

CMCS 611-101

Advanced Computer Architecture

Lecture 6

Introduction to Pipelining

September 23, 2009

www.csee.umbc.edu/~younis/CMSC611/CMSC611.htm



Lecture's Overview

□ Previous Lecture:

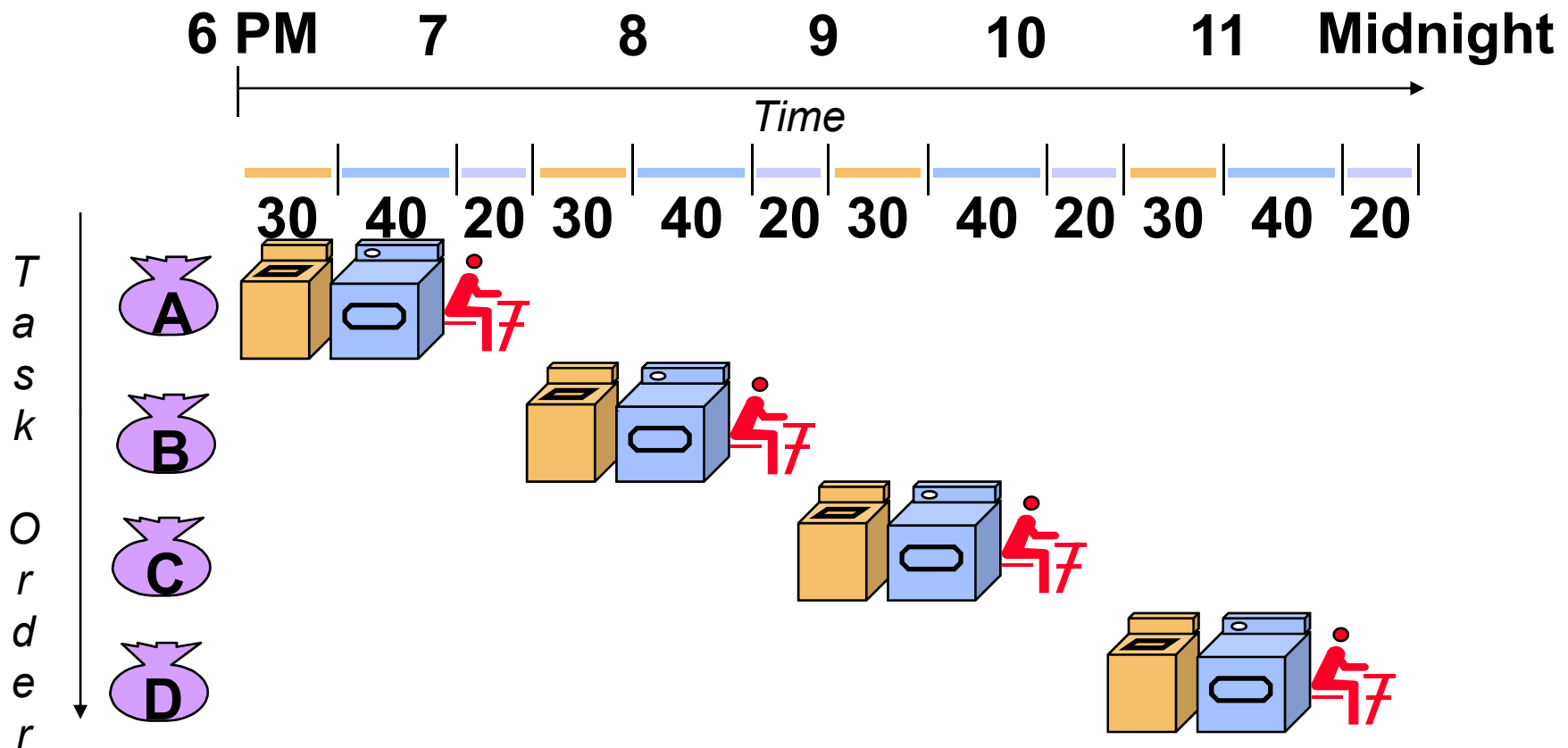
- Type and size of operands
(Famous data types, effect of operand size on design complexity)
- Encoding the instruction set
(Fixed, variable and hybrid encoding, the store program concept)
- The role of the compiler
(Compilation process, compiler optimization, linking and loading)
- Effect of ISA on Compiler Complexity
(Regularity, Primitives, not solutions, Simplify trade-offs, Static binding)

□ This Lecture:

- An overview of pipelining
- Pipeline performance
- Pipelined hazards



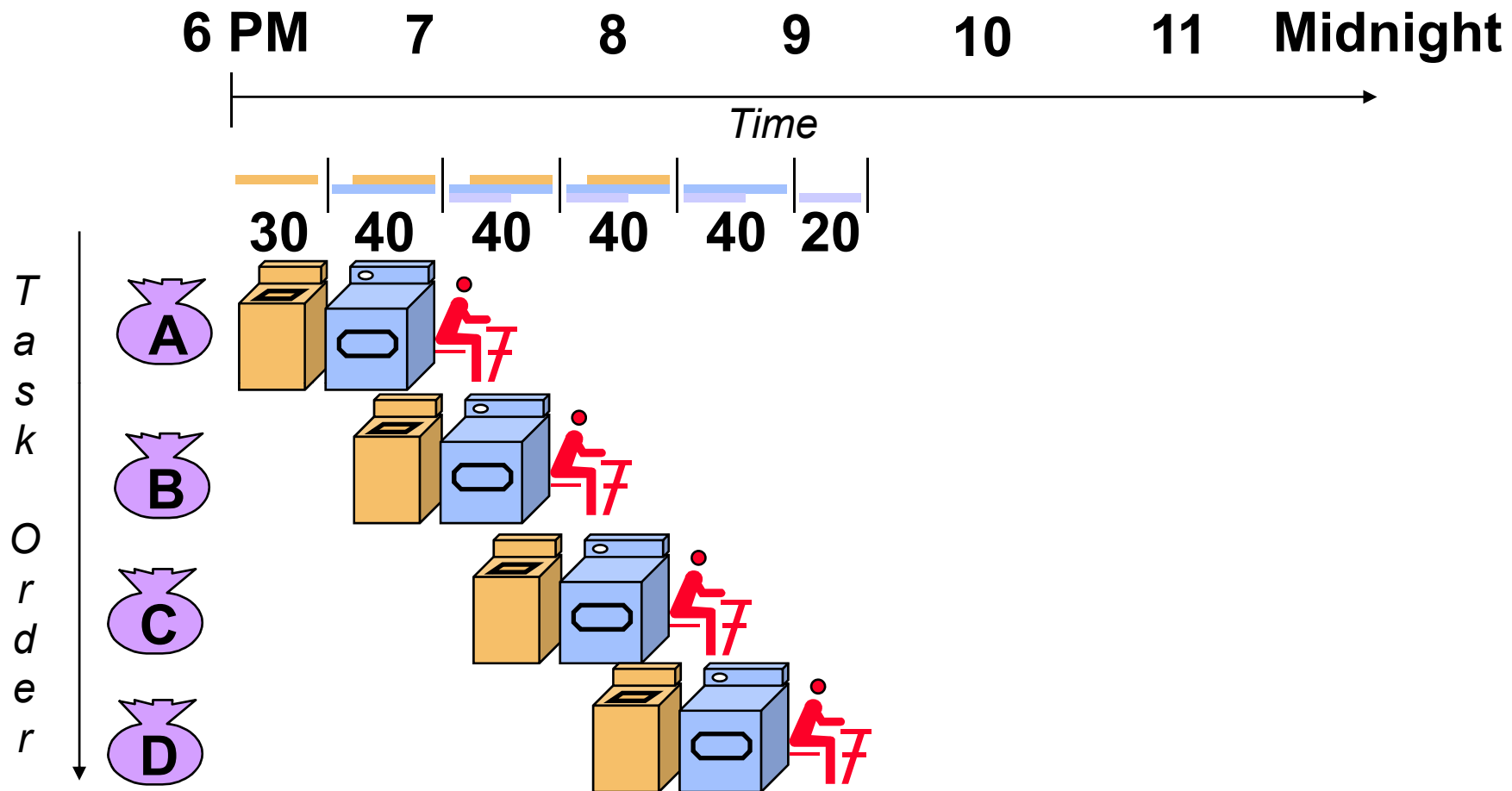
Sequential Laundry



- Washer takes 30 min, Dryer takes 40 min, folding takes 20 min
- Sequential laundry takes 6 hours for 4 loads
- If they learned pipelining, how long would laundry take?



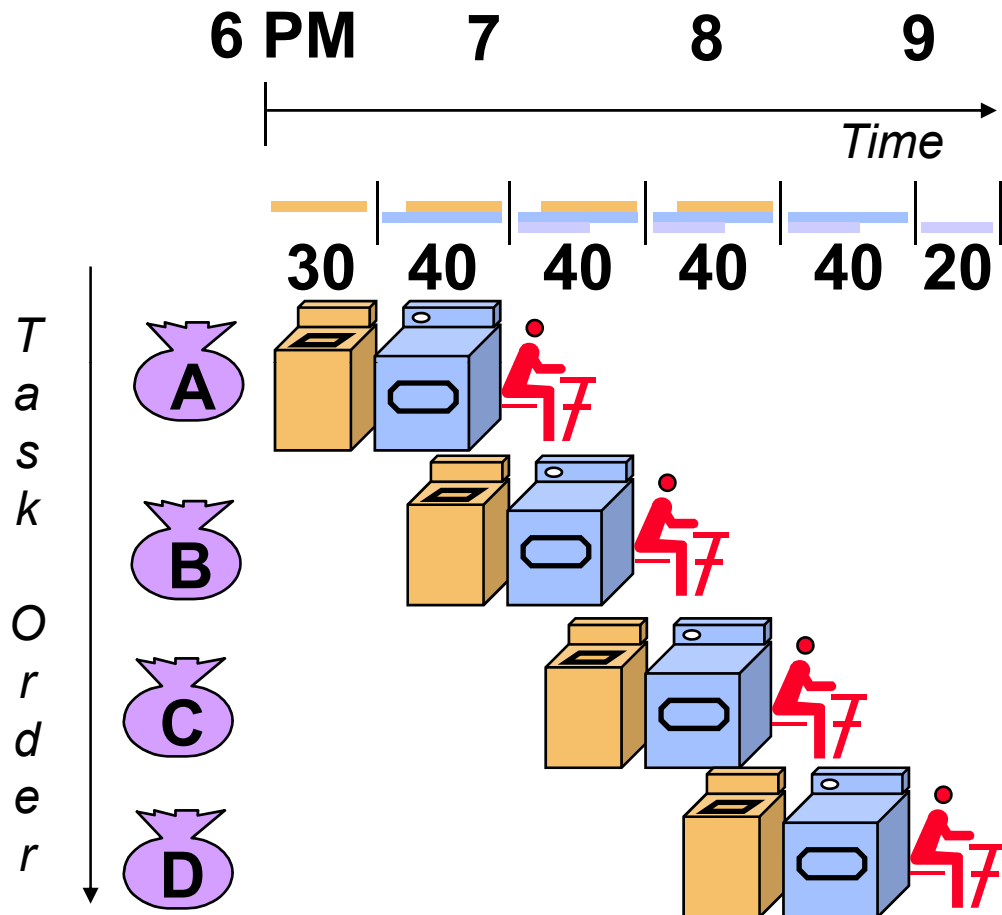
Pipelined Laundry



- ❑ Pipelining means start work as soon as possible
- ❑ Pipelined laundry takes 3.5 hours for 4 loads



Pipelining Lessons



- Pipelining doesn't help **latency** of single task, it helps **throughput** of entire workload
- Pipeline rate limited by **slowest** pipeline stage
- Multiple** tasks operating simultaneously using different resources
- Potential speedup = **Number pipe stages**
- Unbalanced lengths of pipe stages reduces speedup
- Time to “**fill**” pipeline and time to “**drain**” it reduce speedup
- Stall for Dependencies



Basics of a RISC Instruction Set

- ❑ RISC architectures are characterized by the following features that dramatically simplifies the implementation:
 1. All ALU operations apply only on data in registers
 2. Memory is affected only by load and store operations
 3. Instructions follow very few formats and typically are of the same size

❑ All MIPS instructions are 32 bits, following one of three formats:



MIPS Instruction format

① Register-format instructions:

| op | rs | rt | rd | shamt | funct |
|-----------|-----------|-----------|-----------|--------------|--------------|
| 6 bits | 5 bits | 5 bits | 5 bits | 5 bits | 6 bits |

- op*: Basic operation of the instruction, traditionally called opcode
rs: The first register source operand
rt: The second register source operand
rd: The register destination operand, it gets the result of the operation
shamt: Shift amount (explained in future lectures)
funct: This field selects the specific variant of the operation of the op field

- ❑ MIPS assembly language includes two conditional branching instructions using PC -relative addressing:

`beq register1, register2, L1 # go to L1 if (register1) = (register2)`

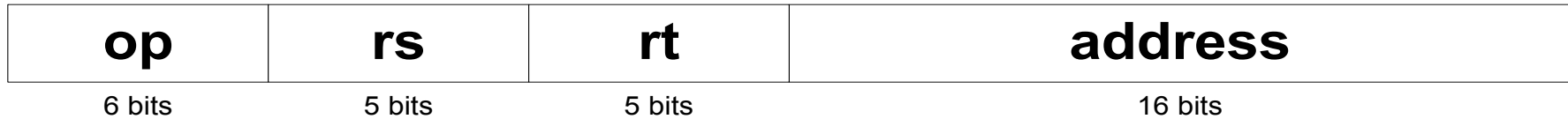
`bne register1, register2, L1 # go to L1 if (register1) ≠ (register2)`

- ❑ Examples:
- | | | |
|------------------|-------------------------------|--|
| <code>add</code> | <code>\$t2, \$t1, \$t1</code> | <code># Temp reg \$t2 = 2 \$t1</code> |
| <code>sub</code> | <code>\$t1, \$s3, \$s4</code> | <code># Temp reg \$t1 = \$s3 - \$s4</code> |
| <code>and</code> | <code>\$t1, \$t2, \$t3</code> | <code># Temp reg \$t1 = \$t2 . \$t3</code> |
| <code>bne</code> | <code>\$s3, \$s4, Else</code> | <code># if \$s3 ≠ \$s4 jump to Else</code> |

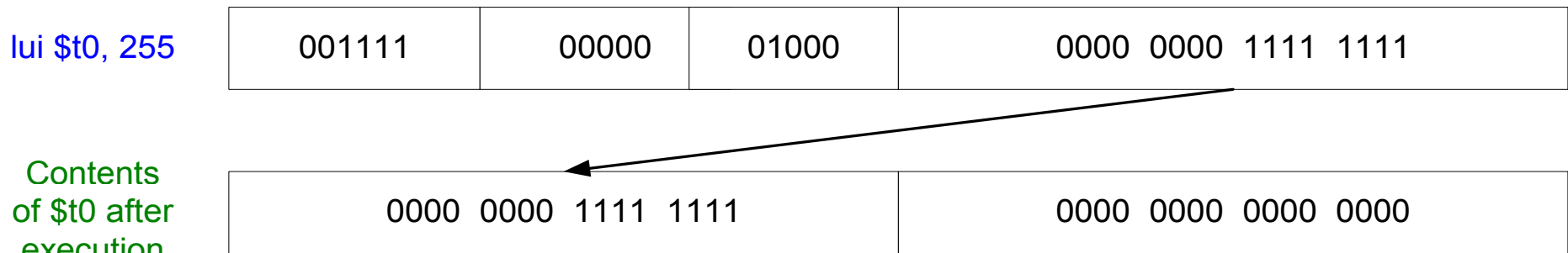


MIPS Instruction format

② Immediate-type instructions:

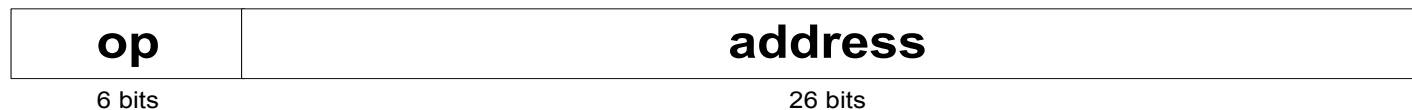


- ❑ The 16-bit address means a load word instruction can load a word within a region of $\pm 2^{15}$ bytes of the address in the base register
- ❑ Examples: `lw $t0, 32($s3)` , `sw $t1, 128($s3)`
- ❑ MIPS handle 16-bit constant efficiently by including the constant value in the address field of an I-type instruction (Immediate-type)
 - `addi $sp, $sp, 4` $\#\$sp = \$sp + 4$
- ❑ For large constants that need more than 16 bits, a load upper-immediate (*lui*) instruction is used to concatenate the second part



Addressing in Branches & Jumps

- ❑ I-type instructions leaves only 16 bits for address reference limiting the size of the jump
- ❑ MIPS branch instructions use the address as an increment to the PC allowing the program to be as large as 2^{32} (called *PC-relative addressing*)
- ❑ Since the program counter gets incremented prior to instruction execution, the branch address is actually relative to (PC + 4)
- ❑ MIPS also supports an J-type instruction format for large jump instructions

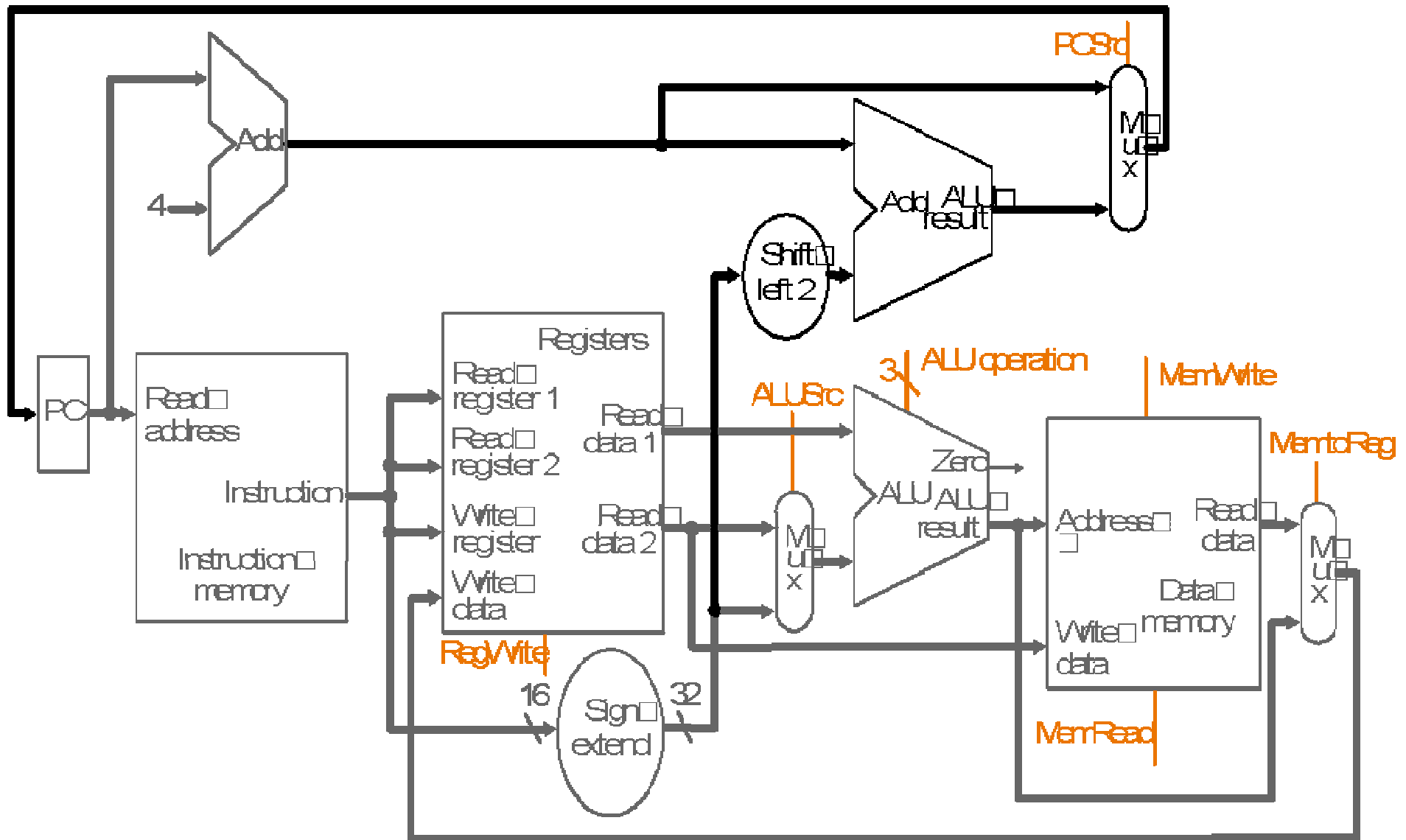


- ❑ The 26-bit address in a J-type instruct. is concatenated to upper 8 bits of PC

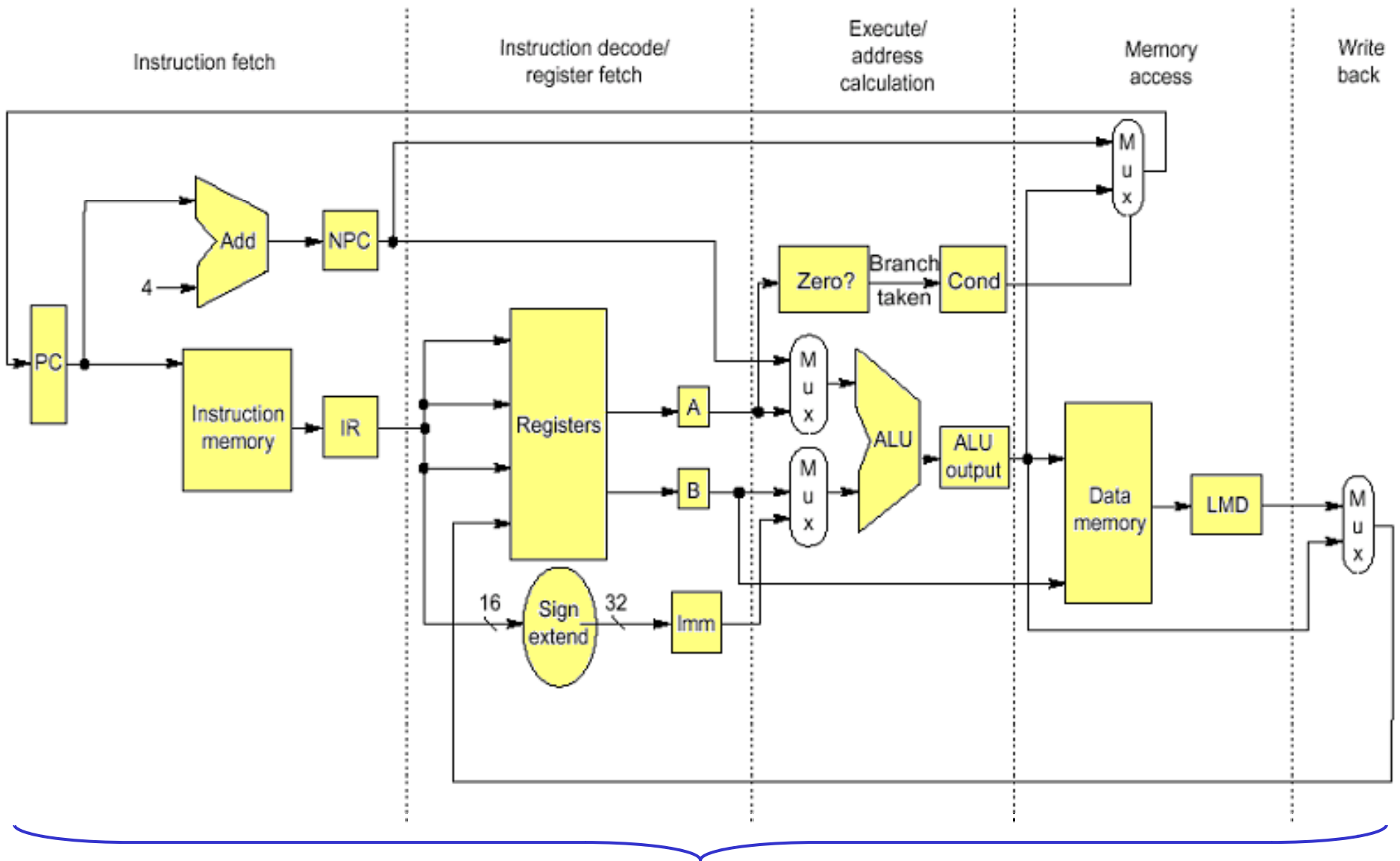
| | | | | | | | | |
|--------------|-----------------------------|-------|-----|-------|----|----|---|----|
| Loop: | add \$t1, \$s3, \$s3 | 80000 | 0 | 19 | 19 | 9 | 0 | 32 |
| | add \$t1, \$t1, \$t1 | 80004 | 0 | 9 | 9 | 9 | 0 | 32 |
| | add \$t1, \$t1, \$s6 | 80008 | 0 | 9 | 22 | 9 | 0 | 32 |
| | lw \$t0, 0(\$t1) | 80012 | 35 | 9 | 8 | 0 | | |
| | bne \$t0, \$s5, Exit | 80016 | 5 | 8 | 21 | 8 | | |
| | add \$s3, \$s3, \$s4 | 80020 | 0 | 19 | 20 | 19 | 0 | 32 |
| | j Loop | 80024 | 2 | 80000 | | | | |
| | | 80028 | ... | | | | | |
| Exit: | | 80012 | 35 | 9 | 8 | 0 | | |



A Simple Implementation of MIPS



Single-cycle Instruction Execution



1



Multi-Cycle Implementation of MIPS

① Instruction fetch cycle (IF)

$IR \leftarrow \text{Mem}[PC]; \quad NPC \leftarrow PC + 4$

② Instruction decode/register fetch cycle (ID)

$A \leftarrow \text{Regs}[IR_{6..10}]; \quad B \leftarrow \text{Regs}[IR_{11..15}]; \quad \text{Imm} \leftarrow ((IR_{16})^{16} \# \# IR_{16..31})$

③ Execution/effective address cycle (EX)

Memory ref: $\text{ALUOutput} \leftarrow A + \text{Imm};$

Reg-Reg ALU: $\text{ALUOutput} \leftarrow A \text{ func } B;$

Reg-Imm ALU: $\text{ALUOutput} \leftarrow A \text{ op } \text{Imm};$

Branch: $\text{ALUOutput} \leftarrow NPC + \text{Imm}; \quad \text{Cond} \leftarrow (A \text{ op } 0)$

④ Memory access/branch completion cycle (MEM)

Memory ref: $\text{LMD} \leftarrow \text{Mem}[\text{ALUOutput}] \quad \text{or} \quad \text{Mem}[\text{ALUOutput}] \leftarrow B;$

Branch: $\text{if (cond)} \text{ PC} \leftarrow \text{ALUOutput};$

⑤ Write-back cycle (WB)

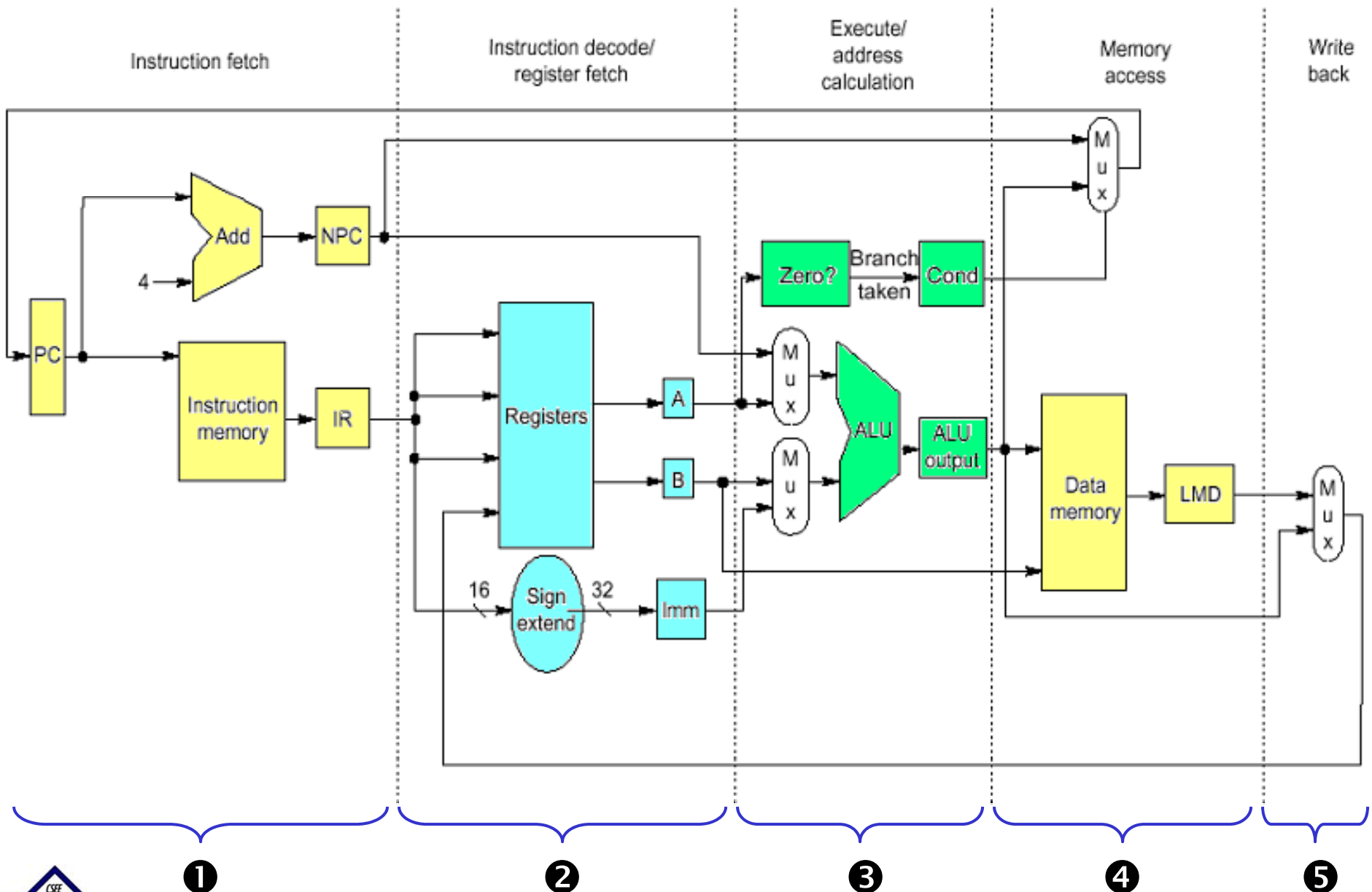
Reg-Reg ALU: $\text{Regs}[IR_{16..20}] \leftarrow \text{ALUOutput};$

Reg-Imm ALU: $\text{Regs}[IR_{11..15}] \leftarrow \text{ALUOutput};$

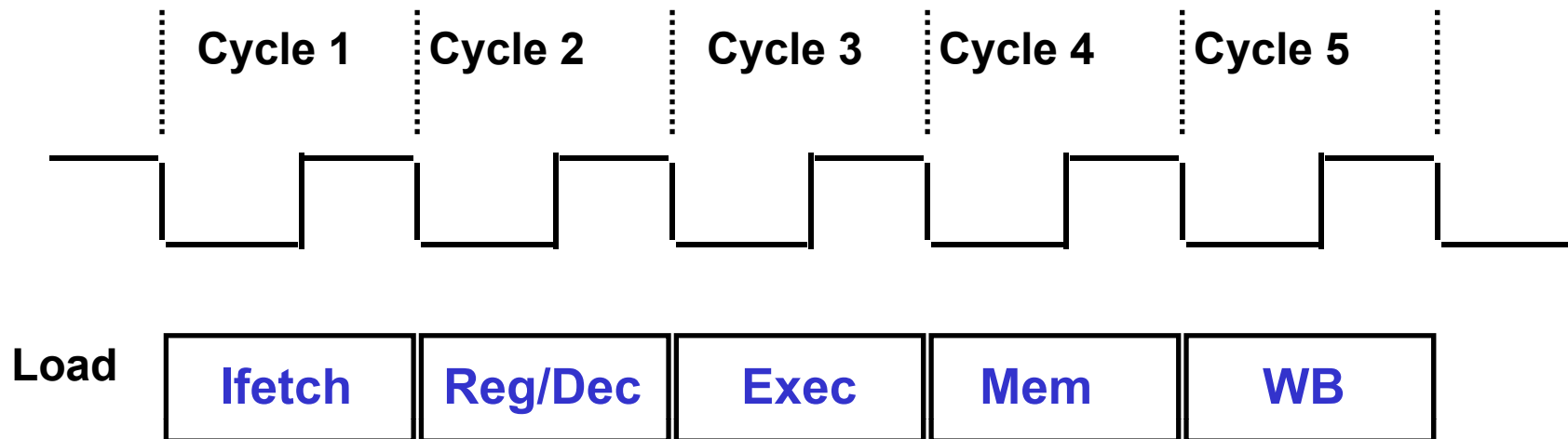
Load: $\text{Regs}[IR_{11..15}] \leftarrow \text{LMD};$



Multi-cycle Instruction Execution



Stages of Instruction Execution

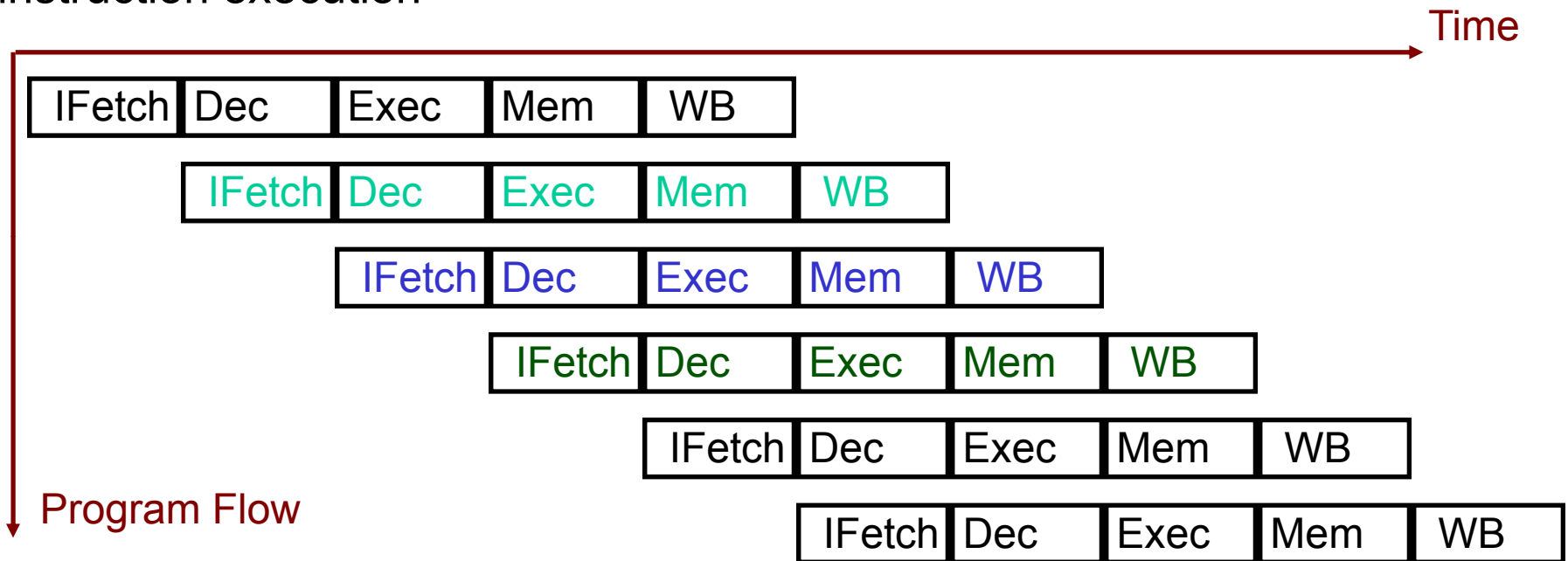


- ❑ The load instruction is the longest
- ❑ All instructions follows at most the following five steps:
 - ➔ **ifetch:** Instruction Fetch
 - Fetch the instruction from the Instruction Memory and update PC
 - ➔ **Reg/Dec:** Registers Fetch and Instruction Decode
 - ➔ **Exec:** Calculate the memory address
 - ➔ **Mem:** Read the data from the Data Memory
 - ➔ **WB:** Write the data back to the register file



Instruction Pipelining

- ❑ Start handling of next instruction while the current instruction is in progress
- ❑ Pipelining is feasible when different devices are used at different stages of instruction execution

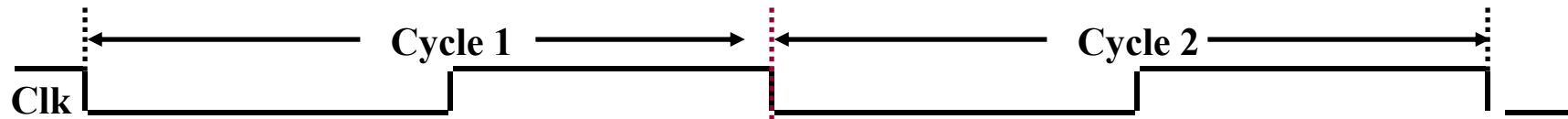


$$\text{Time between instructions}_{\text{pipelined}} = \frac{\text{Time between instructions}_{\text{nonpipelined}}}{\text{Number of pipe stages}}$$

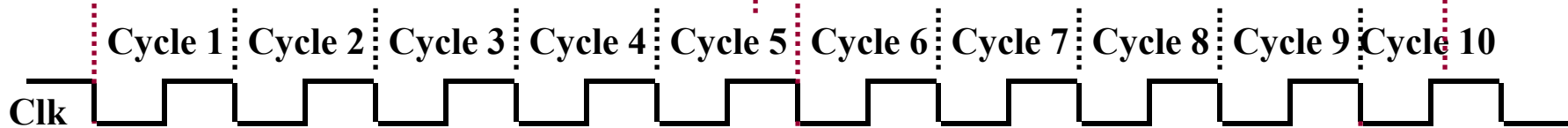
Pipelining improves performance by increasing instruction throughput



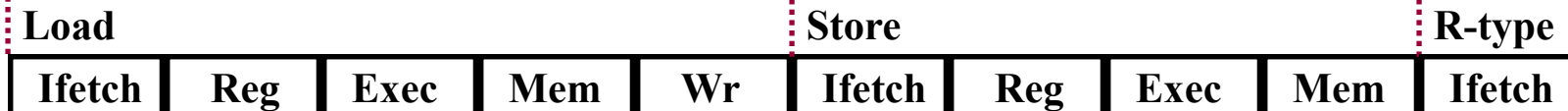
Single Cycle, Multiple Cycle, vs. Pipeline



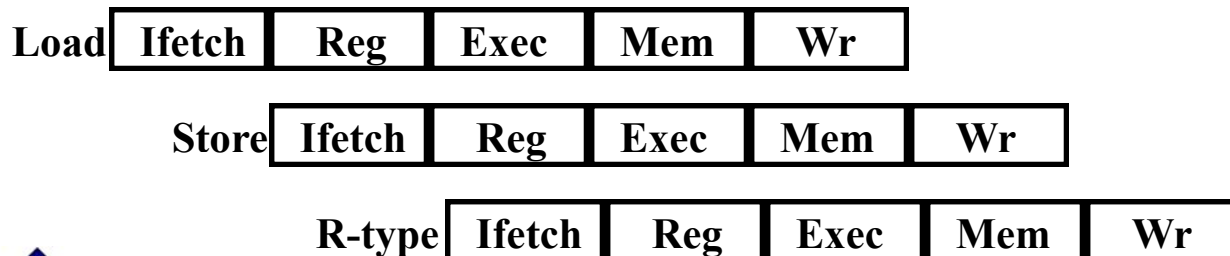
Single Cycle Implementation:



Multiple Cycle Implementation:



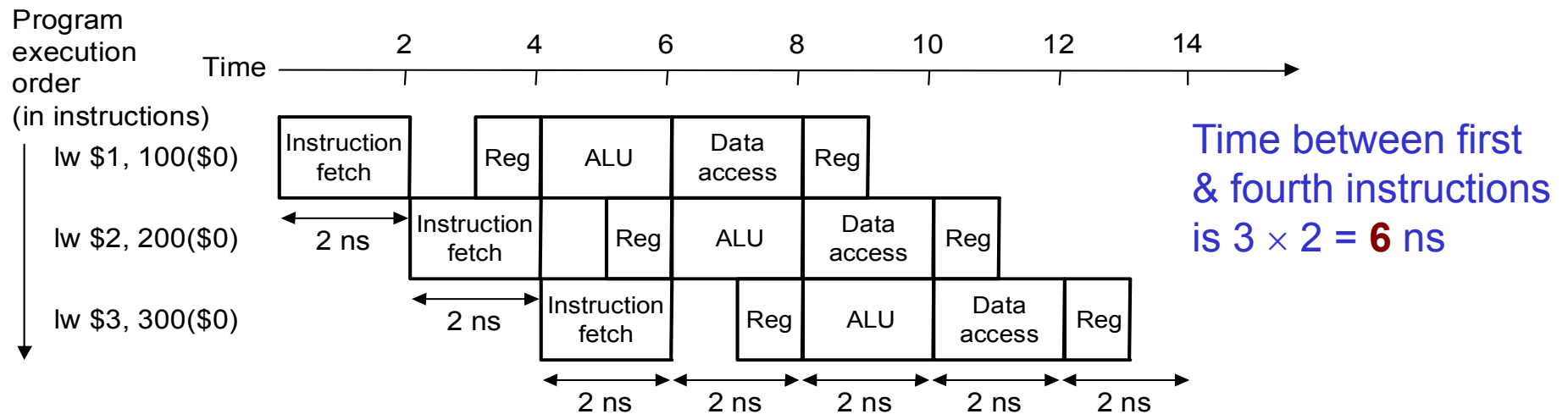
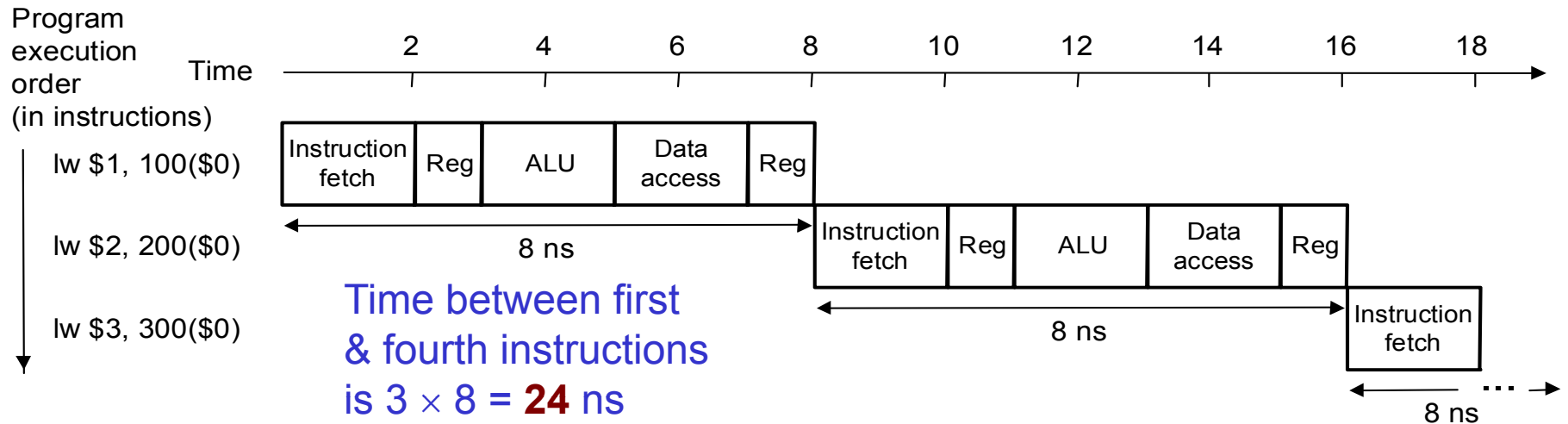
Pipeline Implementation:



* Slide is courtesy of Dave Patterson



Example of Instruction Pipelining



Ideal and upper bound for speedup is number of stages in the pipeline



Pipeline Performance

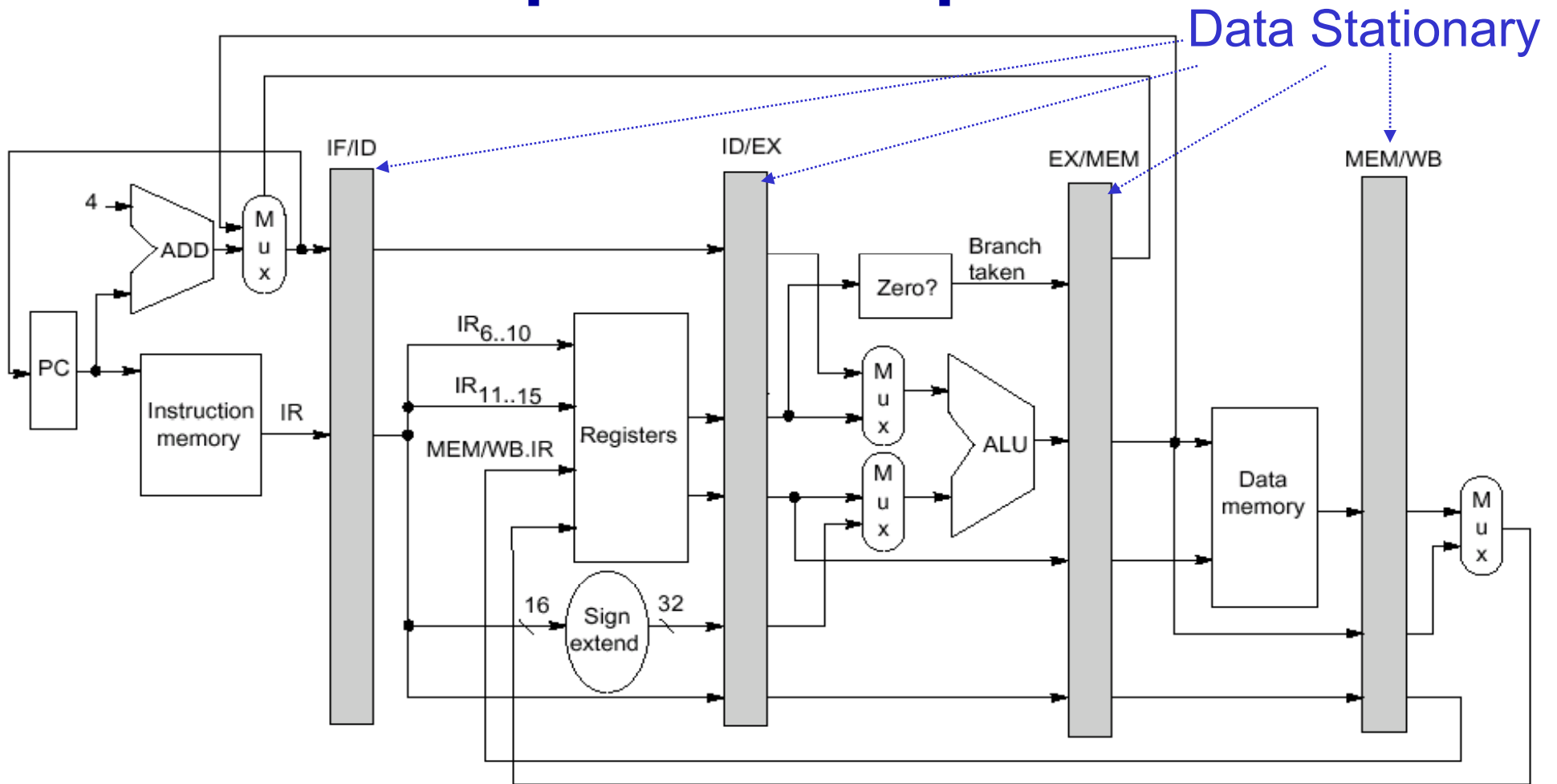
- ❑ Pipeline increases the instruction throughput but does not reduce the execution time of the individual instruction
- ❑ Execution time of the individual instruction in pipeline can be slower due:
 - ➔ Additional pipeline control compared to none pipeline execution
 - ➔ Imbalance among the different pipeline stages
- ❑ Suppose we execute 100 instructions:
 - ➔ **Single Cycle Machine**
 - 45 ns/cycle x 1 CPI x 100 inst = 4500 ns
 - ➔ **Multi-cycle Machine**
 - 10 ns/cycle x 4.2 CPI (due to inst mix) x 100 inst = 4200 ns
 - ➔ **Ideal 5 stages pipelined machine**
 - 10 ns/cycle x (1 CPI x 100 inst + 4 cycle drain) = 1040 ns
- ❑ Due to fill and drain effects of a pipeline ideal performance can be achieved only for very large instructions

Example: a sequence of 1000 load instructions would take 5000 cycles on a multi-cycle machine while taking 1004 on a pipeline machine

$$\Rightarrow \text{speedup} = 5000/1004 \cong 5$$



Pipeline Datapath



- Every stage must be completed in one clock cycle to avoid stalls
- Values must be latched to ensure correct execution of instructions
- The PC multiplexer has moved to the IF stage to prevent two instructions from updating the PC simultaneously (in case of branch instruction)

Pipeline Stage Interface

| Stage | Any Instruction | | |
|------------|---|---|--|
| IF | IF/ID.IR \leftarrow MEM[PC] ; IF/ID.NPC,PC \leftarrow (if ((EX/MEM.opcode == branch) & EX/MEM.cond) {EX/MEM.ALUOutput} else { PC + 4 }) ; | | |
| ID | ID/EX.A = Regs[IF/ID.IR _{6..10}]; ID/EX.B \leftarrow Regs[IF/ID.IR _{11..15}]; ID/EX.NPC \leftarrow IF/ID.NPC ; ID/EX.IR \leftarrow IF/ID.IR; ID/EX.Imm \leftarrow (IF/ID.IR ₁₆) ¹⁶ ## IF/ID.IR _{16..31} ; | | |
| | ALU | Load or Store | Branch |
| EX | EX/MEM.IR = ID/EX.IR; EX/MEM.ALUOutput \leftarrow ID/EX.A func ID/EX.B; Or EX/MEM.ALUOutput \leftarrow ID/EX.A op ID/EX.Imm; EX/MEM.cond \leftarrow 0; | EX/MEM.IR \leftarrow ID/EX.IR; EX/MEM.ALUOutput \leftarrow ID/EX.A + ID/EX.Imm; EX/MEM.cond \leftarrow 0; EX/MEM.B \leftarrow ID/EX.B; | EX/MEM.ALUOutput \leftarrow ID/EX.NPC + ID/EX.Imm; EX/MEM.cond \leftarrow (ID/EX.A op 0); |
| MEM | MEM/WB.IR \leftarrow EX/MEM.IR; MEM/WB.ALUOutput \leftarrow EX/MEM.ALUOutput; | MEM/WB.IR \leftarrow EX/MEM.IR; MEM/WB.LMD \leftarrow Mem[EX/MEM.ALUOutput] ; Or Mem[EX/MEM.ALUOutput] \leftarrow EX/MEM.B ; | |
| WB | Regs[MEM/WB.IR _{16..20}] \leftarrow MEM/WB.ALUOutput; Or Regs[MEM/WB.IR _{11..15}] \leftarrow MEM/WB.ALUOutput ; | For load only: Regs[MEM/WB.IR _{11..15}] \leftarrow MEM/WB.LMD; | |



Pipeline Hazards

- ❑ Pipeline hazards are cases that affect instruction execution semantics and thus need to be detected and corrected
- ❑ Hazards types

Structural hazard: attempt to use a resource two different ways at same time

- E.g., combined washer/dryer would be a structural hazard or folder busy doing something else (watching TV)
- Single memory for instruction and data

Data hazard: attempt to use item before it is ready

- E.g., one sock of pair in dryer and one in washer; can't fold until get sock from washer through dryer
- instruction depends on result of prior instruction still in the pipeline

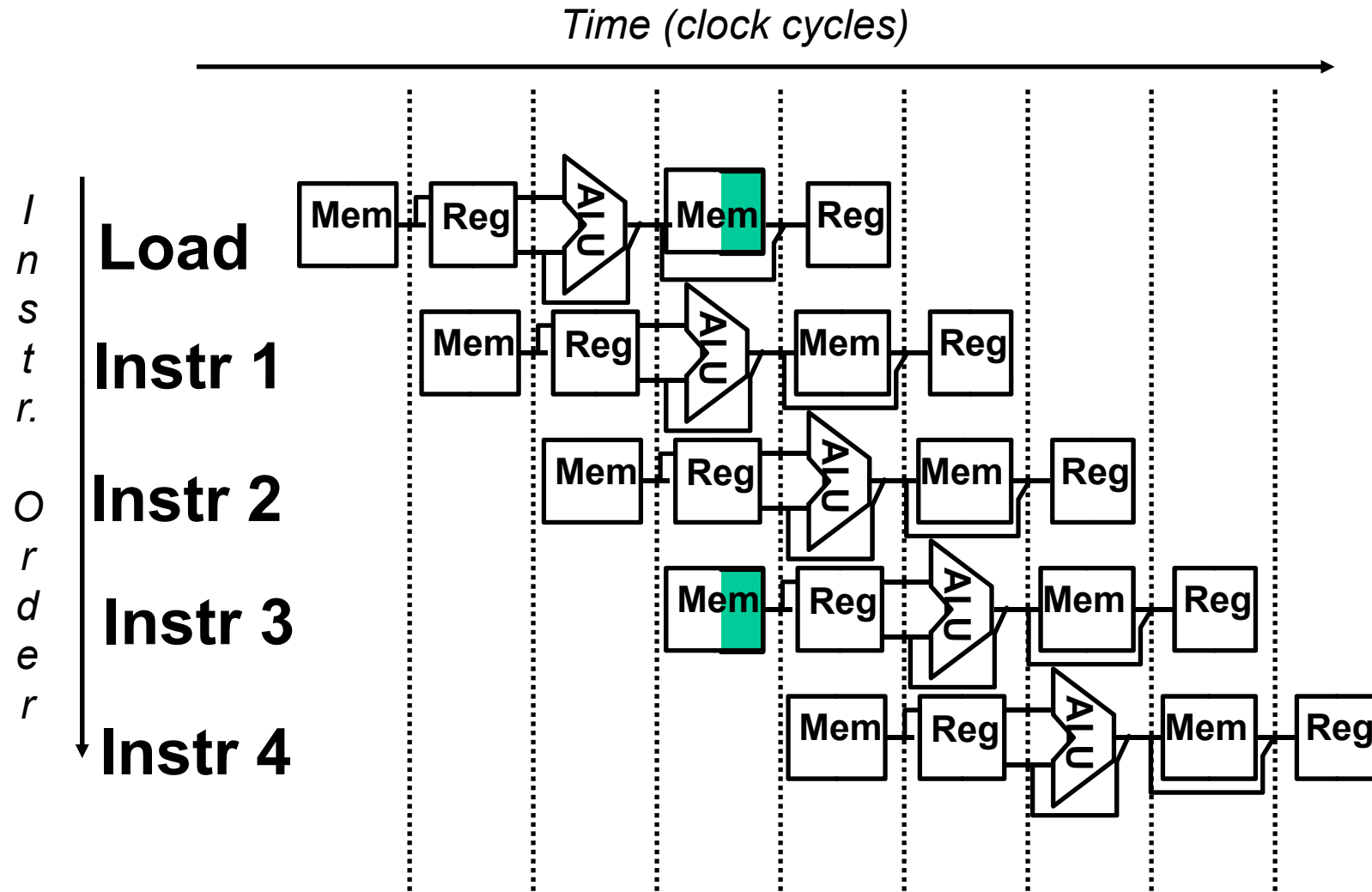
Control hazard: attempt to make a decision before condition is evaluated

- E.g., washing football uniforms and need to get proper detergent level; need to see after dryer before next load in
- branch instructions

- ❑ Hazards can always be resolved by **waiting**



Single Memory is a Structural Hazard



Can be easily detected

Resolved by inserting idle cycles



Stalls & Pipeline Performance

$$\begin{aligned}\text{Speedup from pipelining} &= \frac{\text{Average instruction time unpipelined}}{\text{Average instruction time pipelined}} \\ &= \frac{\text{CPI unpipelined}}{\text{CPI pipelined}} \times \frac{\text{Clock cycle unpipelined}}{\text{Clock cycle pipelined}}\end{aligned}$$

□ Ideally the CPI of the pipeline execution is 1 (after fill-up), thus

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

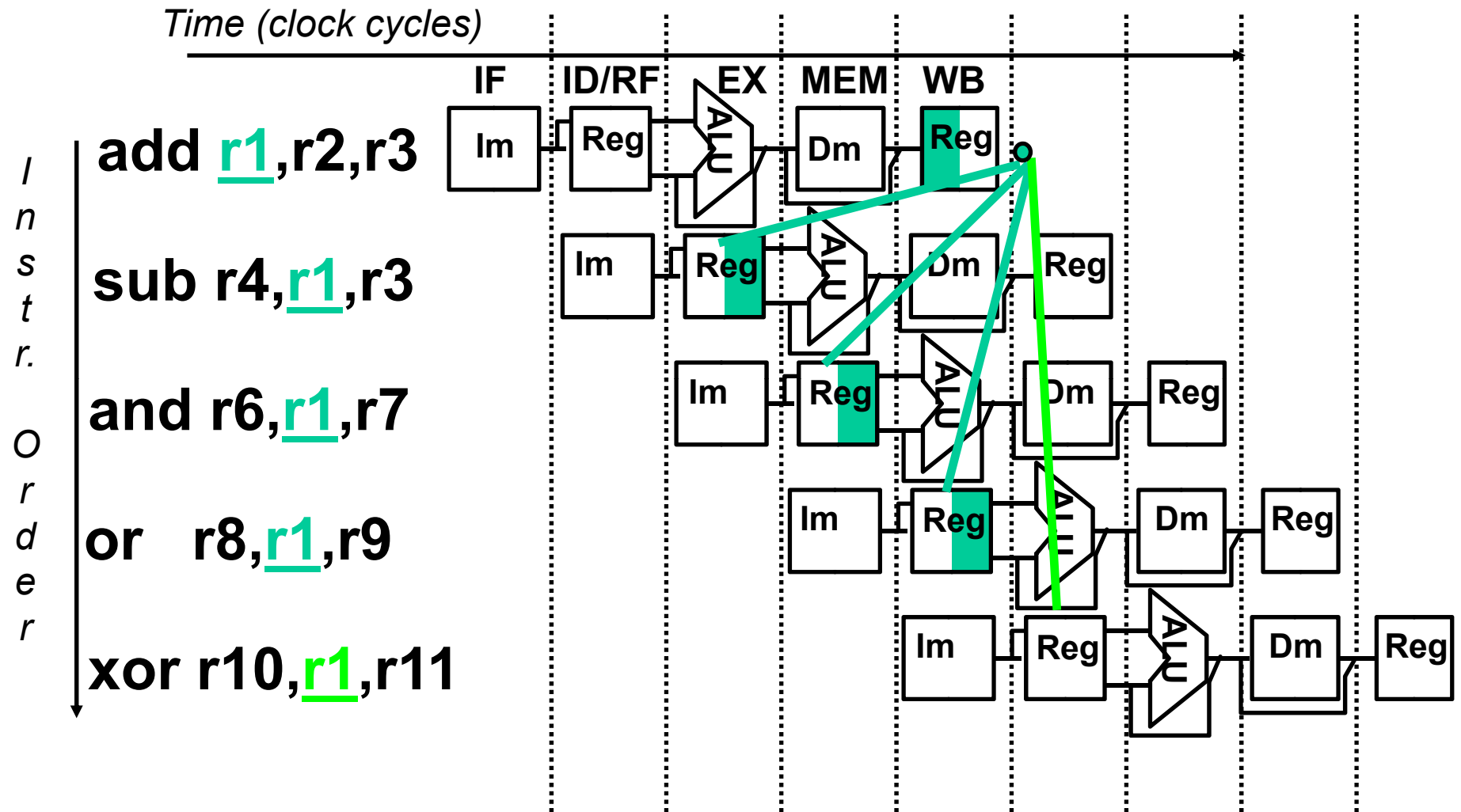
$$\text{Speedup} = \frac{\text{CPI unpipelined}}{1 + \text{Pipeline stall cycles per instruction}} \times \frac{\text{Clock cycle unpipelined}}{\text{Clock cycle pipelined}}$$

□ Assuming all pipeline stages are balanced, then

$$\text{Speedup} = \frac{1}{1 + \text{Pipeline stall cycles per instruction}} \times \text{Pipeline depth}$$



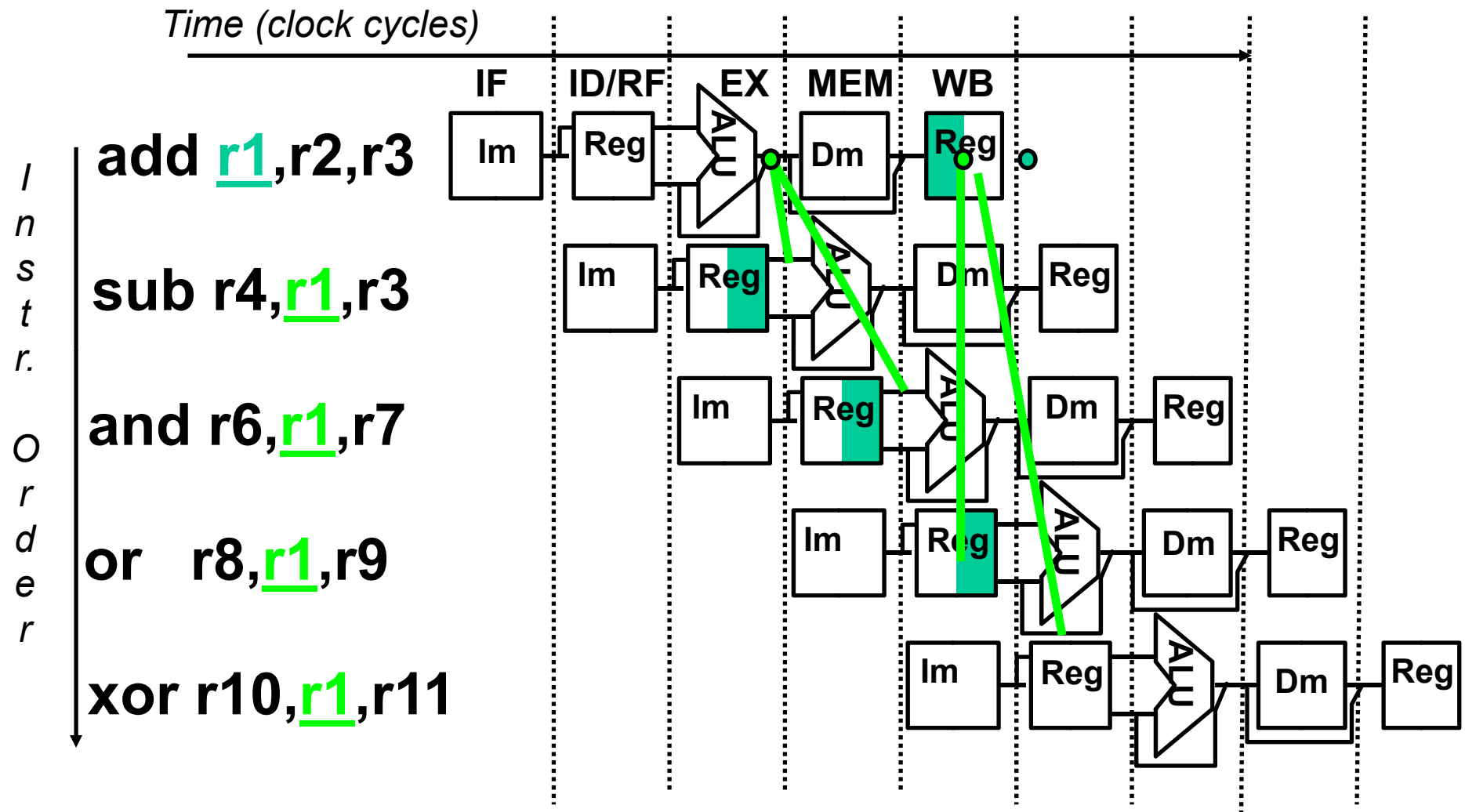
Data Hazard



Dependencies backwards in time are hazards



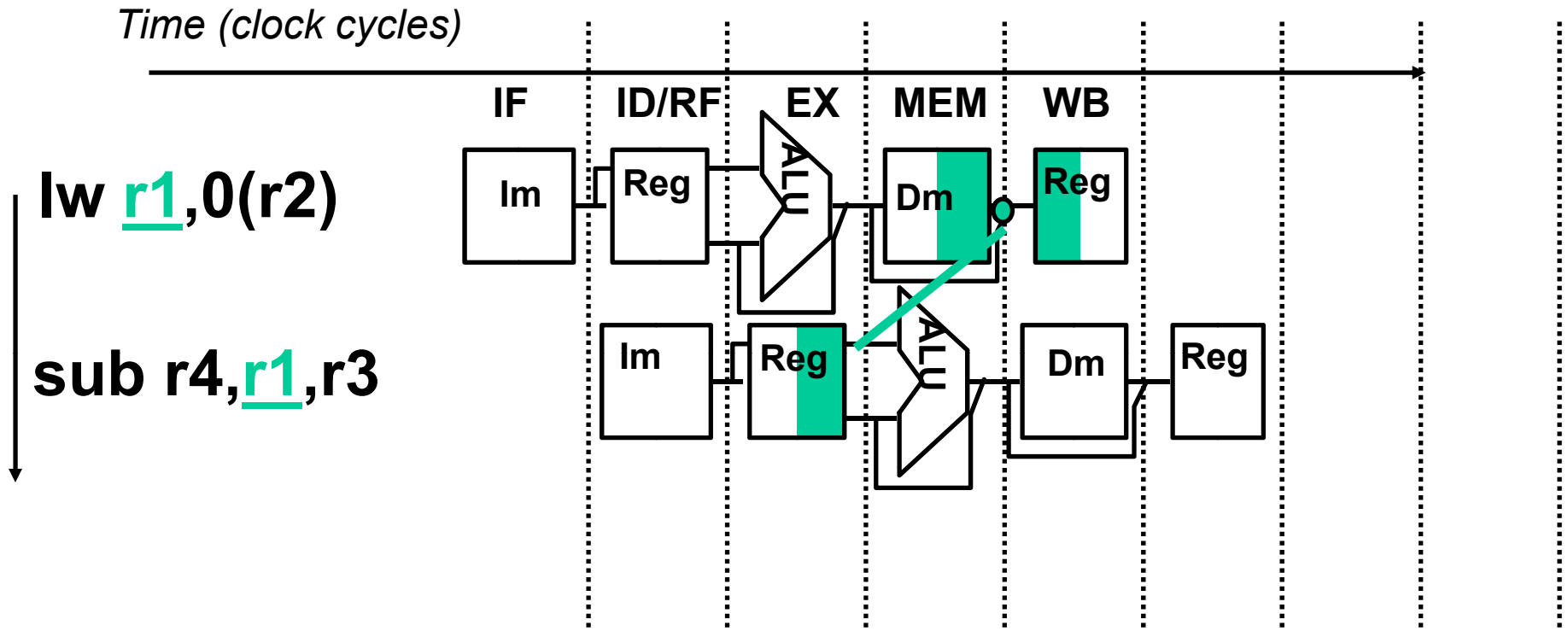
Data Hazard Solution



“Forward” result from one stage to another



Resolving Data Hazards for Loads



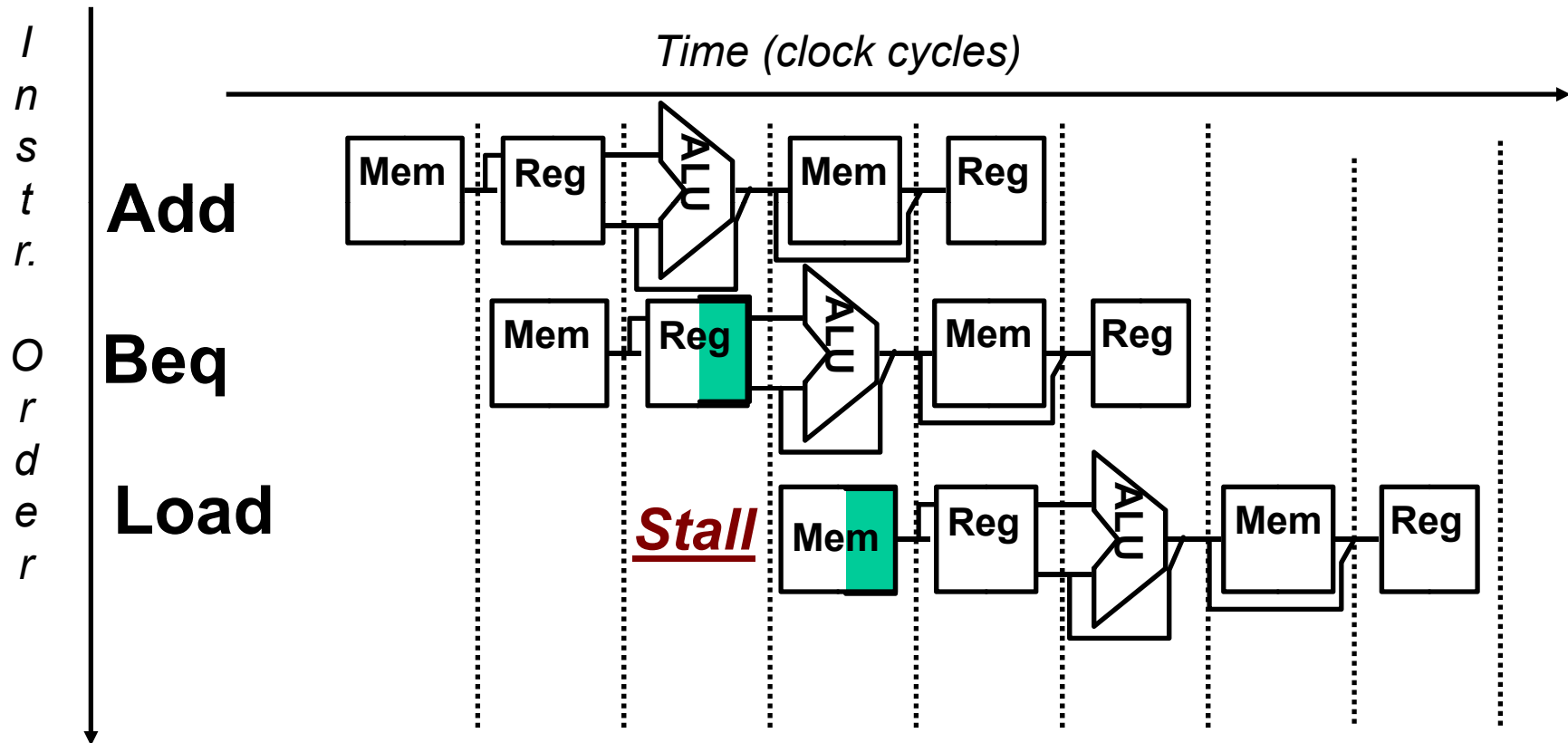
- Dependencies backwards in time are hazards
- Cannot solve with forwarding
- Must delay/stall instruction dependent on loads



Control Hazard

❑ **Stall**: wait until decision is clear

➔ Its possible to move up decision to 2nd stage by adding hardware to check registers as being read



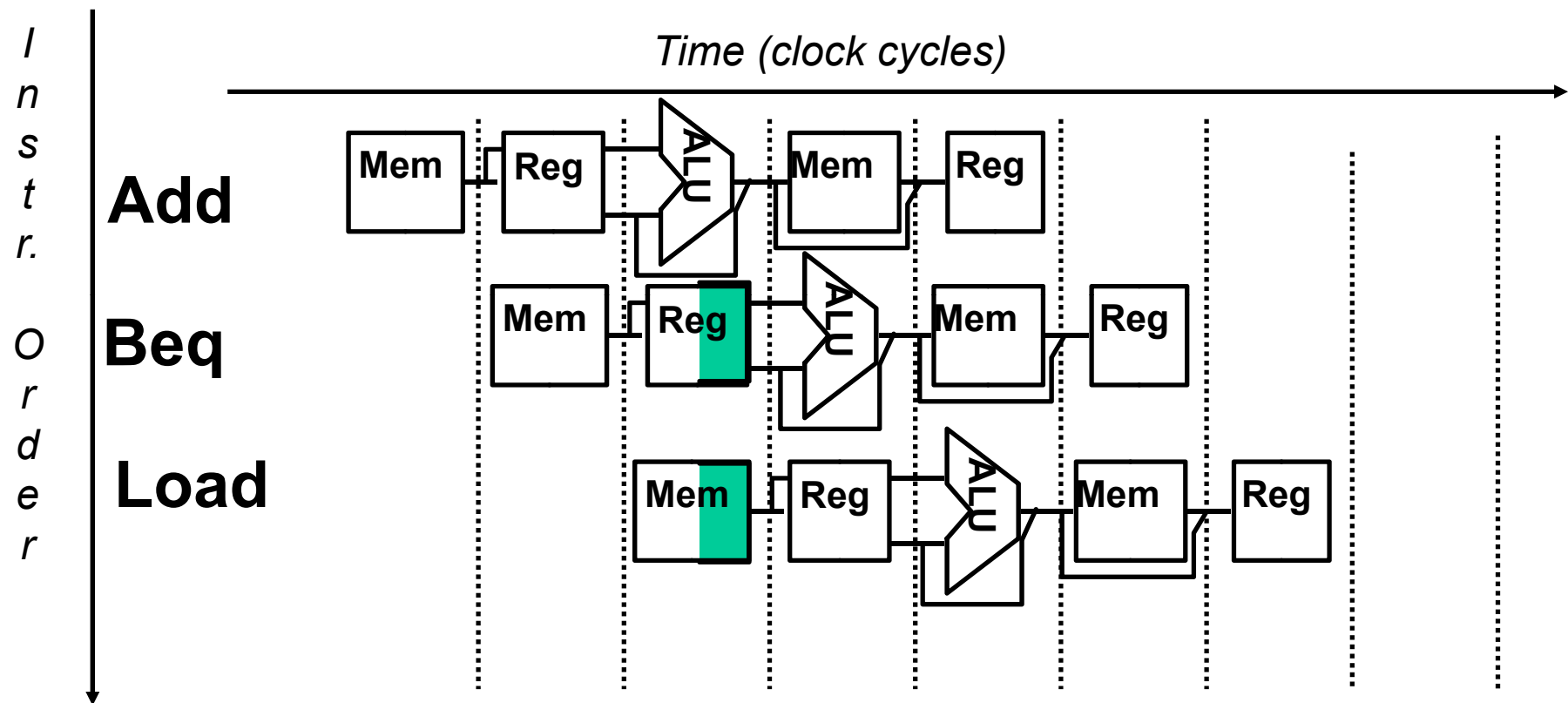
❑ Impact: 2 clock cycles per branch instruction ⇒ **slow**



Control Hazard Solution

- ❑ **Predict:** guess one direction then back up if wrong

➔ Predict not taken

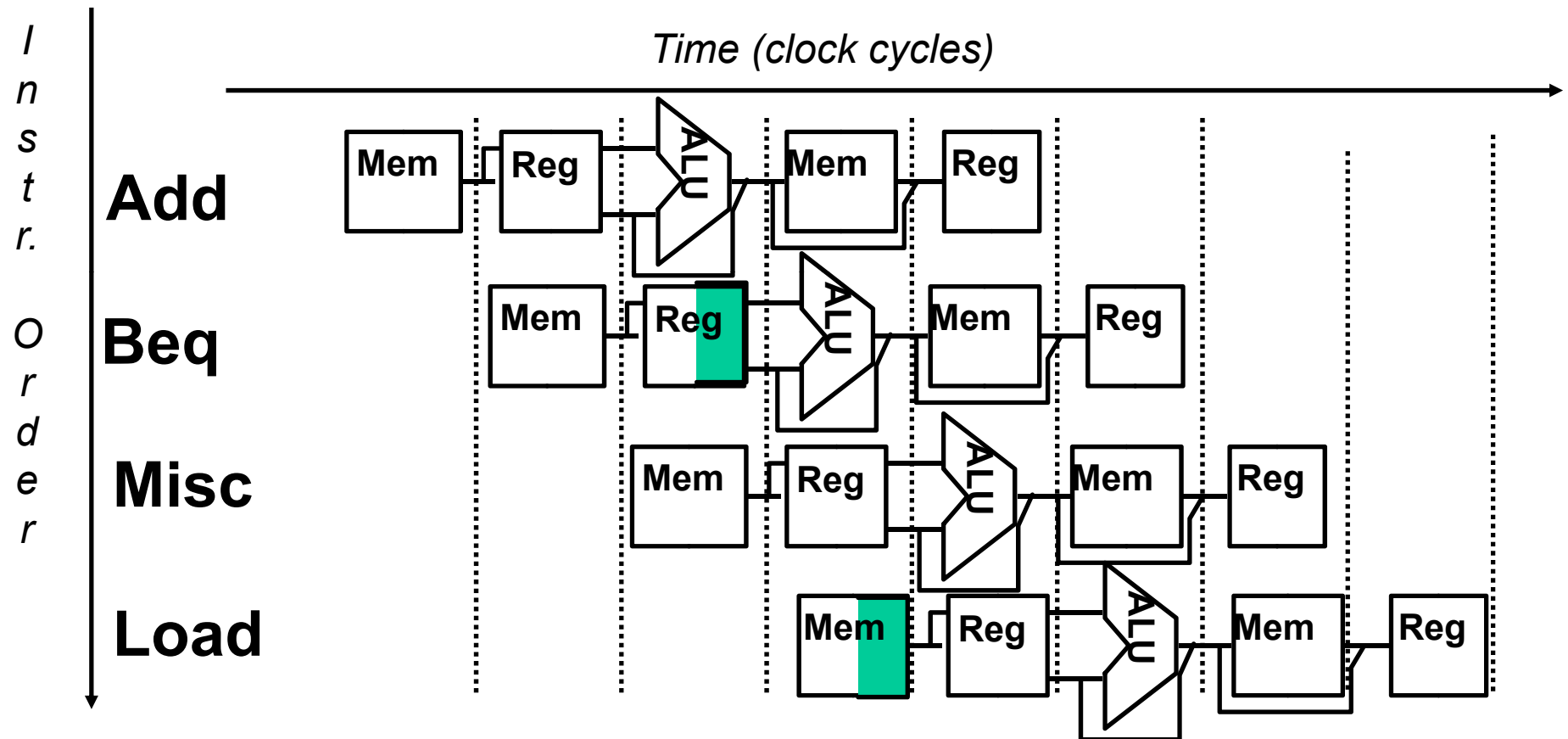


- ❑ **Impact:** 1 clock cycles per branch instruction if right, 2 if wrong (right - 50% of time)
- ❑ **More dynamic scheme:** history of 1 branch (- 90%)



Control Hazard Solution

- Redefine branch behavior (takes place after next instruction)
“delayed branch”



- Impact: 0 clock cycles per branch instruction if can find instruction to put in “slot” (- 50% of time)



Conclusion

□ Summary

- An overview of Pipelining
 - Pipelining concept is natural
 - Start handling of next instruction while current one is in progress
- Pipeline performance
 - Performance improvement by increasing instruction throughput
 - Ideal and upper bound for speedup is number of stages in pipeline
- Pipelined hazards
 - Structural, data and control hazards
 - Hazard resolution techniques

□ Next Lecture

- Data and control Hazards
- Pipelined control

Reading assignment includes Appendix A.1 & A.2 in the textbook

