

# COMP4611: Design and Analysis of Computer Architectures

## Basic Pipelining

Lin Gu  
CSE, HKUST

1

### Introduction to Pipelining – Boeing Everett



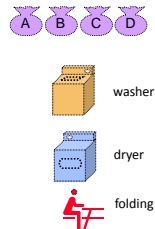
Courtesy photos by Boeing

### Introduction to Pipelining

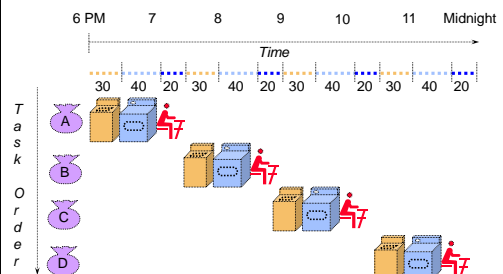
- Pipelining: An implementation technique that overlaps the execution of multiple instructions. It is a **key** technique in achieving high-performance

#### Laundry Example

- Ann, Brian, Cathy, Dave each has one load of clothes to wash, dry, and fold
- Washing takes 30 minutes
- Drying takes 40 minutes
- "Folding" takes 20 minutes

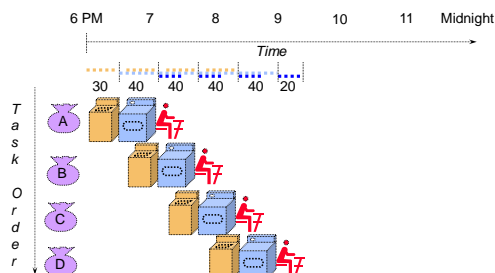


### Sequential Laundry



- Sequential laundry takes 6 hours for 4 loads

### Pipelined Laundry: start work ASAP



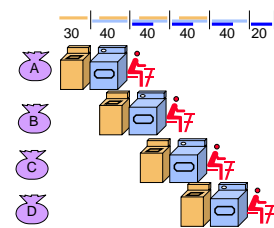
- Pipelined laundry takes 3.5 hours for 4 loads
- Speedup =  $6/3.5 = 1.7$

### Pipelining Lessons

#### Latency vs. Throughput

- Question
  - What is the latency in both cases?
  - What is the throughput in both cases?

- Pipelining doesn't reduce **latency** of single task,
- It increases **throughput** of entire workload

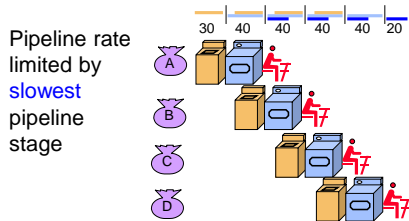


6

## Pipelining Lessons [cont'd...]

### Question

- What is the fastest operation in the example ?
- What is the slowest operation in the example



7

## Pipelining Lessons [cont'd...]

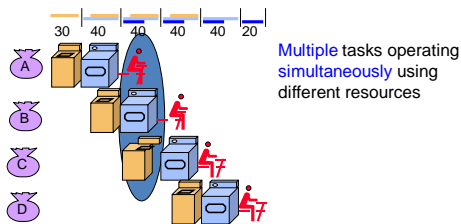
- Washing takes 30 minutes
- Drying takes 40 minutes
- "Folding" takes 20 minutes
- Question
  - Would it matter if "folding" also took 40 minutes



Unbalanced lengths of pipe stages reduce speedup

8

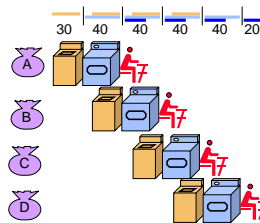
## Pipelining Lessons [cont'd...]



9

## Pipelining Lessons [contd...]

- Question
  - Would the speedup increase if we had more steps ?

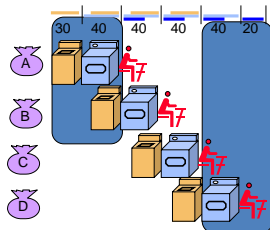


Potential Speedup = Number of pipe stages

More steps = Deeper pipeline

10

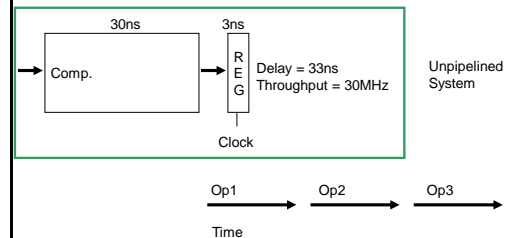
## Pipelining Lessons [cont'd...]



Time to "fill" pipeline and time to "drain" it reduce speedup

11

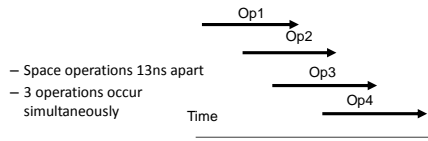
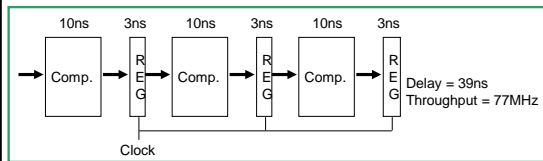
## Pipelining in Computation



- One operation must complete before next can begin
- Operations spaced 33ns apart

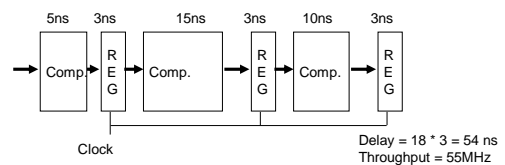
12

### 3 Stage Pipelining



13

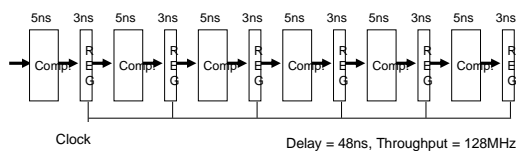
### Limitation: Non-uniform Pipelining



- Throughput limited by slowest stage
- Delay determined by clock period \* number of stages
- Must attempt to balance stages

14

### Limitation: Deep Pipelines



- Diminishing returns as add more pipeline stages
- Register delays become limiting factor
  - Increased latency
  - Small throughput gains
  - More hazards

15

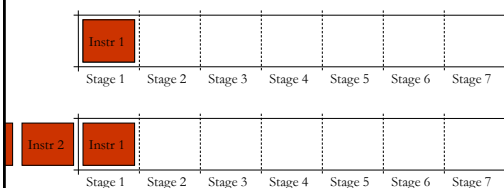
### Computer (Processor) Pipelining

- It is one **KEY** method of achieving High-Performance in modern microprocessors
- A major advantage of pipelining over "parallel processing" is that it is **not visible** to the programmer
- An **instruction** execution pipeline involves a number of steps, where each step completes a part of an instruction.
  - Each step is called a **pipe stage** or a **pipe segment**.

16

### Pipelining

- Multiple instructions overlapped in execution
- Throughput increased
- Doesn't reduce time for individual instructions



17

### Computer Pipelining

- The stages or steps are connected one to the next to form a pipe -- instructions enter at one end and progress through the stage and exit at the other end.
- **Throughput** of an instruction pipeline is determined by how often an instruction exits the pipeline.
- An instruction moves one step down the pipeline after every time interval **C**, where **C** equals to **the cycle time or machine cycle** ( $1/\text{Clock Rate}$ ) and is determined by the stage with the longest processing delay (slowest pipeline stage).

18

### Pipelining: Design Goals

- 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):

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

19

### Pipelining: Design Goals

- Under these ideal conditions:
  - Speedup from pipelining equals the number of pipeline stages:  $n$
  - One instruction is completed every cycle,  $CPI = 1$ .
  - This is an asymptote of course, but +10% is commonly achieved
  - Difference is due to difficulty in achieving balanced stage design
- Two ways to view the performance mechanism
  - Reduced  $CPI$  (i.e. non-piped to piped change)
    - Close to 1 instruction/cycle if you're lucky
  - Reduced cycle-time (i.e. increasing pipeline depth)
    - Work split into more stages
    - Simpler stages result in faster clock cycles

20

### Implementation of MIPS

- We use the MIPS processor as an example to demonstrate the concepts of computer pipelining.
- MIPS ISA design reflects careful measurements and informed architectural decisions.
- The design of its pipeline needs to answer the following questions

How are instructions executed?

How to pipeline them?

21

### How to Execute an Instruction?

*"Vēnī, vīdī, vīcī"*

– "I came, I saw, I conquered"

Bring the instruction to the processor – "Fetch"

Examine the instruction – "Decode"

Do the work – "Execute, Memory-Access, Writeback"



IF, ID, EX, MEM, WB

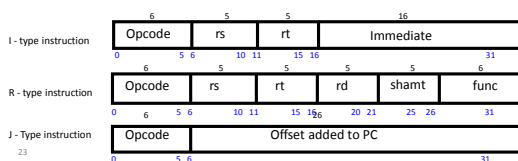
22

### A Basic Multi-Cycle Implementation of MIPS

#### 1 Instruction fetch cycle (IF):

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

Note: IR (instruction register), NPC (next sequential program counter register)



23

### A Basic Multi-Cycle Implementation of MIPS

#### 2 Instruction decode/register fetch cycle (ID):

Parse instruction

Read:

$A \leftarrow \text{Regs}[rs];$

$B \leftarrow \text{Regs}[rt];$

$\text{Imm} \leftarrow ((IR_{16})^{16} \# \# IR_{16..31})$  sign-extended immediate field of IR

Note: A, B, Imm are temporary registers



24

### 3 Execution/Effective address cycle (EX):

- **Memory reference:**

```
ALUOutput ← A + Imm;
```

- **Register-Register ALU instruction:**

```
ALUOutput ← A op B;
```

- **Register-Immediate ALU instruction:**

ALUOutput  $\leftarrow$  A op Imm;

- Branch (simplification – only consider BEQZ):

```
ALUOutput ← NPC + Imm;  
cond ← (A == 0)
```

25

## A Basic Implementation of MIPS (continued)

#### 4 Memory access/branch completion cycle (MEM):

- **Memory reference:**

LMD  $\leftarrow$  Mem[ALUOutput] or  
Mem[ALUOutput]  $\leftarrow$  B;

- **Branch:**

```
if (cond) PC ← ALUOutput else PC ← NPC
```

**Note: LMD (load memory data) register**

26

## A Basic Implementation of MIPS (continued)

### 5 Write-back cycle (WB):

- Register-Register ALU instruction:

```
Regs[rd] ← ALUOutput;
```

- Register-Immediate ALU instruction:

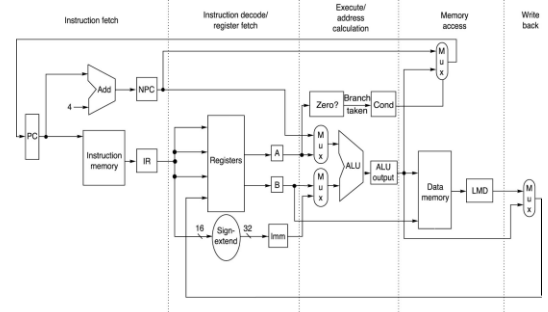
```
Regs[rt] ← ALUOutput;
```

- Load instruction:

$$\text{Regs}[\text{rt}] \leftarrow \text{LMD};$$

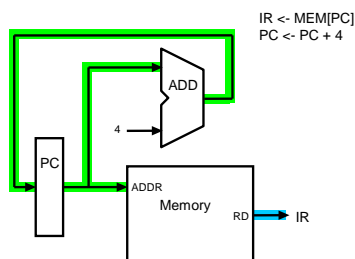
Note: LMD (load memory data) register

## Basic MIPS Integer Datapath Implementation



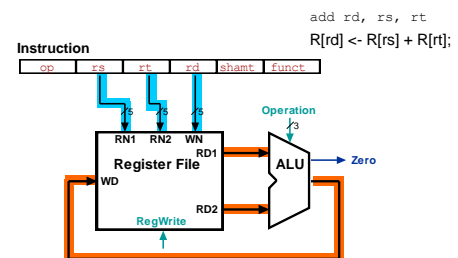
28

## Datapath for Instruction Fetch



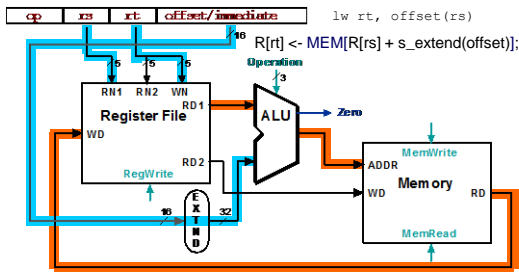
29

### Datapath for R-Type Instructions



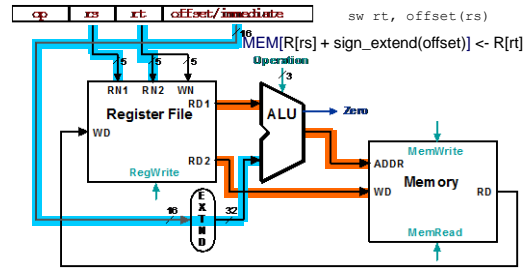
30

## Datapath for Load/Store Instructions



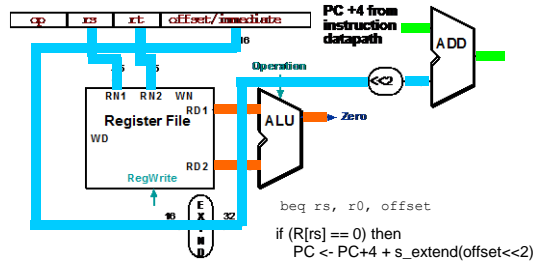
31

## Datapath for Load/Store Instructions



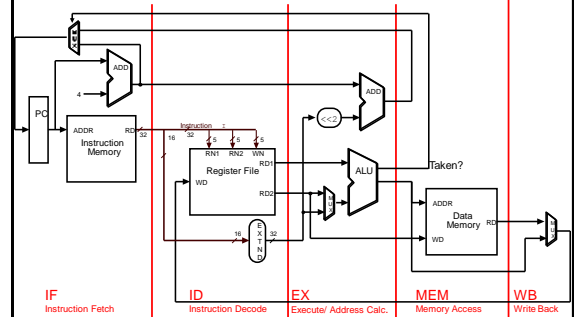
32

## Datapath for Branch Instructions



33

## Basic MIPS Processor (More Details)



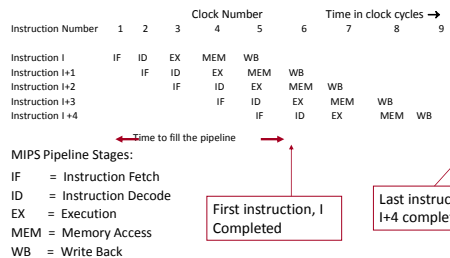
34

## Pipelining - Key Idea

- Question:** What happens if we break execution into multiple cycles?
- Answer:** in the best case, we can start executing a new instruction on each clock cycle - this is pipelining
- Pipelining stages:**
  - IF - Instruction Fetch
  - ID - Instruction Decode
  - EX - Execute / Address Calculation
  - MEM - Memory Access (read / write)
  - WB - Write Back (results into register file)

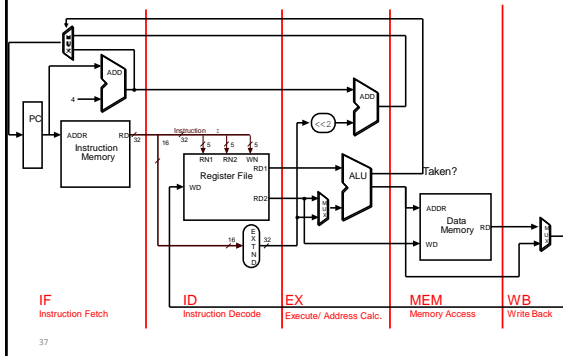
35

## Simple MIPS Pipelined Integer Instruction Processing



36

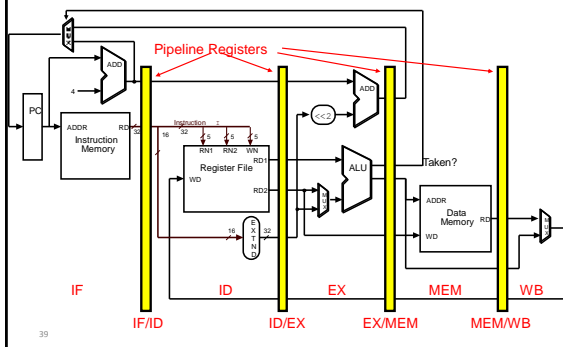
### Basic MIPS Processor (More Details)



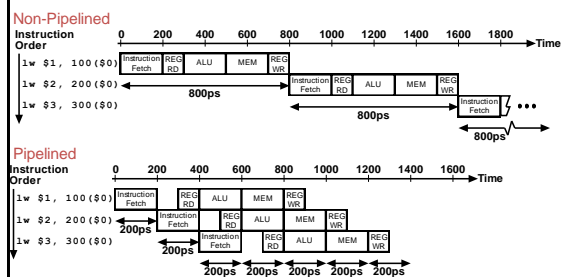
### Pipeline Registers

- Need pipeline registers (latches) between stages to hold data for instructions
- Pipeline registers are named with 2 stages (the stages that the register is "between.")
- ANY information needed in a later pipeline stage **MUST** be passed via a pipeline register
  - Example: IF/ID register gets
    - instruction
    - PC+4
- No register is needed after WB. Results from the WB stage are already stored in the register file.

### Basic Pipelined Processor



### Non-Pipelined vs. Pipelined Execution

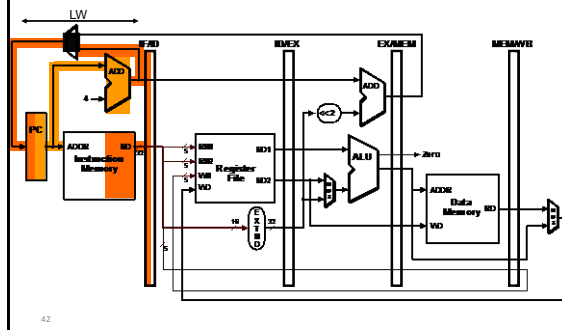


### Pipelined Example - Executing Multiple Instructions

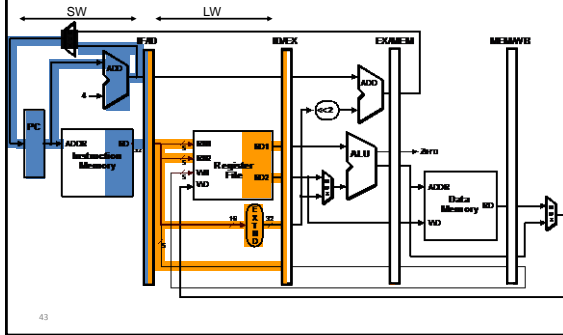
- Consider the following instruction sequence:

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

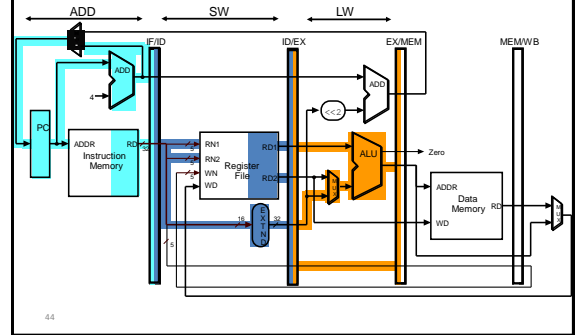
### Executing Multiple Instructions Clock Cycle 1



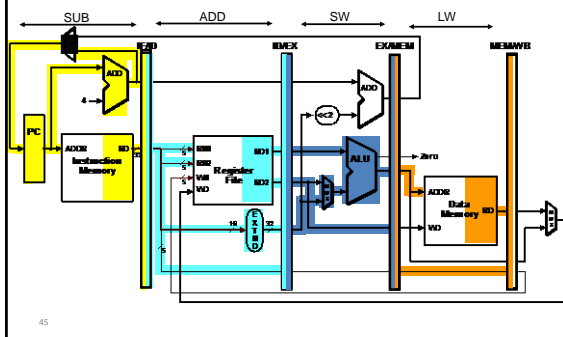
Executing Multiple Instructions  
Clock Cycle 2



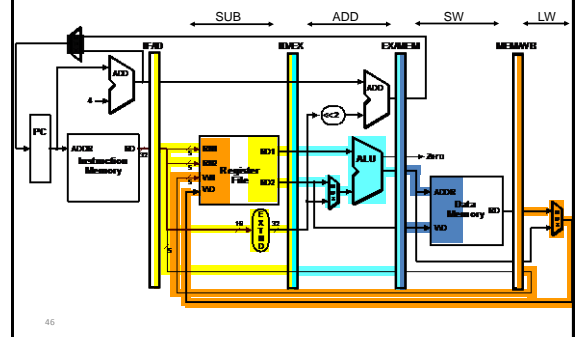
Executing Multiple Instructions  
Clock Cycle 3



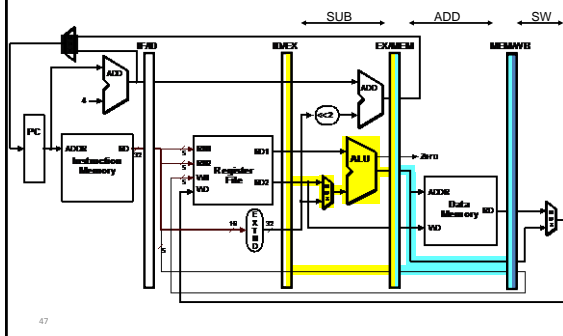
Executing Multiple Instructions  
Clock Cycle 4



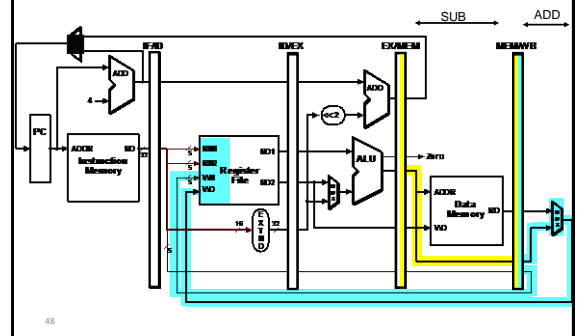
Executing Multiple Instructions  
Clock Cycle 5



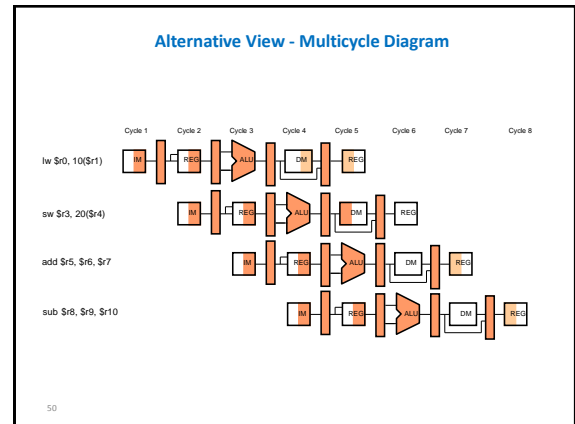
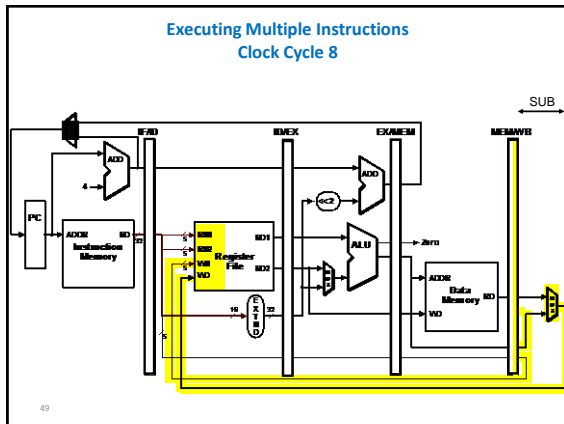
Executing Multiple Instructions  
Clock Cycle 6



Executing Multiple Instructions  
Clock Cycle 7







## Pipeline Performance

- How much does pipelining improve the performance? How to calculate throughput and latency of a pipelined processor?
- How to design the pipeline to increase speedup?
- ...

## Throughput, latency, hazards

## Pipelining Performance Example

- Example: For an unpipelined CPU:
  - Clock cycle = 1ns, 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 0.2 ns to the machine clock cycle then the speedup in instruction execution from pipelining is:

Unpipelined Average instruction execution time = Clock cycle x Average CPI  
 $= 1 \text{ ns} \times ((40\% + 20\%) \times 4 + 40\% \times 5) = 1 \text{ ns} \times 4.4 = 4.4 \text{ ns}$

In the 5-stage pipelined implementation, the latency is 6ns, but the average instruction execution time is:  $1 \text{ ns} + 0.2 \text{ ns} = 1.2 \text{ ns}$

Speedup from pipelining =  $\frac{\text{Instruction time unpipelined}}{\text{Instruction time pipelined}}$   
 $= 4.4 \text{ ns} / 1.2 \text{ ns} = 3.7 \text{ times faster}$

## Pipeline Throughput and Latency: A More realistic Examples



Consider the pipeline above with the indicated delays. We want to know what the *pipeline throughput* and the *pipeline latency* are.

Pipeline throughput: instructions completed per second.

Pipeline latency: how long does it take to execute a single instruction in the pipeline.

## Pipeline Throughput and Latency

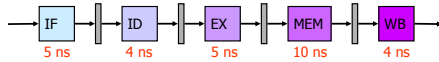


Pipeline latency: how long does it take to execute an instruction in the pipeline.

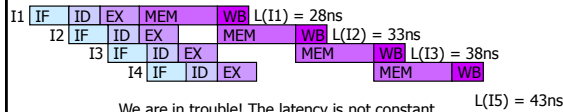
$$L = \text{lat}(\text{IF}) + \text{lat}(\text{ID}) + \text{lat}(\text{EX}) + \text{lat}(\text{MEM}) + \text{lat}(\text{WB})$$

$$= 5\text{ns} + 4\text{ns} + 5\text{ns} + 10\text{ns} + 4\text{ns} = 28\text{ns}$$

### Pipeline Throughput and Latency

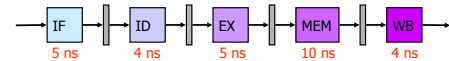


Simply adding the latencies to compute the pipeline latency only would work for an isolated instruction

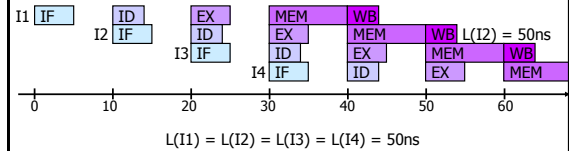


We are in trouble! The latency is not constant. This happens because this is an unbalanced pipeline. The solution is to make every state the same length as the longest one.

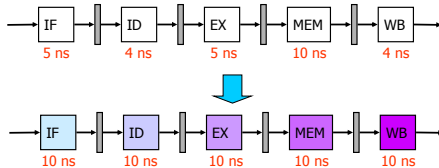
### Synchronous Pipeline Throughput and Latency



The slowest pipeline stage also limits the latency!!



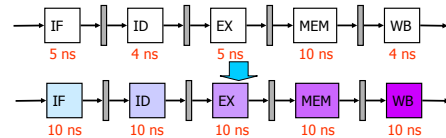
### Pipeline Throughput and Latency



Pipeline throughput: how often an instruction is completed.

$$\begin{aligned}
 &= \text{instr} / \max[\text{lat}(\text{IF}), \text{lat}(\text{ID}), \text{lat}(\text{EX}), \text{lat}(\text{MEM}), \text{lat}(\text{WB})] \\
 &= \text{instr} / \max[5\text{ns}, 4\text{ns}, 5\text{ns}, 10\text{ns}, 4\text{ns}] \\
 &= \text{instr} / 10\text{ns} \quad (\text{ignoring pipeline register overhead})
 \end{aligned}$$

### Pipeline Throughput and Latency



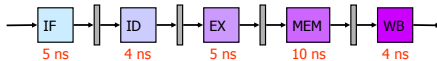
How long does it take to execute (issue) 20000 instructions in this pipeline? (disregard latency, bubbles caused by branches, cache misses, hazards, fill, drain)

$$\text{ExecTime}_{\text{pipe}} = 20000 \times 10\text{ns} = 200000\text{ns} = 200\mu\text{s}$$

How long would it take using the same modules without pipelining?

$$\text{ExecTime}_{\text{non-pipe}} = 20000 \times 28\text{ns} = 560000\text{ns} = 560\mu\text{s}$$

### Pipeline Throughput and Latency



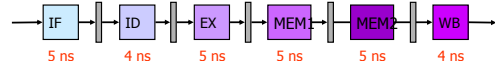
Thus the speedup that we got from the pipeline is:

$$\text{Speedup}_{\text{pipe}} = \frac{\text{ExecTime}_{\text{non-pipe}}}{\text{ExecTime}_{\text{pipe}}} = \frac{560\mu\text{s}}{200\mu\text{s}} = 2.8$$

How can we improve this pipeline design?

We need to reduce the unbalance to increase the clock speed.

### Pipeline Throughput and Latency



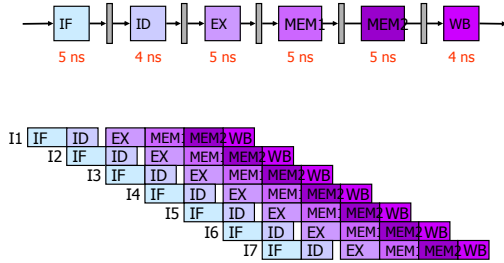
Now we have one more pipeline stage, but the maximum latency of a single stage is reduced in half.

$$\begin{aligned}
 T &= \text{instr} / \max[\text{lat}(\text{IF}), \text{lat}(\text{ID}), \text{lat}(\text{EX}), \text{lat}(\text{MEM}1), \text{lat}(\text{MEM}2), \text{lat}(\text{WB})] \\
 &= \text{instr} / \max[5\text{ns}, 4\text{ns}, 5\text{ns}, 5\text{ns}, 5\text{ns}, 4\text{ns}] \\
 &= \text{instr} / 5\text{ns}
 \end{aligned}$$

The new latency for a single instruction is:

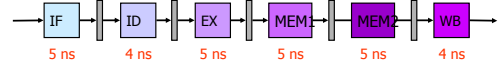
$$L = 6 \times 5\text{ns} = 30\text{ns}$$

### Pipeline Throughput and Latency



61

### Pipeline Throughput and Latency



How long does it take to execute 20000 instructions in this pipeline? (disregard bubbles caused by branches, cache misses, etc, for now)

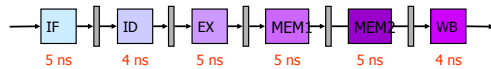
$$ExecTime_{pipe} = 20000 \times 5ns = 100000ns = 100\mu s$$

Thus the speedup that we get from the pipeline is:

$$Speedup_{pipe} = \frac{ExecTime_{non-pipe}}{ExecTime_{pipe}} = \frac{560\mu s}{100\mu s} = 5.6$$

62

### Pipeline Throughput and Latency



What have we learned from this example?

1. It is important to balance the delays in the stages of the pipeline
2. The throughput of a pipeline is  $1/\max(\text{stage\_delay})$ .
3. The latency is  $N \times \max(\text{stage\_delay})$ , where N is the number of stages in the pipeline.

63