarrow IF$
 $S_2 \rightarrow \text{Op and fetch}$
 $S_3 \rightarrow ALU$
 $S_4 \rightarrow WB$.

Read after
I_i write

I _i	S ₁	S ₂	S ₃	S ₄	
I _j	S ₁	S ₂	X	X	S ₂ S ₃ S ₄

Ex-2 I_i: $R_1 \leftarrow R_2 + R_3$
I_j: Load m, R₂

Here, value of m is need to be loaded in R₂. Means I_j need to execute before I_i.

③ WAW : $I_i : R_1 \leftarrow R_2 + R_3$ Say order of execution is
 $I_j : R_4 \leftarrow R_5 + R_6$ $I_i \rightarrow I_j$

But, pipeline get designed in such a way that order of execution becomes $I_j \rightarrow I_i$. So it will cause an error because we are writing to same register.

We can identify the dependency as -

RAW: Ex: $I_i : R_1 \leftarrow R_2 + R_3$ such that I_i occurs before
 $I_j : R_4 \leftarrow R_1 + R_5$ I_j

$$D(I_i) \rightarrow R_2, R_3 \quad R(I_i) \rightarrow R_1$$

$$D(I_j) \rightarrow R_1, R_5 \quad R(I_j) \rightarrow R_4$$

If $R(I_i) \cap D(I_j) \neq \emptyset$
then there is a dependency.

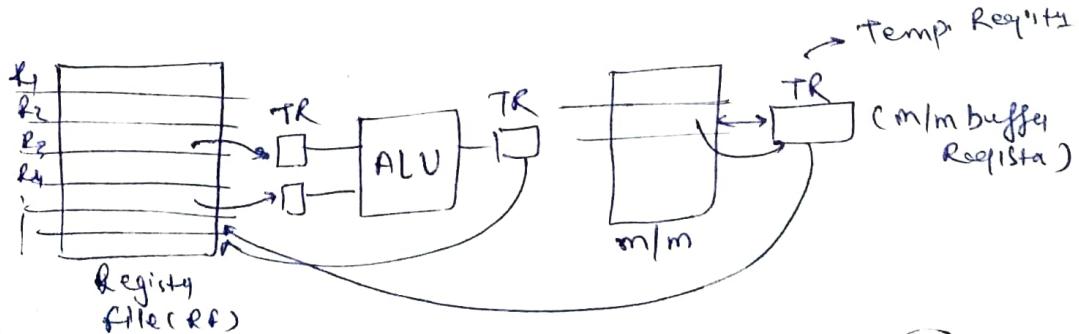
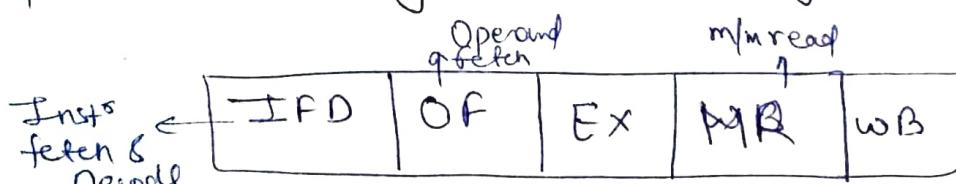
WAR : $I_i : R_1 \leftarrow R_2 + R_3$
 $I_j : \begin{array}{l} \text{Load } m, R_2 \\ \text{or} \\ R_2 \leftarrow R_5 + R_6 \end{array}$

Next instruction is try to update the value. So read the update value of first then perform I_i .

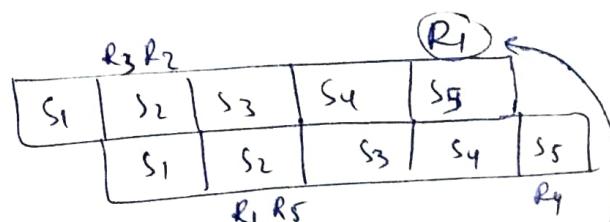
$I_i \rightarrow I_j \rightarrow$ Order

WAW : $(R(I_i) \cap R(I_j)) \neq \emptyset$

Operand forwarding : (Short-circuiting)

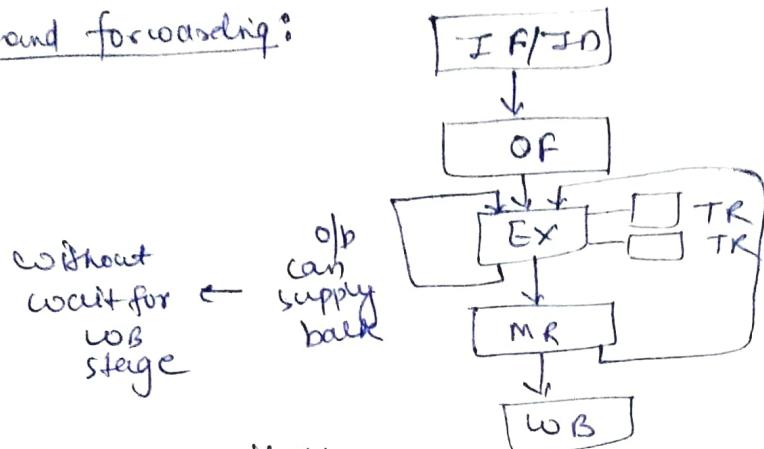


Say $I_1 : R_1 \leftarrow R_2 + R_3$
 $I_2 : R_4 \leftarrow R_1 + R_5$



so it creates a delay bcoz $S_{2,3,4}$ of I_2 gets hazard so we need R_1 here but we will get after S_5

So in operand forwarding:



So, R_4 will be available just after S_3 instead of S_5 in previous examples.

So, we can reduce stalls upto 100% using this technique.

gate-2010: Solution given in videos

Operand forwarding used to resolve RAW hazards,

To resolve WAW, WAR hazards we have concept of register renaming.

$$I_1: R_1 \leftarrow R_2 + R_3$$

$$I_2: R_4 \leftarrow R_4 + R_5$$

	Temporary register	
	200	T_1
	100	T_2
		T_3
		1
		1
		1
		1

say I_1 set value $R_1 \rightarrow 200$
 I_2 set value $R_2 \rightarrow 100$
Time taken (I_1) > Time (I_2)

But I_2 needs to write the final value and in this case I_2 will write final value. So resolve this we have temporary register. We will rename them as $T_1 \leftarrow R_1$, $T_2 \leftarrow R_2$ and keep the value in temp. reg. instead of keeping them in R_4 and when both value will write inside temp. register then we will upload the desired result in R_4 .