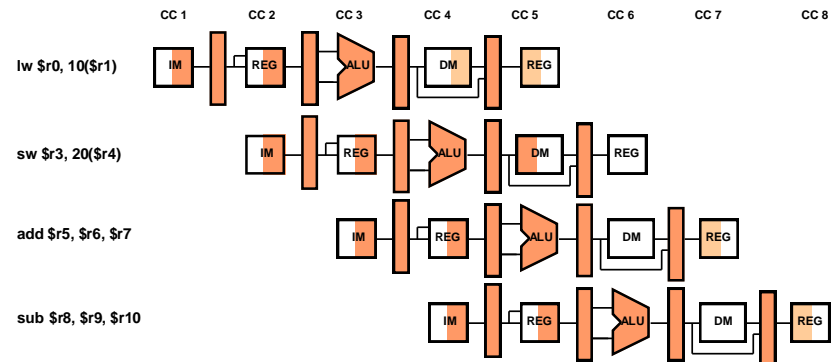# Computer Organization

## Lecture 17 - Pipelined Processor Design 1

**Reading: 4.5-4.6**

# Pipelining Outline

- **Introduction**
  - **Defining Pipelining** ◀
  - **Pipelining Instructions**
  - **Hazards**
- **Pipelined Processor Design**
  - **Datapath**
  - **Control**
- **Advanced Pipelining**
  - **Superscalar**
  - **Dynamic Pipelining**
  - **Examples**
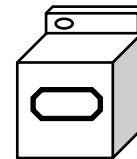
# What is Pipelining?

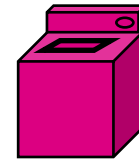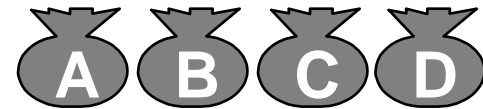▶ **A way of speeding up execution of instructions**

▶ **Key idea: <u>overlap</u> execution of <u>multiple</u> instructions**
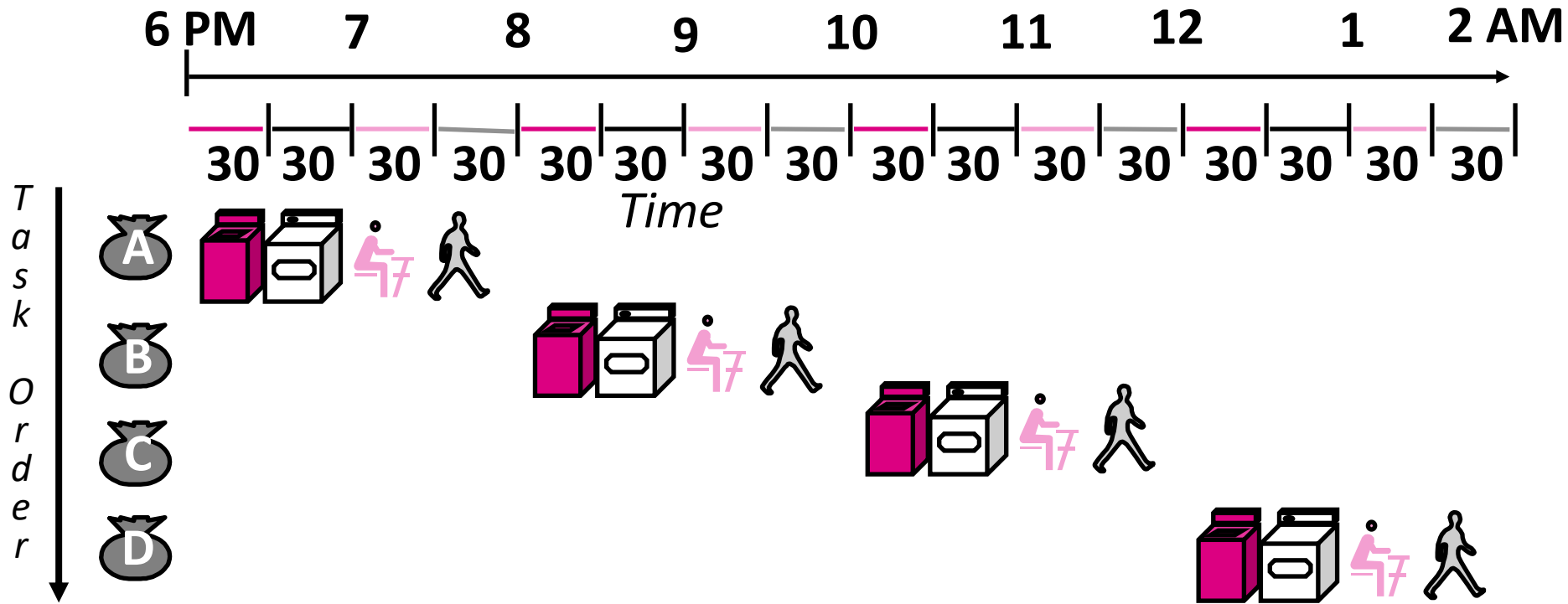
▶ **Analogy: doing your laundry**

# The Laundry Analogy

▸ **Ann, Brian, Cathy, Dave each has one load of clothes to wash, dry, and fold**

▸ **Washer takes 30 minutes**

▸ **Dryer takes 30 minutes**

▸ **"Folder" takes 30 minutes**

▸ **"Stasher" takes 30 minutes to put clothes into drawers**
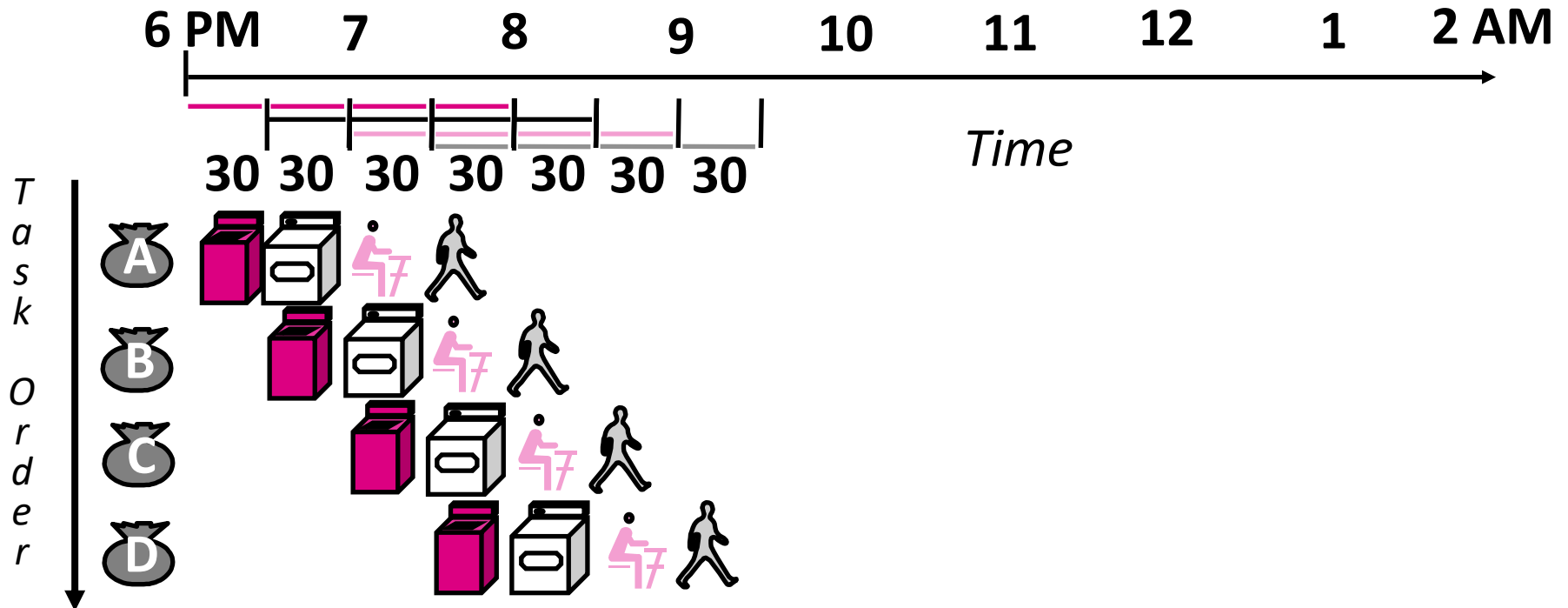
# If we do laundry sequentially...

6 PM    7    8    9    10    11    12    1    2 AM

30 30 30 30 30 30 30 30 30 30 30 30 30 30 30 30

*Time*

T
a
s
k

O
r
d
e
r

A

B

C

D

▸ **Time Required: 8 hours for 4 loads**

# To Pipeline, We Overlap Tasks



6 PM   7   8   9   10   11   12   1   2 AM

*Time*

30 30 30 30 30 30 30

*Task Order*
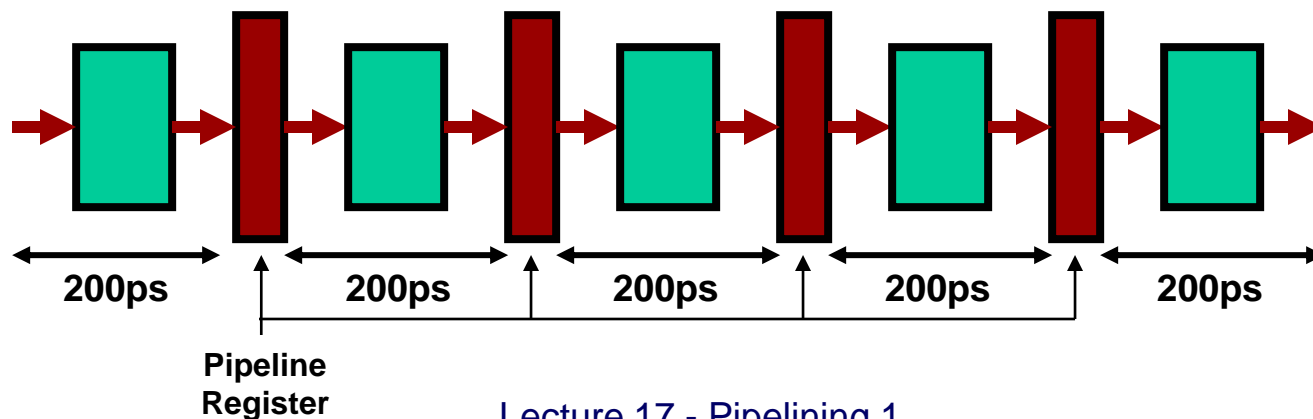
A

B

C

D

▸ **Time Required: 3.5 Hours for 4 Loads**

  ▸ **Latency remains 2 hours**

  ▸ **Throughput improves by factor of 2.3 (decreases for more loads)**

# Pipelining a Digital System

▶ **Key idea: break big computation up into pieces**

1ns

▶ **Separate each piece with a pipeline register**

200ps    200ps    200ps    200ps    200ps

Pipeline
Register

# Pipelining a Digital System

▶ **Why do this?  Because it's faster for repeated computations**

**Non-pipelined:
1 operation finishes
every 1ns**

**1ns**

**Pipelined:
1 operation finishes
every 200ps**

**200ps**      **200ps**      **200ps**      **200ps**      **200ps**

# Comments about pipelining

- **Pipelining increases throughput, but not latency**
  - **Answer available every 200ps, BUT**
  - **A single computation still takes 1ns**
- **Limitations:**
  - **Computations must be divisible into stage size**
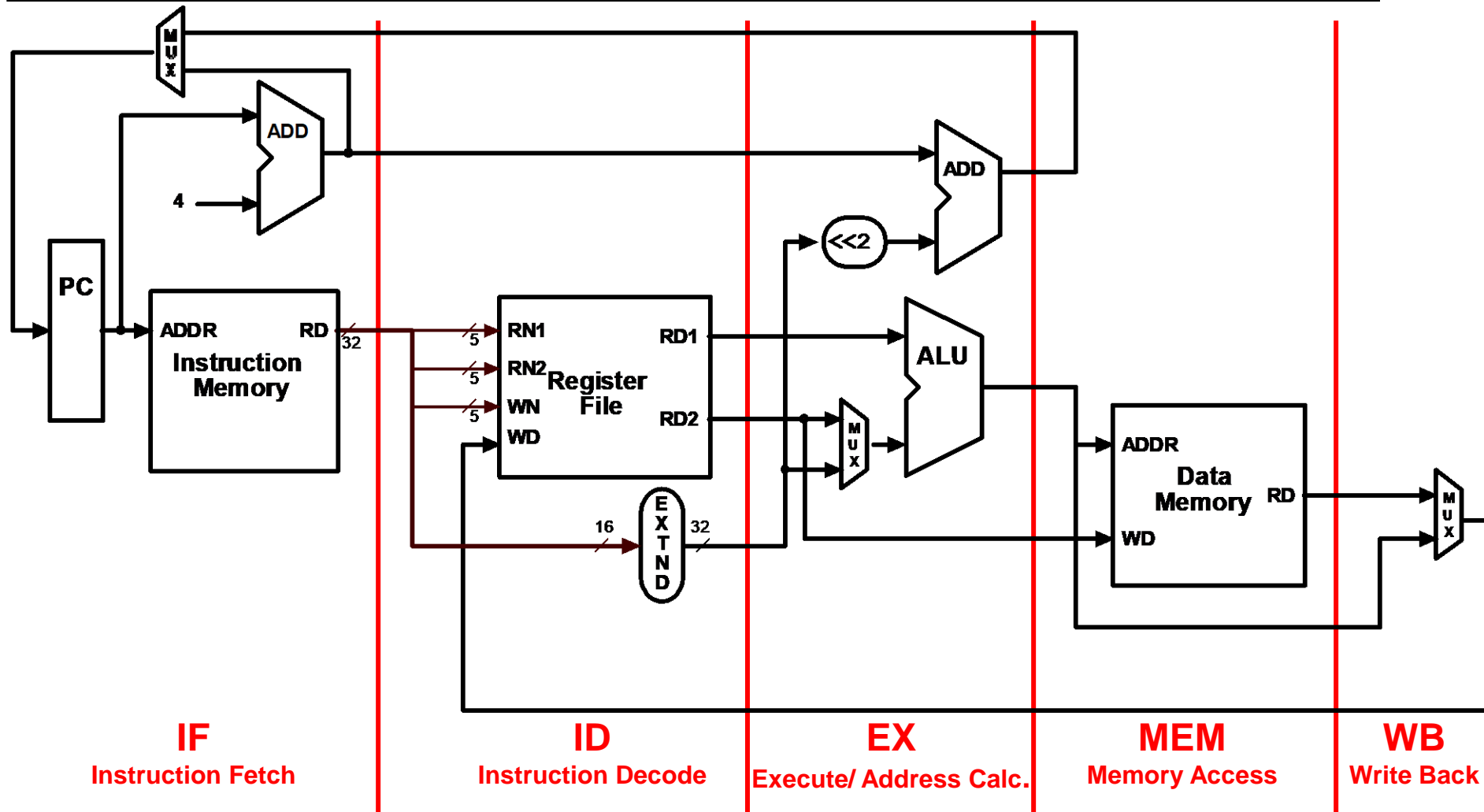  - **Pipeline registers add overhead**

# Pipelining Outline

- ▶ **Introduction**
  - ▶ **Defining Pipelining**
  - ▶ **Pipelining Instructions** ◀
  - ▶ **Hazards**
- ▶ **Pipelined Processor Design**
  - ▶ **Datapath**
  - ▶ **Control**
- ▶ **Advanced Pipelining**
  - ▶ **Superscalar**
  - ▶ **Dynamic Pipelining**
  - ▶ **Examples**

# Pipelining a Processor

▸ **Recall the 5 steps in instruction execution:**

  **1. Instruction Fetch**

  **2. Instruction Decode and Register Read**

  **3. Execution operation or calculate address**

  **4. Memory access**

  **5. Write result into register**

▸ **Review: Single-Cycle Processor**

  ▸ **All 5 steps done in a single clock cycle**

  ▸ **Dedicated hardware required for each step**

▸ **What happens if we break execution into multiple cycles, but keep the extra hardware?**
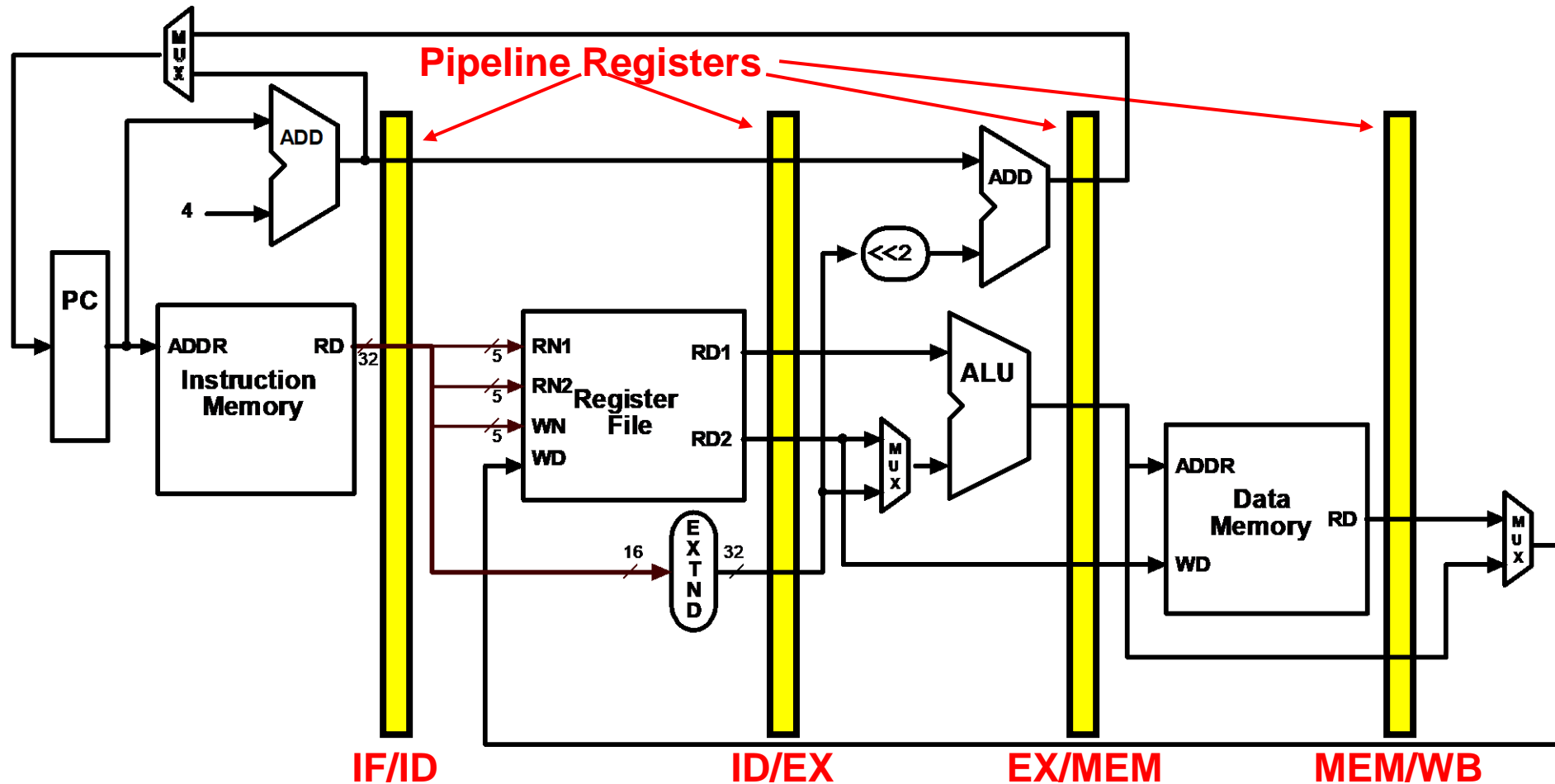
# Review - Single-Cycle Processor



| IF | ID | EX | MEM | WB |
|---|---|---|---|---|
| **Instruction Fetch** | **Instruction Decode** | **Execute/ Address Calc.** | **Memory Access** | **Write Back** |

# Pipelining - Key Idea
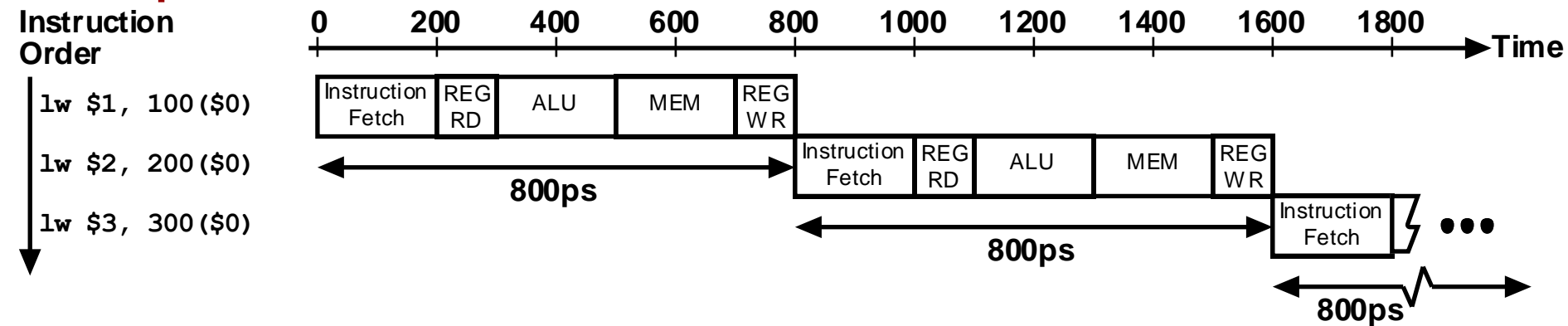
▸ **Question: What happens if we break execution into multiple cycles, but keep the extra hardware?**

▸ **Answer: in the <u>best case</u>, we can start executing a new instruction on each clock cycle - this is <u>pipelining</u>**

▸ **Pipelining <u>stages</u>:**
- ▸ **IF - Instruction Fetch**
- ▸ **ID - Instruction Decode**
- ▸ **EX - Execute / Address Calculation**
- ▸ **MEM - Memory Access (read / write)**
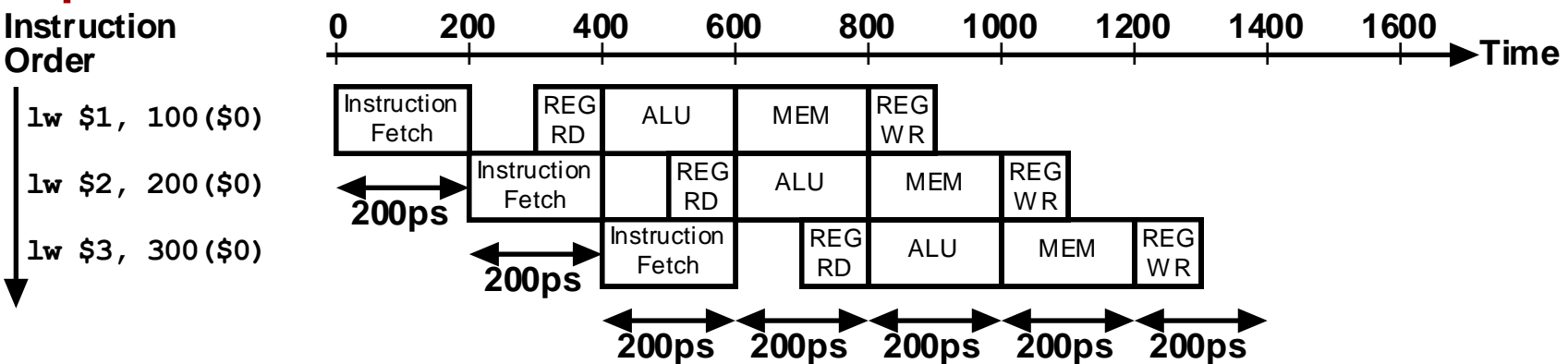- ▸ **WB - Write Back (results into register file)**

# Basic Pipelined Processor

# Single-Cycle vs. Pipelined Execution

**Non-Pipelined**

Instruction Order

lw $1, 100($0)

lw $2, 200($0)

lw $3, 300($0)

Time: 0  200  400  600  800  1000  1200  1400  1600  1800

| Instruction Fetch | REG RD | ALU | MEM | REG WR |

800ps

| Instruction Fetch | REG RD | ALU | MEM | REG WR |

800ps

| Instruction Fetch |

•••

800ps

**Pipelined**

Instruction Order

lw $1, 100($0)

lw $2, 200($0)

lw $3, 300($0)

Time: 0  200  400  600  800  1000  1200  1400  1600

| Instruction Fetch | REG RD | ALU | MEM | REG WR |

200ps

| Instruction Fetch | REG RD | ALU | MEM | REG WR |

200ps

| Instruction Fetch | REG RD | ALU | MEM | REG WR |

200ps  200ps  200ps  200ps  200ps

# Comments about Pipelining

- **The good news**
  - **Multiple instructions are being processed at same time**
  - **This works because stages are <u>isolated</u> by registers**
  - **For N stages, best case speedup of N**
- **The bad news**
  - **Instructions interfere with each other - <u>hazards</u>**
    - **Example: different instructions may need the same piece of hardware (e.g., memory) in same clock cycle**
    - **Example: instruction may require a result produced by an earlier instruction that is not yet complete**
  - **Worst case: must suspend execution - <u>stall</u>**
    - **Wait until needed hardware is available OR**
    - **Wait until needed data is available**

# Pipelined Example - Executing Multiple Instructions

▸ **Consider the following instruction sequence:**

```
lw $r0, 10($r1)
sw $r3, 20($r4)
add $r5, $r6, $r7
sub $r8, $r9, $r10
```
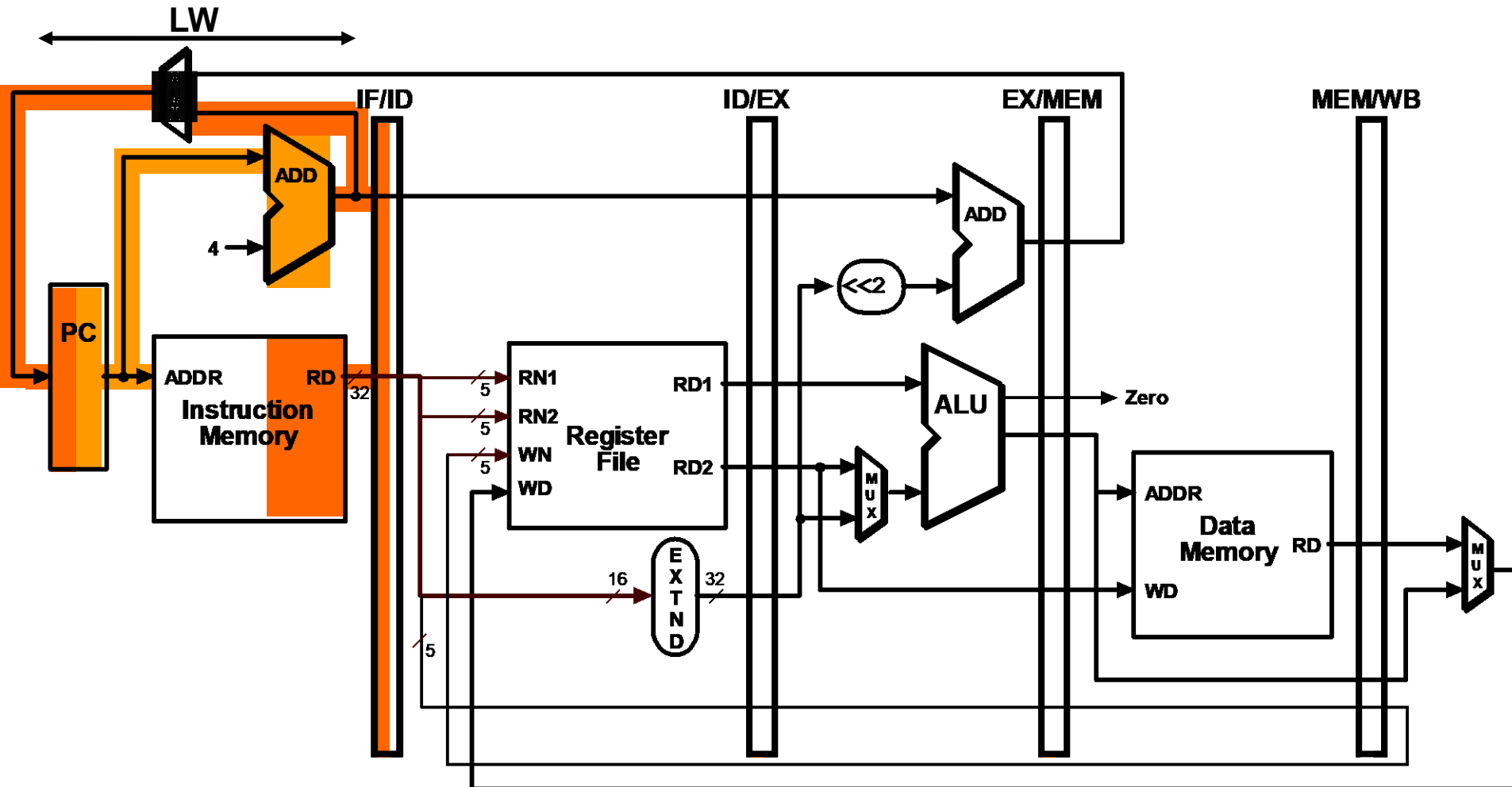
# Executing Multiple Instructions Clock Cycle 1
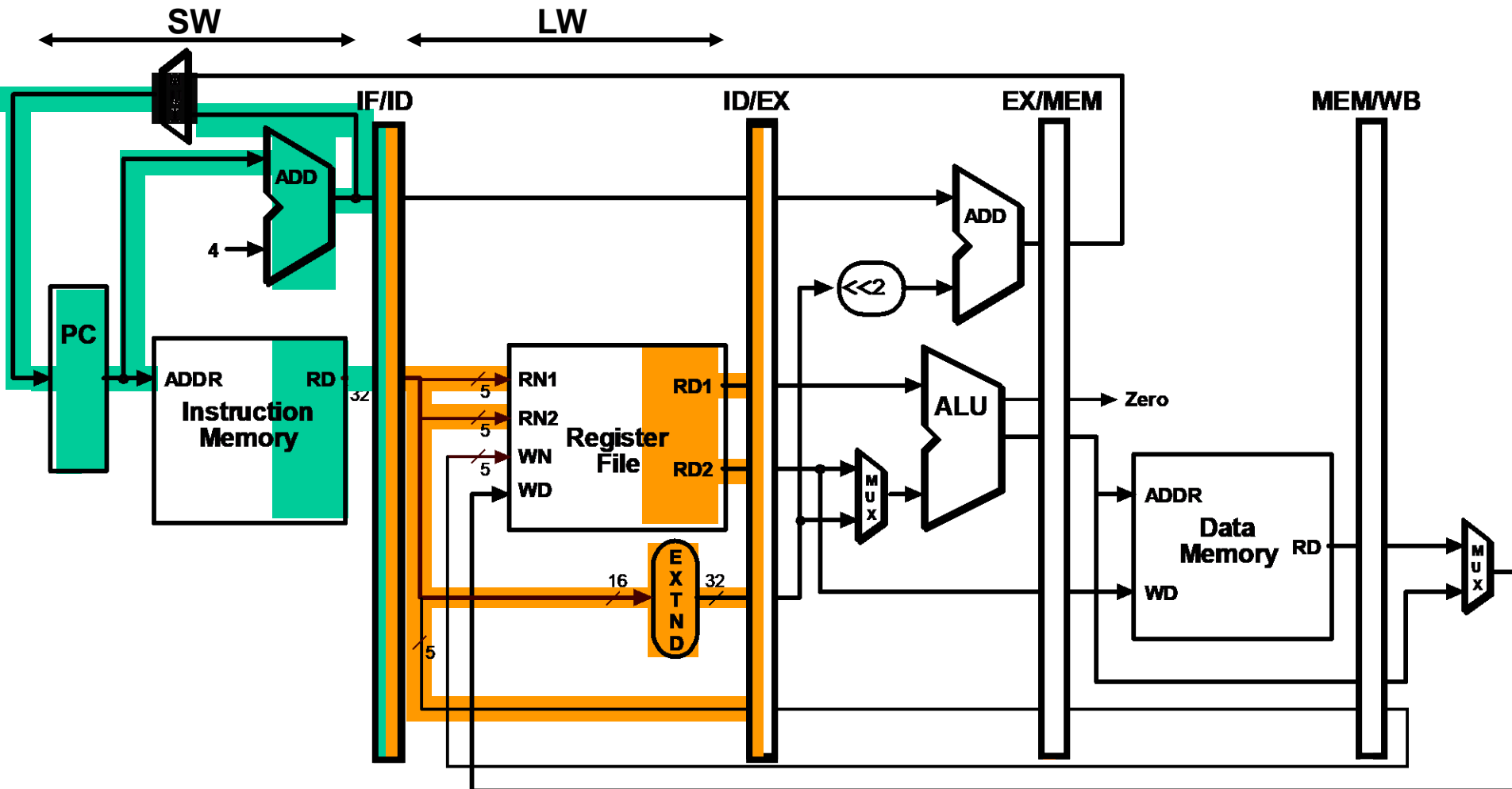
```
lw $r0, 10($r1)
sw $r3, 20($r4)
add $r5, $r6, $r7
sub $r8, $r9, $r10
```
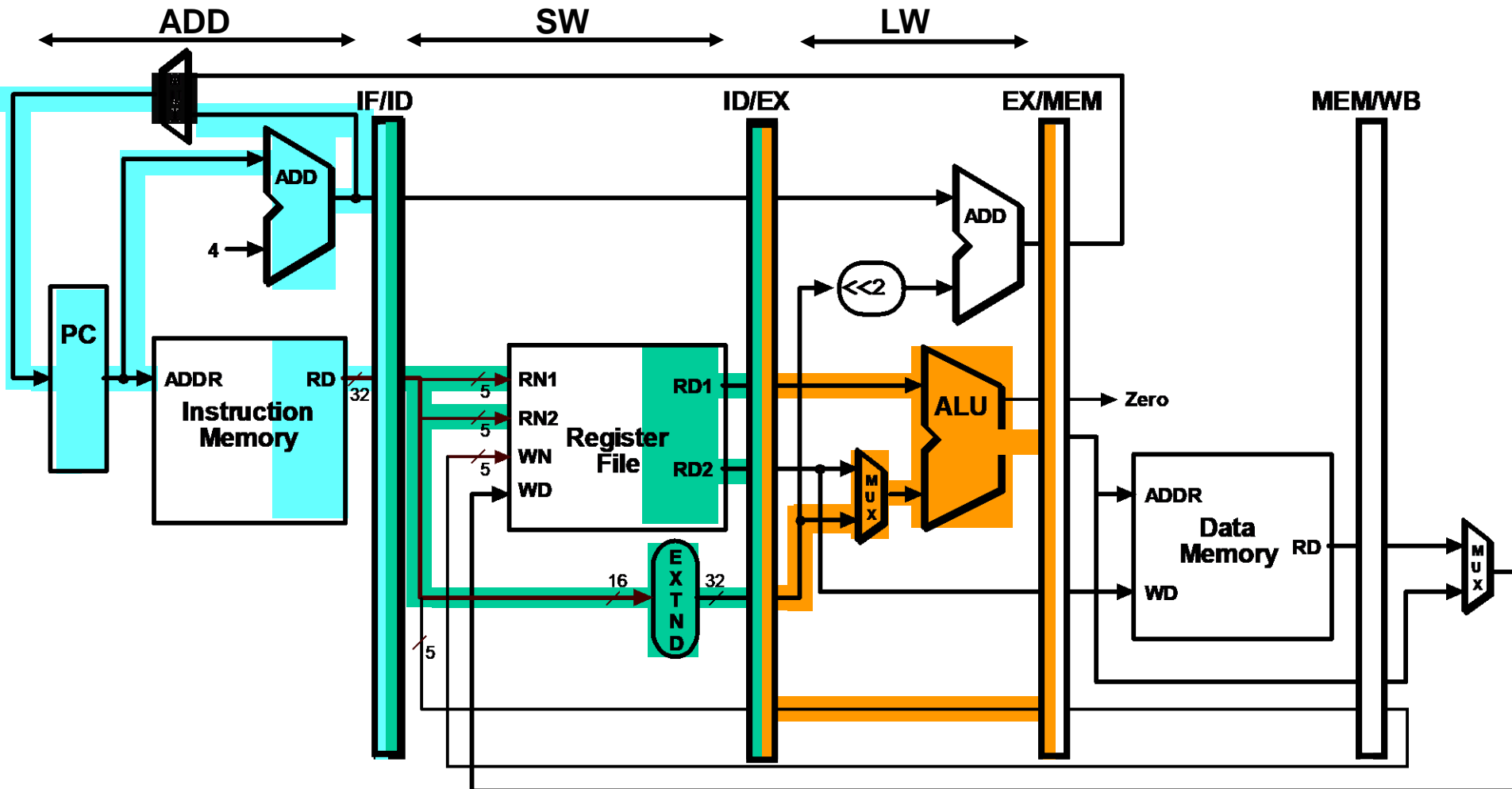
# Executing Multiple Instructions Clock Cycle 2

```
lw $r0, 10($r1)
sw $r3, 20($r4)
add $r5, $r6, $r7
sub $r8, $r9, $r10
```
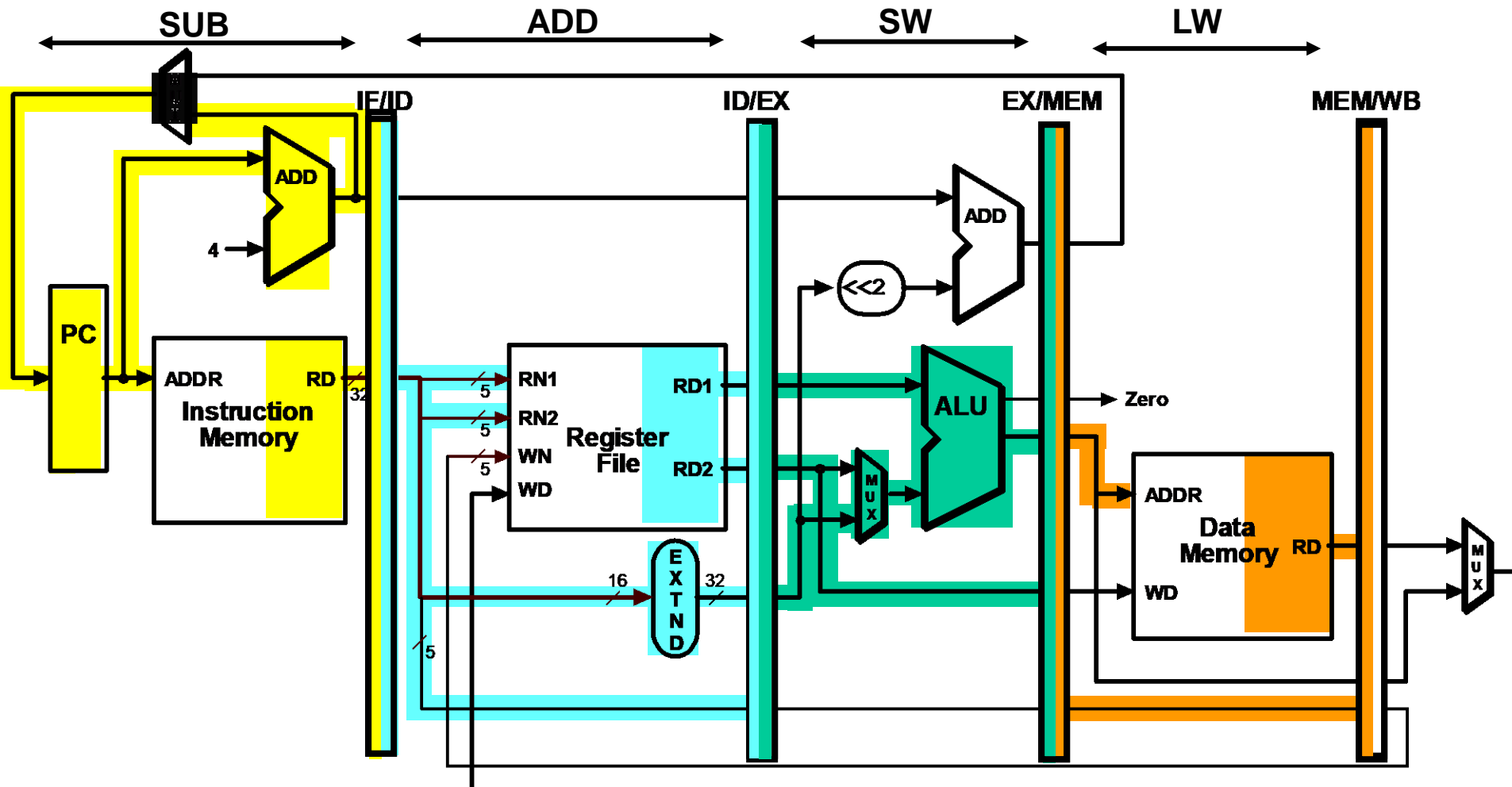
# Executing Multiple Instructions Clock Cycle 3

```
lw $r0, 10($r1)
sw $r3, 20($r4)
add $r5, $r6, $r7
sub $r8, $r9, $r10
```

# Executing Multiple Instructions Clock Cycle 4
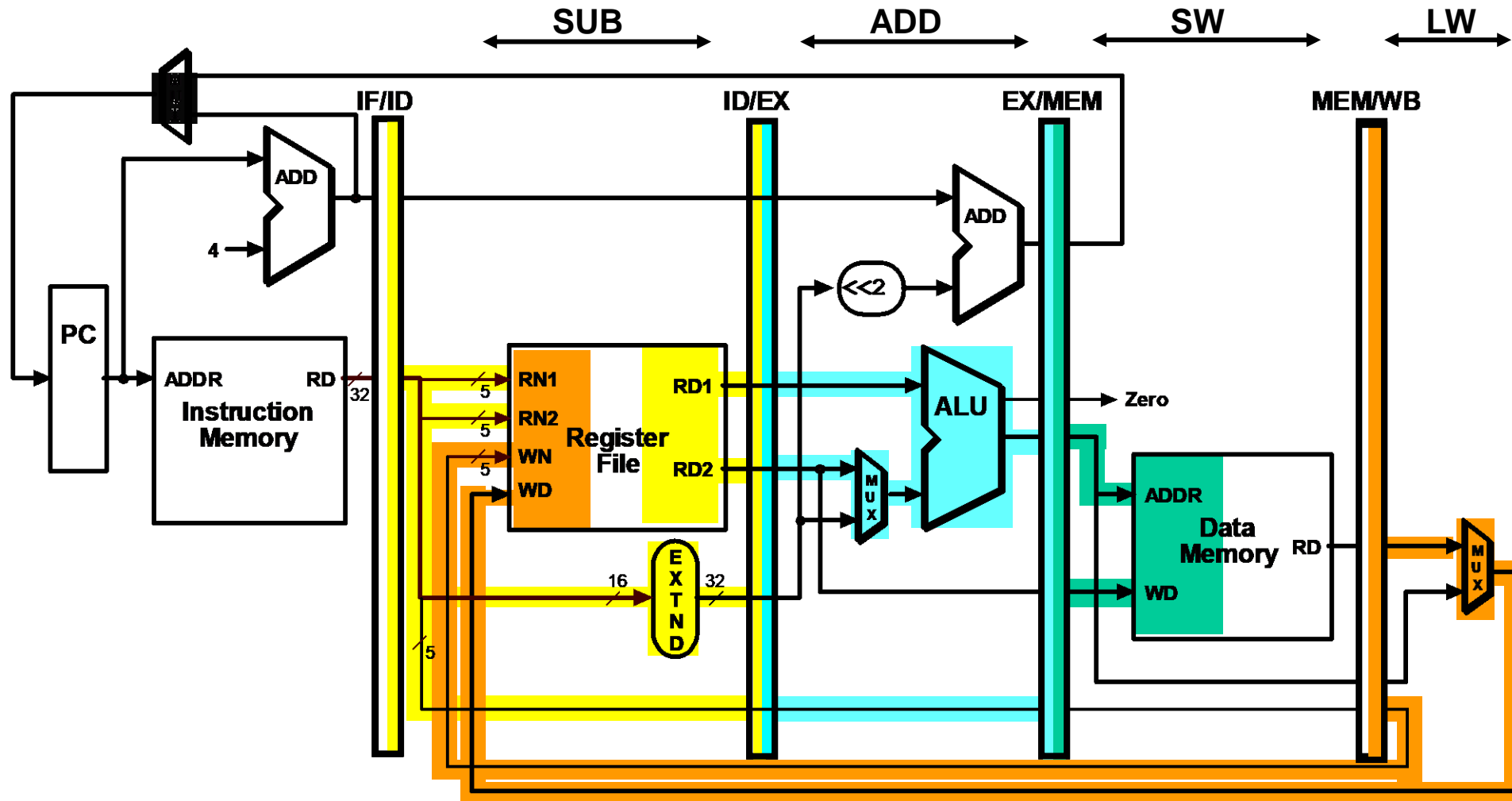
```
lw $r0, 10($r1)
sw $r3, 20($r4)
add $r5, $r6, $r7
sub $r8, $r9, $r10
```

# Executing Multiple Instructions
# Clock Cycle 5

lw $r0, 10($r1)
sw $r3, 20($r4)
add $r5, $r6, $r7
sub $r8, $r9, $r10
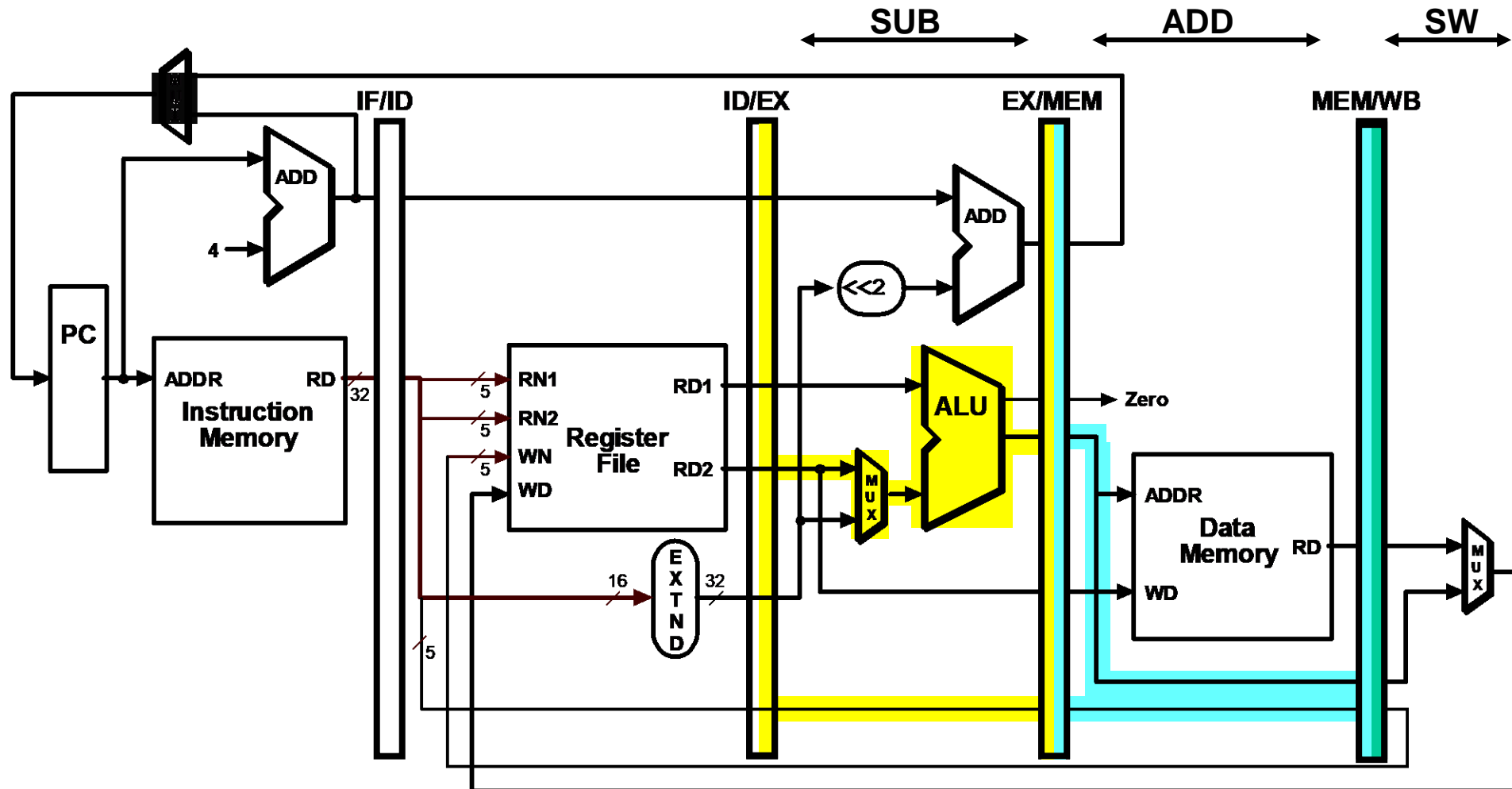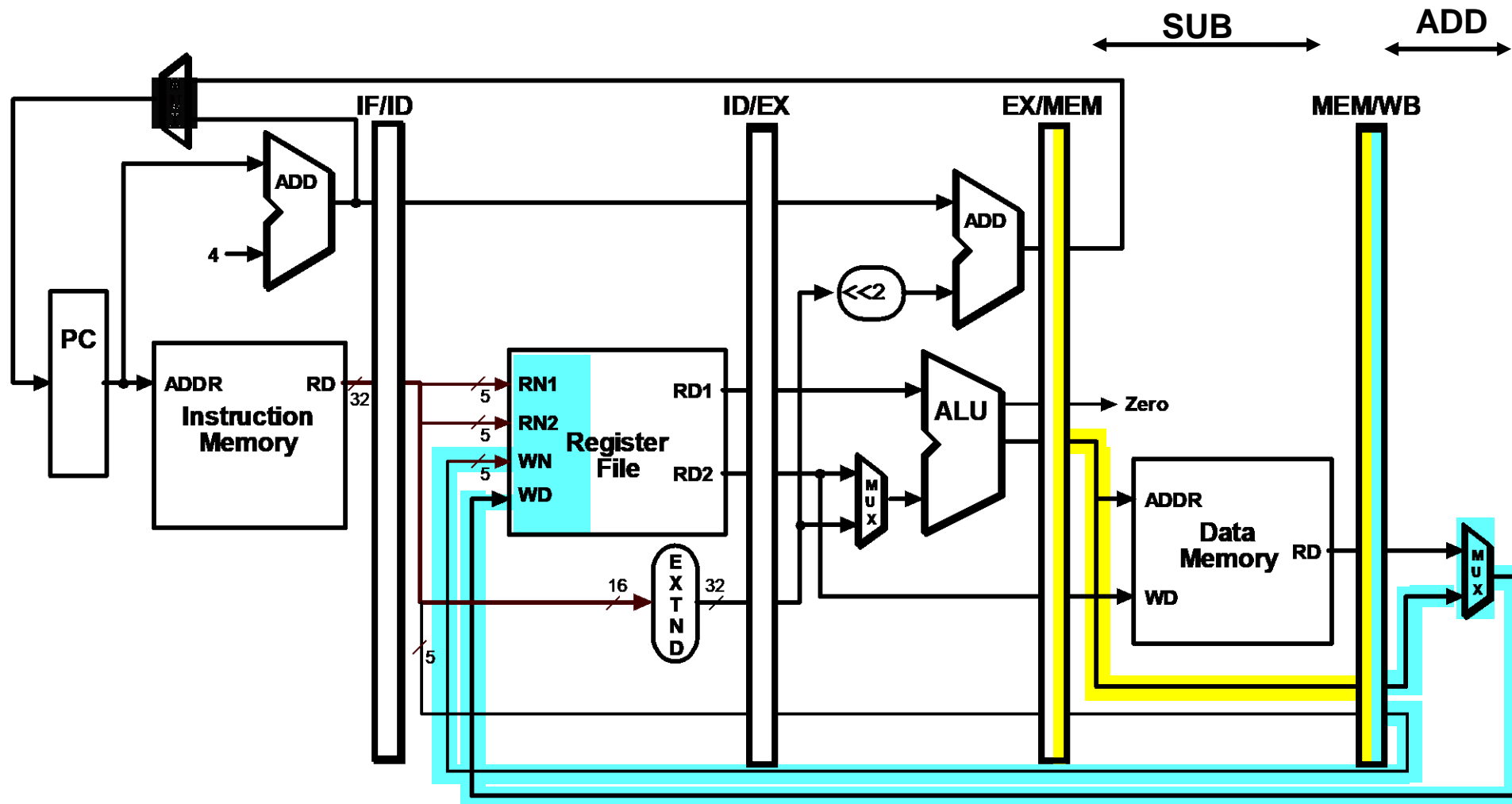
# Executing Multiple Instructions
# Clock Cycle 6

```
lw $r0, 10($r1)
sw $r3, 20($r4)
add $r5, $r6, $r7
sub $r8, $r9, $r10
```

# Executing Multiple Instructions Clock Cycle 7

```
lw $r0, 10($r1)
sw $r3, 20($r4)
add $r5, $r6, $r7
sub $r8, $r9, $r10
```
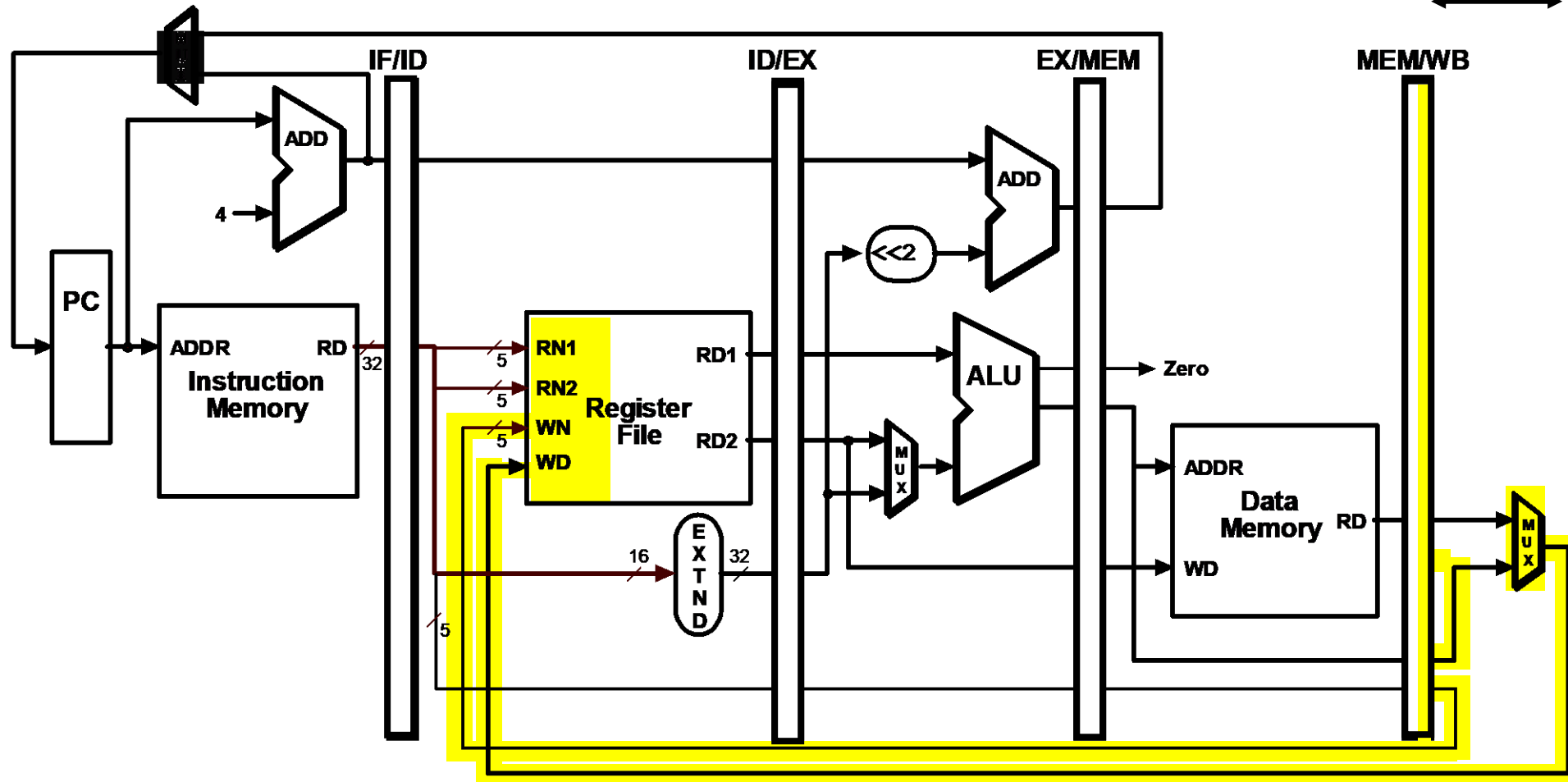
# Executing Multiple Instructions Clock Cycle 8

```
lw $r0, 10($r1)
sw $r3, 20($r4)
add $r5, $r6, $r7
sub $r8, $r9, $r10
```

# Alternative View - Multicycle Diagram

# Hazards

# Pipelining Outline

- ▸ **Introduction**
  - ▸ **Defining Pipelining**
  - ▸ **Pipelining Instructions**
  - ▸ **Hazards**                                                   ◂
- ▸ **Pipelined Processor Design**
  - ▸ **Datapath**
  - ▸ **Control**
- ▸ **Advanced Pipelining**
  - ▸ **Superscalar**
  - ▸ **Dynamic Pipelining**
  - ▸ **Examples**

# Pipeline Hazards

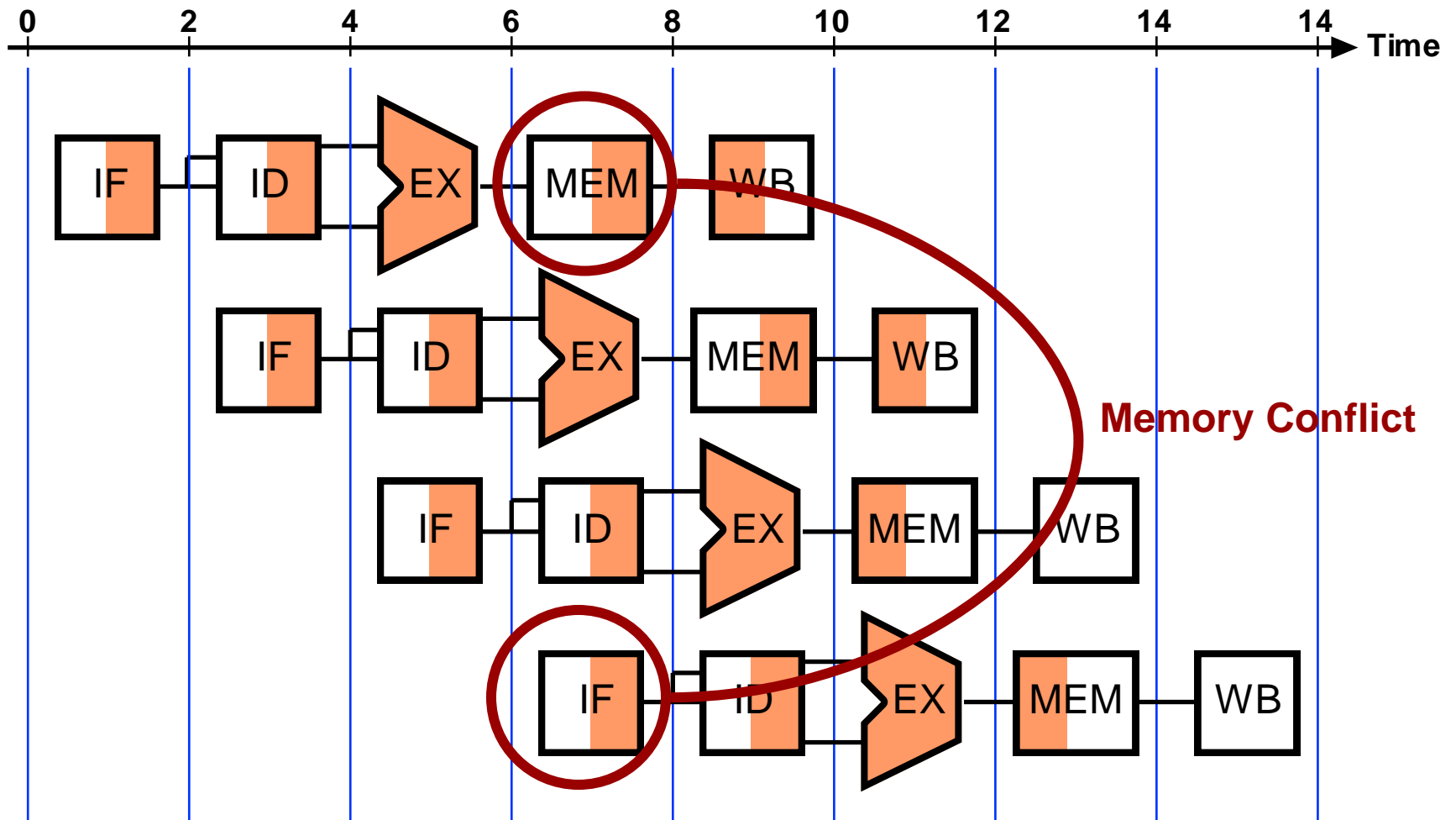▸ **Where one instruction cannot <u>immediately</u> follow another**

▸ **Types of hazards**

   ▸ <u>**Structural hazards**</u> **- attempt to use same resource twice**

   ▸ <u>**Control hazards**</u> **- attempt to make decision before condition is evaluated**

   ▸ <u>**Data hazards**</u> **- attempt to use data before it is ready**

▸ **Hazards can be resolved by waiting ("stalling")**

# Structural Hazards

- **Attempt to use same resource twice at same time**
- **Example: Single Memory for instructions and data**
  - **Accessed by IF stage**
  - **Accessed at the same time by MEM stage**
- **Solutions**
  - **Delay second access by one clock cycle, OR**
  - **Provide separate memories for instructions and data**
    - **This is what the book does**
    - **This is called a "Harvard Architecture"**
    - **Real pipelined processors have separate caches**

# Example Structural Hazard - Single Memory



Memory Conflict

# Control Hazards

▸ **Attempt to make a decision before condition is evaluated**

▸ **Example: `beq $s0, $s1, offset`**

▸ **Assume we add hardware to second stage to:**

 ▸ **Compare fetched registers for equality**

 ▸ **Compute branch target**

▸ **This allows branch to be taken at <u>end</u> of second clock cycle**

▸ **But, this still means result is not ready when we want to load the next instruction!**
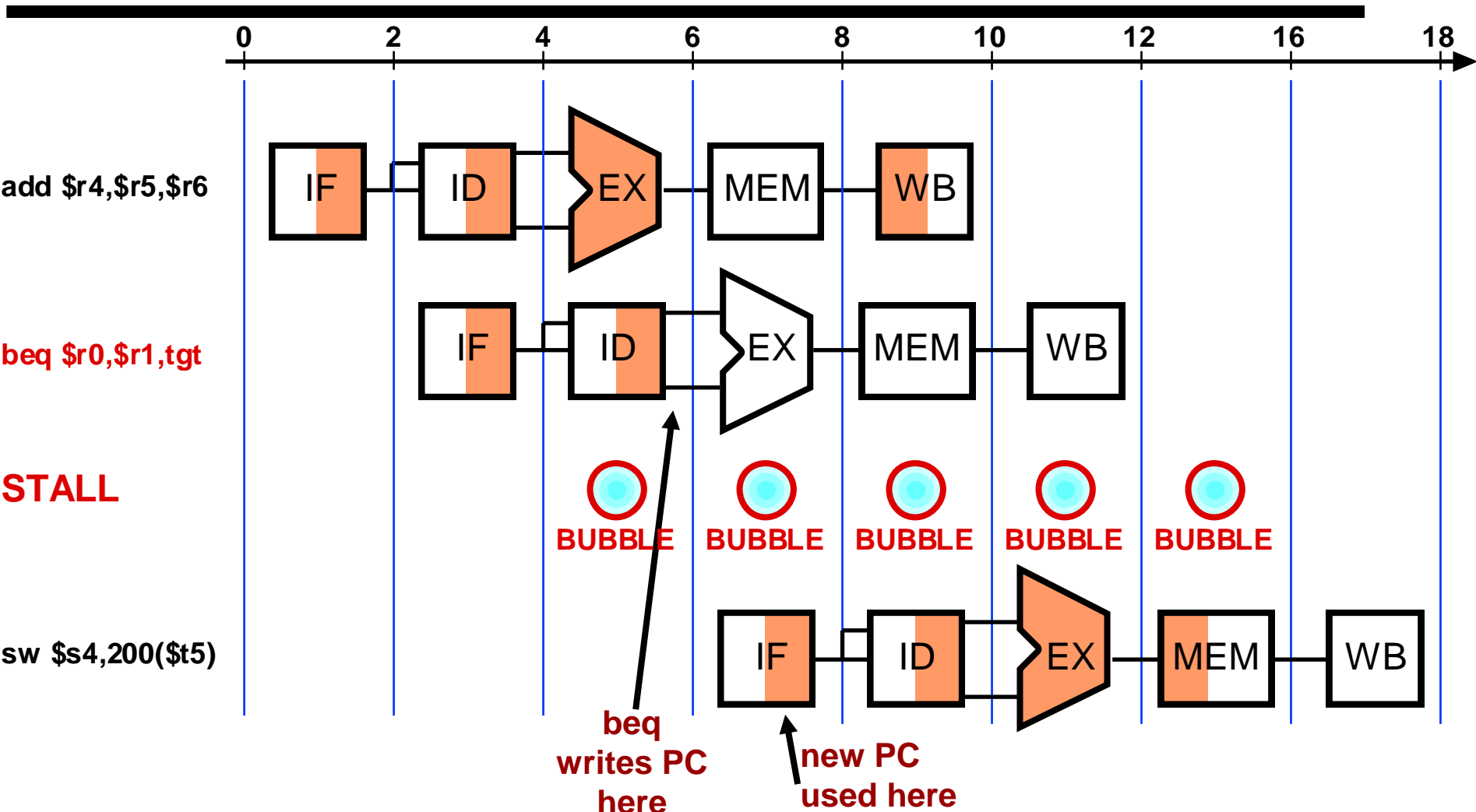
# Control Hazard Solutions
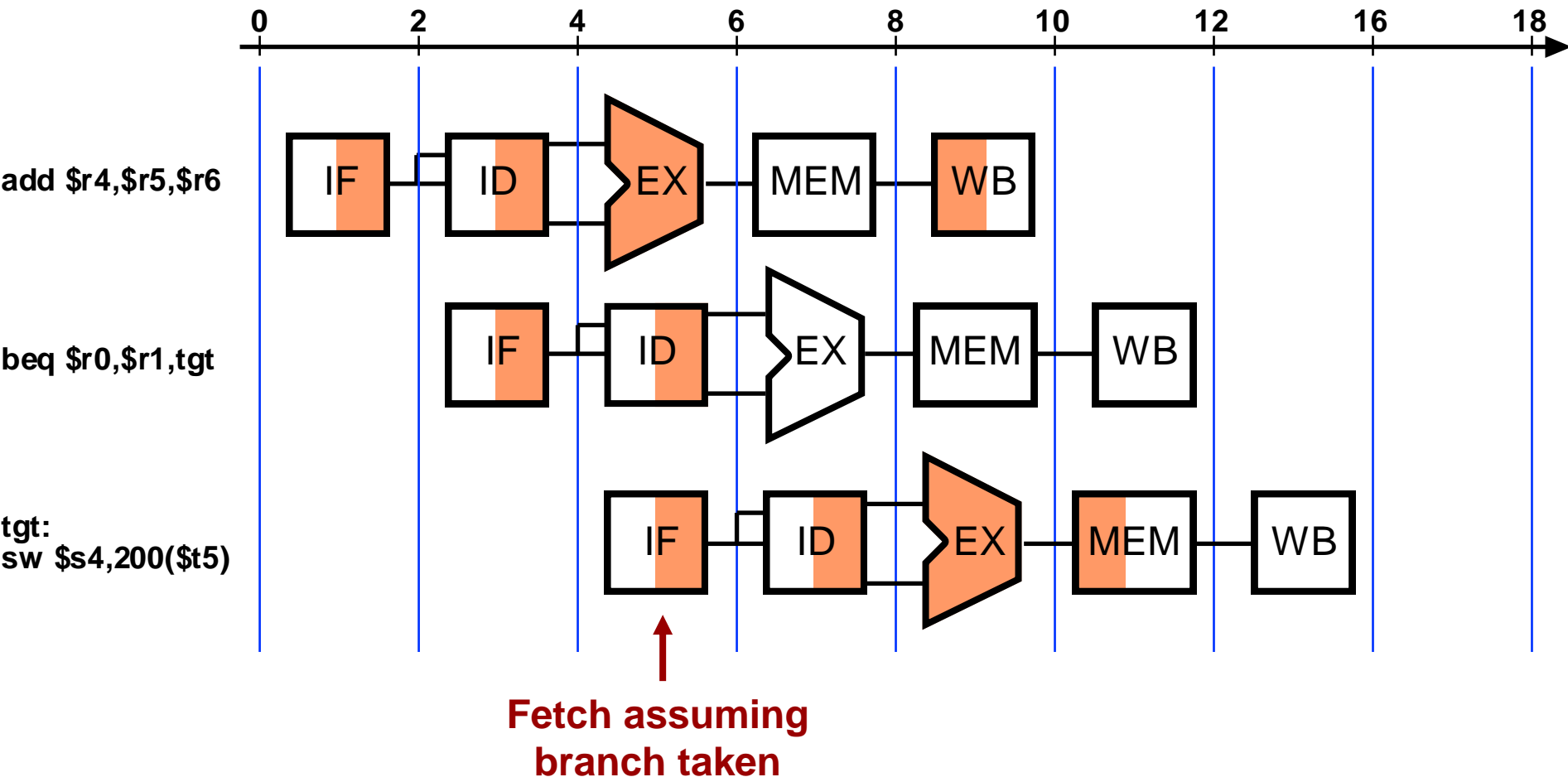
▸ **Stall** - stop loading instructions until result is available

▸ **Predict** - assume an outcome and continue fetching (undo if prediction is wrong)

▸ **Delayed branch** - specify in architecture that following instruction is always executed

# Control Hazard - Stall

| | 0 | 2 | 4 | 6 | 8 | 10 | 12 | 16 | 18 |
|---|---|---|---|---|---|---|---|---|---|

**add $r4,$r5,$r6**  
IF  ID  EX  MEM  WB

**beq $r0,$r1,tgt**  
IF  ID  EX  MEM  WB

**STALL**  
BUBBLE  BUBBLE  BUBBLE  BUBBLE  BUBBLE

**sw $s4,200($t5)**  
IF  ID  EX  MEM  WB

**beq writes PC here**

**new PC used here**

# Control Hazard - Correct Prediction



add $r4,$r5,$r6

beq $r0,$r1,tgt

tgt:
sw $s4,200($t5)

Fetch assuming branch taken

# Control Hazard - Incorrect Prediction



add $r4,$r5,$r6

beq $r0,$r1,tgt

tgt:
sw $s4,200($t5)
(incorrect - STALL)

or $r8,$r8,$r9

IF ID EX MEM WB

BUBBLE BUBBLE BUBBLE BUBBLE

"Squashed" instruction

# Control Hazard - Delayed Branch

add $r4, $r5, $r6
and $s6, $s6, $r7
beq $r0, $r1, tgt



add $r4,$r5,$r6

beq $r0,$r1,tgt

Branch SLOT:
and $r6,$r6,$r7

**always executes**

tgt:
sw $s4,200($t5)

**correct PC avail. here**

tgt:
Sw $s4, 200($t5)

# Summary - Control Hazard Solutions

▸ **Stall** - stop fetching instr. until result is available
  - ▸ **Significant performance penalty**
  - ▸ <u>**Hardware required to stall**</u>

▸ **Predict** - assume an outcome and continue fetching (undo if prediction is wrong)
  - ▸ **Performance penalty only when guess is wrong**
  - ▸ <u>**Hardware required to "squash" instructions**</u>

▸ **Delayed branch** - specify in architecture that following instruction is always executed
  - ▸ **Compiler <u>re-orders instructions into delay slot</u>**
  - ▸ **Insert "NOP" (no-op) operations when can't use (~50%)**
    - • **This is how original MIPS worked**
  - ▸ **But, rationale (이론근거) may not hold up as technology changes**
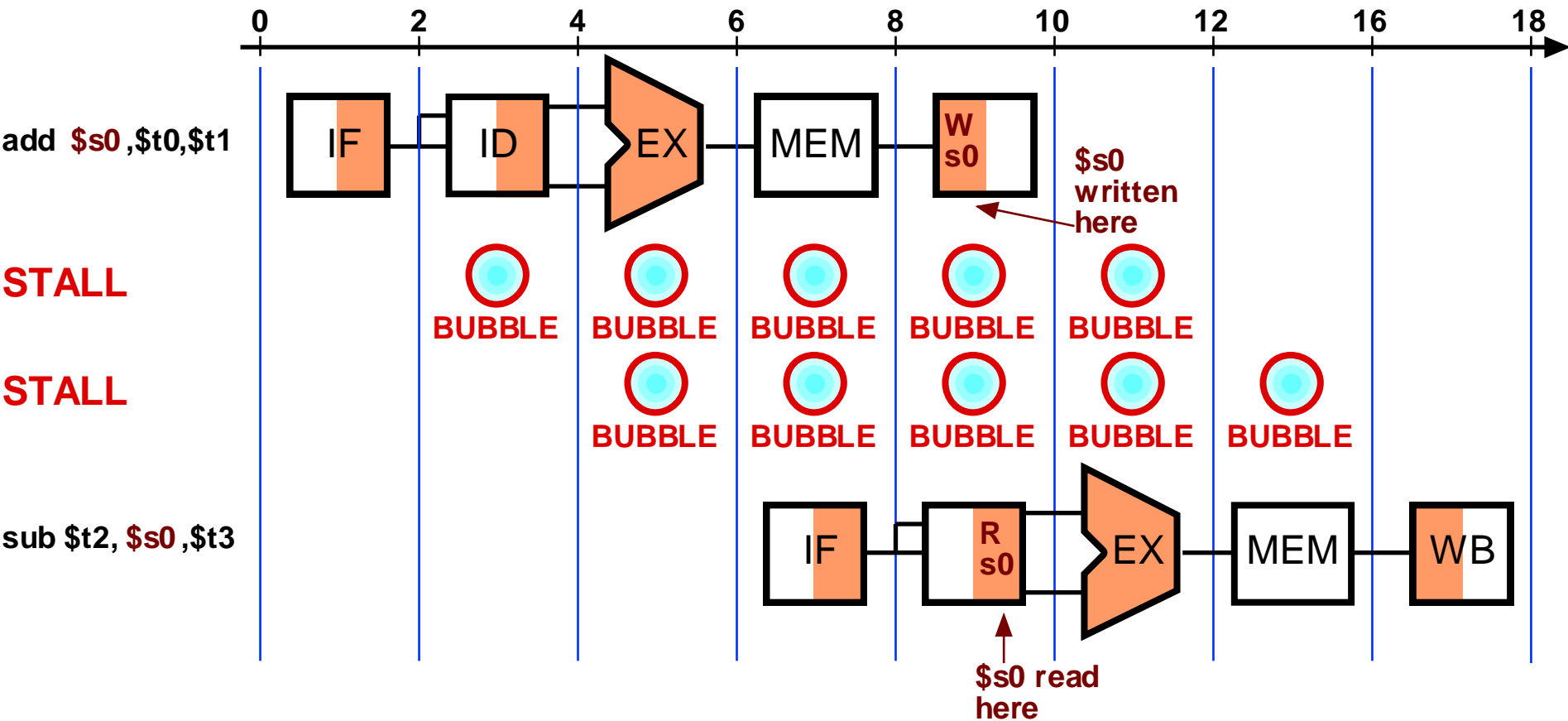
# Data Hazards

▸ **Attempt to use data before it is ready**

▸ **Solutions**

  ▸ **Stalling** - wait until result is available

  ▸ **Forwarding** - make data available inside datapath

  ▸ **Reordering instructions** - use compiler to avoid hazards

▸ **Examples:**

```
add $s0, $t0, $t1 ; $s0 = $t0+$t1
sub $t2, $s0, $t3 ; $t2 = $s0-$t2


lw $s0, 0($t0)     ; $s0 = MEM[$t0]
sub $t2, $s0, $t3 ; $t2 = $s0-$t2
```
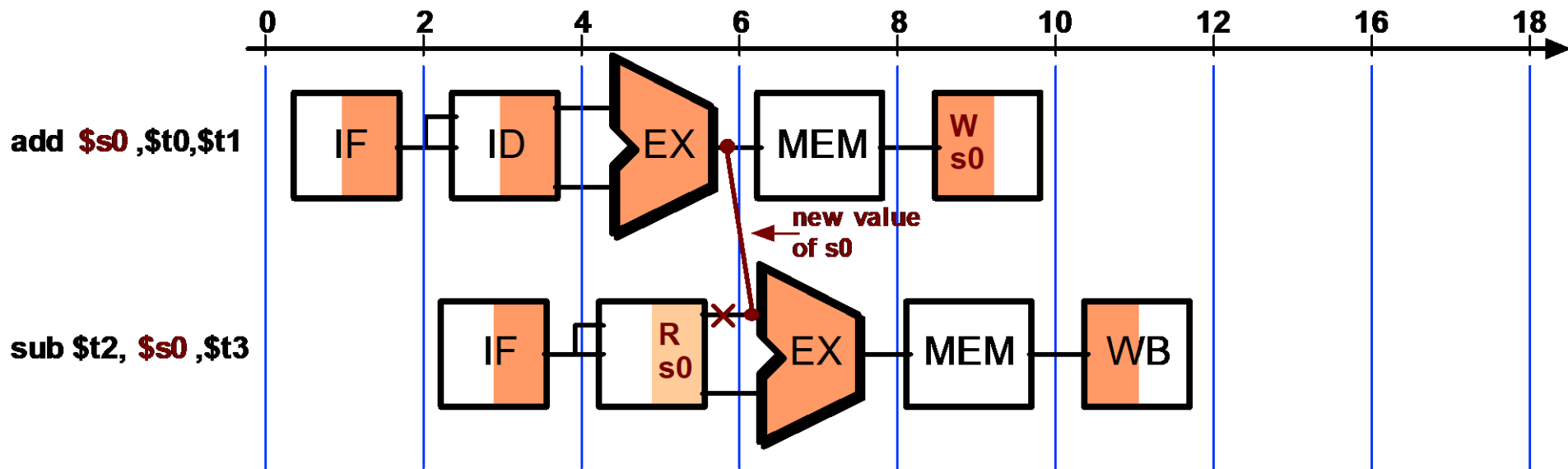
# Data Hazard - Stalling
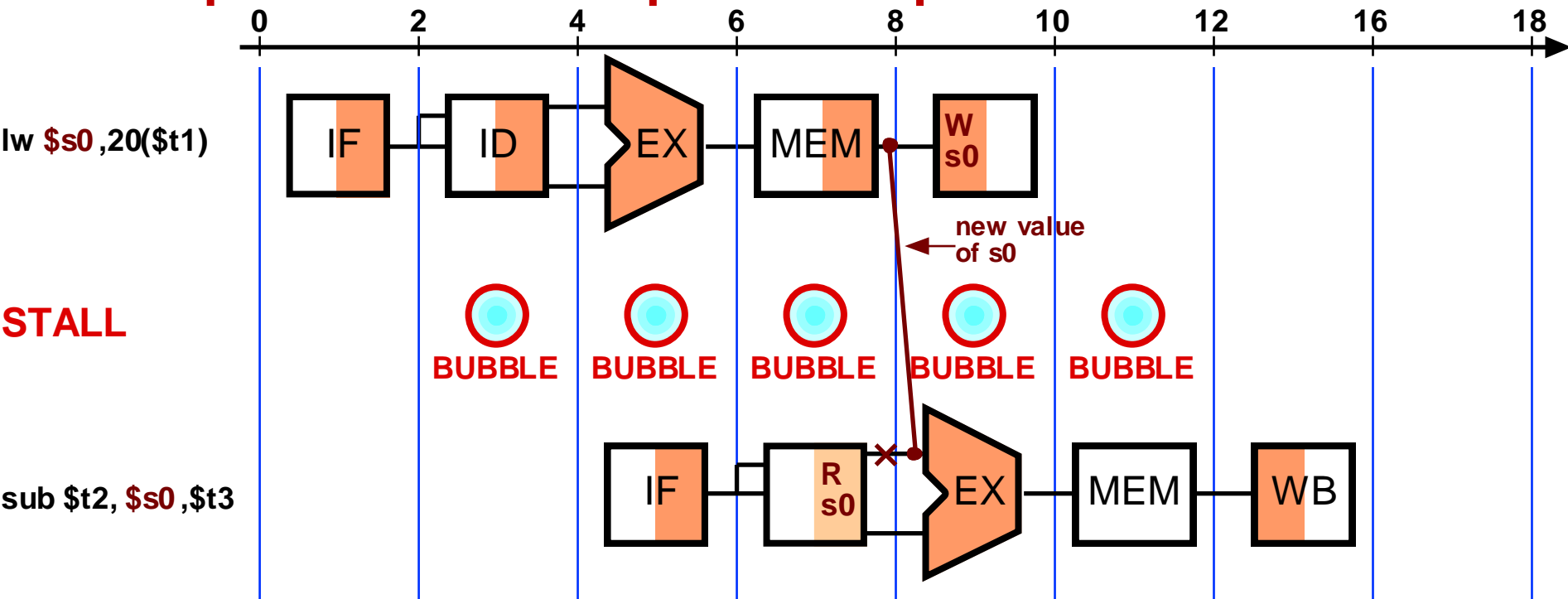
# Data Hazards - **Forwarding**

- **Key idea: connect new value directly to next stage**
- **Still read s0, but ignore in favor of new result**



- **Problem: what about load instructions?**

# Data Hazards - Forwarding

▸ **STALL <u>still</u> required for load - data avail. after MEM**

▸ **MIPS architecture calls this <u>delayed load</u>, initial implementations required compiler to deal with this**
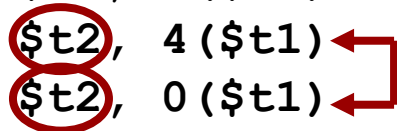
# Data Hazards - Reordering Instructions

▸ **Assuming we have data forwarding, what are the hazards in this code?**

```
lw $t0, 0($t1)
lw $t2, 4($t1)
sw $t2, 0($t1)
sw $t0, 4($t1)
```

▸ **Reorder instructions to remove hazard:**

```
lw $t0, 0($t1)
lw $t2, 4($t1)
sw $t0, 4($t1)
sw $t2, 0($t1)
```

# Summary - Pipelining Overview

▸ **Pipelining increase <u>throughput</u> (but not latency)**

▸ **Hazards limit performance**

  ▸ **Structural hazards**

  ▸ **Control hazards**

  ▸ **Data hazards**

# Pipelining Outline - Coming Up

▸ **Introduction**

  ▸ **Defining Pipelining**

  ▸ **Pipelining Instructions**

  ▸ **Hazards**

▸ **Pipelined Processor Design** ◂

  ▸ **Datapath**

  ▸ **Control**

▸ **Advanced Pipelining**

  ▸ **Superscalar**

  ▸ **Dynamic Pipelining**

  ▸ **Examples**