

Single Cycle :

$$T_c = t_{pcq\_PC} + 2t_{mem} + t_{RRead} + t_{ALU} + t_{mux} + t_{RWrite}$$

$CPI = 1 \rightarrow$  The cycles per instruction  
decide by LW  
Throughput =  $1 = \frac{1}{CPI}$

ALUOp	Meaning
00	add
01	subtract
10	Look at <i>funct</i> field
11	n/a

Instr.	Jump	RegDst	RegWrite	ALUSrc	Branch	ALUOp1	ALUOp0	MemRead	MemWrite	MemtoReg
R-type	0	1	1	0	0	1	0	0	0	0
lw	0	0	1	1	0	0	0	1	0	1
sw	0	x	0	1	0	0	0	0	1	x
addi	0	0	1	1	0	0	0	0	0	0
B-type	0	x	0	0	1	0	1	0	0	x
J-type	1	x	0	x	x	x	x	0	0	x

Multi-cycle :

Fetch

Machine state	Operation	Control signals
T0	$InsR \leftarrow M[PC];$ $PC \leftarrow PC+4$	$lorD=0, IRWrite=1,$ $ALUSrcA=0, ALUSrc=01, ALUOp=00, PCsrc=0, PCWrite=1$

Decode

Machine state	Operation	Control signals
T1	$(PC+4) + SigExt(offset)$	$ALUSrcA=0, ALUSrcB_{1:0} = 11, ALUOp=00$

LW

Machine state	Operation	Control signals
T2	$A + sigEx(offset)$	$ALUSrcA=1, ALUSrcB_{1:0} = 10, ALUOp=00$
T3	$Data \leftarrow M[A+sigEx(off)]$	$lorD=1$
T4	$RF[dest] \leftarrow Data$	$RegDst=0, MemtoReg=1, RegWrite=1, \text{TO}$

SW

Machine state	Operation	Control signals
T2	$A + sigEx(offset)$	$ALUSrcA=1, ALUSrcB_{1:0} = 10, ALUOp=00$
T3	$M[A+sigEx(offset)] \leftarrow B$	$lorD=1, MemWrite=1, \text{TO}$

R -

Machine state	Operation	Control signals
T2	$A \text{ Op } B$	$ALUSrcA=1, ALUSrcB_{1:0} = 00, ALUOp=00$
T3	$RF[dstn] \leftarrow A \text{ Op } B$	$RegDst=1, MemtoReg=0, RegWrite=1, \text{TO}$

B -

Machine state	Operation	Control signals
T1	$(PC+4) + SigExt(offset)$	$ALUSrcA=0, ALUSrcB_{1:0} = 11, ALUOp=00$
T2	$A-B$	$ALUSrcA=1, ALUSrcB_{1:0} = 00, ALUOp=01, Branch=1, \text{TO}$

Decode stage

ADDI -

Machine state	Operation	Control signals
T2	$A \text{ OP } SigExt(offset)$	$ALUSrcA=1, ALUSrcB_{1:0} = 10, ALUOp=00$
T3	$RF[Destn] \leftarrow A \text{ OP } SigExt(offset)$	$RegDst=0, MemtoReg=0, \text{TO}$

Jump -

Machine state	Operation	Control signals
T2	$A \text{ OP } SigExt(offset)$	$Jump=1, \text{TO}$

$T_c$  is going to decide the clock time period

$$T_c = t_{pcq\_PC} + t_{mux} + \max\{t_{ALU} + t_{mux}, t_{mem}\} + t_{registerRead}$$

Multi-cycle processor is less expensive  
It has 5-nonarchitectural elements

# Pipeline :

- Cost/performance ratio

$$\frac{C}{P} = \frac{G + k * L}{\frac{1}{(\frac{T}{k} + S)}}$$

$$= LT + GS + LS k + \frac{GT}{k}$$

$$k_{opt} = \sqrt{\frac{GT}{LS}}$$

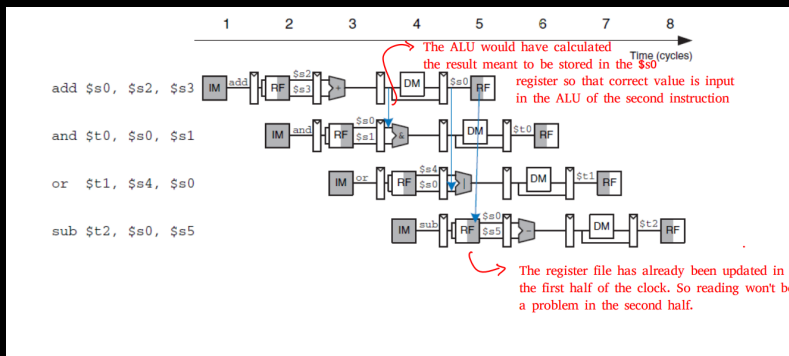
→ Optimal no. of latches

- Clock period (T):  $T_2 - T_1 > T_M - T_m + T_L$

Data Hazard →

a) RAW without LW →

2 types of Forward — (i) MEM to execute stage  
(ii) WB stage to execute stage



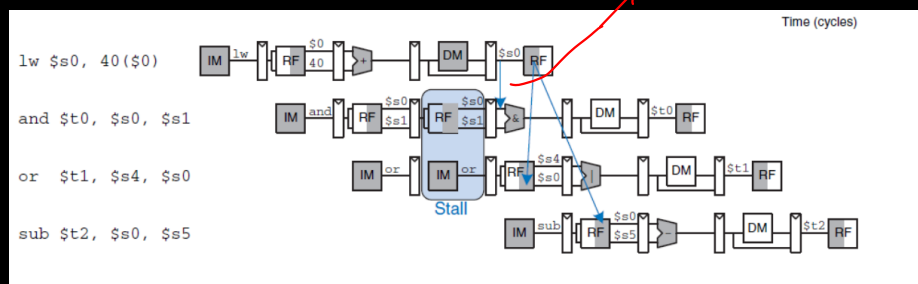
## Data Hazards: H/W-based Solution

- Hazard unit generates control signals for
  - Mux SrcA [ForwardAE]
  - Mux SrcB [ForwardBE]
- The control logic for ForwardAE
  - if  $((rsE \neq 0) \text{ AND } (rsE == WriteRegM))$  and  $RegWriteM$  then ForwardAE = 10 → ALUOut goes in
  - else if  $((rsE \neq 0) \text{ AND } (rsE == WriteRegW))$  and  $RegWriteW$  then ForwardAE = 01 → Output from Data memory goes in.
  - else ForwardAE = 00 → rs, rb fwd.

b) Raw with LW —

Stalling is required

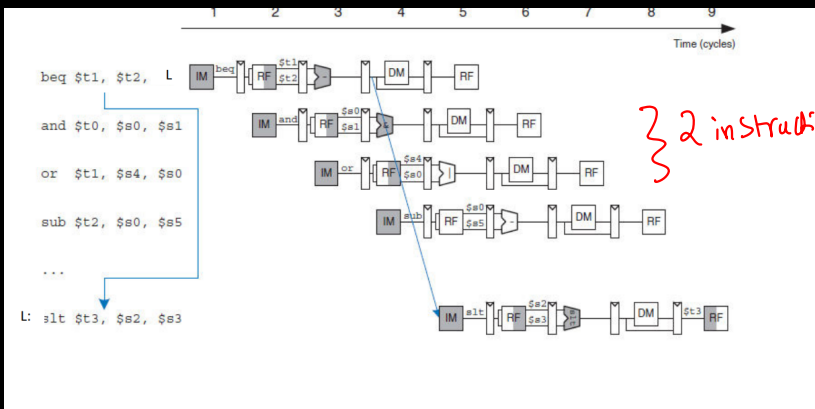
Appropriate forwarding also needs to be done



- $lwstall = ((rsD == rtE) \text{ OR } (rtD == rtE)) \text{ AND } MemtoRegE$
- $StallF = lwstall$  Switch off the program counter
- $StallD = lwstall$  Switch off the latch after instruction stage
- $FlushE = lwstall$  Clear the values in the latch after decode stage

# Control Hazard: Happens due to Branch Instructions.

1)



2 instructions needed to be flushed in this case if Branch is taken.

## 2) Modification Options —

stall : Waste of cycle

Add a Comparator in ID stage : Branch

Instruction exits here itself.

We use this from now on.

Both ZERO, Br. Address are calc. in ID stage

Both causes only 1 instrcn to be flushed

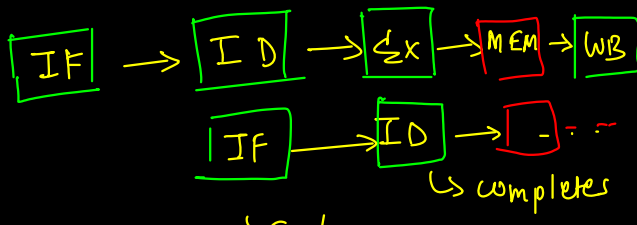
$$\text{Speedup} = \frac{\text{CPI unpipelined}}{1 + \text{pipeline stall cycle from branches}} = \frac{\text{CPI unpipelined}}{1 + \text{Branch frequency} \times \text{Branch penalty}}$$

$$= \frac{T_{\text{unpipe}}}{T_{\text{pipe}}}$$

## Another Problem: —

ADD R1, R2, R3

BEQ R1, R2, offset



But since <sup>the R1=R2+R3</sup> value is not available until end of 'EX',

forwarding of this value can't be done at the start of 'EX' state.

Also, we know BEQ would be making use of this value in the ID stage itself. =>

Solution: Stall the branch instruction. After that Forward it.

- ForwardAD = (rsD != 0) AND (rsD == WriteRegM) AND RegWriteM
- ForwardBD = (rtD != 0) AND (rtD == WriteRegM) AND RegWriteM

What if one of the source operand is in memory stage (rw-type instr.)

- Branchstall=
- BranchD AND RegWriteE AND (WriteRegE == rsD OR WriteRegE == rtD)
- OR
- BranchD AND MemtoRegM AND (WriteRegM == rsD OR WriteRegM == rtD)

- StallF = StallD = FlushE = lwstall OR branchstall

Same kind  
→ of stall  
in both  
cases

$$\text{Cycle period, } T_c = \max \left( \begin{array}{l} t_{pcq} + t_{mem} + t_{setup} \\ 2(t_{RFread} + t_{mux} + t_{eq} + t_{AND} + t_{mux} + t_{setup}) \\ t_{pcq} + t_{mux} + t_{mux} + t_{ALU} + t_{setup} \\ t_{pcq} + t_{memwrite} + t_{setup} \\ 2(t_{pcq} + t_{mux} + t_{RFWrite}) \end{array} \right)$$

Predictions : Whenever a B-type instr<sup>n</sup> actually branches, an instruction has to be flushed. To avoid this, the compiler can take steps based on the prediction it makes.

If prediction is — [Branch untaken] : Would execute normally. (Suitable for forward type branch)

If pred<sup>n</sup> is — [Branch Taken] : —

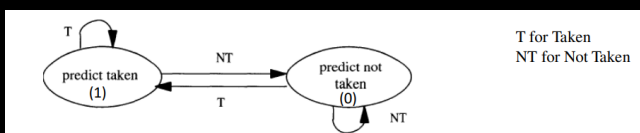
o Rearrange the code such that the next instruction is the one that is not dependent on the Branch block. Like: `if ( -- ) { ... } else { ... }`

Delayed Branch Technique.

This instruction → x y z  
would surely take place  
So, execute it after branch itself.

o There is a prediction bit in the Branch Target Buffer (in the cache) when this bit could be changed by the compiler.

This Buffer is accessed in the IF stage itself and the prediction is made.



## Memory Hierarchy →

1. Spatial : Neighbors being accessed  
Temporal : Same element being accessed
2. → CPU provides address for the memory block required to be accessed along with the offset.

→ Each block would contain Multiple words.

→ 4 bytes = 32 bits.

No. of words = block size =  $b$ .

No. of blocks = No. of cache lines present in the cache =  $S$ .

So, cache capacity =  $C = b * S$

• A cache has the following parameters:  $b$ , block size given in numbers of words;  $S$ , number of blocks;  $N$ , number of ways, and  $A$ , number of address bits.

• In terms of the parameters described, what is the cache capacity,  $C$ ?

•  $C = b * S$

• In terms of the parameters described, what is the total number of bits required to store the tags?

• No of bits for Tag =  $A - \log_2(b) - \log_2(S/N)$

• What are  $S$  and  $N$  for a direct mapped and fully associative cache of size  $C$  words and block size  $b$ ?

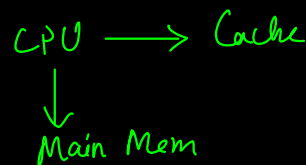
•  $S = C/b$ ,  $N = S$  for fully associative cache

•  $S = C/b$ ,  $N = 1$  for direct mapped cache

## Interaction Policies of Cache : —

1. Read : —

a) Look Aside (Parallel) — Both CPU and Cache



see the Address Bus at the same time.

→ If cache hit, then CPU takes info from here and bus cycle is terminated

→ If miss, take info from Main mem

and write it in cache for the next time.

It is less complex and less expensive

Provide better response to a cache miss

→ Drawback is the processor cannot access cache while another bus master is accessing main memory

b) Look Through (Series) —

if (hit) : { take info from cache ;  
don't send any request to Main mem }

else : { send req. to main mem ;  
store info to cache also after reading }

# Write Policies :

Write hit policy	Write miss policy
Write Through	Write Allocate
<b>Write Through</b>	<b>No Write Allocate</b>
<b>Write Back</b>	<b>Write Allocate</b>
Write Back	No Write Allocate

• **Write Through** - the information is written to both the block in the cache and to the block in the lower-level memory.

*Advantage:*

- read miss never results in writes to main memory
- easy to implement
- main memory always has the most current copy of the data (consistent)

*Disadvantage:*

- write is slower
- every write needs a main memory access
- as a result uses more memory bandwidth

• **Write back** - the information is written only to the block in the cache. The modified cache block is written to main memory only when it is replaced. To reduce the frequency of writing back blocks on replacement, a **dirty bit** is commonly used. This status bit indicates whether the block is dirty (modified while in the cache) or clean (not modified). If it is clean the block is not written on a miss.

*Advantage:*

- writes occur at the speed of the cache memory
- multiple writes within a block require only one write to main memory
- as a result uses less memory bandwidth

*Disadvantage:*

- harder to implement
- main memory is not always consistent with cache
- reads that result in replacement may cause writes of dirty blocks to main memory

Related

if cache is full and some line is being replaced, then dirty blocks must be written.

**Write Through with Write Allocate:**

- on hits it writes to cache and main memory
- on misses it updates the block in main memory and brings the block to the cache
- Bringing the block to cache on a miss does not make a lot of sense in this combination because the next hit to this block will generate a write to main memory anyway (according to Write Through policy)

**Write Through with No Write Allocate:**

- on hits it writes to cache and main memory;
- on misses it updates the block in main memory not bringing that block to the cache;
- Subsequent writes to the block will update main memory because Write Through policy is employed. So, some time is saved not bringing the block in the cache on a miss because it appears useless anyway.

**Write Back with Write Allocate:**

- on hits it writes to cache setting **dirty** bit for the block, main memory is not updated;
- on misses it updates the block in main memory and brings the block to the cache;
- Subsequent writes to the same block, if the block originally caused a miss, will hit in the cache next time, setting dirty bit for the block. That will eliminate extra memory accesses and result in very efficient execution compared with Write Through with Write Allocate combination.

**Write Back with No Write Allocate:**

- on hits it writes to cache setting **dirty** bit for the block, main memory is not updated;
- on misses it updates the block in main memory not bringing that block to the cache;
- Subsequent writes to the same block, if the block originally caused a miss, will generate misses all the way and result in very inefficient execution.