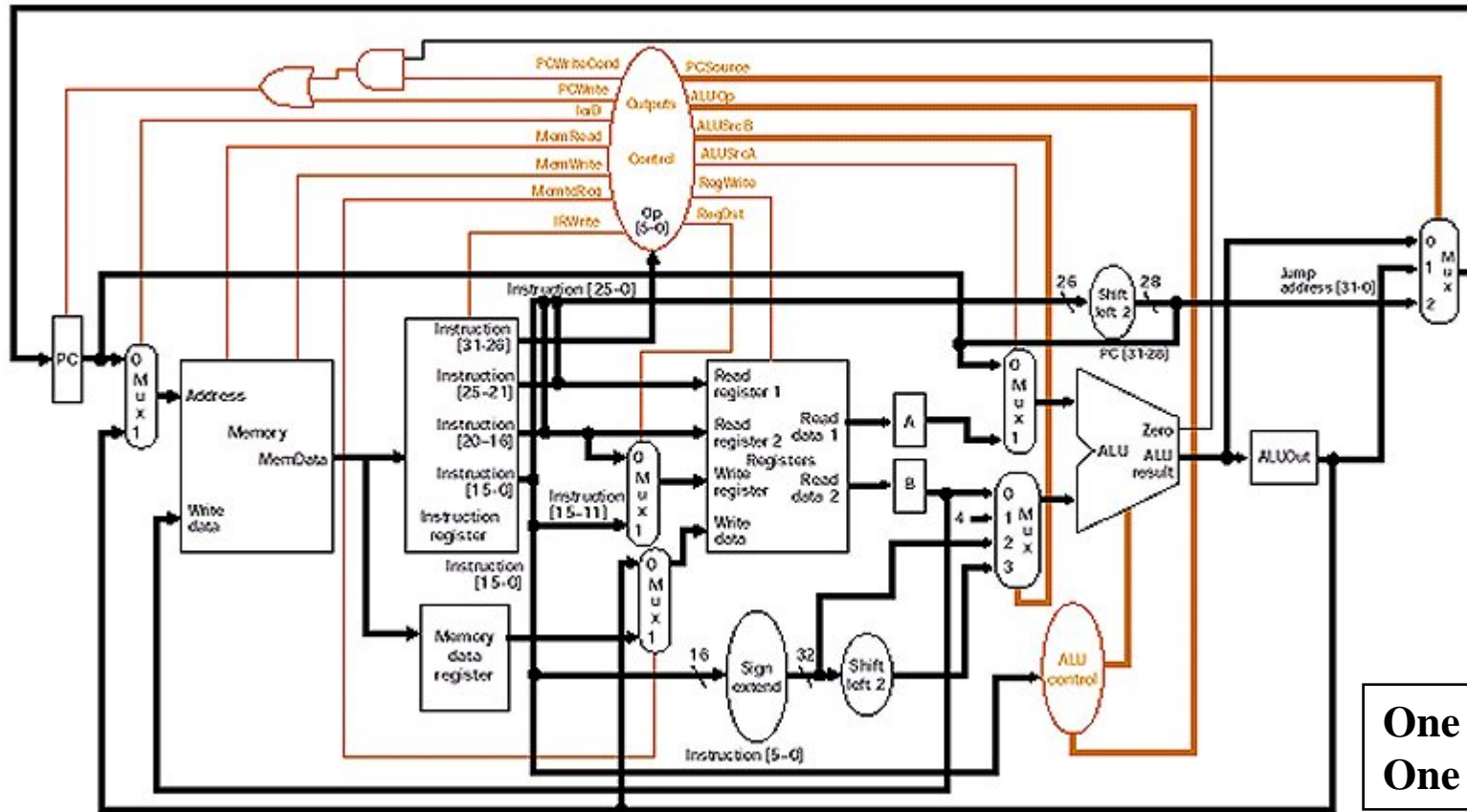


# EECC550 Exam Review

## 4 out of 6 questions

- 1 **Multicycle CPU performance vs. Pipelined CPU performance**
- 2 **Given MIPS code, MIPS pipeline (similar to questions 4, 5 of HW#4)**
  - **Performance of code as is on a given CPU**
  - **Schedule the code to reduce stalls + resulting performance**
- 3 **Cache Operation: Given a series of word memory address references, cache capacity and organization: (similar to question #1 of HW #5)**
  - **Find Hits/misses, Hit rate, Final content of cache**
- 4 **Pipelined CPU performance with non-ideal memory and unified or split cache**
  - **Find AMAT, CPI, performance ...**
- 5 **For a cache level with given characteristics find:**
  - **Address fields, mapping function, storage requirements etc.**
- 6 **Performance evaluation of non-ideal pipelined CPUs using non ideal memory + cache:**
  - **Desired performance maybe given: Find missing parameter**

# MIPS CPU Design: Multi-Cycle Datapath (Textbook Version)



**CPI: R-Type = 4, Load = 5, Store 4, Jump/Branch = 3**  
**Only one instruction being processed in datapath**

*How to lower CPI further without increasing CPU clock cycle time, C?*

$$T = I \times \text{CPI} \times C$$

**EECC550 - Shaaban**

Processing an instruction starts when the previous instruction is completed

# Operations In Each Cycle

	R-Type	Load	Store	Branch	Jump
IF	Instruction Fetch $IR \leftarrow Mem[PC]$ $PC \leftarrow PC + 4$	Instruction Fetch $IR \leftarrow Mem[PC]$ $PC \leftarrow PC + 4$	Instruction Fetch $IR \leftarrow Mem[PC]$ $PC \leftarrow PC + 4$	Instruction Fetch $IR \leftarrow Mem[PC]$ $PC \leftarrow PC + 4$	Instruction Fetch $IR \leftarrow Mem[PC]$ $PC \leftarrow PC + 4$
ID	Instruction Decode $A \leftarrow R[rs]$ $B \leftarrow R[rt]$ $ALUout \leftarrow PC + (SignExt(imm16) \times 4)$	Instruction Decode $A \leftarrow R[rs]$ $B \leftarrow R[rt]$ $ALUout \leftarrow PC + (SignExt(imm16) \times 4)$	Instruction Decode $A \leftarrow R[rs]$ $B \leftarrow R[rt]$ $ALUout \leftarrow PC + (SignExt(imm16) \times 4)$	Instruction Decode $A \leftarrow R[rs]$ $B \leftarrow R[rt]$ $ALUout \leftarrow PC + (SignExt(imm16) \times 4)$	Instruction Decode $A \leftarrow R[rs]$ $B \leftarrow R[rt]$ $ALUout \leftarrow PC + (SignExt(imm16) \times 4)$
EX	Execution $ALUout \leftarrow A \text{ funct } B$	Execution $ALUout \leftarrow A + SignEx(Im16)$	Execution $ALUout \leftarrow A + SignEx(Im16)$	Execution $Zero \leftarrow A - B$ $Zero: PC \leftarrow ALUout$	Execution $PC \leftarrow \text{Jump Address}$
MEM	Memory	Memory $M \leftarrow Mem[ALUout]$	Memory $Mem[ALUout] \leftarrow B$		
WB	Write Back $R[rd] \leftarrow ALUout$	Write Back $R[rt] \leftarrow M$			

$$T = I \times CPI \times C$$

Reducing the CPI by combining cycles increases CPU clock cycle

Instruction Fetch (IF) & Instruction Decode (ID) cycles are common for all instructions

**EECC550 - Shaaban**

# Multi-cycle Datapath Instruction CPI

- **R-Type/Immediate: Require four cycles, CPI = 4**
  - IF, ID, EX, WB
- **Loads: Require five cycles, CPI = 5**
  - IF, ID, EX, MEM, WB
- **Stores: Require four cycles, CPI = 4**
  - IF, ID, EX, MEM
- **Branches: Require three cycles, CPI = 3**
  - IF, ID, EX
- **Average program  $3 \leq \text{CPI} \leq 5$  depending on program profile (instruction mix).**



**Non-overlapping Instruction Processing:**

Processing an instruction starts when the previous instruction is completed.

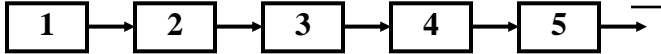
# MIPS Multi-cycle Datapath Performance Evaluation

- What is the average CPI?
  - State diagram gives CPI for each instruction type.
  - Workload (program) below gives frequency of each type.

Type	CPI <sub>i</sub> for type	Frequency	CPI <sub>i</sub> x frequ <sub>i</sub>
Arith/Logic	4	40%	1.6
Load	5	30%	1.5
Store	4	10%	0.4
branch	3	20%	0.6
Average CPI:			4.1

**Better than CPI = 5 if all instructions took the same number of clock cycles (5).**

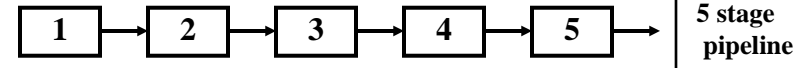
# Instruction Pipelining

- Instruction pipelining is a CPU implementation technique where multiple operations on a number of instructions are overlapped.
  - For Example: The next instruction is fetched in the next cycle without waiting for the current instruction to complete.
- An instruction execution pipeline involves a number of steps, where each step completes a part of an instruction. Each step is called *a pipeline stage* or *a pipeline segment*.
- The stages or steps are connected in a linear fashion: one stage to the next to form the pipeline (or pipelined CPU datapath) -- instructions enter at one end and progress through the stages and exit at the other end.  5 stage pipeline
- The time to move an instruction one step down the pipeline is equal to *the machine (CPU) cycle* and is determined by the stage with the longest processing delay.
- Pipelining increases the CPU instruction throughput: The number of instructions completed per cycle.
  - Instruction Pipeline Throughput : The instruction completion rate of the pipeline and is determined by how often an instruction exists the pipeline.
  - Under ideal conditions (no stall cycles), instruction throughput is one instruction per machine cycle, or ideal effective CPI = 1 Or ideal IPC = 1
- Pipelining does not reduce the execution time of an individual instruction: The time needed to complete all processing steps of an instruction (also called instruction completion latency).
  - Minimum instruction latency =  $n$  cycles, where  $n$  is the number of pipeline stages

# Pipelining: Design Goals

- The length of the machine clock cycle is determined by the time required for the slowest pipeline stage. Similar to non-pipelined multi-cycle CPU

- An important pipeline design consideration is to balance the length of each pipeline stage.



- If all stages are perfectly balanced, then the time per instruction on a pipelined machine (assuming ideal conditions with no stalls):

$$\frac{\text{Time per instruction on unpipelined machine}}{\text{Number of pipeline stages}}$$

- Under these ideal conditions:
  - Speedup from pipelining = the number of pipeline stages =  $n$
  - Goal: One instruction is completed every cycle: **CPI = 1**.

$$T = I \times \text{CPI} \times C$$

**EECC550 - Shaaban**

# Ideal Pipelined Instruction Processing

(i.e no stall cycles)

## Timing Representation

**n = 5 stage pipeline**

Fill Cycles = number of stages -1

Clock cycle Number

Time in clock cycles →

Instruction Number      1      2      3      4      5      6      7      8      9

Instruction I	IF	ID	EX	MEM	WB				
Instruction I+1		IF	ID	EX	MEM	WB			
Instruction I+2			IF	ID	EX	MEM	WB		
Instruction I+3				IF	ID	EX	MEM	WB	
Instruction I +4					IF	ID	EX	MEM	WB

**Ideal CPI = 1**  
(or IPC =1)

← 4 cycles = n -1 = 5 -1 Time to fill the pipeline →

Program Order

**n= 5 Pipeline Stages:**

- IF** = Instruction Fetch
- ID** = Instruction Decode
- EX** = Execution
- MEM** = Memory Access
- WB** = Write Back

Any individual instruction goes through all five pipeline stages taking 5 cycles to complete  
Thus instruction latency= 5 cycles

**First instruction, I Completed**

**Instruction, I+4 completed**

**Pipeline Fill Cycles:** No instructions completed yet  
Number of fill cycles = Number of pipeline stages - 1  
Here 5 - 1 = 4 fill cycles

**Ideal pipeline operation:** After fill cycles, one instruction is completed per cycle giving the ideal pipeline CPI = 1 (ignoring fill cycles)

**EECC550 - Shaaban**

Ideal pipeline operation without any stall cycles

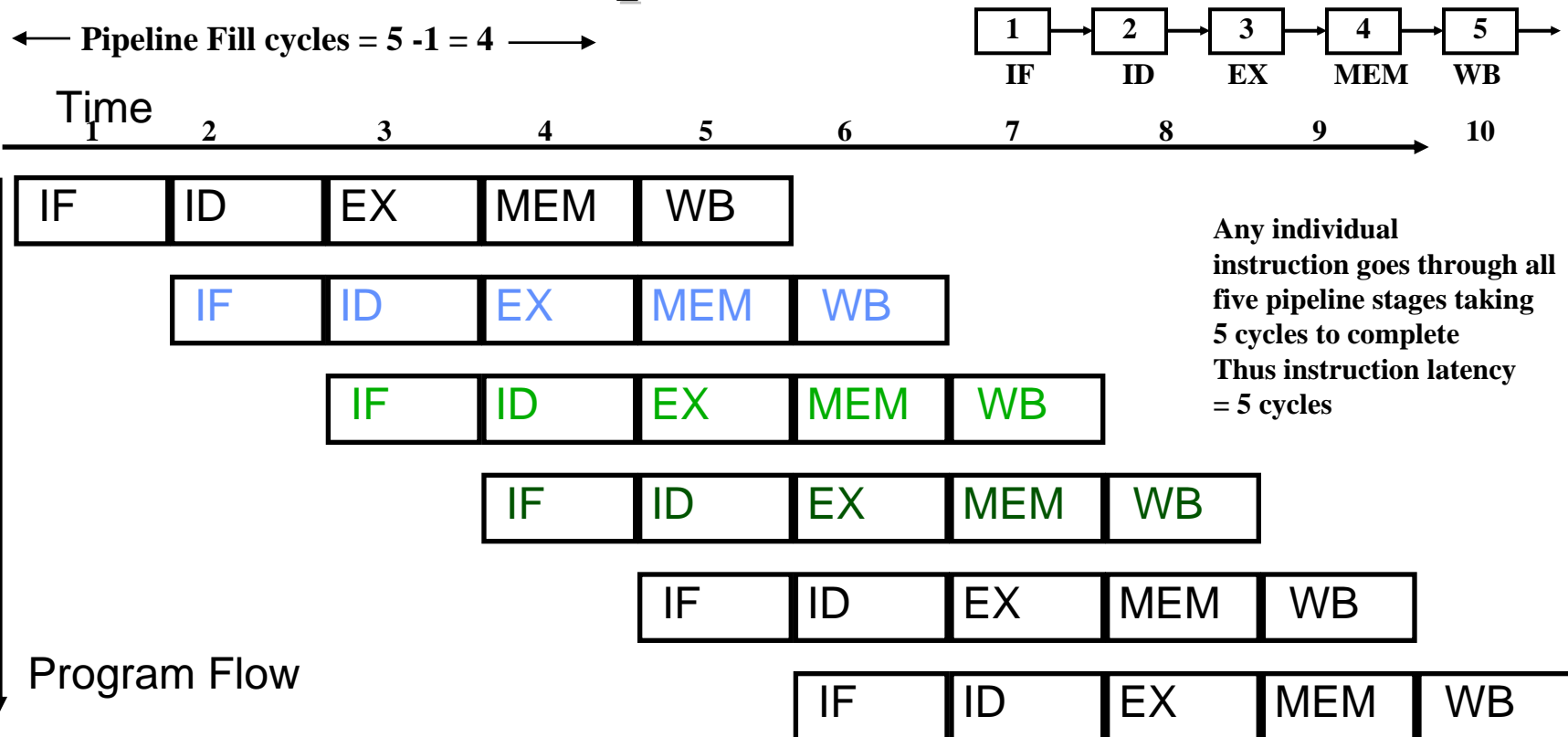


# Ideal Pipelined Instruction Processing

(i.e no stall cycles)

## Representation

5 Stage Pipeline



Here  $n = 5$  pipeline stages or steps

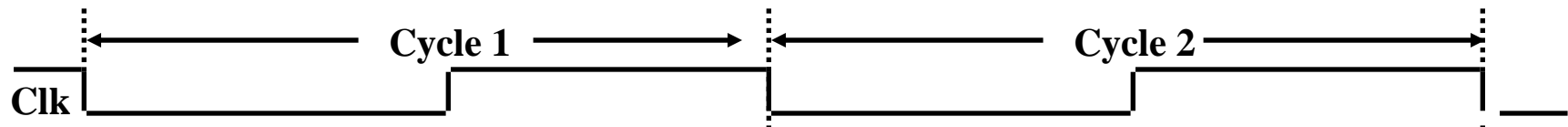
Number of pipeline fill cycles = Number of stages - 1 Here  $5 - 1 = 4$

After fill cycles: One instruction is completed every cycle (Effective CPI = 1)  
(ideally)

Ideal pipeline operation without any stall cycles

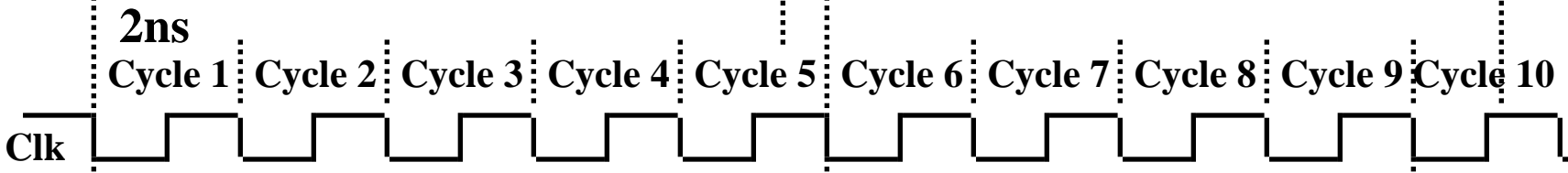
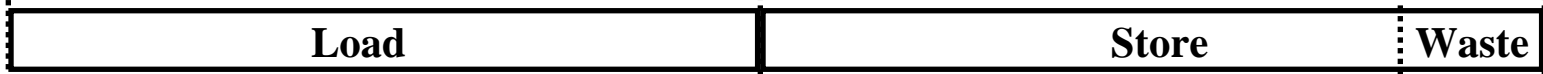
**EECC550 - Shaaban**

# Single Cycle, Multi-Cycle, Vs. Pipelined CPU

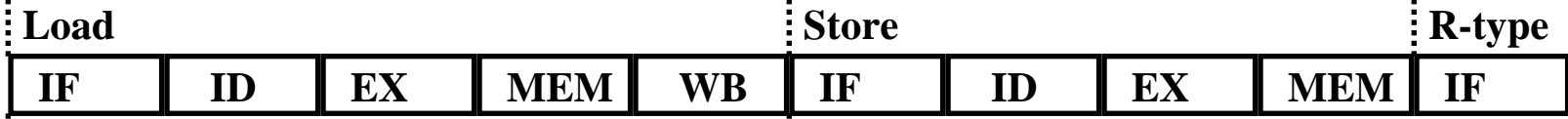


Single Cycle Implementation:

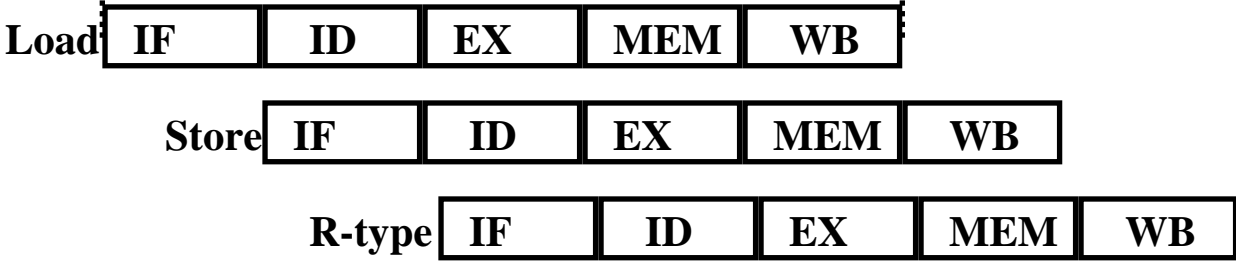
8 ns



Multiple Cycle Implementation:



Pipeline Implementation:



Assuming the following datapath/control hardware components delays:  
 Memory Units: 2 ns    ALU and adders: 2 ns  
 Register File: 1 ns    Control Unit < 1 ns

**EECC550 - Shaaban**

# Single Cycle, Multi-Cycle, Pipeline: Performance Comparison Example

For 1000 instructions, execution time:

$$T = I \times \text{CPI} \times C$$

- Single Cycle Machine:

- 8 ns/cycle x 1 CPI x 1000 inst = 8000 ns

- Multi-cycle Machine:

- 2 ns/cycle x 4.6 CPI (due to inst mix) x 1000 inst = 9200 ns

Depends on program instruction mix

- Ideal pipelined machine, 5-stages (effective CPI = 1):

- 2 ns/cycle x (1 CPI x 1000 inst + 4 cycle fill) = 2008 ns

- Speedup =  $8000/2008 = 3.98$  faster than single cycle CPU

- Speedup =  $9200/2008 = 4.58$  times faster than multi cycle CPU

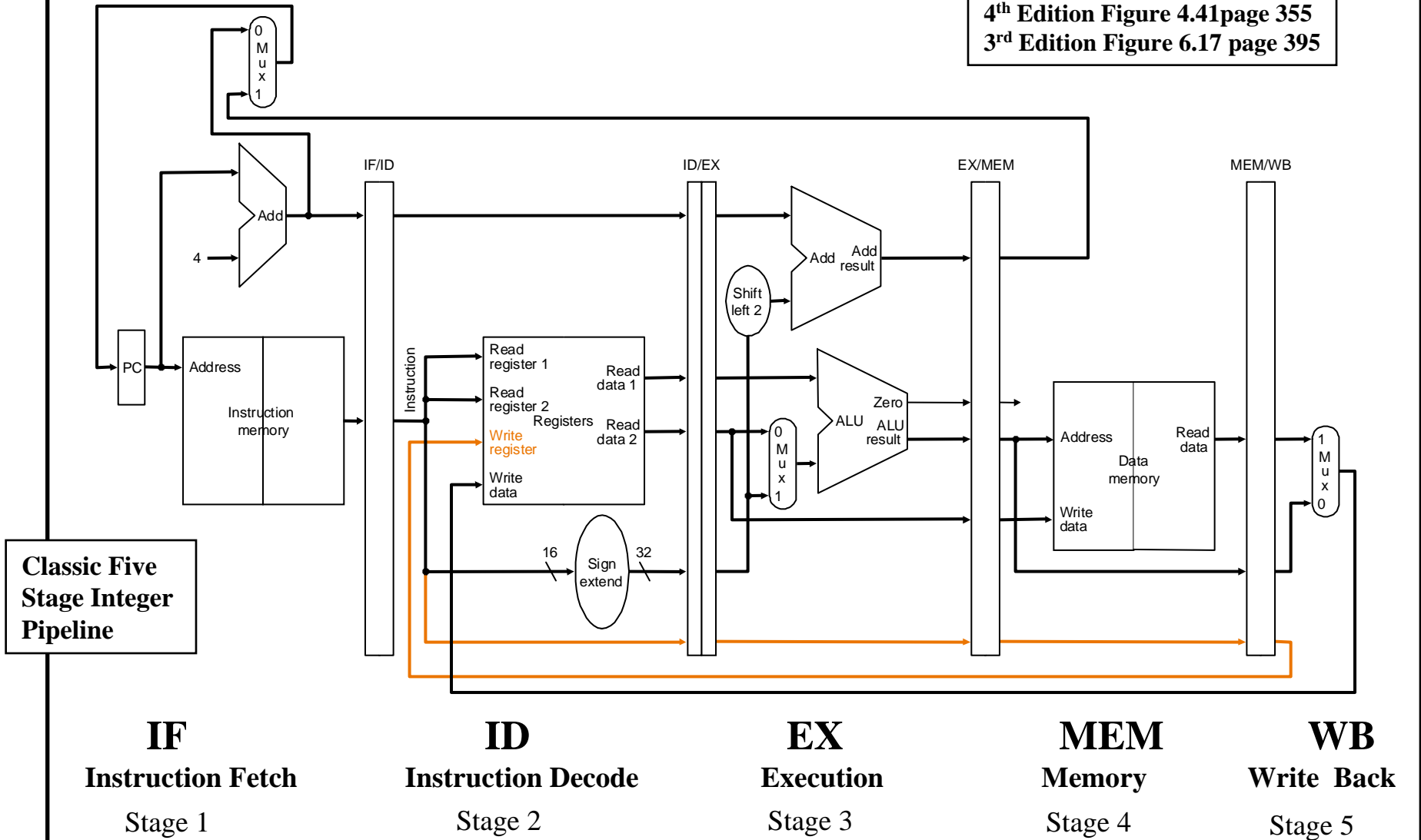
# Basic Pipelined CPU Design Steps

1. Analyze instruction set operations using independent RTN => datapath requirements.
2. Select required datapath components and connections.
3. Assemble an initial datapath meeting the ISA requirements.
4. Identify pipeline stages based on operation, balancing stage delays, and ensuring no hardware conflicts exist when common hardware is used by two or more stages simultaneously in the same cycle.
5. Divide the datapath into the stages identified above by adding buffers between the stages of sufficient width to hold:
  - Instruction fields.
  - Remaining control lines needed for remaining pipeline stages.
  - All results produced by a stage and any unused results of previous stages.
6. Analyze implementation of each instruction to determine setting of control points that effects the register transfer taking pipeline hazard conditions into account . (More on this a bit later)
7. Assemble the control logic.

i.e registers

# A Basic Pipelined Datapath

4<sup>th</sup> Edition Figure 4.41 page 355  
3<sup>rd</sup> Edition Figure 6.17 page 395



Version 1: No forwarding, Branch resolved in MEM stage

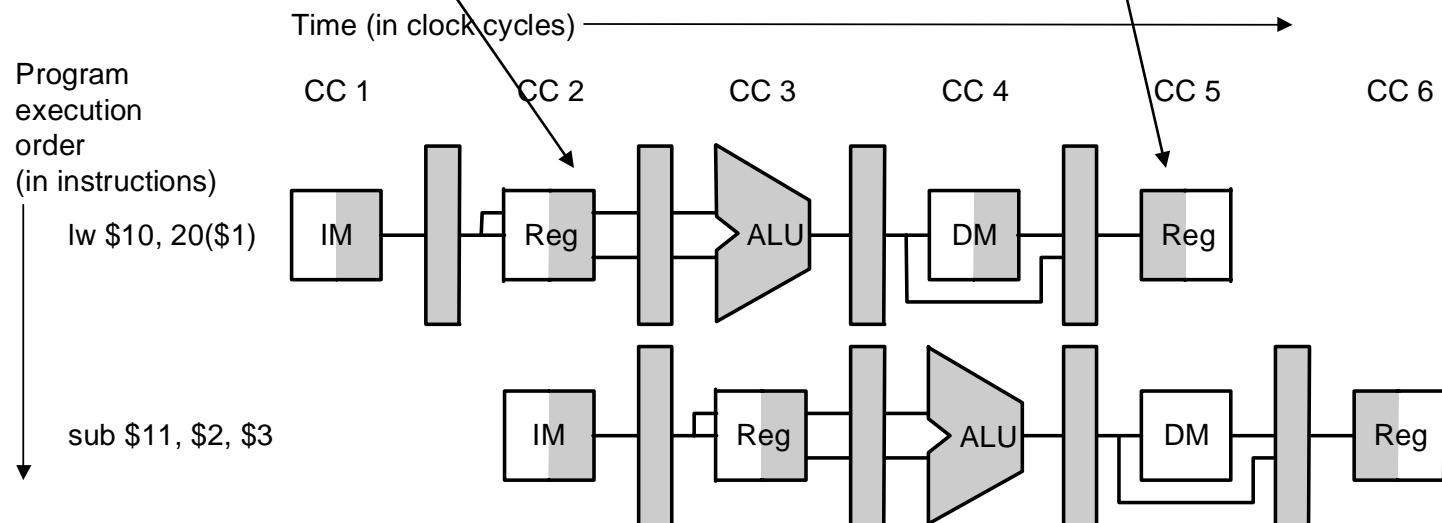
**EECC550 - Shaaban**

# Read/Write Access To Register Bank

- Two instructions need to access the register bank in the same cycle:
  - One instruction to read operands in its instruction decode (ID) cycle.
  - The other instruction to write to a destination register in its Write Back (WB) cycle.
- This represents a potential hardware conflict over access to the register bank.
- Solution: Coordinate register reads and write in the same cycle as follows:

• Operand register reads in Instruction Decode ID cycle occur in the second half of the cycle (indicated here by the dark shading of the second half of the cycle)

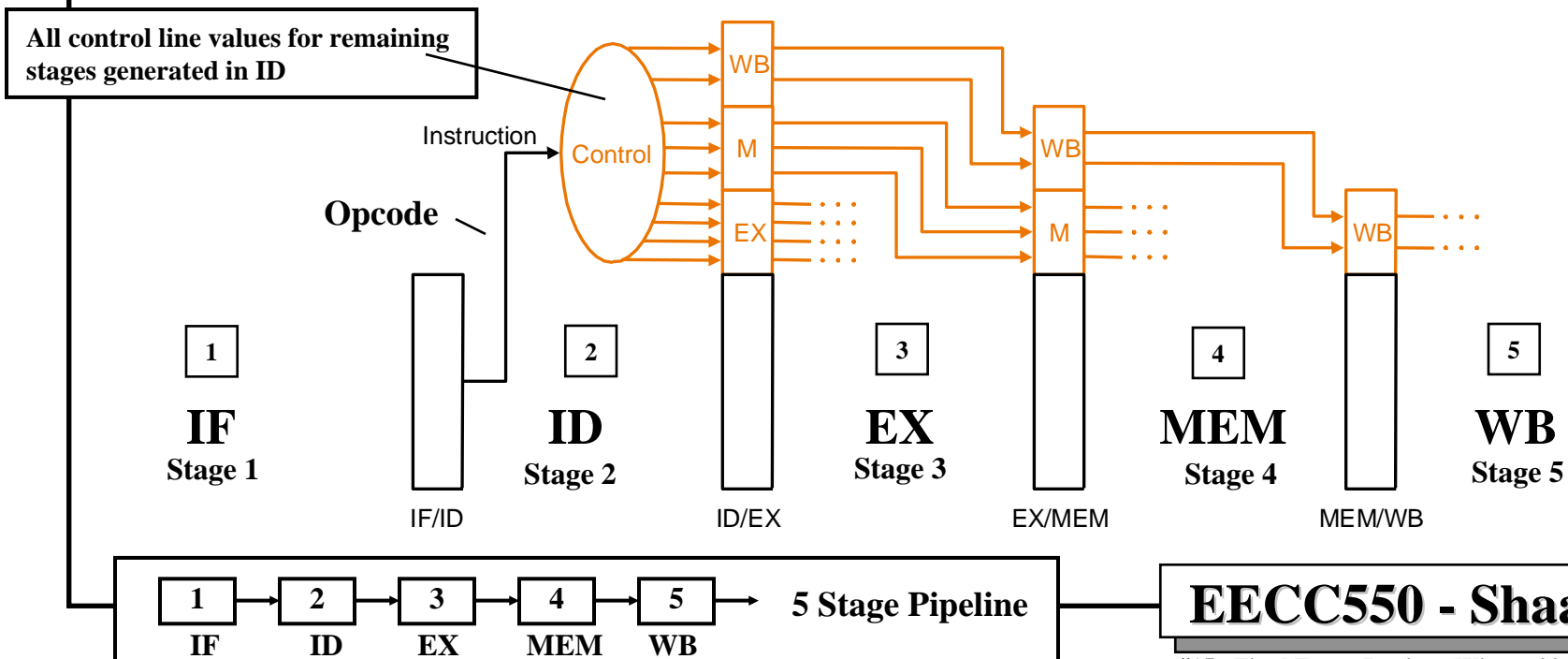
• Register write in Write Back WB cycle occur in the first half of the cycle. (indicated here by the dark shading of the first half of the WB cycle)



# Pipeline Control

- Pass needed control signals along from one stage to the next as the instruction travels through the pipeline just like the needed data

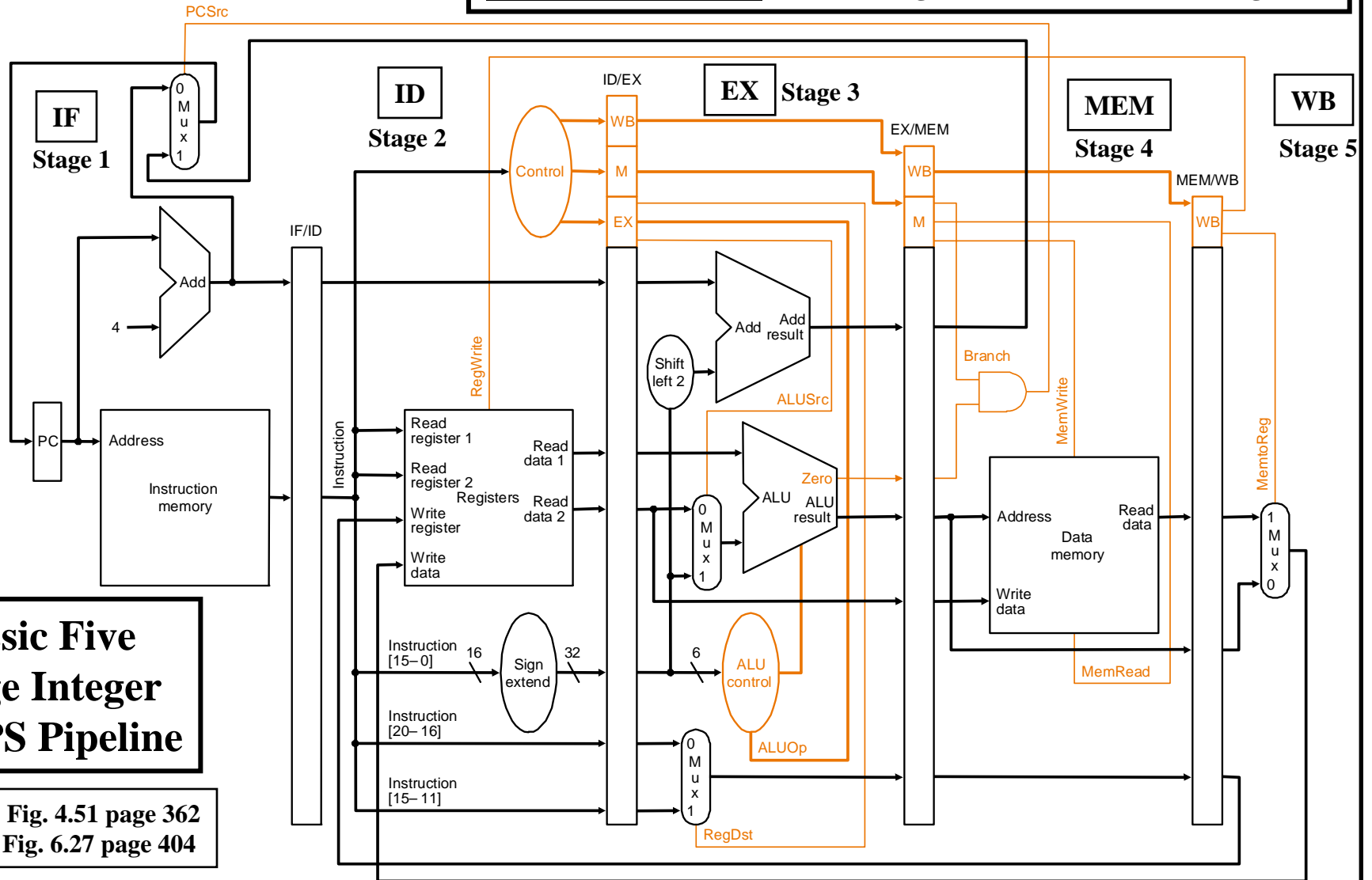
Instruction	EX				MEM			WB	
	Execution/Address Calculation stage control lines				Memory access stage control lines			stage control lines	
	Reg Dst	ALU Op1	ALU Op0	ALU Src	Branch	Mem Read	Mem Write	Reg write	Mem to Reg
R-format	1	1	0	0	0	0	0	1	0
lw	0	0	0	1	0	1	0	1	1
sw	X	0	0	1	0	0	1	0	X
beq	X	0	1	0	1	0	0	0	X



# Pipelined Datapath with Control Added

MIPS Pipeline Version #1

MIPS Pipeline Version 1: No forwarding, branch resolved in MEM stage



**Classic Five Stage Integer MIPS Pipeline**

4<sup>th</sup> Ed. Fig. 4.51 page 362  
3<sup>rd</sup> Ed. Fig. 6.27 page 404

Target address of branch determined in EX but PC is updated in MEM stage (i.e branch is resolved in MEM, stage 4)

**EECC550 - Shaaban**



# Basic Performance Issues In Pipelining

- Pipelining increases the CPU instruction throughput:

The number of instructions completed per unit time.

Under ideal conditions (i.e. No stall cycles):

$$T = I \times \text{CPI} \times C$$

- Pipelined CPU instruction throughput is one instruction completed per machine cycle, or  $\text{CPI} = 1$

(ignoring pipeline fill cycles)

Or Instruction throughput: Instructions Per Cycle =  $\text{IPC} = 1$

- Pipelining does not reduce the execution time of an individual instruction: The time needed to complete all processing steps of an instruction (also called instruction completion latency).

- It usually slightly increases the execution time of individual instructions over unpipelined CPU implementations due to:

- The increased control overhead of the pipeline and pipeline stage registers delays +
- Every instruction goes through every stage in the pipeline even if the stage is not needed. (i.e MEM pipeline stage in the case of R-Type instructions)

Here  $n = 5$  stages

**EECC550 - Shaaban**

# Pipelining Performance Example

- **Example: For an unpipelined multicycle CPU:**
  - Clock cycle = 10ns, 4 cycles for ALU operations and branches and 5 cycles for memory operations with instruction frequencies of 40%, 20% and 40%, respectively.
  - If pipelining adds 1ns to the CPU clock cycle then the speedup in instruction execution from pipelining is:

$$\begin{aligned} \text{Non-pipelined Average execution time/instruction} &= \text{Clock cycle} \times \text{Average CPI} \\ &= 10 \text{ ns} \times ((40\% + 20\%) \times 4 + 40\% \times 5) = 10 \text{ ns} \times 4.4 = 44 \text{ ns} \\ &\quad \text{CPI} = 4.4 \end{aligned}$$

In the pipelined CPU implementation, ideal CPI = 1

CPI = 1

$$\begin{aligned} \text{Pipelined execution time/instruction} &= \text{Clock cycle} \times \text{CPI} \\ &= (10 \text{ ns} + 1 \text{ ns}) \times 1 = 11 \text{ ns} \times 1 = 11 \text{ ns} \end{aligned}$$

$$\begin{aligned} \text{Speedup from pipelining} &= \frac{\text{Time Per Instruction time unpipelined}}{\text{Time per Instruction time pipelined}} \\ &= 44 \text{ ns} / 11 \text{ ns} = 4 \text{ times faster} \end{aligned}$$

$T = I \times \text{CPI} \times C$  here I did not change

**EECC550 - Shaaban**

# Pipeline Hazards

$$\text{CPI} = 1 + \text{Average Stalls Per Instruction}$$

- Hazards are situations in pipelined CPUs which prevent the next instruction in the instruction stream from executing during the designated clock cycle possibly resulting in one or more stall (or wait) cycles.
- Hazards reduce the ideal speedup (increase  $\text{CPI} > 1$ ) gained from pipelining and are classified into three classes:

i.e A resource the instruction requires for correct execution is not available in the cycle needed

Resource  
Not available:

Hardware  
Component

Correct  
Operand  
(data) value

Correct  
PC

– **Structural hazards**: Arise from hardware resource conflicts when the available hardware cannot support all possible combinations of instructions.

Hardware structure (component) conflict

– **Data hazards**: Arise when an instruction depends on the results of a previous instruction in a way that is exposed by the overlapping of instructions in the pipeline.

Operand not ready yet when needed in EX

– **Control hazards**: Arise from the pipelining of conditional branches and other instructions that change the PC.

Correct PC not available when needed in IF

**EECC550 - Shaaban**

# Performance of Pipelines with Stalls

- Hazard conditions in pipelines may make it necessary to stall the pipeline by a number of cycles degrading performance from the ideal pipelined CPU CPI of 1.

$$\begin{aligned} \text{CPI pipelined} &= \text{Ideal CPI} + \text{Pipeline stall clock cycles per instruction} \\ &= 1 + \text{Pipeline stall clock cycles per instruction} \end{aligned}$$

- If pipelining overhead is ignored and we assume that the stages are perfectly balanced then speedup from pipelining is given by:

$$\begin{aligned} \text{Speedup} &= \text{CPI unpipelined} / \text{CPI pipelined} \\ &= \text{CPI unpipelined} / (1 + \text{Pipeline stall cycles per instruction}) \end{aligned}$$

- When all instructions in the multicycle CPU take the same number of cycles equal to the number of pipeline stages then:

$$\text{Speedup} = \text{Pipeline depth} / (1 + \text{Pipeline stall cycles per instruction})$$

# Structural (or Hardware) Hazards

- In pipelined machines overlapped instruction execution requires pipelining of functional units and duplication of resources to allow all possible combinations of instructions in the pipeline. *To prevent hardware structures conflicts*

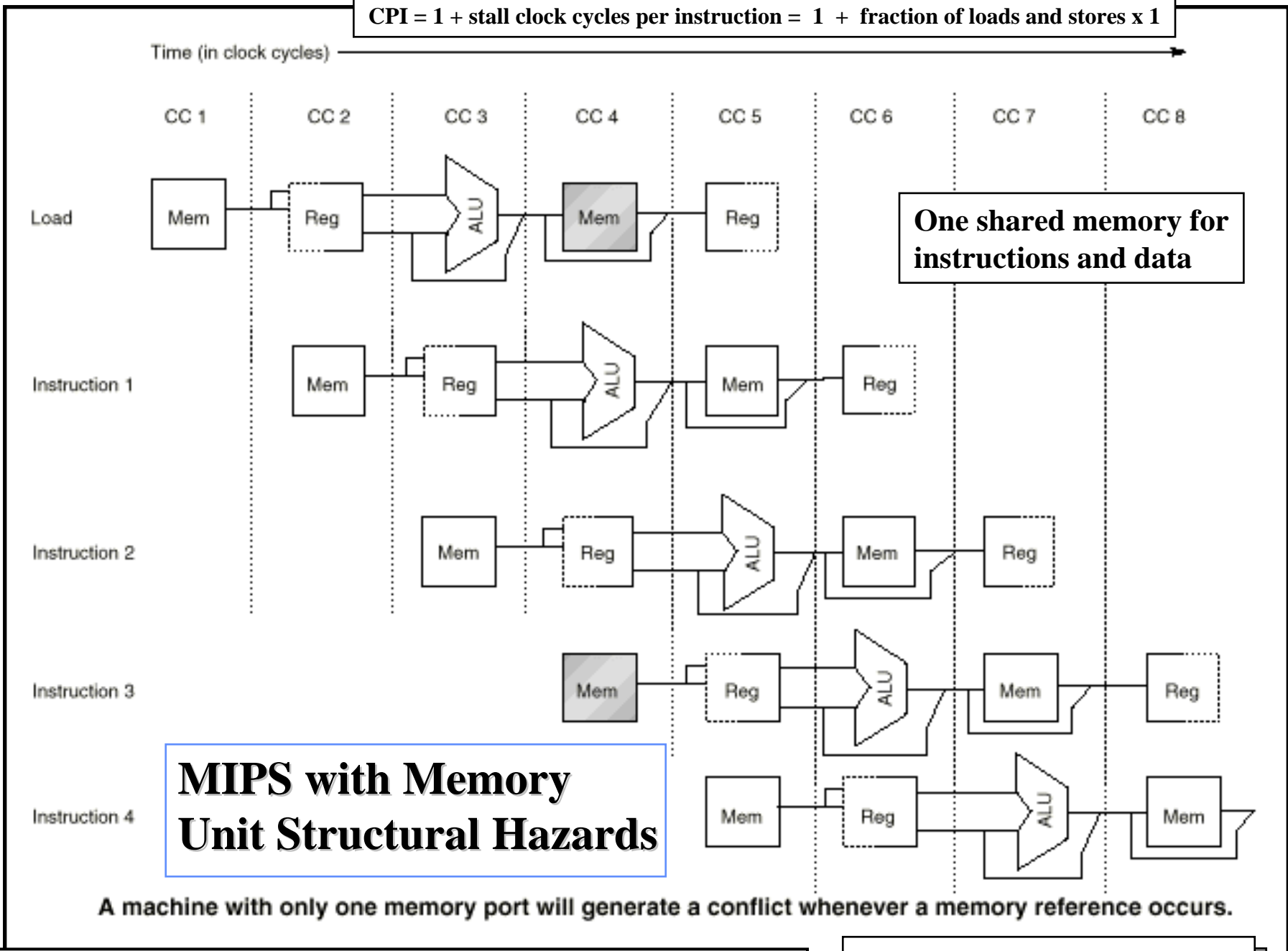
- If a resource conflict arises due to a hardware resource being required by more than one instruction in a single cycle, and one or more such instructions cannot be accommodated, then a structural hazard has occurred, for example:

- e.g. – When a pipelined machine has a shared single-memory for both data and instructions.  
→ stall the pipeline for one cycle for memory data access

i.e A hardware component the instruction requires for correct execution is not available in the cycle needed

**EECC550 - Shaaban**

$$\text{CPI} = 1 + \text{stall clock cycles per instruction} = 1 + \text{fraction of loads and stores} \times 1$$



**One shared memory for instructions and data**

**MIPS with Memory Unit Structural Hazards**

Instructions 1-4 above are assumed to be instructions other than loads/stores

**EECC550 - Shaaban**

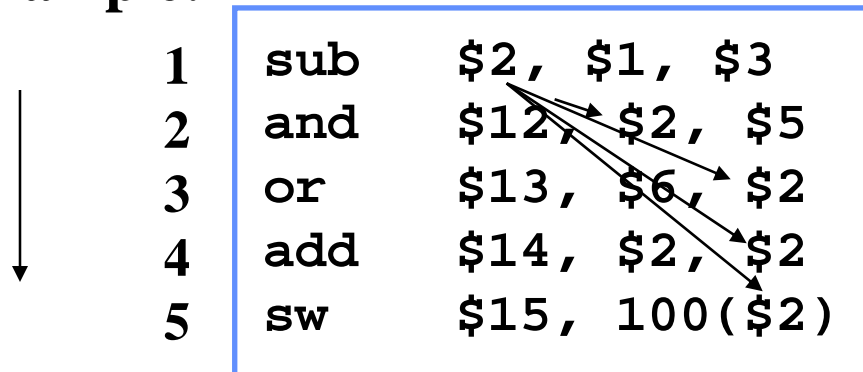
i.e Operands

# Data Hazards

- Data hazards occur when the pipeline changes the order of read/write accesses to instruction operands in such a way that the resulting access order differs from the original sequential instruction operand access order of the unpipelined CPU resulting in incorrect execution.
- Data hazards may require one or more instructions to be stalled in the pipeline to ensure correct execution.

- **Example:**

$$\text{CPI} = 1 + \text{stall clock cycles per instruction}$$



Arrows represent data dependencies between instructions

Instructions that have no dependencies among them are said to be parallel or independent

A high degree of Instruction-Level Parallelism (ILP) is present in a given code sequence if it has a large number of parallel instructions

- All the instructions after the sub instruction use its result data in register \$2
- As part of pipelining, these instruction are started before sub is completed:
  - Due to this data hazard instructions need to be stalled for correct execution.

(As shown next)

i.e Correct operand data not ready yet when needed in EX cycle

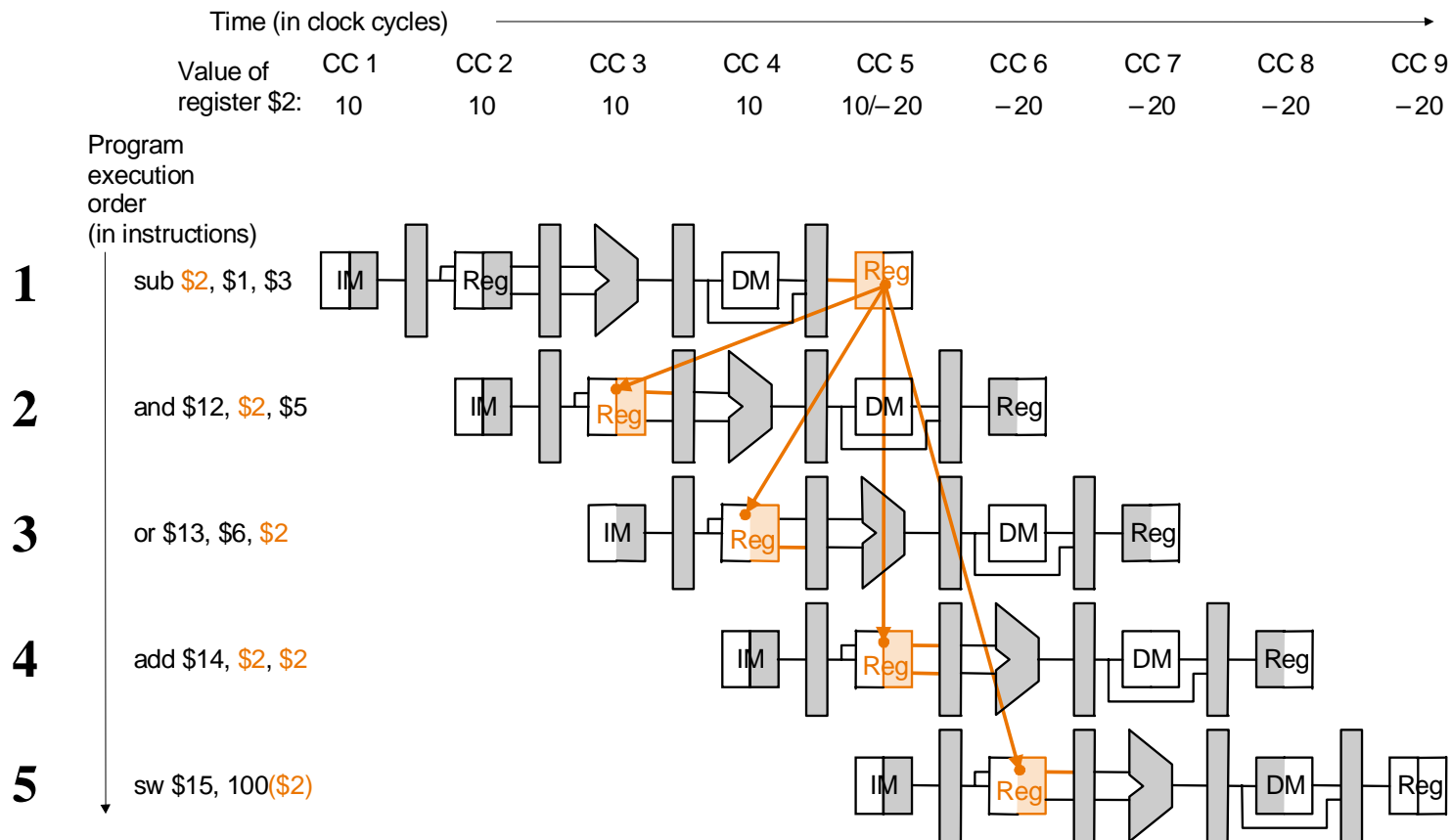
**EECC550 - Shaaban**

# Data Hazards Example

- Problem with starting next instruction before first is finished

- Data dependencies here that “go backward in time” create data hazards.

1	sub	\$2, \$1, \$3
2	and	\$12, \$2, \$5
3	or	\$13, \$6, \$2
4	add	\$14, \$2, \$2
5	sw	\$15, 100(\$2)





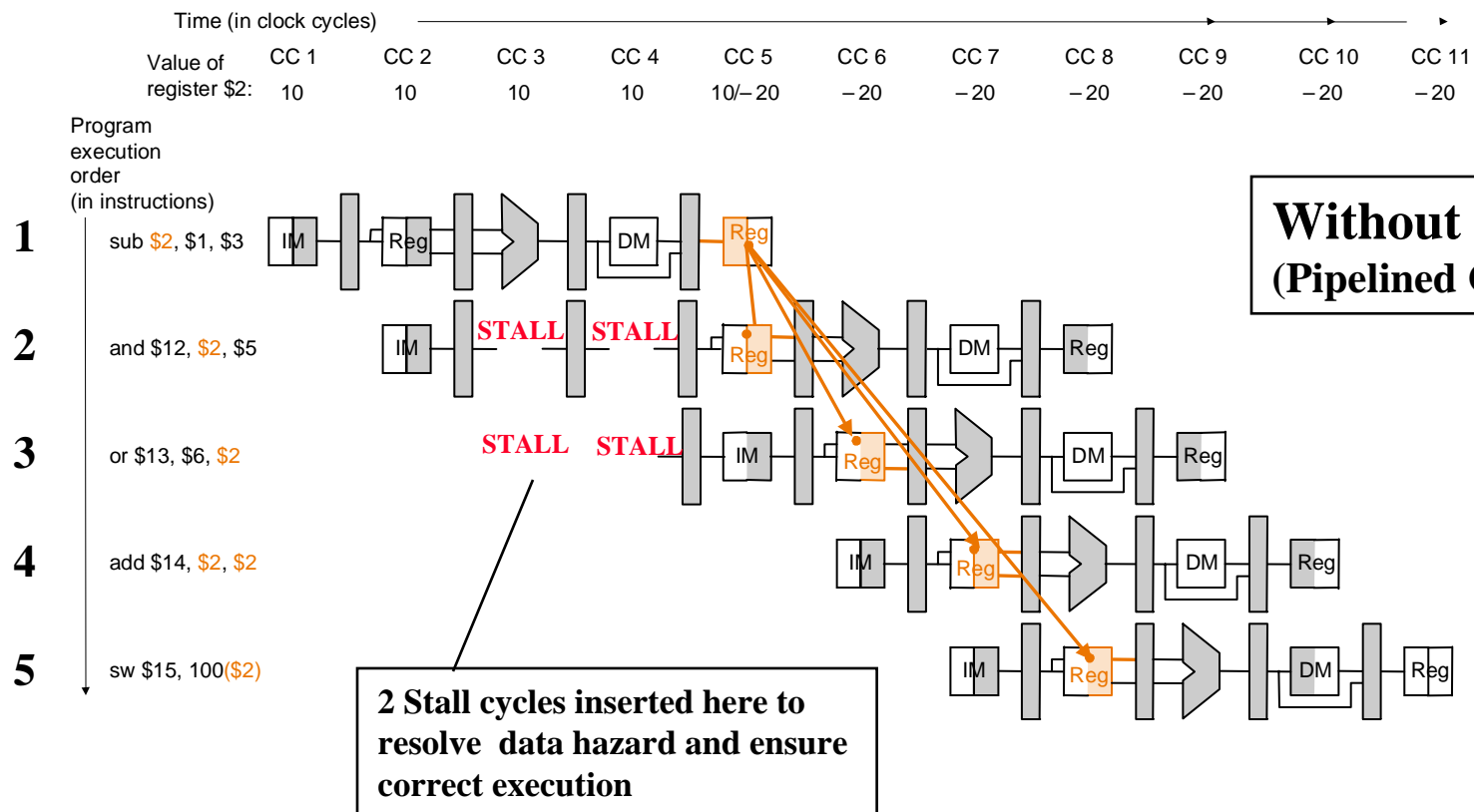
# Data Hazard Resolution: Stall Cycles

Stall the pipeline by a number of cycles.

The control unit must detect the need to insert stall cycles.

In this case two stall cycles are needed.

$$\text{CPI} = 1 + \text{stall clock cycles per instruction}$$



# Data Hazard Resolution/Stall Reduction: Data Forwarding

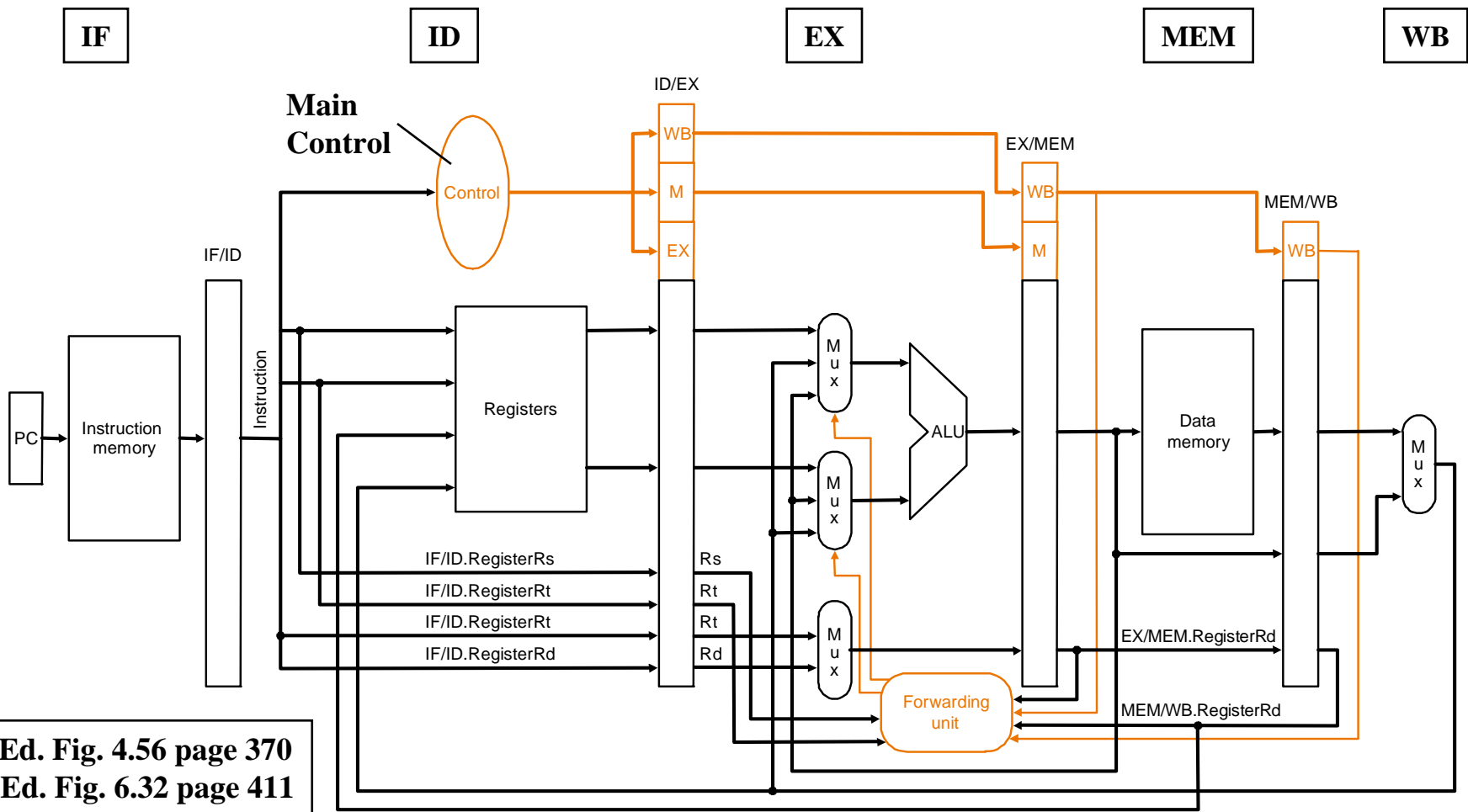
- **Observation:**

Why not use temporary results produced by memory/ALU and not wait for them to be written back in the register bank.

- **Data Forwarding** is a hardware-based technique (also called register bypassing or register short-circuiting) used to eliminate or minimize data hazard stalls that makes use of this observation.
- Using forwarding hardware, the result of an instruction is copied directly (i.e. forwarded) from where it is produced (ALU, memory read port etc.), to where subsequent instructions need it (ALU input register, memory write port etc.)

# Pipelined Datapath With Forwarding

(Pipelined CPU Version 2: With forwarding, Branches still resolved in MEM Stage)

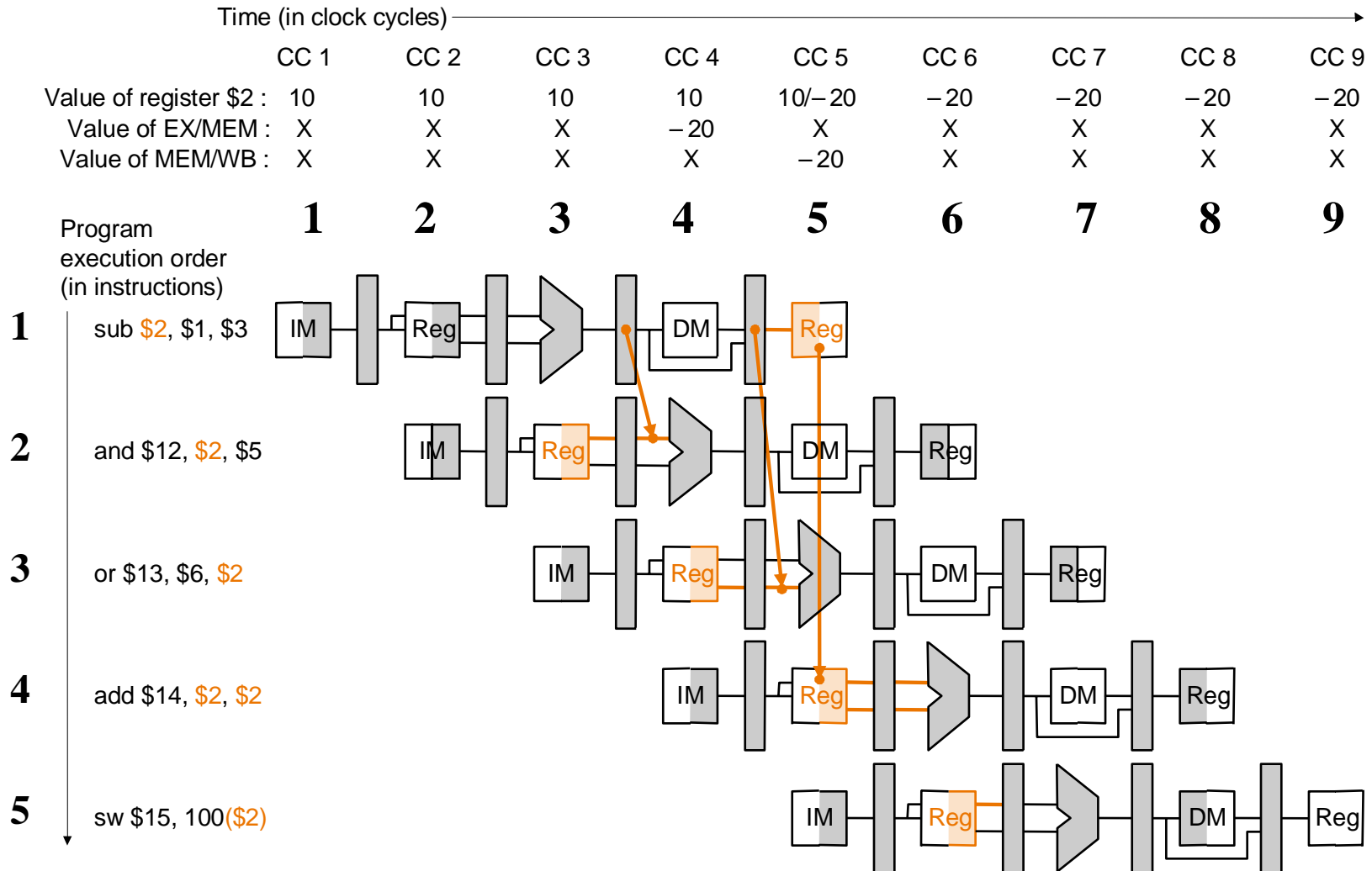


4<sup>th</sup> Ed. Fig. 4.56 page 370  
3<sup>rd</sup> Ed. Fig. 6.32 page 411

- The forwarding unit compares operand registers of the instruction in EX stage with destination registers of the previous two instructions in MEM and WB
- If there is a match one or both operands will be obtained from forwarding paths bypassing the registers

**EECC550 - Shaaban**

# Data Hazard Example With Forwarding



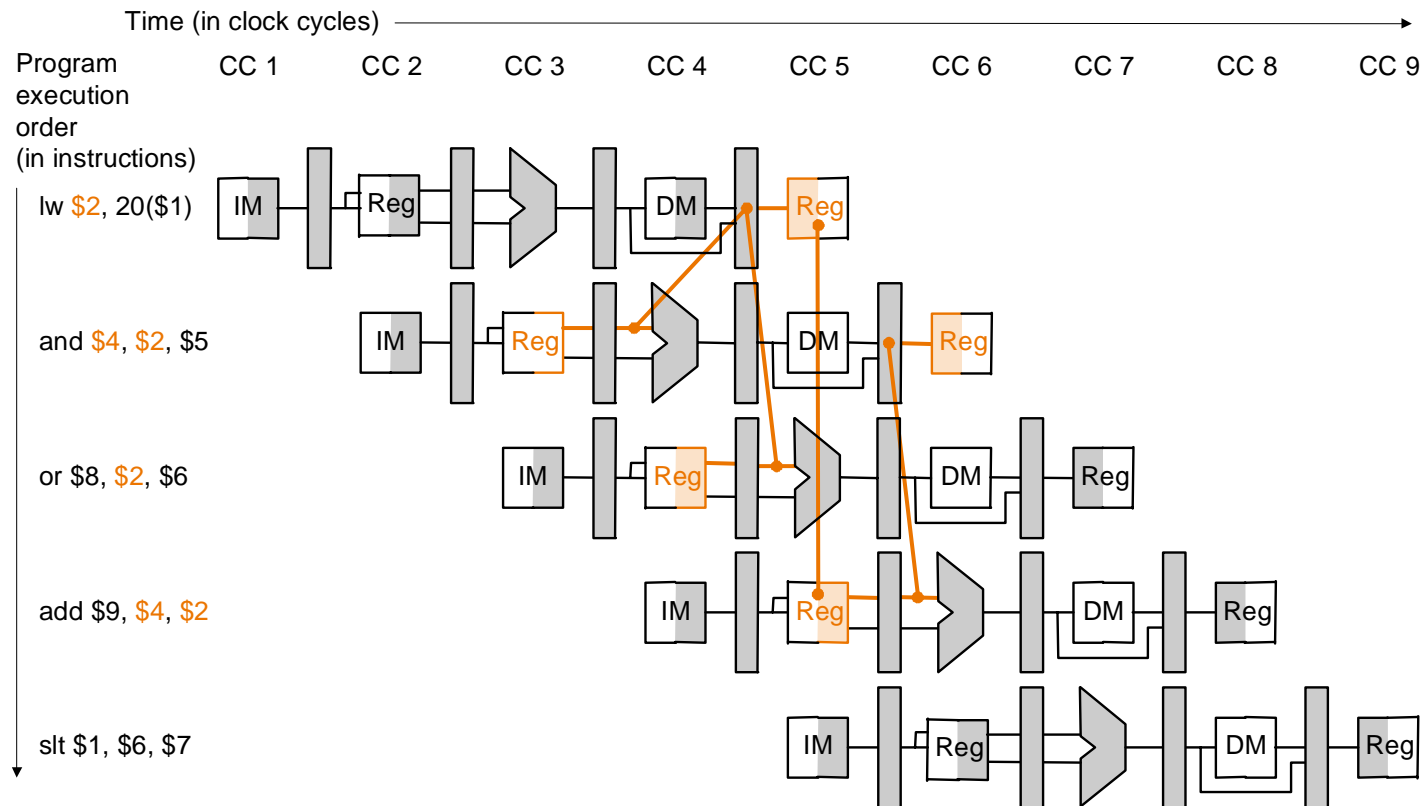
What registers numbers are being compared by the forwarding unit during cycle 5? What about in Cycle 6?

**EECC550 - Shaaban**

# A Data Hazard Requiring A Stall

A load followed by an R-type instruction that uses the loaded value

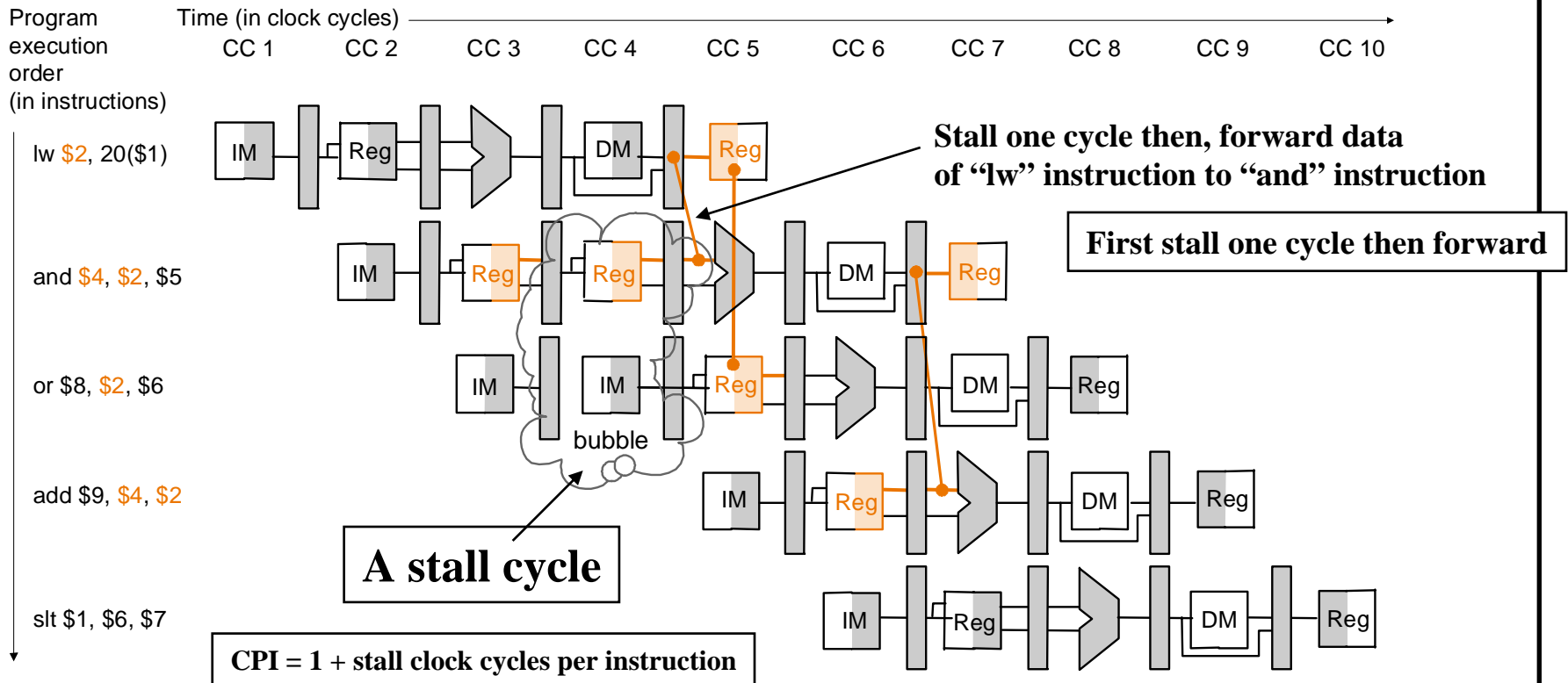
(or any other type of instruction that needs loaded value in ex stage)



**Even with forwarding in place a stall cycle is needed (shown next)  
This condition must be detected by hardware**

# A Data Hazard Requiring A Stall

A load followed by an R-type instruction that uses the loaded value results in a single stall cycle even with forwarding as shown:



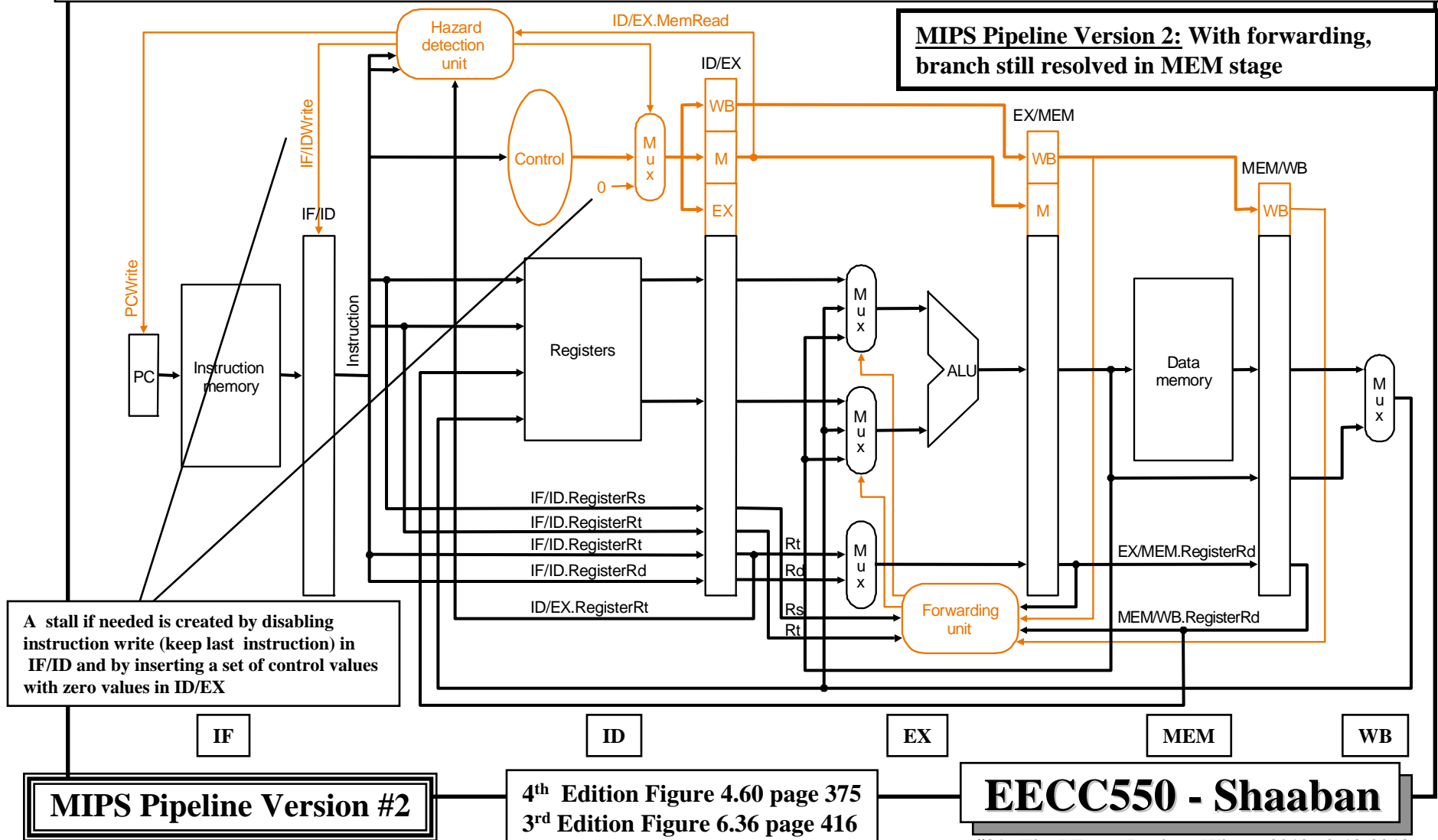
- We can stall the pipeline by keeping all instructions following the "lw" instruction in the same pipeline stage for one cycle

What is the hazard detection unit (shown next slide) doing during cycle 3?

**EECC550 - Shaaban**

# Datapath With Hazard Detection Unit

A load followed by an instruction that uses the loaded value is detected by the hazard detection unit and a stall cycle is inserted.  
The hazard detection unit checks if the instruction in the EX stage is a load by checking its MemRead control line value  
If that instruction is a load it also checks if any of the operand registers of the instruction in the decode stage (ID) match the destination register of the load. In case of a match it inserts a stall cycle (delays decode and fetch by one cycle).



# Compiler Instruction Scheduling (Re-ordering) Example

- Reorder the instructions to avoid as many pipeline stalls as possible:

Stall →

lw	\$15, 0(\$2)	<b>Original Code</b>
lw	\$16, 4(\$2)	
add	\$14, \$5, \$16	
sw	\$16, 4(\$2)	

- The data hazard occurs on register \$16 between the second lw and the add instruction resulting in a stall cycle even with forwarding
- With forwarding we (or the compiler) need to find only one independent instruction to place between them, swapping the lw instructions works:

i.e pipeline version #2	<b>No Stalls</b>	lw	\$16, 4(\$2)	<b>Scheduled Code</b>
i.e pipeline version #1		lw	\$15, 0(\$2)	
	add	\$14, \$5, \$16		
	sw	\$16, 4(\$2)		

- Without forwarding we need two independent instructions to place between them, so in addition a nop is added (or the hardware will insert a stall).

Or stall cycle →

lw	\$16, 4(\$2)
lw	\$15, 0(\$2)
nop	
add	\$14, \$5, \$16
sw	\$16, 4(\$2)



# Control Hazards

- When a conditional branch is executed it may change the PC (when taken) and, without any special measures, leads to stalling the pipeline for a number of cycles until the branch condition is known and PC is updated (branch is resolved).

Here end of stage 4 (MEM)

Versions  
1 and 2

– Otherwise the PC may not be correct when needed in IF

- In current MIPS pipeline, the conditional branch is resolved in stage 4 (MEM stage) resulting in three stall cycles as shown below:

Branch instruction	IF	ID	EX	MEM	WB					
Branch successor		stall	stall	stall	IF	ID	EX	MEM	WB	
Branch successor + 1					IF	ID	EX	MEM	WB	
Branch successor + 2						IF	ID	EX	MEM	
Branch successor + 3							IF	ID	EX	
Branch successor + 4								IF	ID	
Branch successor + 5									IF	

3 stall cycles

Branch Penalty

Correct PC available here  
(end of MEM cycle or stage)

Assuming we stall or flush the pipeline on a branch instruction:

Three clock cycles are wasted for every branch for current MIPS pipeline

Branch Penalty = stage number where branch is resolved - 1

here Branch Penalty = 4 - 1 = 3 Cycles

i.e Correct PC is not available when needed in IF

**EECC550 - Shaaban**

# Basic Branch Handling in Pipelines

- 1 One scheme discussed earlier is to always stall (*flush or freeze*) the pipeline whenever a conditional branch is decoded by holding or deleting any instructions in the pipeline until the branch destination is known (zero pipeline registers, control lines).

Pipeline stall cycles from branches = frequency of branches X branch penalty

- Ex: Branch frequency = 20% branch penalty = 3 cycles

$$\text{CPI} = 1 + .2 \times 3 = 1.6$$

CPI = 1 + stall clock cycles per instruction

- 2 Another method is to assume or predict that the branch is not taken where the state of the machine is not changed until the branch outcome is definitely known. Execution here continues with the next instruction; *stall occurs here when the branch is taken.*

Pipeline stall cycles from branches = frequency of taken branches X branch penalty

- Ex: Branch frequency = 20% of which 45% are taken branch penalty = 3 cycles

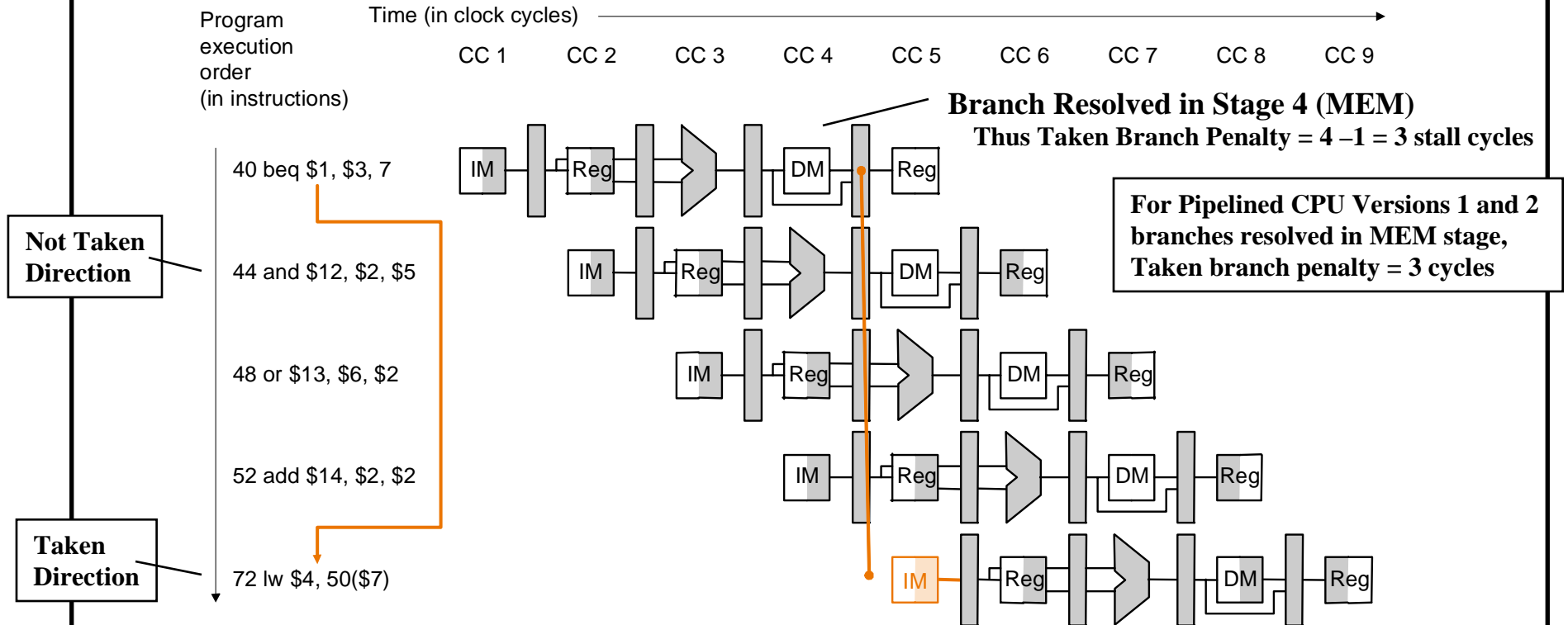
$$\text{CPI} = 1 + .2 \times .45 \times 3 = 1.27$$

CPI = 1 + Average Stalls Per Instruction

EECC550 - Shaaban

# Control Hazards: Example

- Three other instructions are in the pipeline before branch instruction target decision is made when BEQ is in MEM stage.



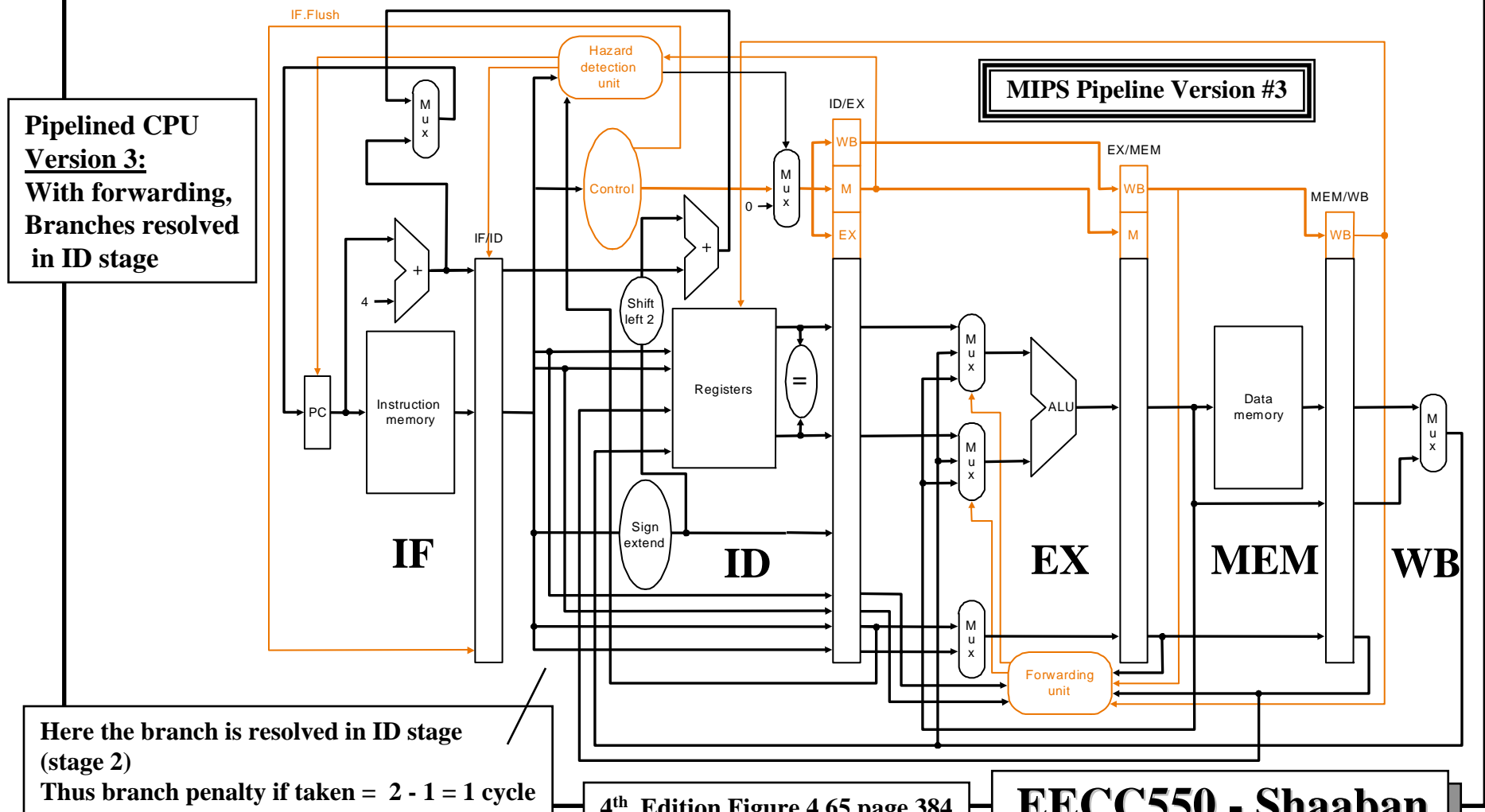
- In the above diagram, we are predicting “branch not taken”
  - Need to add hardware for flushing the three following instructions if we are wrong losing three cycles when the branch is taken.

i.e the branch was resolved as taken in MEM stage

**EECC550 - Shaaban**

# Reducing Delay (Penalty) of Taken Branches

- So far: Next PC of a branch known or resolved in MEM stage: Costs three lost cycles if the branch is taken.
- If next PC of a branch is known or resolved in EX stage, one cycle is saved.
- Branch address calculation can be moved to ID stage (stage 2) using a register comparator, costing only one cycle if branch is taken as shown below. Branch Penalty = stage 2 - 1 = 1 cycle



**Pipelined CPU  
Version 3:  
With forwarding,  
Branches resolved  
in ID stage**

**Here the branch is resolved in ID stage  
(stage 2)  
Thus branch penalty if taken = 2 - 1 = 1 cycle**

**4<sup>th</sup> Edition Figure 4.65 page 384  
3<sup>rd</sup> Edition Figure 6.41 page 427**

**EECC550 - Shaaban**

# Pipeline Performance Example

- Assume the following MIPS instruction mix:

Type	Frequency	
Arith/Logic	40%	
Load	30%	of which 25% are followed immediately by an instruction using the loaded value <span style="border: 1px solid black; padding: 2px;">1 stall</span>
Store	10%	
branch	20%	of which 45% are taken <span style="border: 1px solid black; padding: 2px;">1 stall</span>

- What is the resulting CPI for the pipelined MIPS with forwarding and branch address calculation in ID stage when using the branch not-taken scheme?

i.e Version 3  
Branch Penalty = 1 cycle

- $CPI = \text{Ideal CPI} + \text{Pipeline stall clock cycles per instruction}$

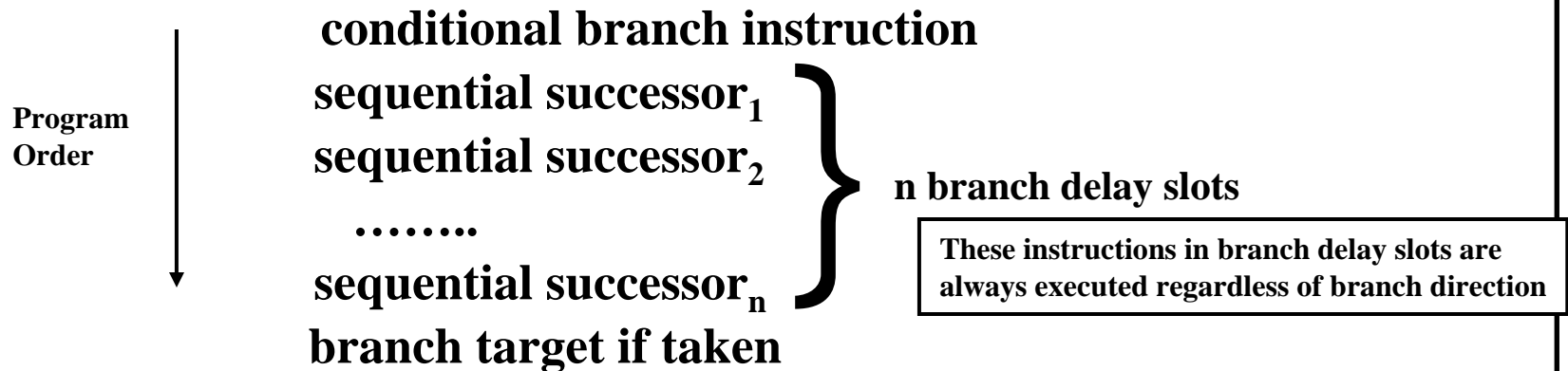
$$\begin{aligned}
 &= 1 + \text{stalls by loads} + \text{stalls by branches} \\
 &= 1 + .3 \times .25 \times 1 + .2 \times .45 \times 1 \\
 &= 1 + .075 + .09 \\
 &= 1.165
 \end{aligned}$$

When the ideal memory assumption is removed this CPI becomes the base CPI with ideal memory or  $CPI_{\text{execution}}$

**EECC550 - Shaaban**

# ISA Reduction of Branch Penalties: Delayed Branch

- When delayed branch is used in an ISA, the branch is delayed by  $n$  cycles (or instructions), following this execution pattern:



- The sequential successor instructions are said to be in the branch delay slots. These instructions are executed whether or not the branch is taken.
- In Practice, all ISAs that utilize delayed branching including MIPS utilize a single instruction branch delay slot. (All RISC ISAs)
  - The job of the compiler is to make the successor instruction in the delay slot a valid and useful instruction.

# Compiler Instruction Scheduling Example With Branch Delay Slot

- Schedule the following MIPS code for the pipelined MIPS CPU with forwarding and reduced branch delay using a single branch delay slot to minimize stall cycles:

```
loop:  lw $1,0($2)           # $1 array element
      add $1, $1, $3        # add constant in $3
      sw $1,0($2)          # store result array element
      addi $2, $2, -4       # decrement address by 4
      bne $2, $4, loop      # branch if $2 != $4
```

- Assuming the initial value of  $\$2 = \$4 + 40$   
(i.e it loops 10 times)
  - What is the CPI and total number of cycles needed to run the code with and without scheduling?

i.e. Pipelined CPU  
Version 3:  
With forwarding,  
Branches resolved  
in ID stage

For MIPS Pipeline Version 3

**EECC550 - Shaaban**

# Compiler Instruction Scheduling Example (With Branch Delay Slot)

- Without compiler scheduling

loop: lw \$1,0(\$2)

Stall

add \$1, \$1, \$3

sw \$1,0(\$2)

addi \$2, \$2, -4

Needed because new value of \$2 is not produced yet →

Stall

bne \$2, \$4, loop

Stall (or NOP)

Ignoring the initial 4 cycles to fill the

pipeline:

Each iteration takes = 8 cycles

CPI =  $8/5 = 1.6$

Total cycles =  $8 \times 10 = 80$  cycles

- With compiler scheduling

loop: lw \$1,0(\$2)

addi \$2, \$2, -4

add \$1, \$1, \$3

bne \$2, \$4, loop

sw \$1, 4(\$2)

Move between lw add

Move to branch delay slot

Adjust address offset

Ignoring the initial 4 cycles to fill the

pipeline:

Each iteration takes = 5 cycles

CPI =  $5/5 = 1$

Total cycles =  $5 \times 10 = 50$  cycles

Speedup =  $80/50 = 1.6$

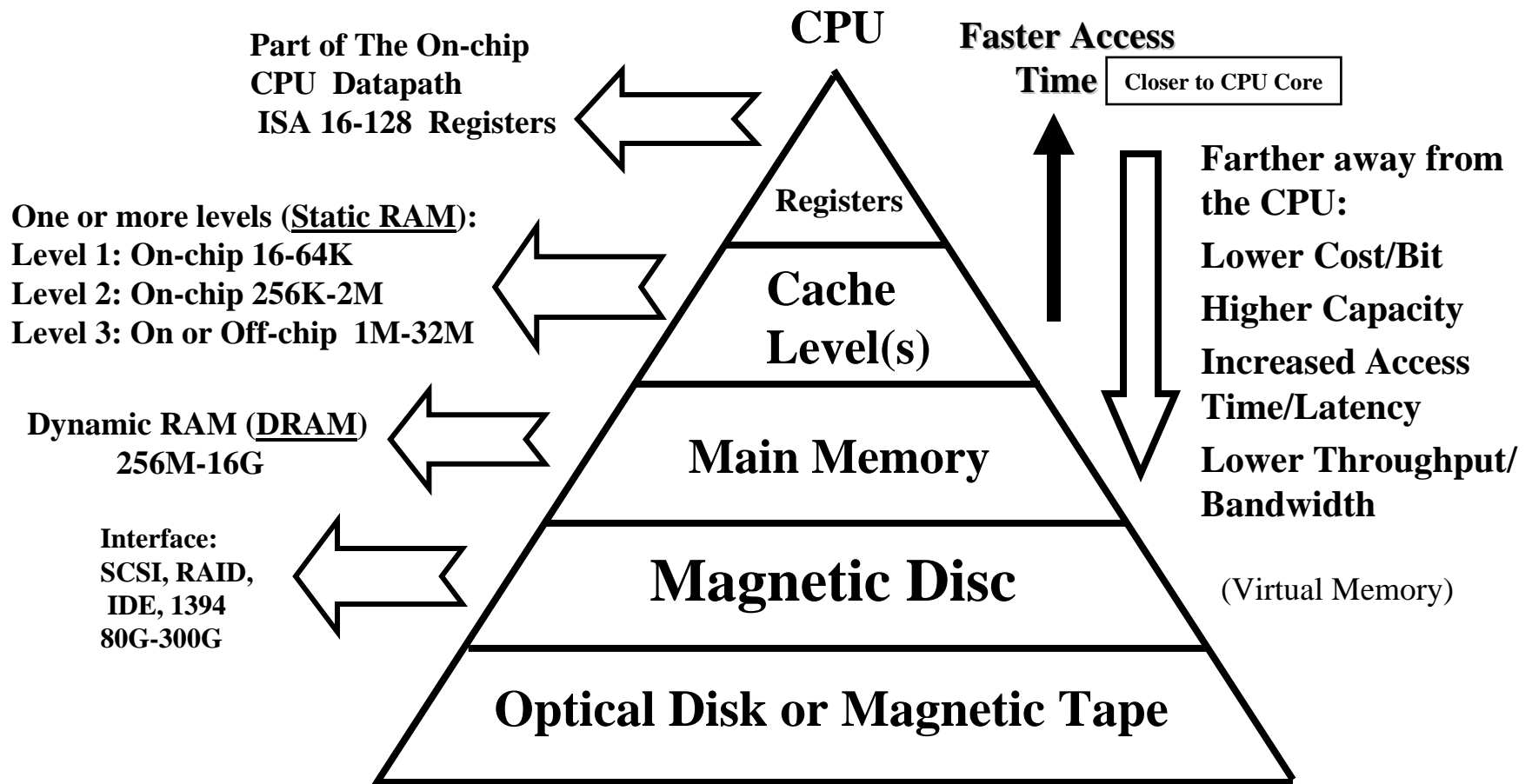
Target CPU: Pipelined CPU Version 3: With forwarding,  
Branches resolved in ID stage

EECC550 - Shaaban



# Levels of The Memory Hierarchy

In this course, we concentrate on the design, operation and performance of a single level of cache L1 (either unified or separate) when using non-ideal main memory



4<sup>th</sup> Edition Chapter 5.1-5.3 - 3<sup>rd</sup> Edition Chapter 7.1-5.3

EECC550 - Shaaban

# Memory Hierarchy Operation

- If an instruction or operand is required by the CPU, the levels of the memory hierarchy are searched for the item starting with the level closest to the CPU (Level 1 cache):

Hit rate for level one cache =  $H_1$

If the item is found, it's delivered to the CPU resulting in a cache hit without searching lower levels.

Hit rate for level one cache =  $H_1$

Cache Miss

– If the item is missing from an upper level, resulting in a cache miss, the level just below is searched.

Miss rate for level one cache =  $1 - \text{Hit rate} = 1 - H_1$

– For systems with several levels of cache, the search continues with cache level 2, 3 etc.

– If all levels of cache report a miss then main memory is accessed for the item.

- CPU ↔ cache ↔ memory: Managed by hardware.

– If the item is not found in main memory resulting in a page fault, then disk (virtual memory), is accessed for the item.

- Memory ↔ disk: Managed by the operating system with hardware support

In this course, we concentrate on the design, operation and performance of a single level of cache L1 (either unified or separate) when using non-ideal main memory

**EECC550 - Shaaban**

# Memory Hierarchy: Terminology

- **A Block:** The smallest unit of information transferred between two levels.
- **Hit:** Item is found in some block in the upper level (example: Block X)

e. g.  $H_1$

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

– **Hit Time:** Time to access the upper level which consists of

Ideally = 1 Cycle

Hit rate for level one cache =  $H_1$

(S)RAM access time + Time to determine hit/miss

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

e. g.  $1 - H_1$

– **Miss Rate** =  $1 - (\text{Hit Rate})$

Miss rate for level one cache =  $1 - \text{Hit rate} = 1 - H_1$

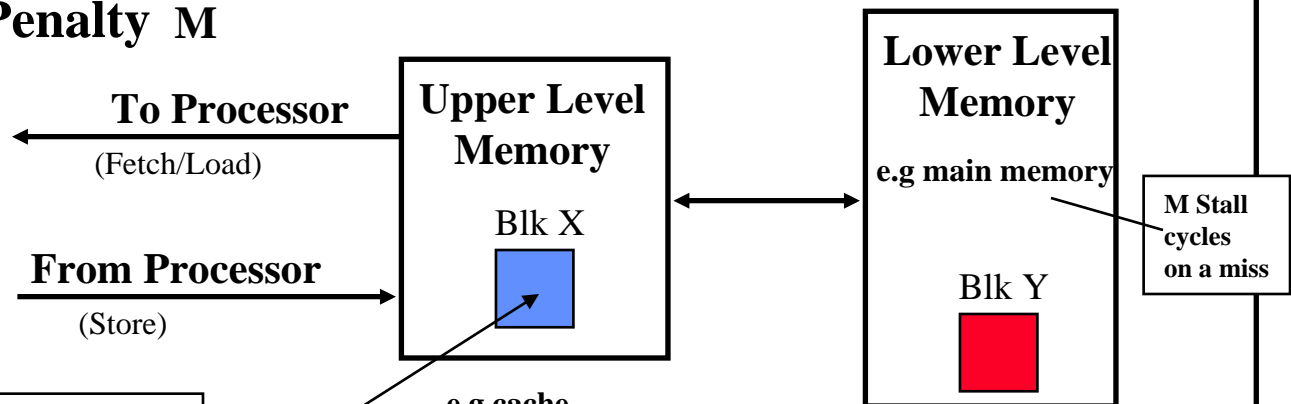
– **Miss Penalty:** Time to replace a block in the upper level +

M

Time to deliver the missed block to the processor

- **Hit Time**  $\ll$  Miss Penalty M

Ideally = 1 Cycle



Typical Cache Block (or line) Size: 16-64 bytes

A block e.g. cache

Hit if block is found in cache

**EECC550 - Shaaban**

# Basic Cache Concepts

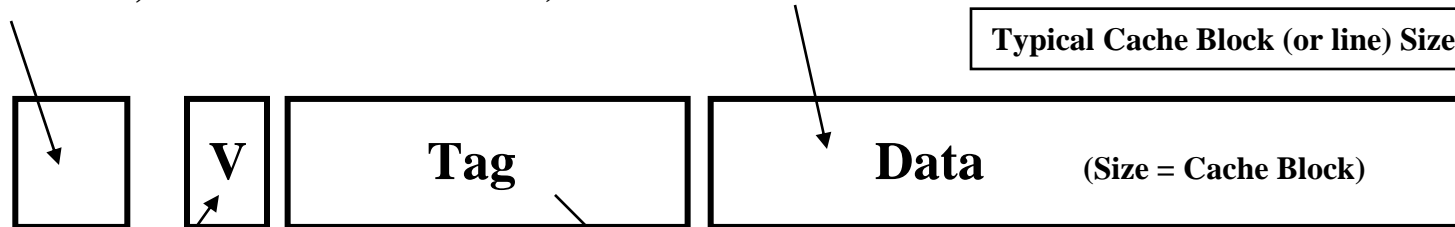
- Cache is the first level of the memory hierarchy once the address leaves the CPU and is searched first for the requested data.
- If the data requested by the CPU is present in the cache, it is retrieved from cache and the data access is a cache hit otherwise a cache miss and data must be read from main memory.
- On a cache miss a block of data must be brought in from main memory to cache to possibly replace an existing cache block.
- The allowed block addresses where blocks can be mapped (placed) into cache from main memory is determined by cache placement strategy.
- Locating a block of data in cache is handled by cache block identification mechanism (tag checking).
- On a cache miss choosing the cache block being removed (replaced) is handled by the block replacement strategy in place.

# Cache Block Frame

Cache is comprised of a number of cache block frames

Other status/access bits:  
(e.g. modified, read/write access bits)

Data Storage: Number of bytes is the size of a cache block or cache line size (Cached instructions or data go here)



Valid Bit: Indicates whether the cache block frame contains valid data

Tag: Used to identify if the address supplied matches the address of the data stored

The tag and valid bit are used to determine whether we have a cache hit or miss

Nominal  
Cache  
Size

Stated nominal cache capacity or size only accounts for space used to store instructions/data and ignores the storage needed for tags and status bits:

**Nominal Cache Capacity = Number of Cache Block Frames x Cache Block Size**

e.g For a cache with block size = 16 bytes and  $1024 = 2^{10} = 1\text{k}$  cache block frames  
Nominal cache capacity =  $16 \times 1\text{k} = 16 \text{ Kbytes}$

Cache utilizes faster memory (SRAM)

**EECC550 - Shaaban**

# Locating A Data Block in Cache

- Each block frame in cache has an address tag.
- The tags of every cache block that might contain the required data are checked or searched in parallel. Tag Matching
- A valid bit is added to the tag to indicate whether this entry contains a valid address.
- The byte address from the CPU to cache is divided into:
  - A block address, further divided into:
    - 1 • An index field to choose/map a block set in cache.  
(no index field when fully associative).
    - 2 • A tag field to search and match addresses in the selected set.
  - A byte block offset to select the data from the block. 3



Index = Mapping

**EECC550 - Shaaban**

# Cache Organization & Placement Strategies

Placement strategies or mapping of a main memory data block onto cache block frame addresses divide cache into three organizations:

1 **Direct mapped cache**: A block can be placed in only one location (cache block frame), given by the mapping function: Least complex to implement

Mapping  
Function

$$\text{index} = (\text{Block address}) \text{ MOD } (\text{Number of blocks in cache})$$

2 **Fully associative cache**: A block can be placed anywhere in cache. (no mapping function). Most complex cache organization to implement = Frame #

3 **Set associative cache**: A block can be placed in a restricted set of places, or cache block frames. A set is a group of block frames in the cache. A block is first mapped onto the set and then it can be placed anywhere within the set. The set in this case is chosen by:

Mapping  
Function

$$\text{index} = (\text{Block address}) \text{ MOD } (\text{Number of sets in cache})$$
= Set #

If there are  $n$  blocks in a set the cache placement is called  $n$ -way set-associative. Most common cache organization

**EECC550 - Shaaban**

# Cache Organization: Direct Mapped Cache



Cache Block Frame

A block in memory can be placed in one location (cache block frame) only, given by: (Block address) MOD (Number of blocks in cache)

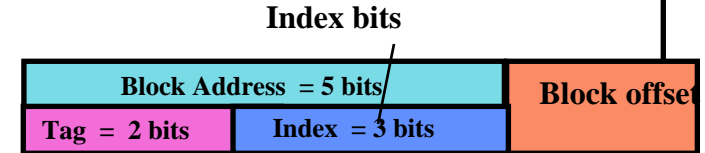
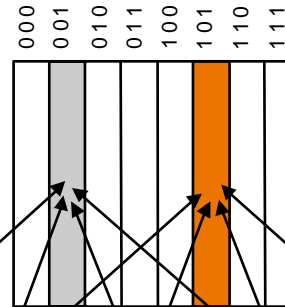
In this case, mapping function: (Block address) MOD (8) = Index

Index

Cache

(i.e low three bits of block address)

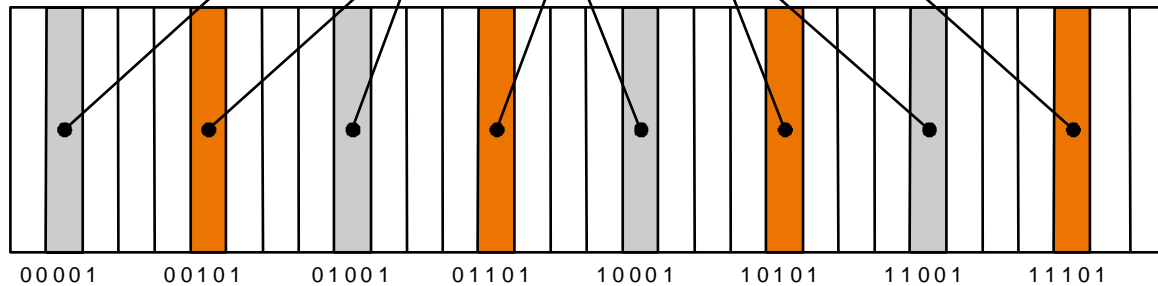
8 cache block frames



Here four blocks in memory map to the same cache block frame

Example:  
29 MOD 8 = 5  
(11101) MOD (1000) = 101

32 memory blocks cacheable



Index size =  $\log_2 8 = 3$  bits

Memory

Limitation of Direct Mapped Cache: Conflicts between memory blocks that map to the same cache block frame may result in conflict cache misses

**EECC550 - Shaaban**



# 4KB Direct Mapped Cache Example

4 Kbytes = Nominal Cache Capacity

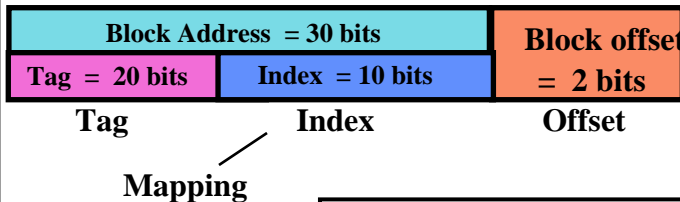
1K =  $2^{10}$  = 1024 Blocks  
 Each block = one word  
 (4 bytes)

Can cache up to  
 $2^{32}$  bytes = 4 GB  
 of memory

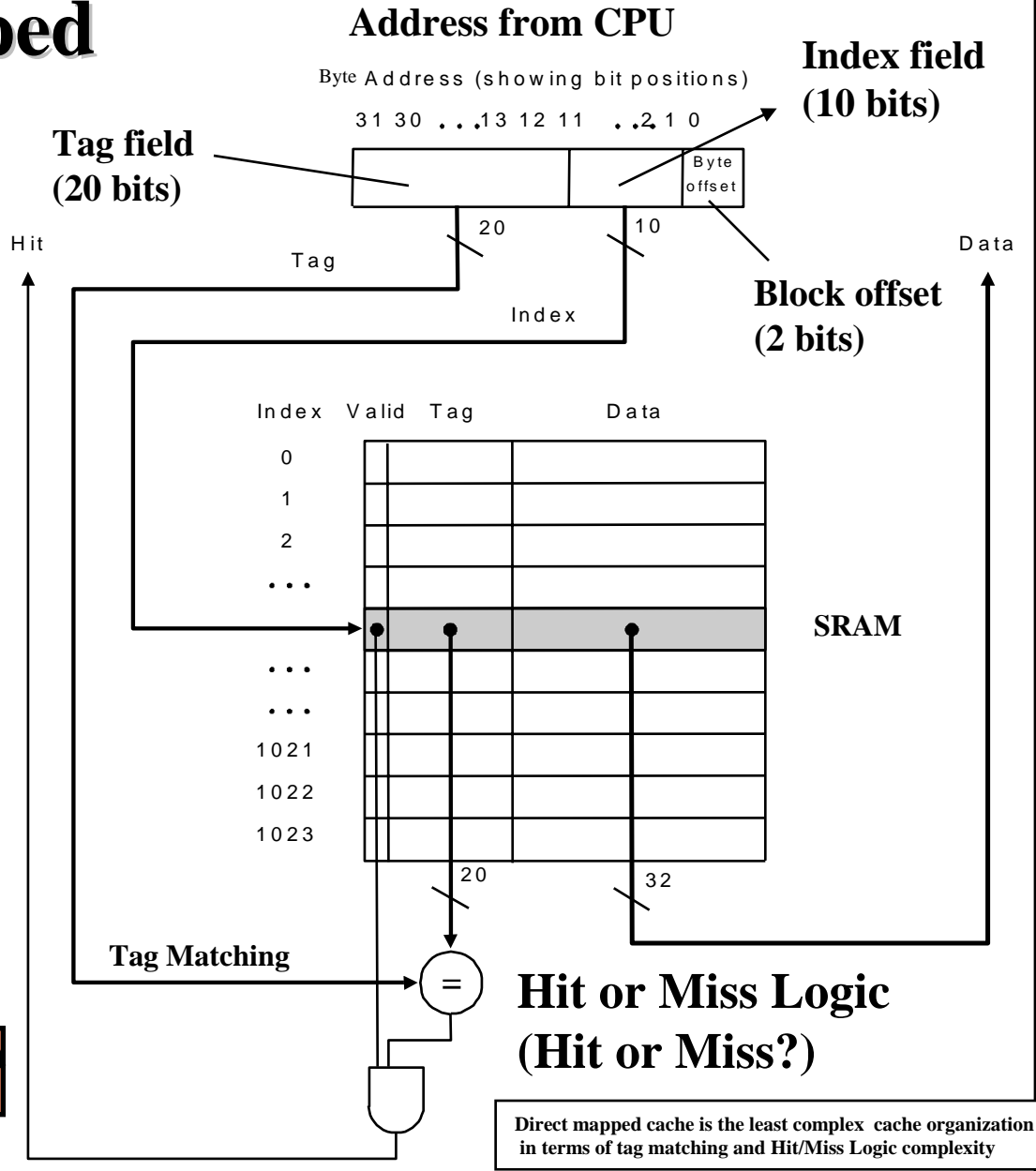
## Mapping function:

Cache Block frame number =  
 (Block address) MOD (1024)

i.e . Index field or 10 low bits of  
 block address



Hit Access Time = SRAM Delay + Hit/Miss Logic Delay



**Hit or Miss Logic  
 (Hit or Miss?)**

**EECC550 - Shaaban**

# Direct Mapped Cache Operation Example

- Given a series of 16 memory address references given as word addresses:

1, 4, 8, 5, 20, 17, 19, 56, 9, 11, 4, 43, 5, 6, 9, 17.

Here:  
Block Address = Word Address

- Assume a direct mapped cache with 16 one-word blocks that is initially empty, label each reference as a hit or miss and show the final content of cache

- Here: Block Address = Word Address      Mapping Function = (Block Address) MOD 16 = Index

Cache Block Frame#	1	4	8	5	20	17	19	56	9	11	4	43	5	6	9	17	Hit/Miss
	Miss	Miss	Miss	Miss	Miss	Miss	Miss	Miss	Miss	Miss	Miss	Miss	Hit	Miss	Hit	Hit	
0																	
1	1	1	1	1	1	17	17	17	17	17	17	17	17	17	17	17	17
2																	
3							19	19	19	19	19	19	19	19	19	19	19
4		4	4	4	20	20	20	20	20	20	4	4	4	4	4	4	4
5				5	5	5	5	5	5	5	5	5	5	5	5	5	5
6														6	6	6	6
7																	
8			8	8	8	8	8	56	56	56	56	56	56	56	56	56	56
9									9	9	9	9	9	9	9	9	9
10																	
11										11	11	43	43	43	43	43	43
12																	
13																	
14																	
15																	

Initial Cache Content (empty)

**Cache Content After Each Reference**

Final Cache Content

**Hit Rate = # of hits / # memory references = 3/16 = 18.75%**

Mapping Function = Index = (Block Address) MOD 16  
i.e 4 low bits of block address

EECC550 - Shaaban

Nominal Capacity

# 64KB Direct Mapped Cache Example

4K =  $2^{12} = 4096$  blocks

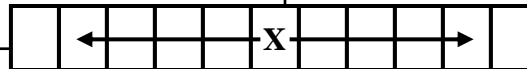
Each block = four words = 16 bytes

Can cache up to  $2^{32}$  bytes = 4 GB of memory

SRAM

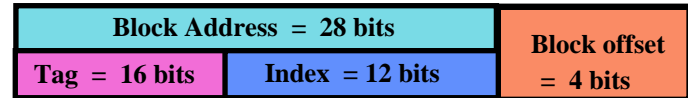
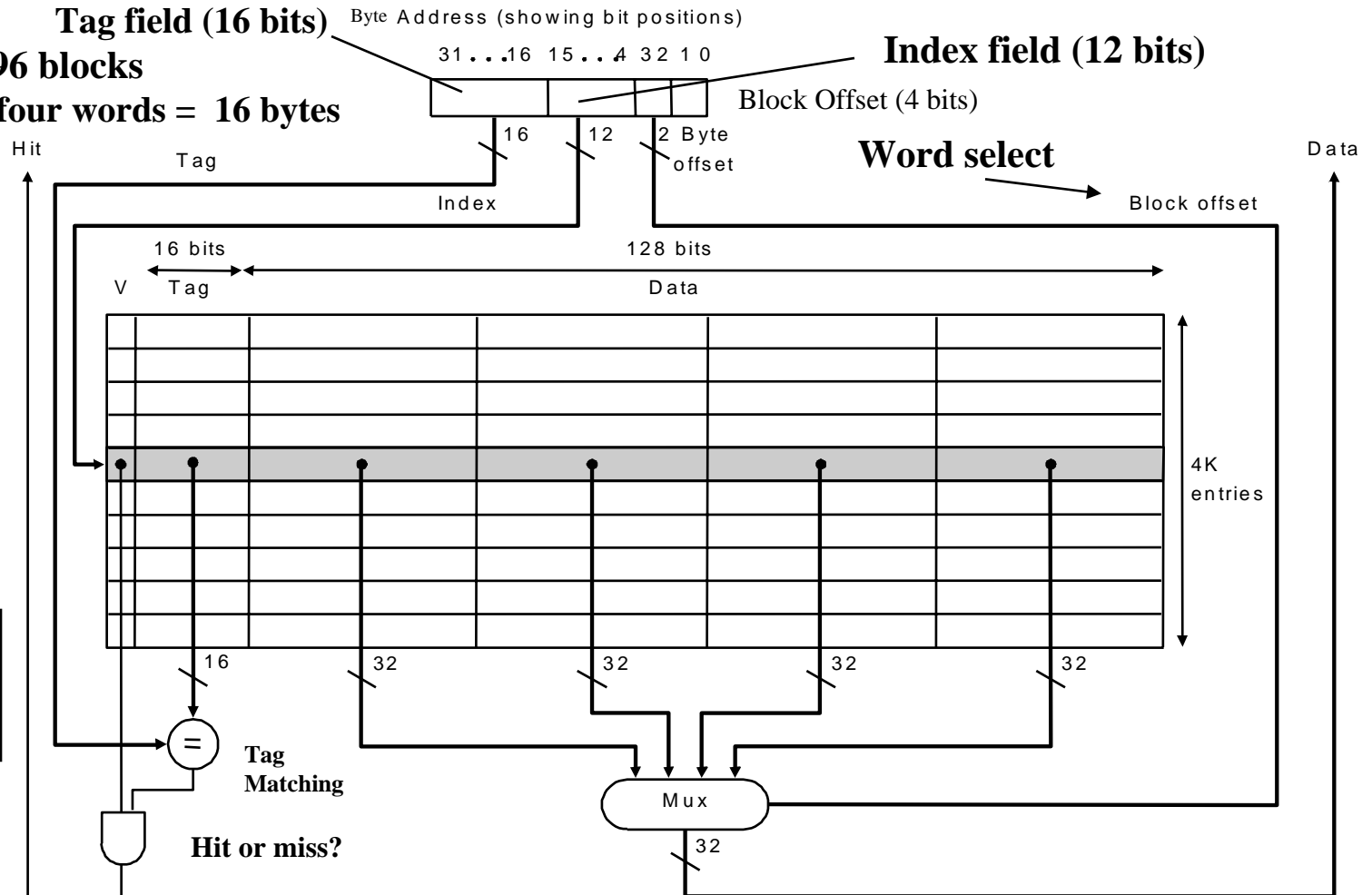
Typical cache Block or line size: 64 bytes

Larger cache blocks take better advantage of spatial locality and thus may result in a lower miss rate



**Mapping Function:** Cache Block frame number = (Block address) MOD (4096)  
i.e. index field or 12 low bit of block address

Hit Access Time = SRAM Delay + Hit/Miss Logic Delay



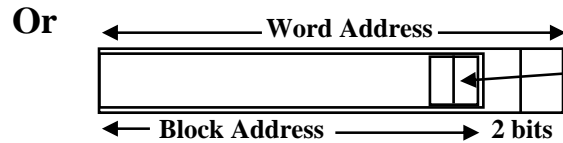
**EECC550 - Shaaban**

# Direct Mapped Cache Operation Example

With Larger Cache Block Frames

- Given the same series of 16 memory address references given as word addresses:  
1, 4, 8, 5, 20, 17, 19, 56, 9, 11, 4, 43, 5, 6, 9, 17.
- Assume a direct mapped cache with four word blocks and a total of 16 words that is initially empty, label each reference as a hit or miss and show the final content of cache
- Cache has  $16/4 = 4$  cache block frames (each has four words)
- Here:  $\text{Block Address} = \text{Integer} (\text{Word Address}/4)$

i.e We need to find block addresses for mapping



Mapping Function = (Block Address) MOD 4

i.e 2 low bits of block address

Cache Block Frame#	0	1	2	1	5	4	4	14	2	2	1	10	1	1	2	4	Word addresses
	1	4	8	5	20	17	19	56	9	11	4	43	5	6	9	17	
	Miss	Miss	Miss	Hit	Miss	Miss	Hit	Miss	Miss	Hit	Miss	Miss	Hit	Hit	Miss	Hit	Hit/Miss
0	0	0	0	0	0	16	16	16	16	16	16	16	16	16	16	16	
1		4	4	4	20	20	20	20	20	20	4	4	4	4	4	4	
2			8	8	8	8	8	56	8	8	8	40	40	40	8	8	
3																	

Initial Cache Content (empty)

Starting word address of Cache Frames Content After Each Reference

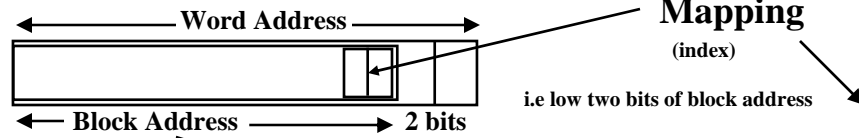
Final Cache Content

Hit Rate = # of hits / # memory references = 6/16 = 37.5%

Here: Block Address  $\neq$  Word Address

EECC550 - Shaaban

Block size = 4 words



**Word Addresses vs. Block Addresses and Frame Content for Previous Example**

Given Word address	Block address	Cache Block Frame # (Block address)mod 4	word address range in frame (4 words)
1	0	0	0-3
4	1	1	4-7
8	2	2	8-11
5	1	1	4-7
20	5	1	20-23
17	4	0	16-19
19	4	0	16-19
56	14	2	56-59
9	2	2	8-11
11	2	2	8-11
4	1	1	4-7
43	10	2	40-43
5	1	1	4-7
6	1	1	4-7
9	2	2	8-11
17	4	0	16-19

**Block Address = Integer (Word Address/4)**

**EECC550 - Shaaban**

# Cache Organization:



Cache Block Frame

## Set Associative Cache

Why set associative?

Set associative cache reduces cache misses by reducing conflicts between blocks that would have been mapped to the same cache block frame in the case of direct mapped cache

One-way set associative (direct mapped)

Block	Tag	Data
0		
1		
2		
3		
4		
5		
6		
7		

1-way set associative: (direct mapped) 1 block frame per set

Two-way set associative

Set	Tag	Data	Tag	Data
0				
1				
2				
3				

2-way set associative: 2 blocks frames per set

Four-way set associative

Set	Tag	Data	Tag	Data	Tag	Data	Tag	Data
0								
1								

4-way set associative: 4 blocks frames per set

8-way set associative: 8 blocks frames per set  
In this case it becomes fully associative since total number of block frames = 8

Eight-way set associative (fully associative)

Tag	Data	Tag	Data	Tag	Data	Tag	Data	Tag	Data	Tag	Data	Tag	Data

A cache with a total of 8 cache block frames shown

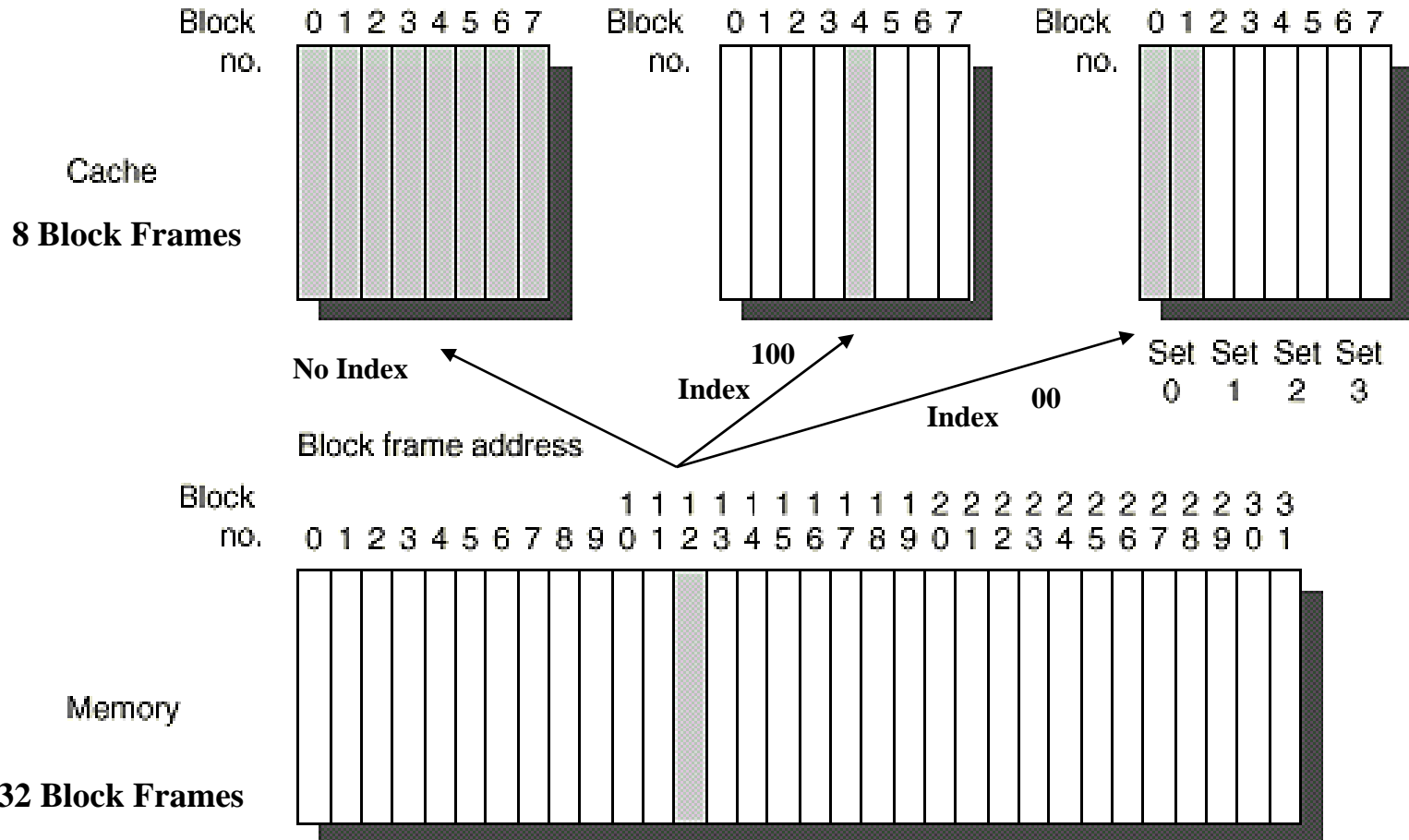
**EECC550 - Shaaban**

# Cache Organization/Mapping Example

Fully associative:  
block 12 can go  
anywhere  
(No mapping function)

Direct mapped:  
block 12 can go  
only into block 4  
(12 mod 8) = index = 100

2-way  
Set associative:  
block 12 can go  
anywhere in set 0  
(12 mod 4) = index = 00



This example cache has eight block frames and memory has 32 blocks.

$$12 = 1100$$

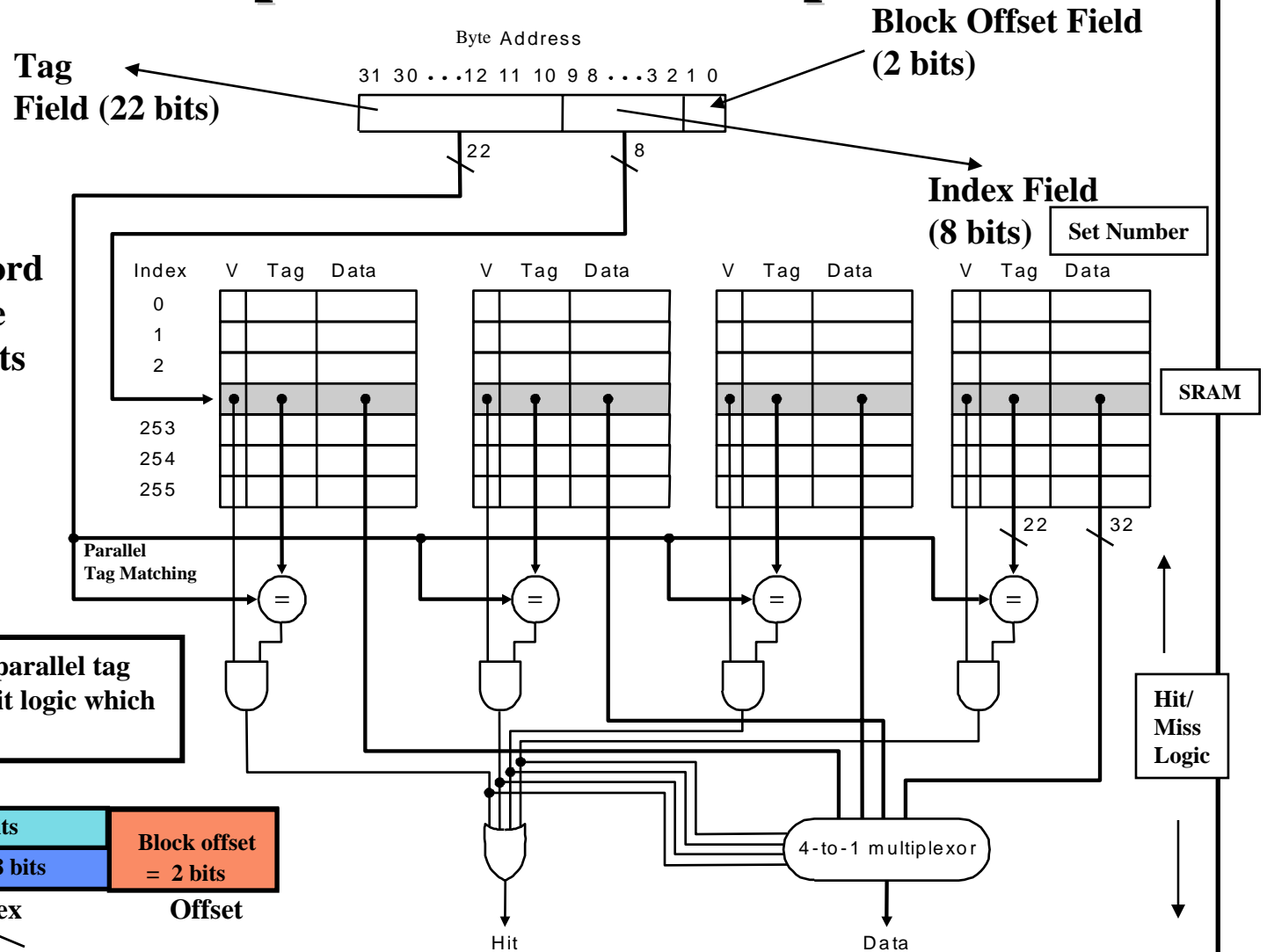
**EECC550 - Shaaban**

# 4K Four-Way Set Associative Cache: MIPS Implementation Example

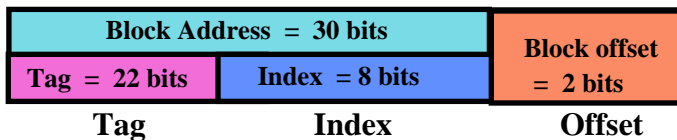
Nominal Capacity

1024 block frames  
Each block = one word  
4-way set associative  
 $1024 / 4 = 2^8 = 256$  sets

Can cache up to  
 $2^{32}$  bytes = 4 GB  
of memory



Set associative cache requires parallel tag matching and more complex hit logic which may increase hit time



Mapping Function: Cache Set Number =  $\text{index} = (\text{Block address}) \text{ MOD } (256)$

Hit Access Time = SRAM Delay + Hit/Miss Logic Delay

**EECC550 - Shaaban**



# Cache Replacement Policy

Which block to replace on a cache miss?

- When a cache miss occurs the cache controller may have to select a block of cache data to be removed from a cache block frame and replaced with the requested data, such a block is selected by one of three methods:

*(No cache replacement policy in direct mapped cache)*

No choice on which block to replace

## 1 – Random:

- Any block is randomly selected for replacement providing uniform allocation.
- Simple to build in hardware. Most widely used cache replacement strategy.

## 2 – Least-recently used (LRU):

- Accesses to blocks are recorded and the block replaced is the one that was not used for the longest period of time.
- Full LRU is *expensive* to implement, as the number of blocks to be tracked increases, and is usually approximated by block usage bits that are cleared at regular time intervals.

## 3 – First In, First Out (FIFO):

- Because LRU can be complicated to implement, this approximates LRU by determining the oldest block rather than LRU

**EECC550 - Shaaban**

# Miss Rates for Caches with Different Size, Associativity & Replacement Algorithm

## Sample Data

Nominal

Associativity: Size	2-way		4-way		8-way	
	LRU	Random	LRU	Random	LRU	Random
16 KB	5.18%	5.69%	4.67%	5.29%	4.39%	4.96%
64 KB	1.88%	2.01%	1.54%	1.66%	1.39%	1.53%
256 KB	1.15%	1.17%	1.13%	1.13%	1.12%	1.12%

Program steady state cache miss rates are given  
Initially cache is empty and miss rates ~ 100%

FIFO replacement miss rates (not shown here) is better than random but worse than LRU

For SPEC92

Miss Rate = 1 – Hit Rate = 1 – H1

EECC550 - Shaaban

# 2-Way Set Associative Cache Operation Example

- Given the same series of 16 memory address references given as word addresses: Here: Block Address = Word Address  
1, 4, 8, 5, 20, 17, 19, 56, 9, 11, 4, 43, 5, 6, 9, 17. (LRU Replacement)
- Assume a two-way set associative cache with one word blocks and a total size of 16 words that is initially empty, label each reference as a hit or miss and show the final content of cache
- Here: Block Address = Word Address      Mapping Function = Set # = (Block Address) MOD 8

Cache Set #	1	4	8	5	20	17	19	56	9	11	4	43	5	6	9	17	
	Miss	Miss	Miss	Miss	Miss	Miss	Miss	Miss	Miss	Miss	Hit	Miss	Hit	Miss	Hit	Hit	Hit/Miss
<b>0</b>			8	8	8	8	8	8	8	8	8	8	8	8	8	8	LRU
								56	56	56	56	56	56	56	56	56	
<b>1</b>	1	1	1	1	1	1	1	1	9	9	9	9	9	9	9	9	LRU
						17	17	17	17	17	17	17	17	17	17	17	
<b>2</b>																	
							19	19	19	19	19	43	43	43	43	43	
<b>3</b>										11	11	11	11	11	11	11	LRU
		4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	
<b>4</b>					20	20	20	20	20	20	20	20	20	20	20	20	LRU
				5	5	5	5	5	5	5	5	5	5	5	5	5	
<b>5</b>																	
														6	6	6	
<b>6</b>																	
<b>7</b>																	

Initial  
Cache  
Content  
(empty)

Cache Content After Each Reference

Final  
Cache  
Content

Hit Rate = # of hits / # memory references = 4/16 = 25%

**EECC550 - Shaaban**

Replacement policy: LRU = Least Recently Used

# Address Field Sizes/Mapping

← Physical Address Generated by CPU →  
(The size of this address depends on amount of cacheable physical main memory)



Block offset size =  $\log_2(\text{block size})$

Index size =  $\log_2(\text{Total number of blocks/associativity})$

Tag size = address size - index size - offset size

Mapping function: (From memory block to cache)

Cache set or block frame number = Index =  
= (Block Address) MOD (Number of Sets)

Number of Sets  
in cache

Fully associative cache has no index field or mapping function  
e.g. no index field

**EECC550 - Shaaban**

# Calculating Number of Cache Bits Needed



Cache Block Frame (or just cache block)

Address Fields

- How many total bits are needed for a direct-mapped cache with 64 KBytes of data and one word blocks, assuming a 32-bit address?

- 64 Kbytes = 16 K words =  $2^{14}$  words =  $2^{14}$  blocks
- Block size = 4 bytes  $\Rightarrow$  offset size =  $\log_2(4) = 2$  bits,
- #sets = #blocks =  $2^{14}$   $\Rightarrow$  index size = 14 bits
- Tag size = address size - index size - offset size =  $32 - 14 - 2 = 16$  bits
- Bits/block = data bits + tag bits + valid bit =  $32 + 16 + 1 = 49$
- Bits in cache = #blocks x bits/block =  $2^{14} \times 49 = 98$  Kbytes

i.e nominal cache Capacity = 64 KB

Number of cache block frames

Actual number of bits in a cache block frame

- How many total bits would be needed for a 4-way set associative cache to store the same amount of data?

- Block size and #blocks does not change.
- #sets = #blocks/4 =  $(2^{14})/4 = 2^{12}$   $\Rightarrow$  index size = 12 bits
- Tag size = address size - index size - offset =  $32 - 12 - 2 = 18$  bits
- Bits/block = data bits + tag bits + valid bit =  $32 + 18 + 1 = 51$
- Bits in cache = #blocks x bits/block =  $2^{14} \times 51 = 102$  Kbytes

- Increase associativity  $\Rightarrow$  increase bits in cache

Word = 4 bytes

More bits in tag

1 k = 1024 =  $2^{10}$

EECC550 - Shaaban

# Calculating Cache Bits Needed



Address Fields



Cache Block Frame (or just cache block)

- How many total bits are needed for a direct- mapped cache with 64 KBytes of data and 8 word (32 byte) blocks, assuming a 32-bit address (it can cache  $2^{32}$  bytes in memory)?
  - 64 Kbytes =  $2^{14}$  words =  $(2^{14})/8 = 2^{11}$  blocks Number of cache block frames
  - block size = 32 bytes  
 => offset size = block offset + byte offset =  $\log_2(32) = 5$  bits,
  - #sets = #blocks =  $2^{11}$  => index size = 11 bits
  - tag size = address size - index size - offset size =  $32 - 11 - 5 = 16$  bits
  - bits/block = data bits + tag bits + valid bit =  $8 \times 32 + 16 + 1 = 273$  bits
  - bits in cache = #blocks x bits/block =  $2^{11} \times 273 = 68.25$  Kbytes Actual number of bits in a cache block frame
- Increase block size => decrease bits in cache. Fewer cache block frames thus fewer tags/valid bits

Word = 4 bytes

1 k = 1024 =  $2^{10}$

**EECC550 - Shaaban**

# Unified vs. Separate Level 1 Cache

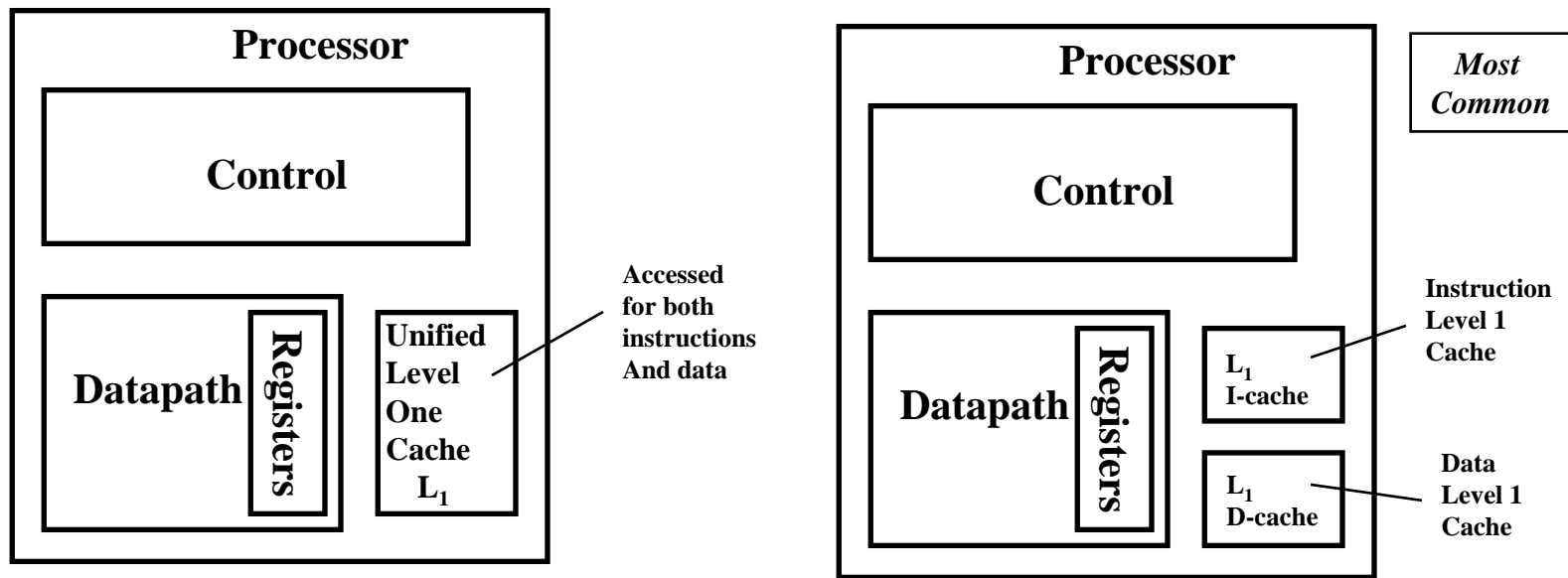
- Unified Level 1 Cache (Princeton Memory Architecture). AKA Shared Cache

A single level 1 ( $L_1$ ) cache is used for both instructions and data.

Or Split

- Separate instruction/data Level 1 caches (Harvard Memory Architecture):

The level 1 ( $L_1$ ) cache is split into two caches, one for instructions (instruction cache,  $L_1$  I-cache) and the other for data (data cache,  $L_1$  D-cache).



AKA shared

Unified Level 1 Cache  
(Princeton Memory Architecture)

Separate (Split) Level 1 Caches  
(Harvard Memory Architecture)

Split Level 1 Cache is more preferred in pipelined CPUs to avoid instruction fetch/Data access structural hazards

**EECC550 - Shaaban**

## ***Memory Hierarchy/Cache Performance:***

### **Average Memory Access Time (AMAT), Memory Stall cycles**

- **The Average Memory Access Time (AMAT):** The number of cycles required to complete an average memory access request by the CPU.
- **Memory stall cycles per memory access:** The number of stall cycles added to CPU execution cycles for one memory access.

- **Memory stall cycles per average memory access = (AMAT -1)**

- **For ideal memory: AMAT = 1 cycle, this results in zero memory stall cycles.**

- **Memory stall cycles per average instruction =**

**Number of memory accesses per instruction**

Instruction Fetch  $\rightarrow$   $\text{Number of memory accesses per instruction} \times \text{Memory stall cycles per average memory access}$   
 $= (1 + \text{fraction of loads/stores}) \times (\text{AMAT} - 1)$

$$\text{Base CPI} = \text{CPI}_{\text{execution}} = \text{CPI with ideal memory}$$

$$\text{CPI} = \text{CPI}_{\text{execution}} + \text{Mem Stall cycles per instruction}$$

cycles = CPU cycles



# Cache Performance: Single Level L1 Princeton (Unified) Memory Architecture

$CPU_{time} = \text{Instruction count} \times CPI \times \text{Clock cycle time}$

$CPI_{execution} = \text{CPI with ideal memory}$

$$CPI = CPI_{execution} + \text{Mem Stall cycles per instruction}$$

Mem Stall cycles per instruction =

Memory accesses per instruction  $\times$  Memory stall cycles per access

i.e No hit penalty

Assuming no stall cycles on a cache hit (cache access time = 1 cycle, stall = 0)

Cache Hit Rate =  $H1$       Miss Rate =  $1 - H1$       Miss Penalty =  $M$

Memory stall cycles per memory access = Miss rate  $\times$  Miss penalty =  $(1 - H1) \times M$

AMAT =  $1 + \text{Miss rate} \times \text{Miss penalty}$

Memory accesses per instruction =  $(1 + \text{fraction of loads/stores})$

Miss Penalty =  $M = \text{the number of stall cycles resulting from missing in cache}$   
 $= \text{Main memory access time} - 1$

Thus for a unified L1 cache with no stalls on a cache hit:

$$CPI = CPI_{execution} + (1 + \text{fraction of loads/stores}) \times (1 - H1) \times M$$

$$AMAT = 1 + (1 - H1) \times M$$

$$CPI = CPI_{execution} + (1 + \text{fraction of loads and stores}) \times \text{stall cycles per access}$$

$$= CPI_{execution} + (1 + \text{fraction of loads and stores}) \times (AMAT - 1)$$

# Memory Access Tree: For Unified Level 1 Cache

Probability to be here

## CPU Memory Access

Unified

L<sub>1</sub>

H1

100%  
or 1

(1-H1)

### L1 Hit:

% = Hit Rate = H1

Hit Access Time = 1

Stall cycles per access = 0

Stall = H1 x 0 = 0

( No Stall)

### L1 Miss:

% = (1- Hit rate) = (1-H1)

Access time = M + 1

Stall cycles per access = M

Stall = M x (1-H1)

Assuming:

Ideal access on a hit

$$AMAT = \overset{\text{Hit Rate}}{H1} \times \overset{\text{Hit Time}}{1} + \overset{\text{Miss Rate}}{(1-H1)} \times \overset{\text{Miss Time}}{(M+1)} = 1 + M \times (1-H1)$$

Stall Cycles Per Access = AMAT - 1 = M x (1 -H1)

CPI = CPI<sub>execution</sub> + (1 + fraction of loads/stores) x M x (1 -H1)

M = Miss Penalty = stall cycles per access resulting from missing in cache

M + 1 = Miss Time = Main memory access time

H1 = Level 1 Hit Rate

1- H1 = Level 1 Miss Rate

AMAT = 1 + Stalls per average memory access

**EECC550 - Shaaban**

# Cache Performance Example

- Suppose a CPU executes at Clock Rate = 200 MHz (5 ns per cycle) with a single level of cache.
- $CPI_{\text{execution}} = 1.1$  (i.e base CPI with ideal memory)
- Instruction mix: 50% arith/logic, 30% load/store, 20% control
- Assume a cache miss rate of 1.5% and a miss penalty of  $M = 50$  cycles.

$$CPI = CPI_{\text{execution}} + \text{mem stalls per instruction}$$

$$\text{Mem Stalls per instruction} = \text{Mem accesses per instruction} \times \text{Miss rate} \times \text{Miss penalty}$$

$$\text{Mem accesses per instruction} = 1 + .3 = 1.3$$

Instruction fetch ↗
← Load/store

$$\text{Mem Stalls per memory access} = (1 - H1) \times M = .015 \times 50 = .75 \text{ cycles}$$

$$AMAT = 1 + .75 = 1.75 \text{ cycles}$$

$$\text{Mem Stalls per instruction} = 1.3 \times .015 \times 50 = 0.975$$

$$CPI = 1.1 + .975 = 2.075$$

The ideal memory CPU with no misses is  $2.075/1.1 = 1.88$  times faster

**EECC550 - Shaaban**

$M$  = Miss Penalty = stall cycles per access resulting from missing in cache

# Cache Performance Example

- Suppose for the previous example we double the clock rate to 400 MHz, how much faster is this machine, assuming similar miss rate, instruction mix?
- Since memory speed is not changed, the miss penalty takes more CPU cycles:

$$\text{Miss penalty} = M = 50 \times 2 = 100 \text{ cycles.}$$

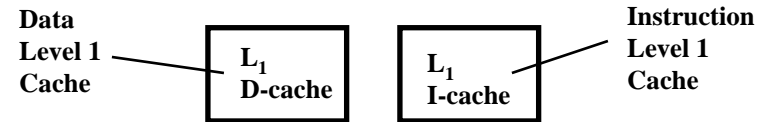
$$\text{CPI} = 1.1 + 1.3 \times .015 \times 100 = 1.1 + 1.95 = 3.05$$

$$\begin{aligned} \text{Speedup} &= (\text{CPI}_{\text{old}} \times C_{\text{old}}) / (\text{CPI}_{\text{new}} \times C_{\text{new}}) \\ &= 2.075 \times 2 / 3.05 = 1.36 \end{aligned}$$

The new machine is only 1.36 times faster rather than 2 times faster due to the increased effect of cache misses.

→ *CPUs with higher clock rate, have more cycles per cache miss and more memory impact on CPI.*

# Cache Performance:



Usually: Data Miss Rate >> Instruction Miss Rate

$$\text{Miss rate} = 1 - \text{data H1}$$

$$\text{Miss rate} = 1 - \text{instruction H1}$$

## Single Level L1 Harvard (Split) Memory Architecture

For a CPU with separate or split level one (L1) caches for instructions and data (Harvard memory architecture) and no stalls for cache hits:

$$\text{CPUtime} = \text{Instruction count} \times \text{CPI} \times \text{Clock cycle time}$$

$$\text{CPI} = \text{CPI}_{\text{execution}} + \text{Mem Stall cycles per instruction}$$

Mem Stall cycles per instruction =

This is one method to find stalls per instruction another method is shown in next slide →

$$\text{Instruction Fetch Miss rate} \times M + \text{Data Memory Accesses Per Instruction} \times \text{Data Miss Rate} \times M$$

1- Instruction H1

Fraction of Loads and Stores

1- Data H1

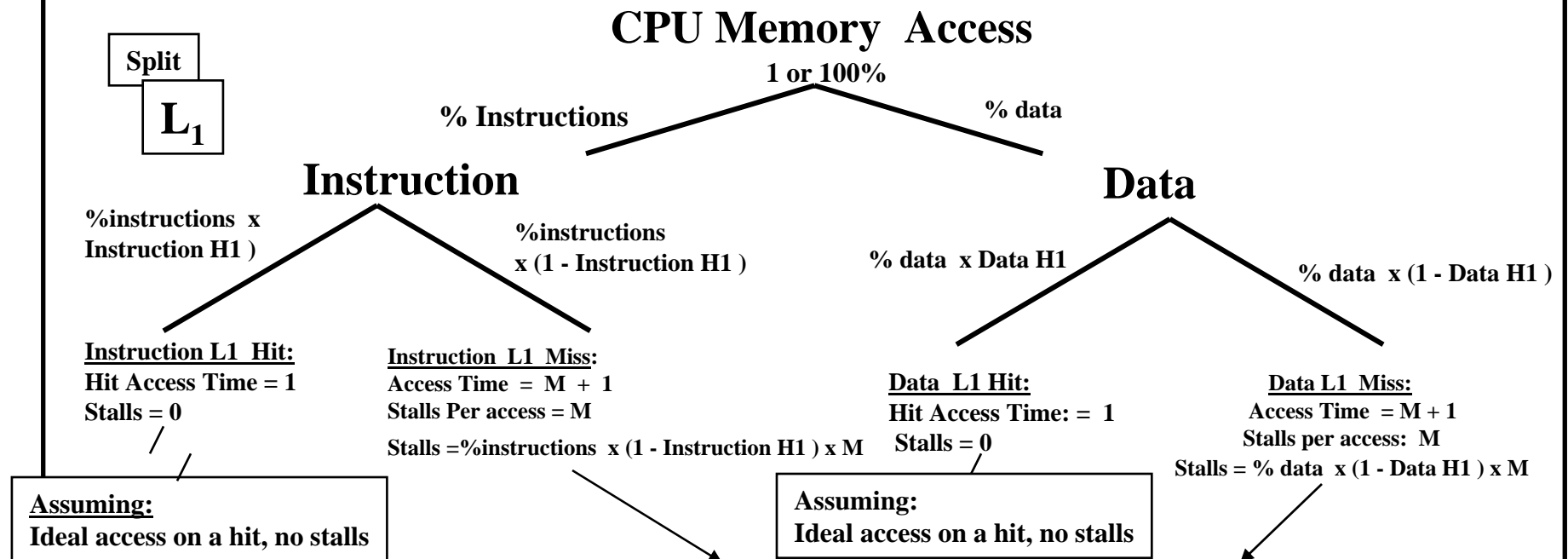
**M = Miss Penalty = stall cycles per access to main memory resulting from missing in cache**

$$\text{CPI}_{\text{execution}} = \text{base CPI with ideal memory}$$

**EECC550 - Shaaban**

# Memory Access Tree

## For Separate Level 1 Caches



$$\text{Stall Cycles Per Access} = \% \text{ Instructions} \times ( 1 - \text{Instruction H1} ) \times M + \% \text{ data} \times ( 1 - \text{Data H1} ) \times M$$

$$\text{AMAT} = 1 + \text{Stall Cycles per access}$$

$$\text{Stall cycles per instruction} = ( 1 + \text{fraction of loads/stores} ) \times \text{Stall Cycles per access}$$

$$\text{CPI} = \text{CPI}_{\text{execution}} + \text{Stall cycles per instruction}$$

$$= \text{CPI}_{\text{execution}} + ( 1 + \text{fraction of loads/stores} ) \times \text{Stall Cycles per access}$$

M = Miss Penalty = stall cycles per access resulting from missing in cache

M + 1 = Miss Time = Main memory access time

Data H1 = Level 1 Data Hit Rate

1 - Data H1 = Level 1 Data Miss Rate

Instruction H1 = Level 1 Instruction Hit Rate

1 - Instruction H1 = Level 1 Instruction Miss Rate

% Instructions = Percentage or fraction of instruction fetches out of all memory accesses

% Data = Percentage or fraction of data accesses out of all memory accesses

**EECC550 - Shaaban**

# Split L1 Cache Performance Example

- Suppose a CPU uses separate level one (L1) caches for instructions and data (Harvard memory architecture) with different miss rates for instruction and data access:
  - A cache hit incurs no stall cycles while a cache miss incurs 200 stall cycles for both memory reads and writes.
  - $CPI_{\text{execution}} = 1.1$  (i.e base CPI with ideal memory)
  - Instruction mix: 50% arith/logic, 30% load/store, 20% control
  - Assume a cache miss rate of 0.5% for instruction fetch and a cache data miss rate of 6%.
  - A cache hit incurs no stall cycles while a cache miss incurs 200 stall cycles for both memory reads and writes.
- Find the resulting stalls per access, AMAT and CPI using this cache? M

$$CPI = CPI_{\text{execution}} + \text{mem stalls per instruction}$$

$$\text{Memory Stall cycles per instruction} = \text{Instruction Fetch Miss rate} \times \text{Miss Penalty} + \text{Data Memory Accesses Per Instruction} \times \text{Data Miss Rate} \times \text{Miss Penalty}$$

$$\text{Memory Stall cycles per instruction} = 0.5/100 \times 200 + 0.3 \times 6/100 \times 200 = 1 + 3.6 = 4.6 \text{ cycles}$$

$$\text{Stall cycles per average memory access} = 4.6/1.3 = 3.54 \text{ cycles}$$

$$\text{AMAT} = 1 + \text{Stall cycles per average memory access} = 1 + 3.54 = 4.54 \text{ cycles}$$

$$CPI = CPI_{\text{execution}} + \text{mem stalls per instruction} = 1.1 + 4.6 = 5.7 \text{ cycles}$$

- What is the miss rate of a single level unified cache that has the same performance?

$$4.6 = 1.3 \times \text{Miss rate} \times 200 \quad \text{which gives a miss rate of } 1.8 \% \text{ for an equivalent unified cache}$$

- How much faster is the CPU with ideal memory?

The CPU with ideal cache (no misses) is  $5.7/1.1 = 5.18$  times faster

With no cache at all the CPI would have been  $= 1.1 + 1.3 \times 200 = 261.1$  cycles !!

# Memory Access Tree For Separate Level 1 Caches Example

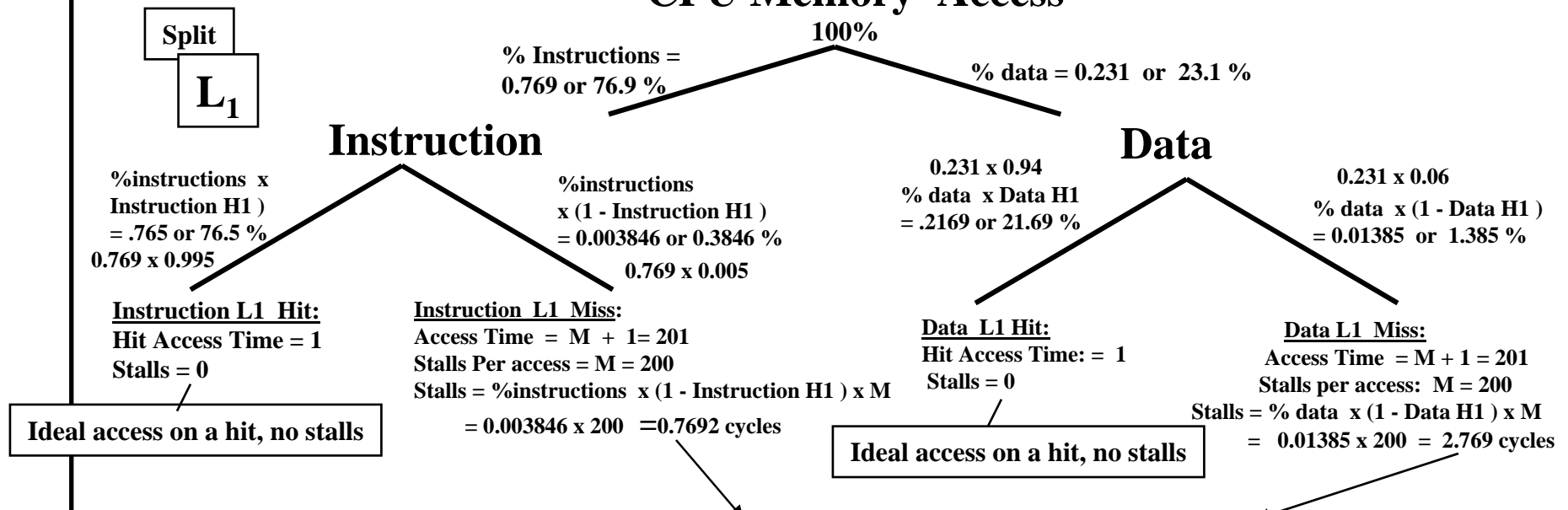
For Last Example

30% of all instructions executed are loads/stores, thus:

Fraction of instruction fetches out of all memory accesses =  $1 / (1+0.3) = 1/1.3 = 0.769$  or 76.9 %

Fraction of data accesses out of all memory accesses =  $0.3 / (1+0.3) = 0.3/1.3 = 0.231$  or 23.1 %

## CPU Memory Access



$$\text{Stall Cycles Per Access} = \% \text{ Instructions} \times (1 - \text{Instruction H1}) \times M + \% \text{ data} \times (1 - \text{Data H1}) \times M$$

$$= 0.7692 + 2.769 = 3.54 \text{ cycles}$$

$$\text{AMAT} = 1 + \text{Stall Cycles per access} = 1 + 3.5 = 4.54 \text{ cycles}$$

$$\text{Stall cycles per instruction} = (1 + \text{fraction of loads/stores}) \times \text{Stall Cycles per access} = 1.3 \times 3.54 = 4.6 \text{ cycles}$$

$$\text{CPI} = \text{CPI}_{\text{execution}} + \text{Stall cycles per instruction} = 1.1 + 4.6 = 5.7$$

Given as 1.1

M = Miss Penalty = stall cycles per access resulting from missing in cache = 200 cycles  
 M + 1 = Miss Time = Main memory access time = 200+1 = 201 cycles    L1 access Time = 1 cycle  
 Data H1 = 0.94 or 94%    1- Data H1 = 0.06 or 6%  
 Instruction H1 = 0.995 or 99.5%    1- Instruction H1 = 0.005 or 0.5 %  
 % Instructions = Percentage or fraction of instruction fetches out of all memory accesses = 76.9 %  
 % Data = Percentage or fraction of data accesses out of all memory accesses = 23.1 %

**EECC550 - Shaaban**