

The Big Picture



Requirements

Algorithms

Prog. Lang./OS

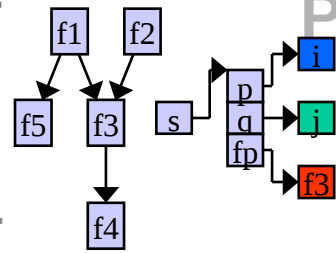
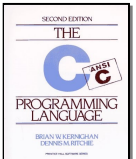
ISA

uArch

Circuit

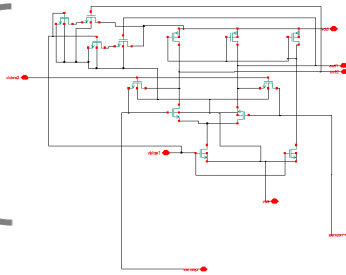
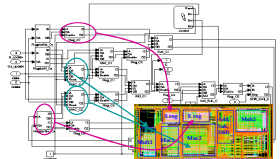
Device

```
f2() {
  f3(s2, &j, &i);
  *s2->p = 10;
  i = *s2->q + i;
}
```

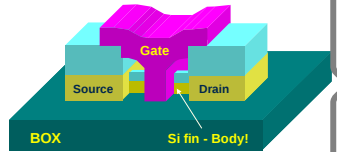


Problem Focus

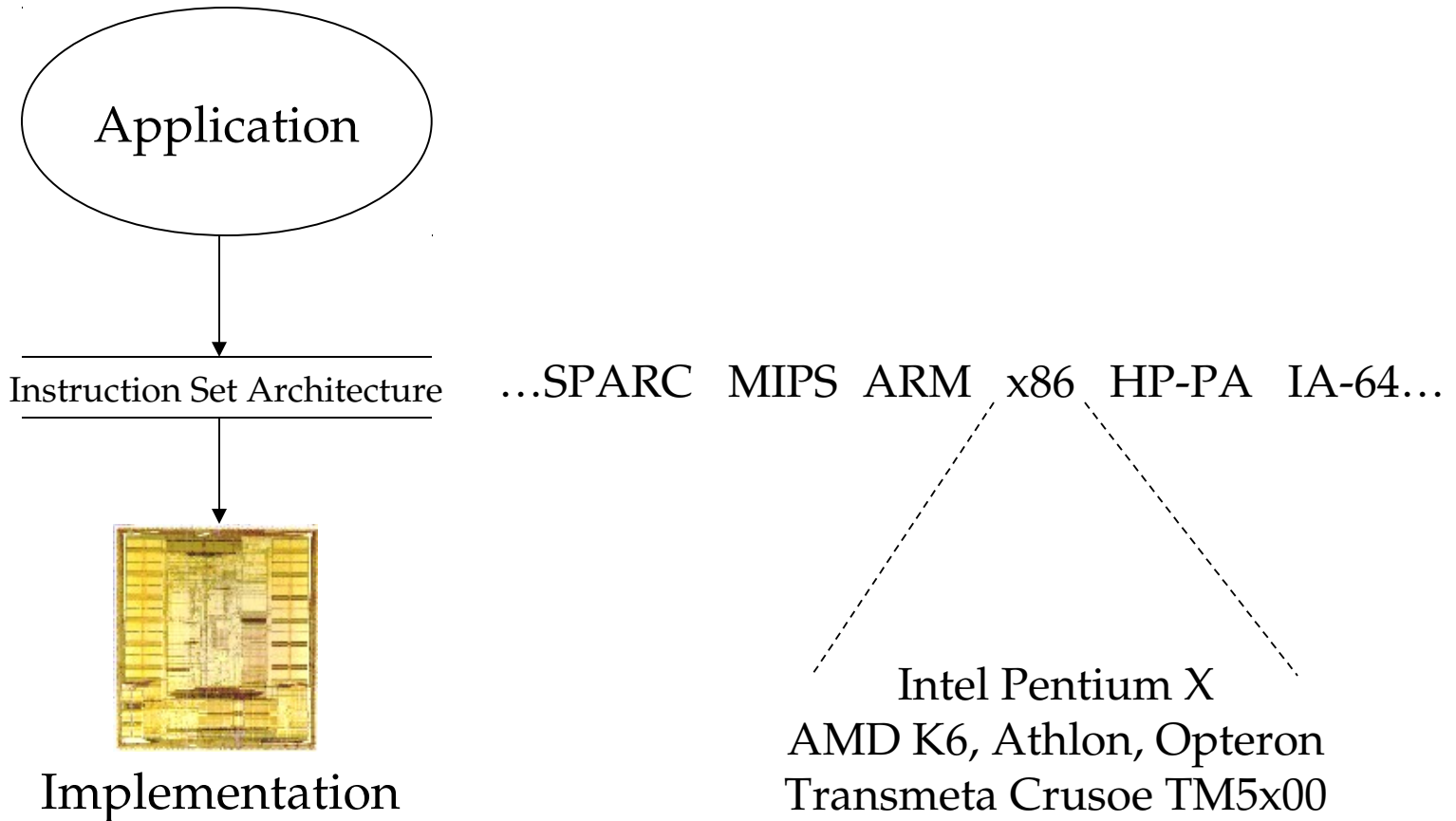
```
i1: ld r1, b      <p1>
i2: ld r2, c      <p1>
i3: ld r5, z      <p3>
i4: mul r6, r5, 3  <p3>
i5: add r3, r1, r2 <p1>
```



Performance Focus



Instruction Set Architecture



Instruction Set Architecture

- **Strong influence on cost/performance**
- **New ISAs are rare, but versions are not**
 - 16-bit, 32-bit and 64-bit X86 versions
- **Longevity is a strong function of marketing prowess**

Traditional Issues

- **Strongly constrained by the number of bits available to instruction encoding**
- **Opcodes/operands**
- **Registers/memory**
- **Addressing modes**
- **Orthogonality**
- **0, 1, 2, 3 address machines**
- **Instruction formats**
- **Decoding uniformity**

Introduction

A.1 What is Pipelining?

A.2 The Major Hurdle of Pipelining-Structural Hazards

- Data Hazards
- Control Hazards

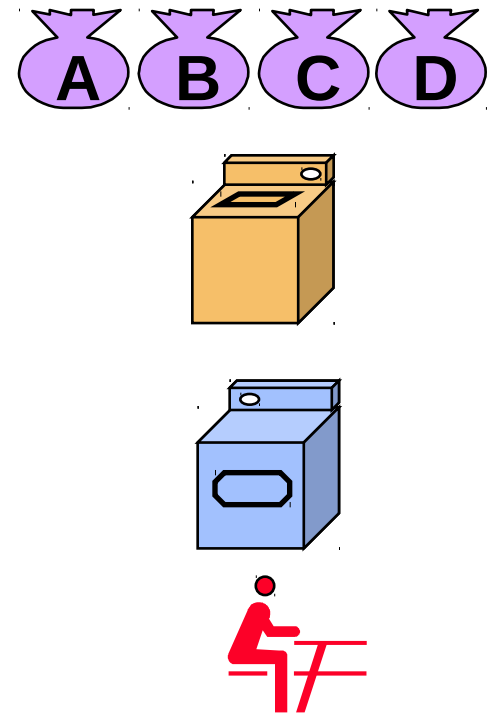
A.3 How is Pipelining Implemented

A.4 What Makes Pipelining Hard to Implement?

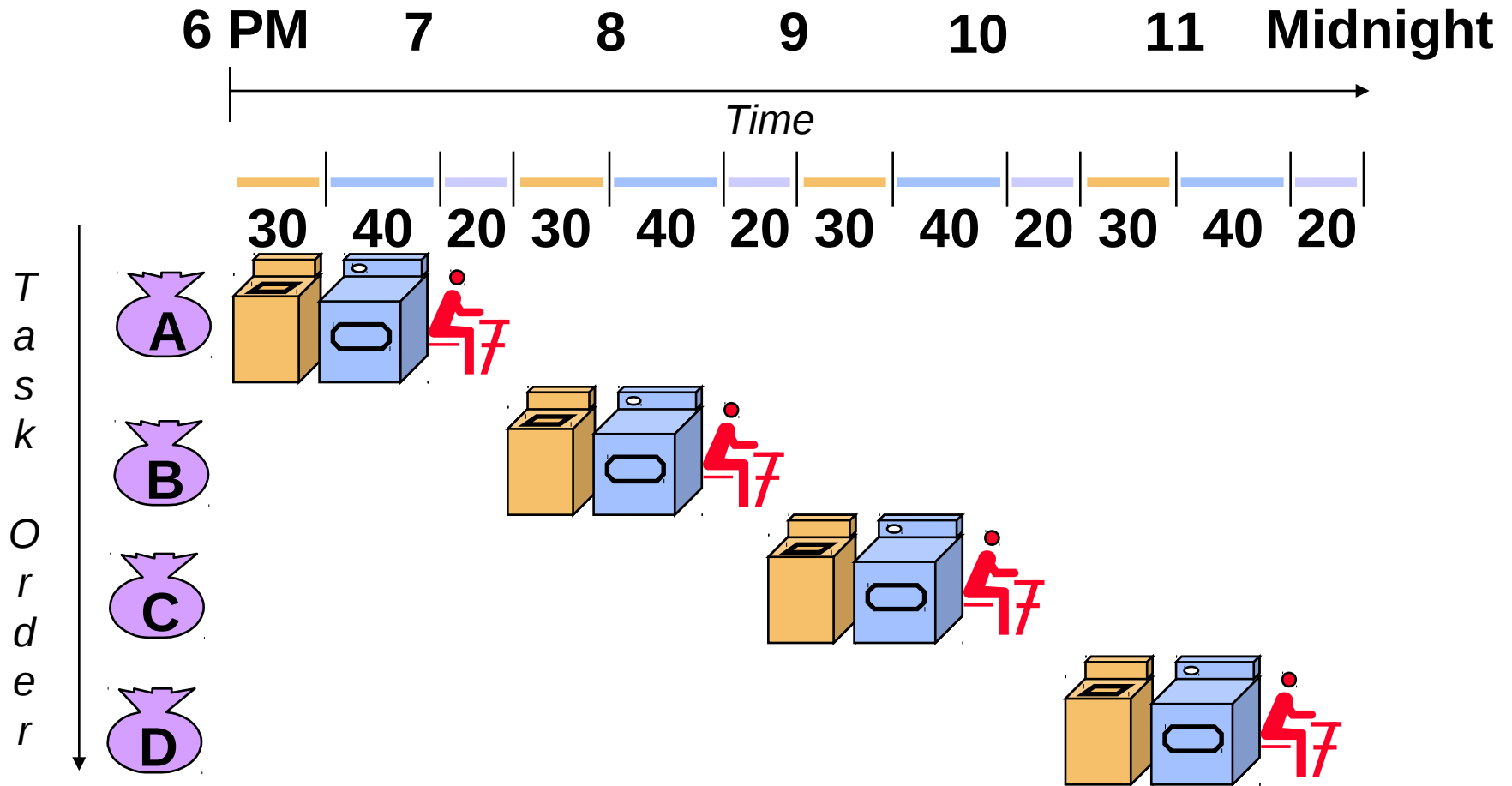
A.5 Extending the MIPS Pipeline to Handle Multi-cycle Operations

What Is Pipelining

- **Laundry Example**
- **Ann, Brian, Cathy, Dave** each have one load of clothes to wash, dry, and fold
- **Washer takes 30 minutes**
- **Dryer takes 40 minutes**
- **“Folder” takes 20 minutes**



What Is Pipelining

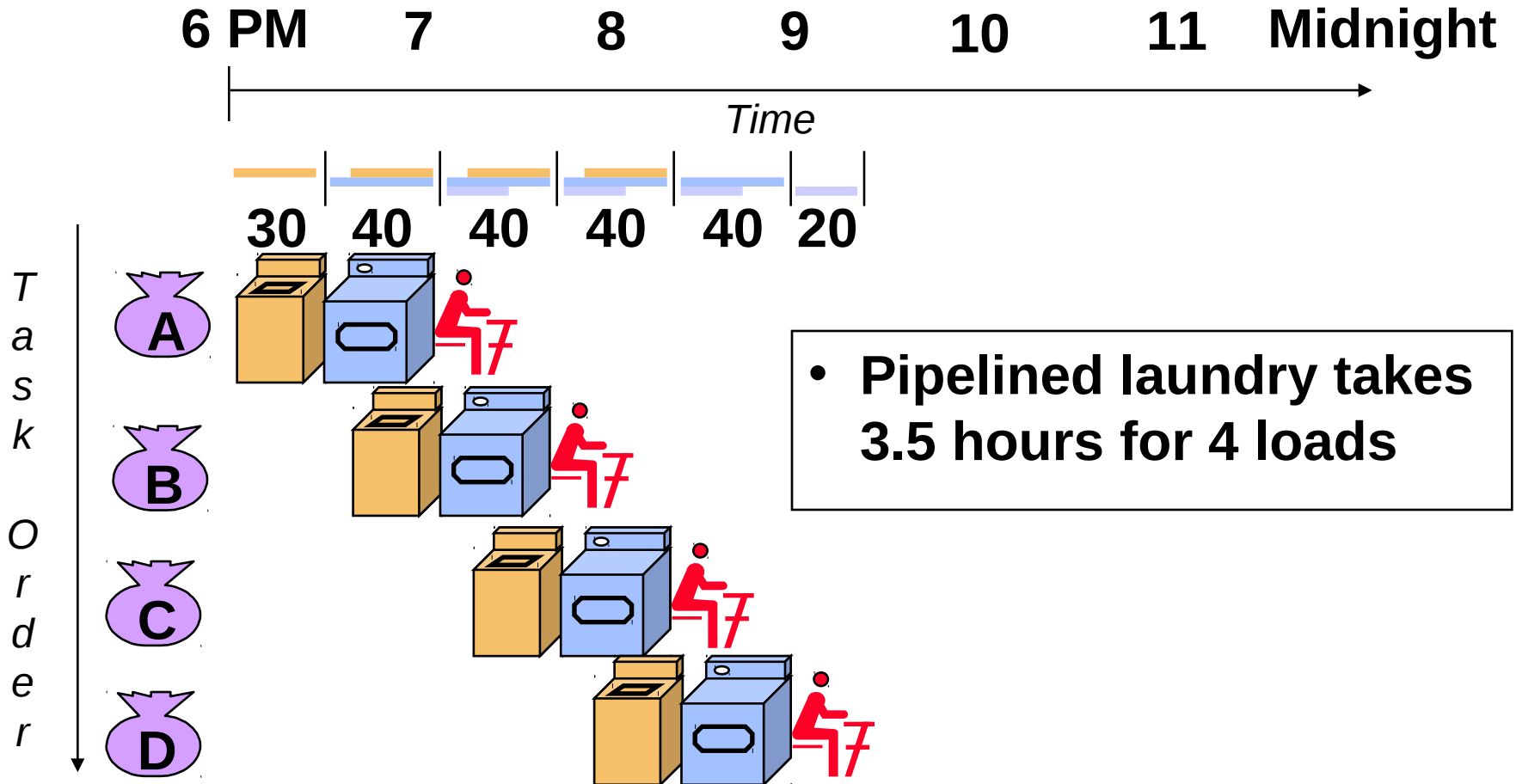


Sequential laundry takes 6 hours for 4 loads

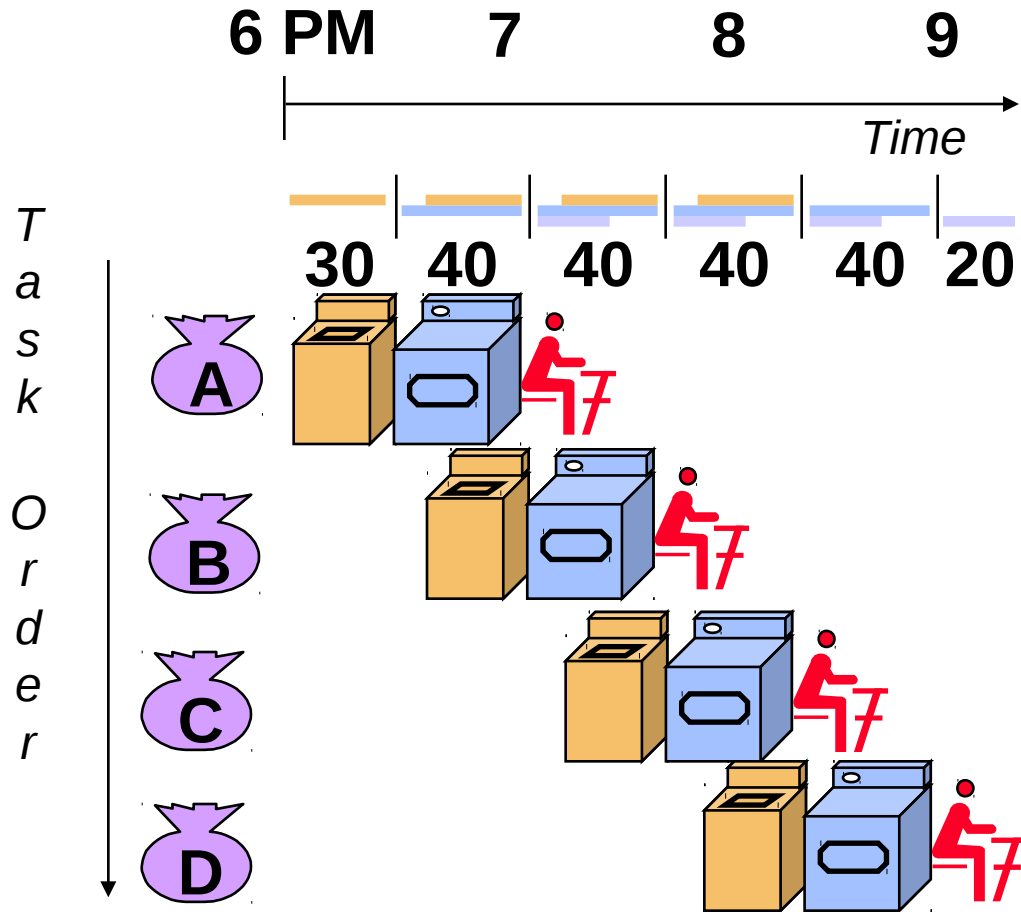
If they learned pipelining, how long would laundry take?

What Is Pipelining

Start work ASAP



What Is Pipelining

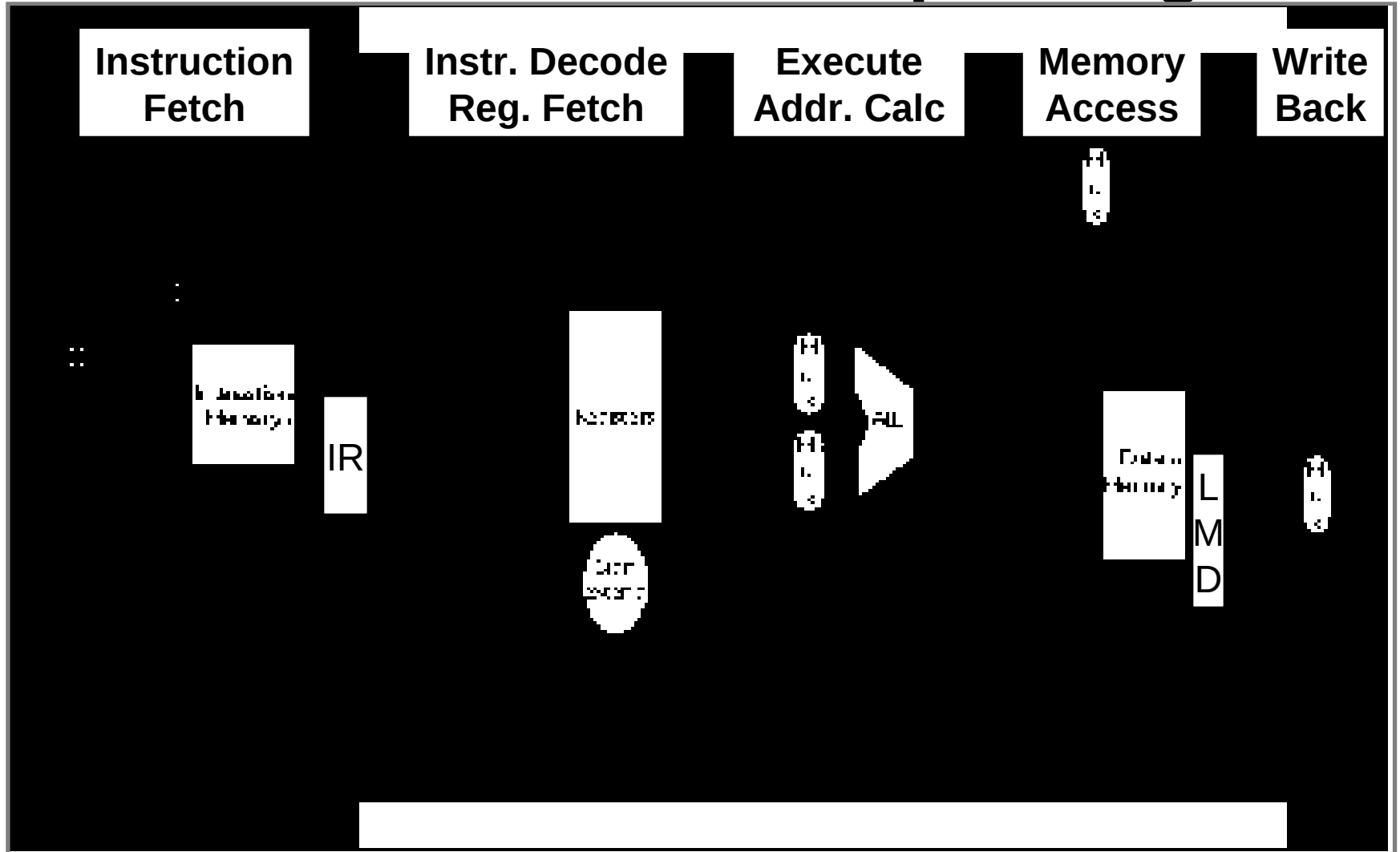


Pipelining Lessons

- 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
- Potential speedup = **Number pipe stages**
- Unbalanced lengths of pipe stages reduces speedup
- Time to “**fill**” pipeline and time to “**drain**” it reduces speedup

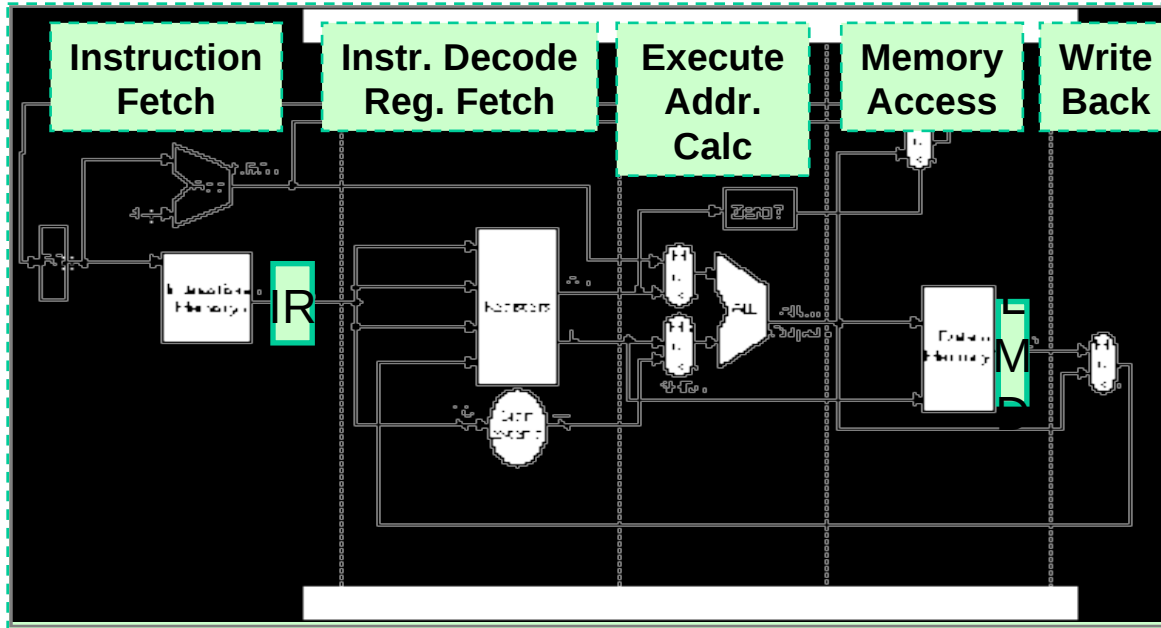
What Is Pipelining

MIPS Without Pipelining



What Is Pipelining

MIPS Functions



Passed To Next Stage
 $IR \leftarrow Mem[PC]$
 $NPC \leftarrow PC + 4$

Instruction Fetch (IF):

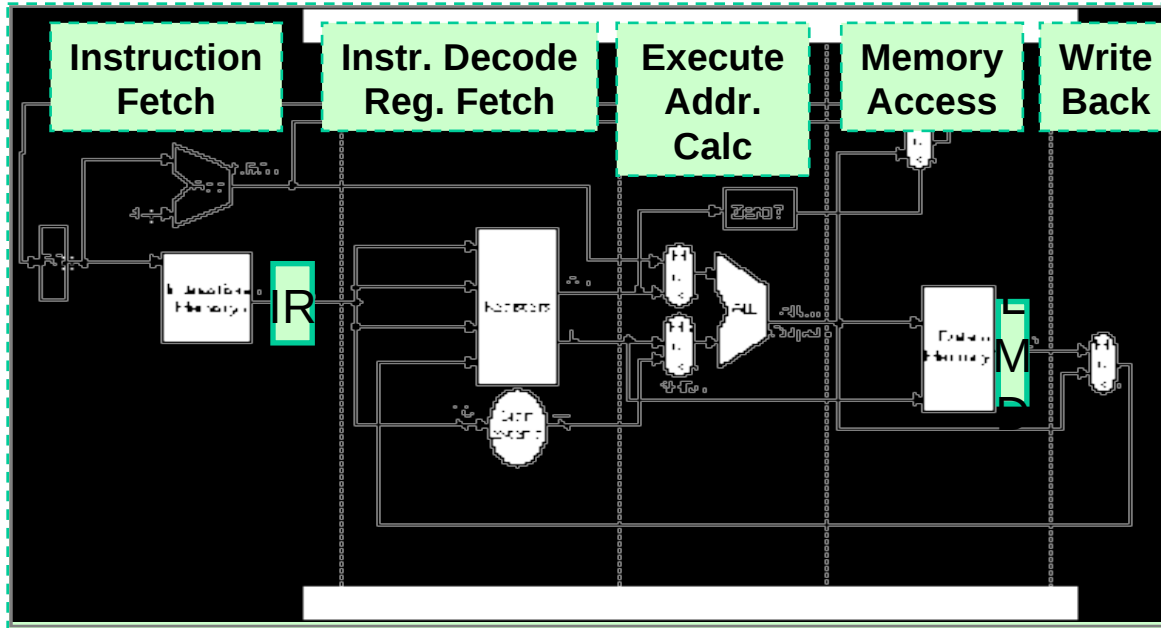
Send out the PC and fetch the instruction from memory into the instruction register (IR); increment the PC by 4 to address the next sequential instruction.

IR holds the instruction that will be used in the next stage.

NPC holds the value of the next PC.

What Is Pipelining

MIPS Functions



Passed To Next Stage
 $A \leftarrow \text{Regs}[\text{IR6}..\text{IR10}];$
 $B \leftarrow \text{Regs}[\text{IR10}..\text{IR15}];$
 $\text{Imm} \leftarrow ((\text{IR16}) \# \# \text{IR16-31})$

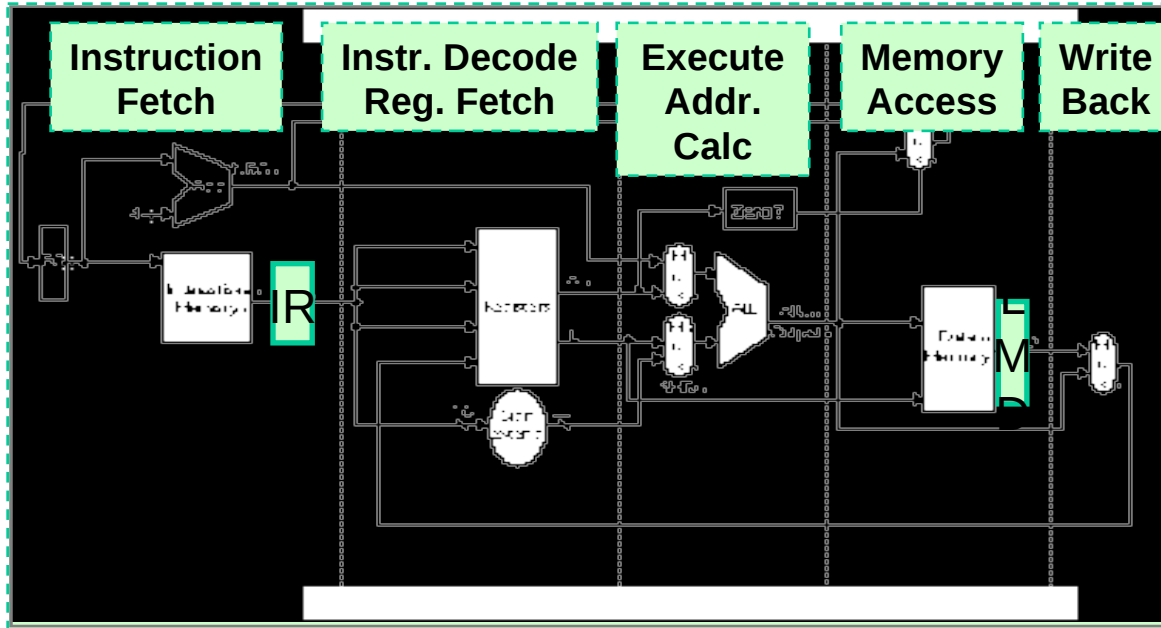
Instruction Decode/Register Fetch Cycle (ID):

Decode the instruction and access the register file to read the registers. The outputs of the general purpose registers are read into two temporary registers (A & B) for use in later clock cycles.

We extend the sign of the lower 16 bits of the Instruction Register.

What Is Pipelining

MIPS Functions



Passed To Next Stage
`A <- A func. B`
`cond = 0;`

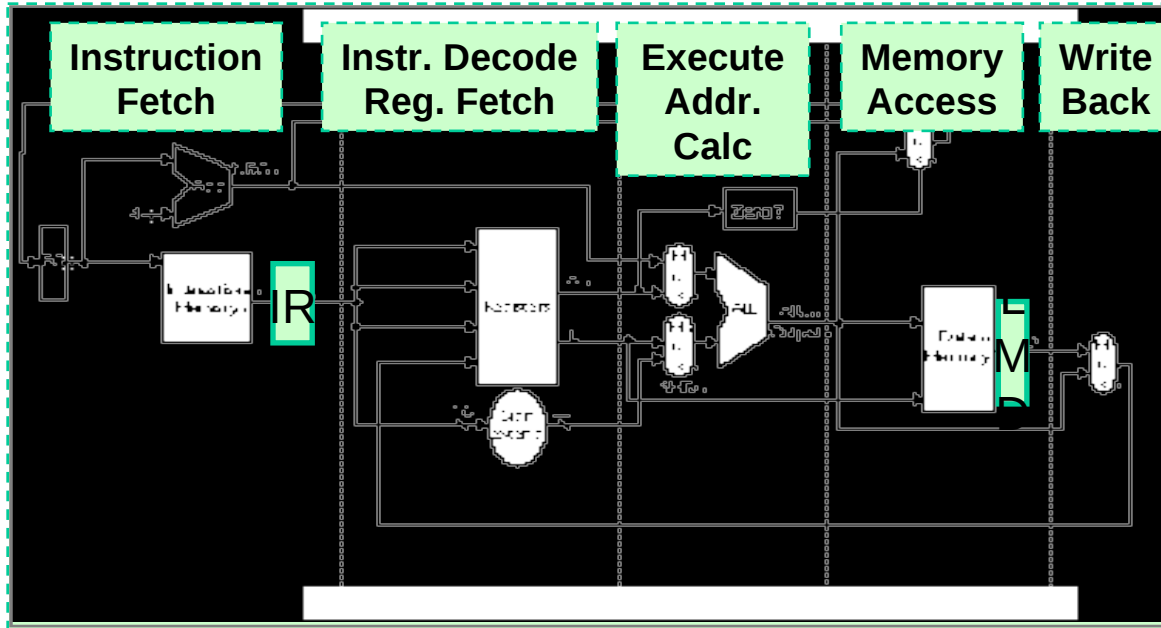
Execute Address Calculation (EX):

We perform an operation (for an ALU) or an address calculation (if it's a load or a Branch).

If an ALU, actually do the operation. If an address calculation, figure out how to obtain the address and stash away the location of that address for the next cycle.

What Is Pipelining

MIPS Functions



Passed To Next Stage
 $A = \text{Mem}[\text{prev. } B]$
or
 $\text{Mem}[\text{prev. } B] = A$

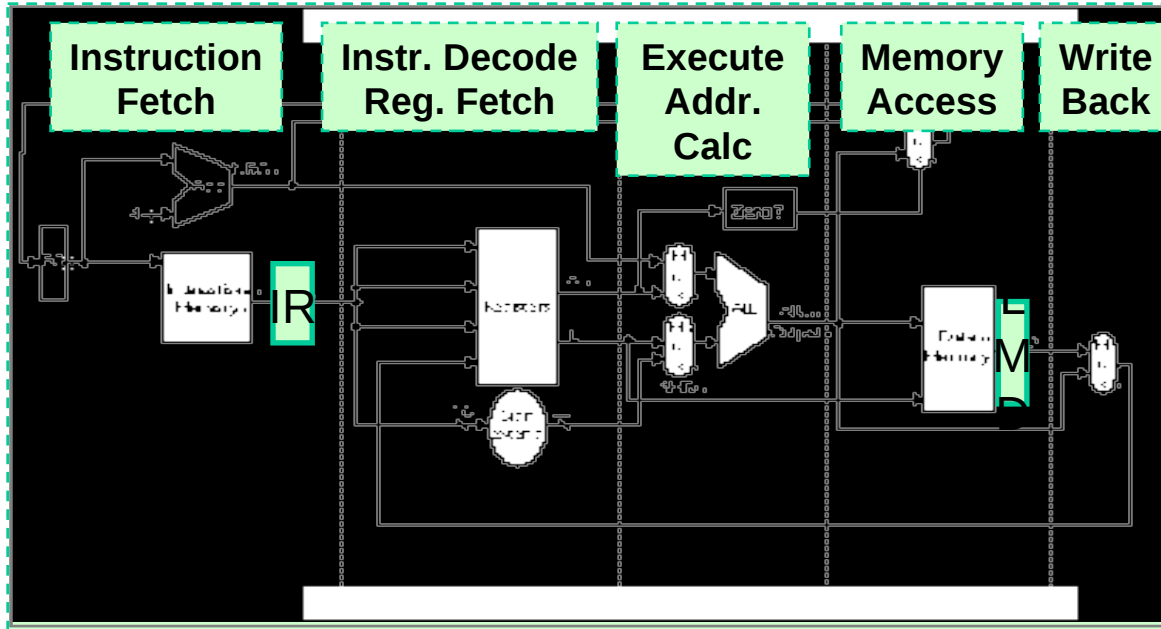
MEMORY ACCESS (MEM):

If this is an ALU, do nothing.

If a load or store, then access memory.

What Is Pipelining

MIPS Functions

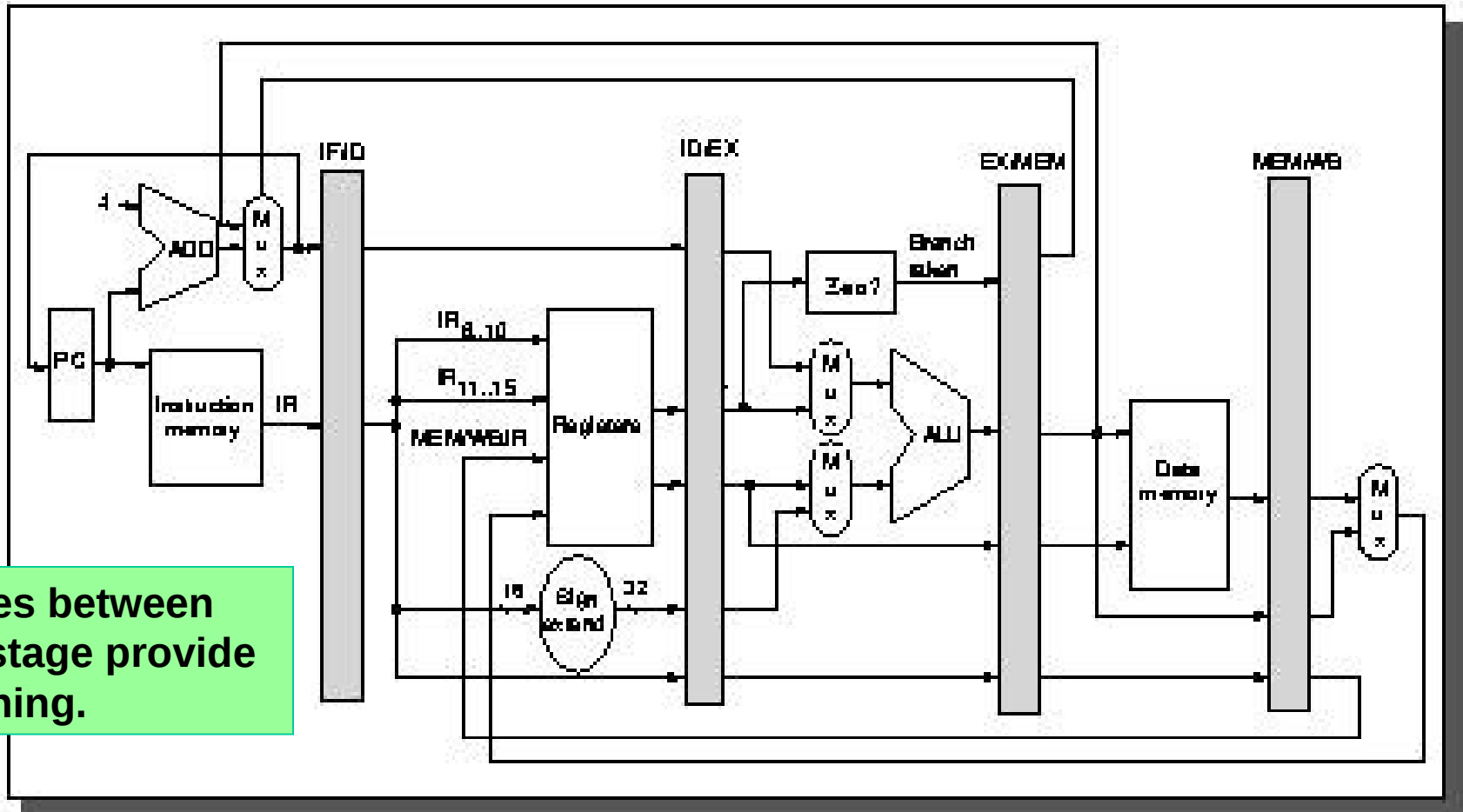


Passed To Next Stage
Regs \leftarrow A, B;

WRITE BACK (WB):

Update the registers from either the ALU or from the data loaded.

The Basic Pipeline For MIPS



Latches between each stage provide pipelining.

FIGURE 3.4 The datapath is pipelined by adding a set of registers, one between each pair of pipe stages.

The Basic Pipeline For MIPS

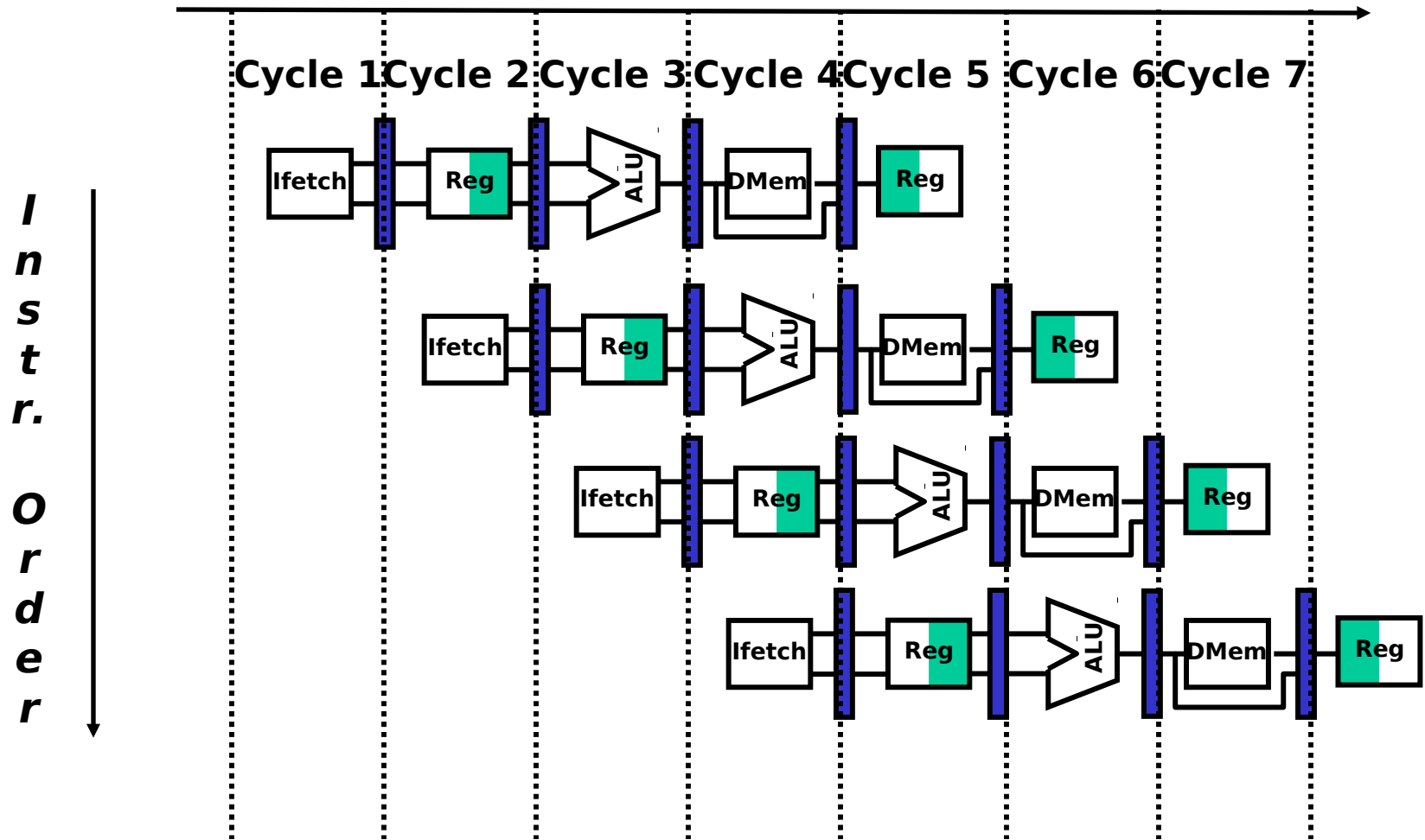


Figure 3.3

Pipeline Hurdles

A.1 What is Pipelining?

A.2 The Major Hurdle of Pipelining- Structural Hazards

- Structural Hazards
- Data Hazards
- Control Hazards

A.3 How is Pipelining Implemented

A.4 What Makes Pipelining Hard to Implement?

A.5 Extending the MIPS Pipeline to Handle Multi-cycle Operations

Limits to pipelining: Hazards prevent next instruction from executing during its designated clock cycle

- **Structural hazards**: HW cannot support this combination of instructions (single person to fold and put clothes away)
- **Data hazards**: Instruction depends on result of prior instruction still in the pipeline (missing sock)
- **Control hazards**: Pipelining of branches & other instructions that change the PC
- Common solution is to **stall** the pipeline until the hazard is resolved, inserting one or more “**bubbles**” in the pipeline

Pipeline Hurdles

Definition

- conditions that lead to incorrect behavior if not fixed
- Structural hazard
 - two different instructions use same h/w in same cycle
- Data hazard
 - two different instructions use same storage
 - must appear as if the instructions execute in correct order
- Control hazard
 - one instruction affects which instruction is next

Resolution

- Pipeline interlock logic detects hazards and fixes them
- simple solution: stall -
- increases CPI, decreases performance
- better solution: partial stall -
- some instruction stall, others proceed better to stall early than late

Structural Hazards

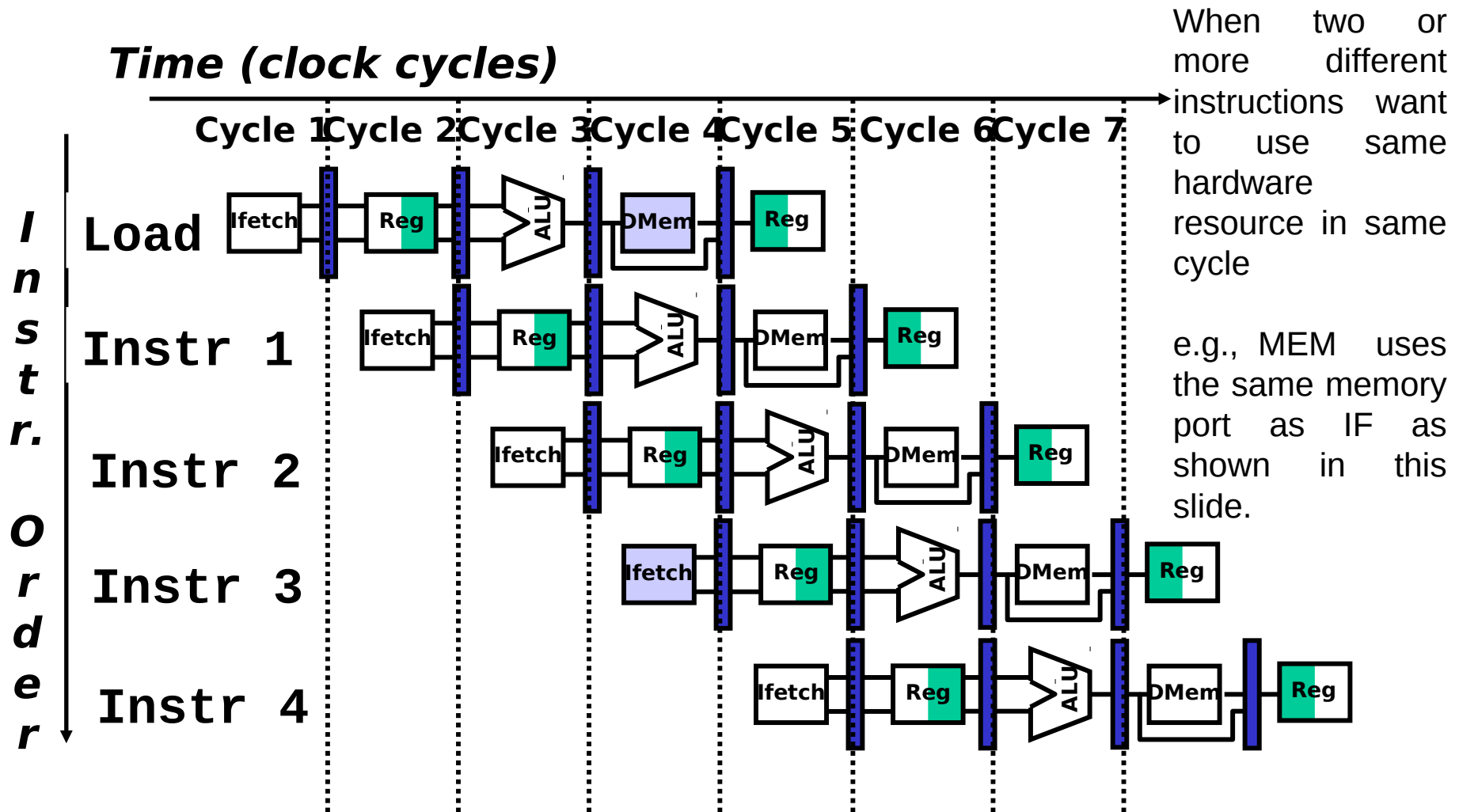


Figure 3.6

Structural Hazards

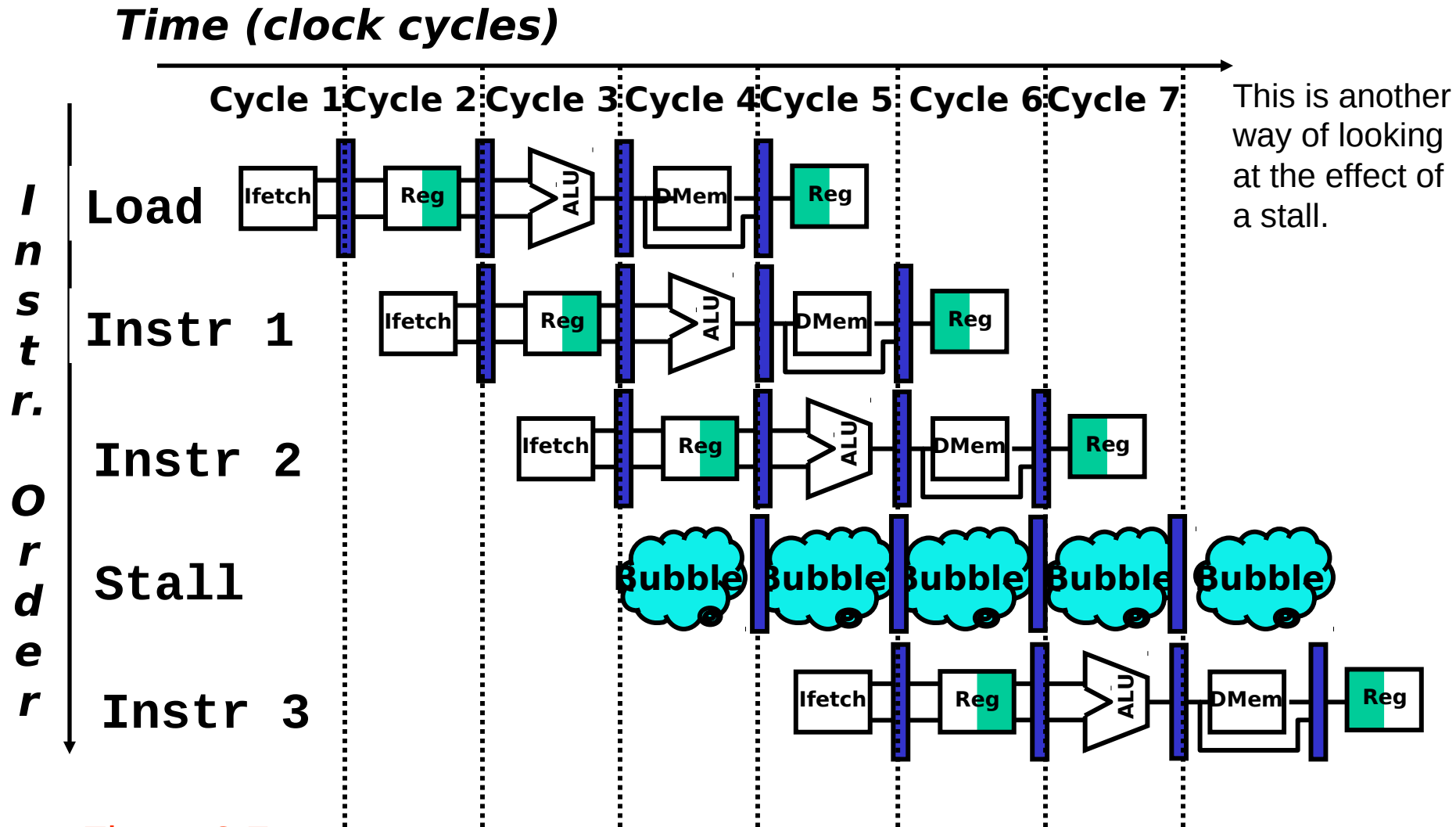


Figure 3.7

Structural Hazards

Instruction	Clock cycle number									
	1	2	3	4	5	6	7	8	9	10
Load instruction	IF	ID	EX	MEM	WB					
Instruction $j + 1$		IF	ID	EX	MEM	WB				
Instruction $j + 2$			IF	ID	EX	MEM	WB			
Instruction $j + 3$				stall	IF	ID	EX	MEM	WB	
Instruction $j + 4$						IF	ID	EX	MEM	WB
Instruction $j + 5$							IF	ID	EX	MEM
Instruction $j + 6$								IF	ID	EX

This is another way to represent the stall we saw on the last few pages.

Structural Hazards

Dealing with Structural Hazards

Stall

- low cost, simple
- Increases CPI
- use for rare case since stalling has performance effect

Pipeline hardware resource

- useful for multi-cycle resources
- good performance
- sometimes complex e.g., RAM

Replicate resource

- good performance
- increases cost (+ maybe interconnect delay)
- useful for cheap or divisible resources

Structural Hazards

Structural hazards are reduced with these rules:

- Each instruction uses a resource at most once
- Always use the resource in the same pipeline stage
- Use the resource for one cycle only

Many RISC ISA's designed with this in mind

Sometimes very complex to do this. For example, memory of necessity is used in the IF and MEM stages.

Some common Structural Hazards:

- Memory - we've already mentioned this one.
- Floating point - Since many floating point instructions require many cycles, it's easy for them to interfere with each other.
- Starting up more of one type of instruction than there are resources. For instance, the PA-8600 can support two ALU + two load/store instructions per cycle - that's how much hardware it has available.

Structural Hazards

This is the example on Page 144.

We want to compare the performance of two machines. Which machine is faster?

- Machine A: Dual ported memory - so there are no memory stalls
- Machine B: Single ported memory, but its pipelined implementation has a 1.05 times faster clock rate

Assume:

- Ideal CPI = 1 for both
- Loads are 40% of instructions executed

$$\begin{aligned}\text{SpeedUp}_A &= \text{Pipeline Depth} / (1 + 0) \times (\text{clock}_{\text{unpipe}} / \text{clock}_{\text{pipe}}) \\ &= \text{Pipeline Depth}\end{aligned}$$

$$\begin{aligned}\text{SpeedUp}_B &= \text{Pipeline Depth} / (1 + 0.4 \times 1) \\ &\quad \times (\text{clock}_{\text{unpipe}} / (\text{clock}_{\text{unpipe}} / 1.05)) \\ &= (\text{Pipeline Depth} / 1.4) \times 1.05 \\ &= 0.75 \times \text{Pipeline Depth}\end{aligned}$$

$$\text{SpeedUp}_A / \text{SpeedUp}_B = \text{Pipeline Depth} / (0.75 \times \text{Pipeline Depth}) = 1.33$$

- Machine A is 1.33 times faster

Data Hazards

A.1 What is Pipelining?

A.2 The Major Hurdle of Pipelining- Structural Hazards

- Structural Hazards
- Data Hazards
- Control Hazards

A.3 How is Pipelining Implemented

A.4 What Makes Pipelining Hard to Implement?

A.5 Extending the MIPS Pipeline to Handle Multi-cycle Operations

These occur when at any time, there are instructions active that need to access the same data (memory or register) locations.

Where there's real trouble is when we have:

**instruction A
instruction B**

and B manipulates (reads or writes) data before A does. This violates the order of the instructions, since the architecture implies that A completes entirely before B is executed.

Data Hazards

Execution Order is:
Instr_i
Instr_j

Read After Write (RAW)

Instr_j tries to read operand before Instr_i writes it

 I: add r1, r2, r3
J: sub r4, r1, r3

- Caused by a “**Dependence**” (in compiler nomenclature). This hazard results from an actual need for communication.

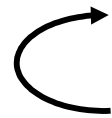
Data Hazards

Execution Order is:
Instr_i
Instr_j

Write After Read (WAR)

Instr_j tries to write operand before Instr_i reads it

- Gets wrong operand

 I: sub r4, r1, r3
J: add r1, r2, r3
K: mul r6, r1, r7

- Called an “anti-dependence” by compiler writers. This results from reuse of the name “r1”.

- **Can't happen in MIPS 5 stage pipeline because:**
 - All instructions take 5 stages, and
 - Reads are always in stage 2, and
 - Writes are always in stage 5

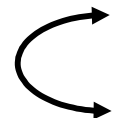
Data Hazards

Execution Order is:
Instr_i
Instr_j

Write After Write (WAW)

Instr_j tries to write operand *before* **Instr_i** writes it

- Leaves wrong result (**Instr_i** not **Instr_j**)



```
I: sub r1, r4, r3  
J: add r1, r2, r3  
K: mul r6, r1, r7
```

- Called an “**output dependence**” by compiler writers
This also results from the reuse of name “**r1**”.
- **Can’t happen in MIPS 5 stage pipeline because:**
 - All instructions take 5 stages, and
 - Writes are always in stage 5
- **Will see WAR and WAW in later more complicated pipes**

Data Hazards

Simple Solution to RAW

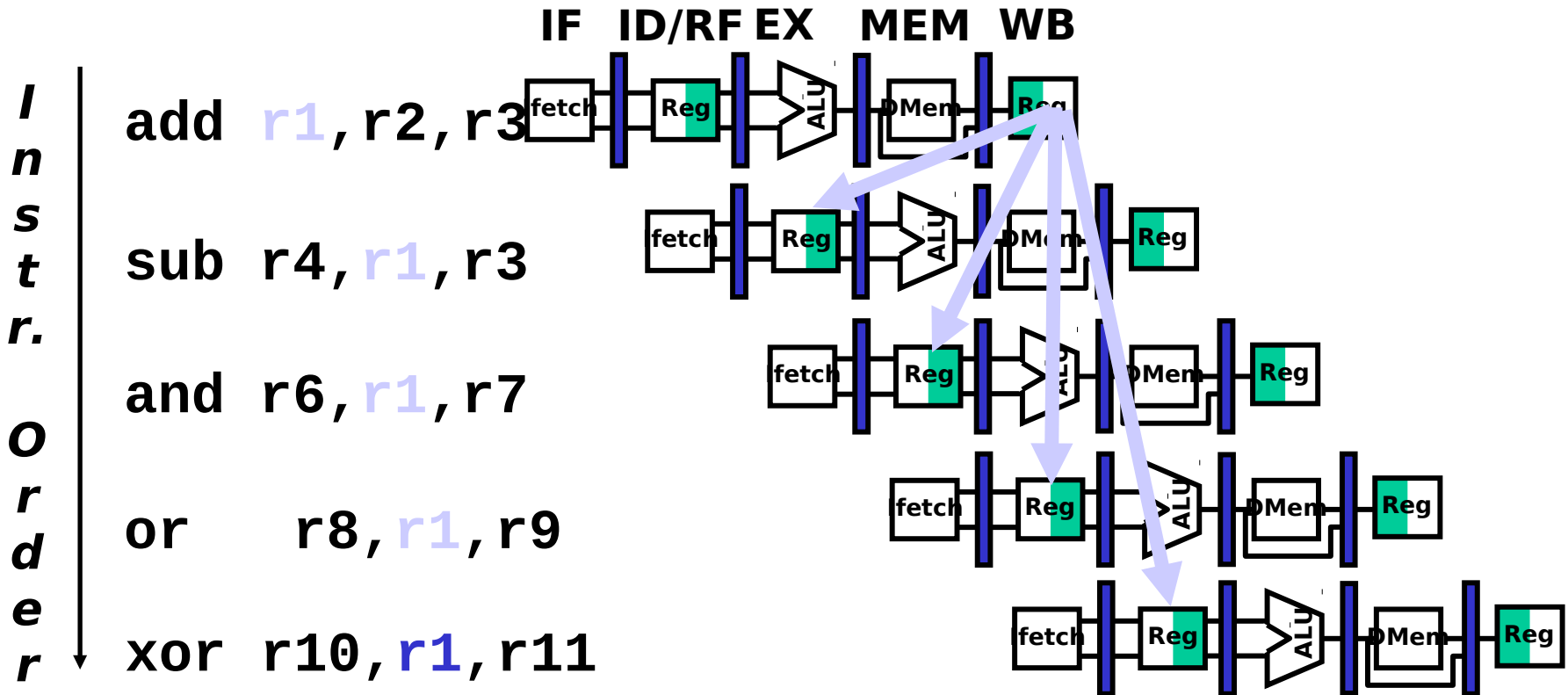
- Hardware detects RAW and stalls
- Assumes register written then read each cycle
 - + low cost to implement, simple
 - reduces IPC
- Try to minimize stalls

Minimizing RAW stalls

- Bypass/forward/short-circuit (We will use the word “forward”)
- Use data before it is in the register
 - + reduces/avoids stalls
 - complex
- Crucial for common RAW hazards

Data Hazards

Time (clock cycles)



The use of the result of the ADD instruction in the next three instructions causes a hazard, since the register is not written until after those instructions read it.

Figure 3.9

Data Hazards

Forwarding To Avoid
Data Hazard

Forwarding is the concept of making data available to the input of the ALU for subsequent instructions, even though the generating instruction hasn't gotten to WB in order to write the memory or registers.

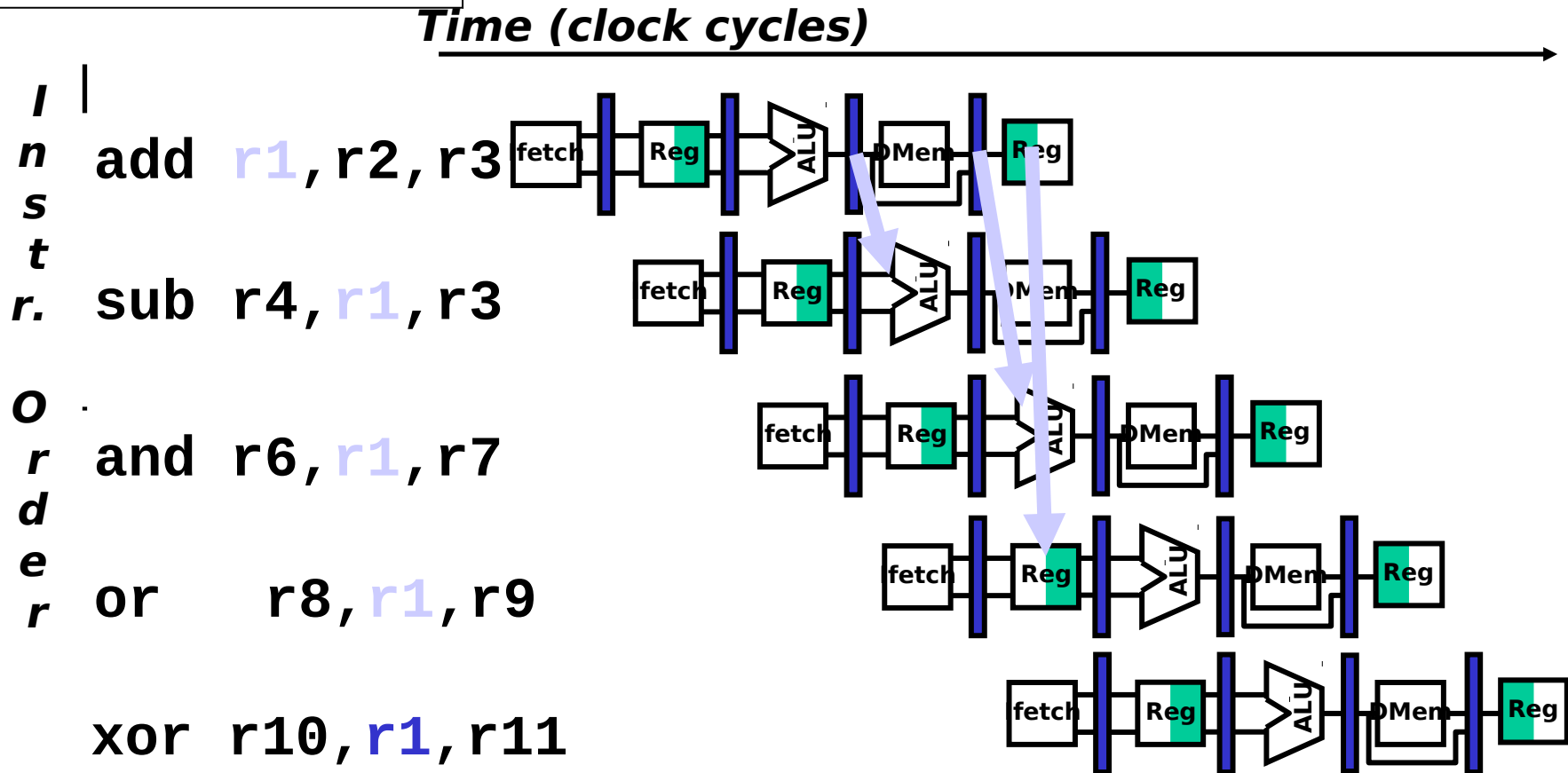
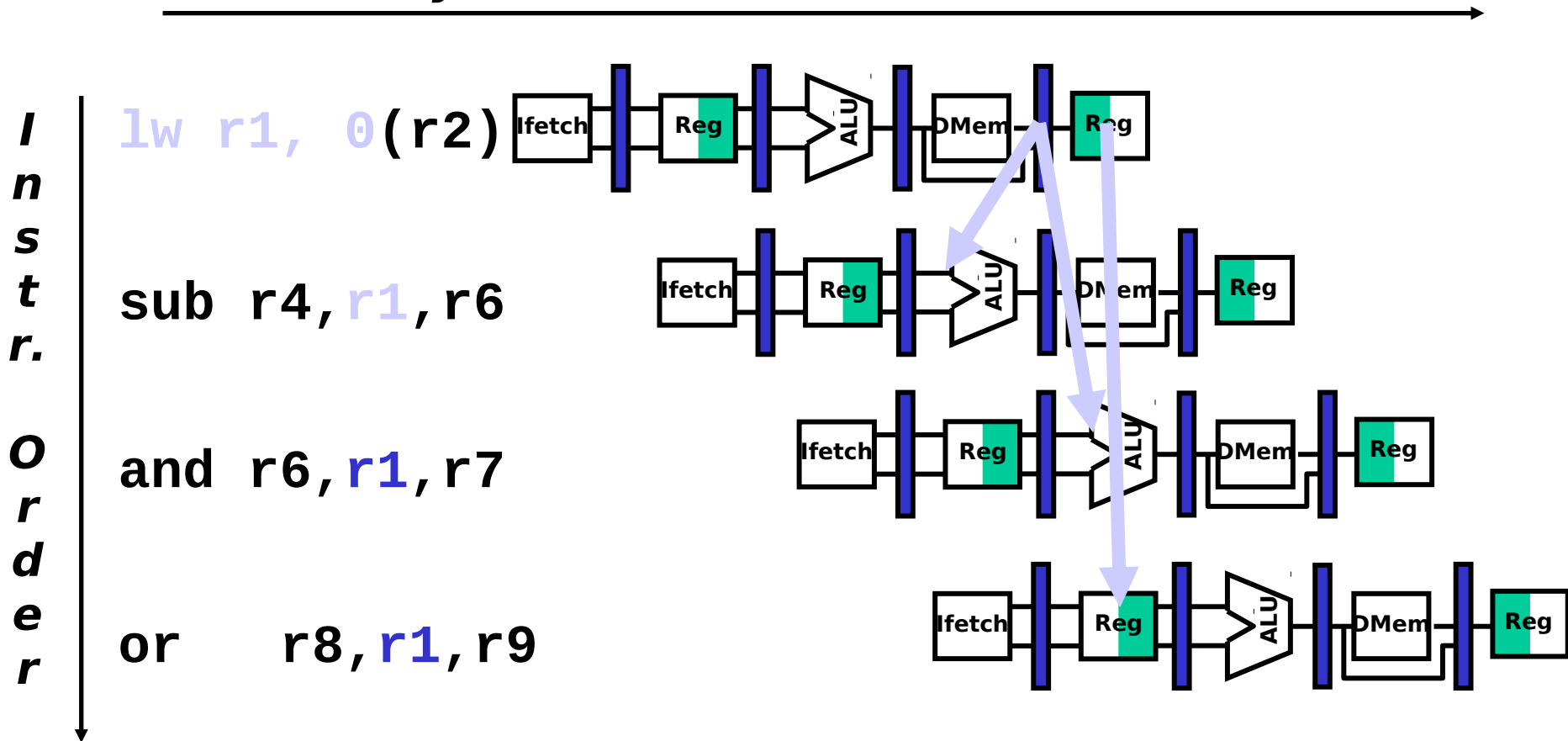


Figure 3.10

Data Hazards

The data isn't loaded until after the MEM stage.

Time (clock cycles)

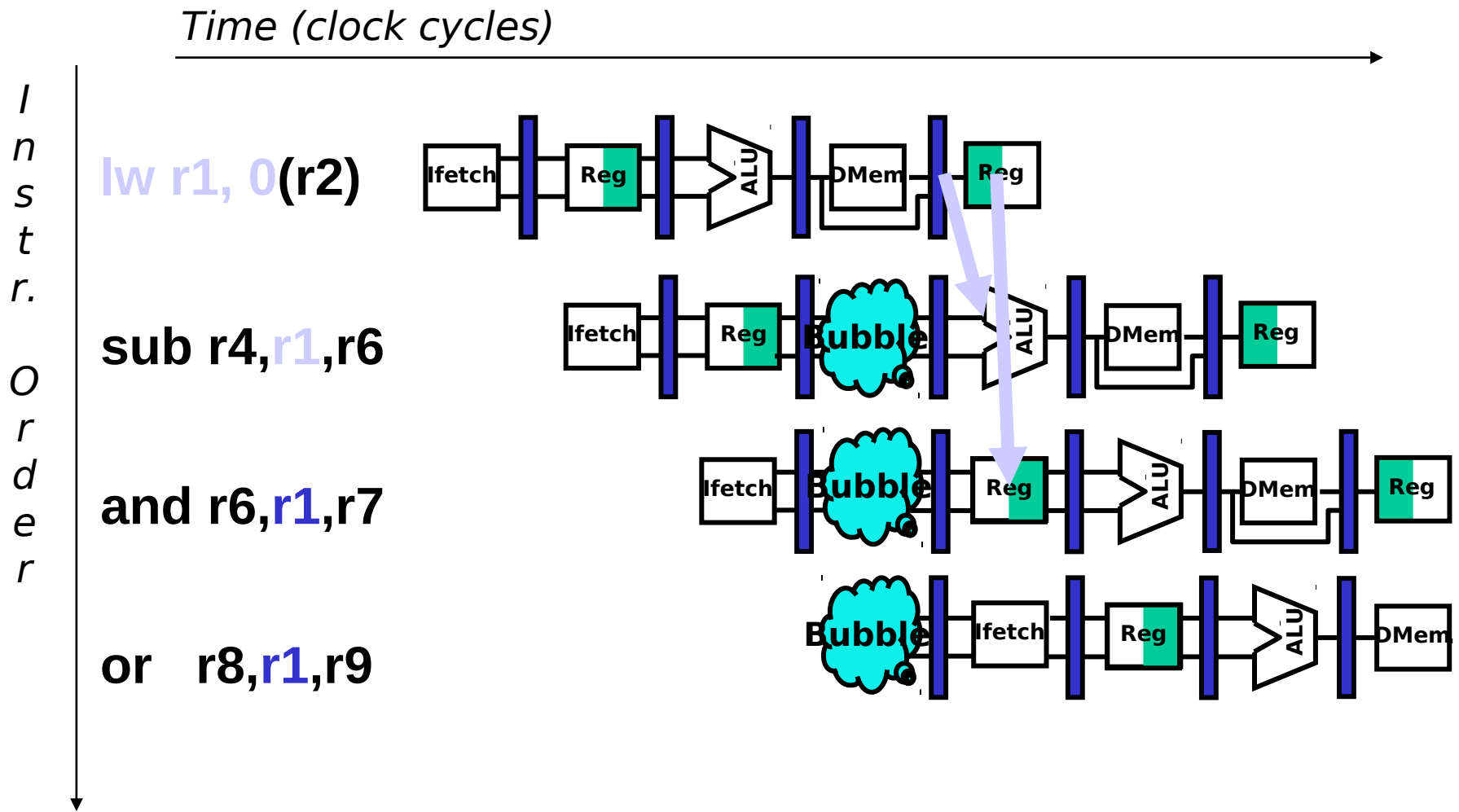


There are some instances where hazards occur, even with forwarding.

Figure 3.12

Data Hazards

The stall is necessary as shown here.



There are some instances where hazards occur, even with forwarding.

Figure 3.13

Data Hazards

This is another representation of the stall.

LW	R1, 0(R2)	IF	ID	EX	MEM	WB			
SUB	R4, R1, R5		IF	ID	EX	MEM	WB		
AND	R6, R1, R7			IF	ID	EX	MEM	WB	
OR	R8, R1, R9				IF	ID	EX	MEM	WB

LW	R1, 0(R2)	IF	ID	EX	MEM	WB				
SUB	R4, R1, R5		IF	ID	stall	EX	MEM	WB		
AND	R6, R1, R7			IF	stall	ID	EX	MEM	WB	
OR	R8, R1, R9				stall	IF	ID	EX	MEM	WB

Data Hazards

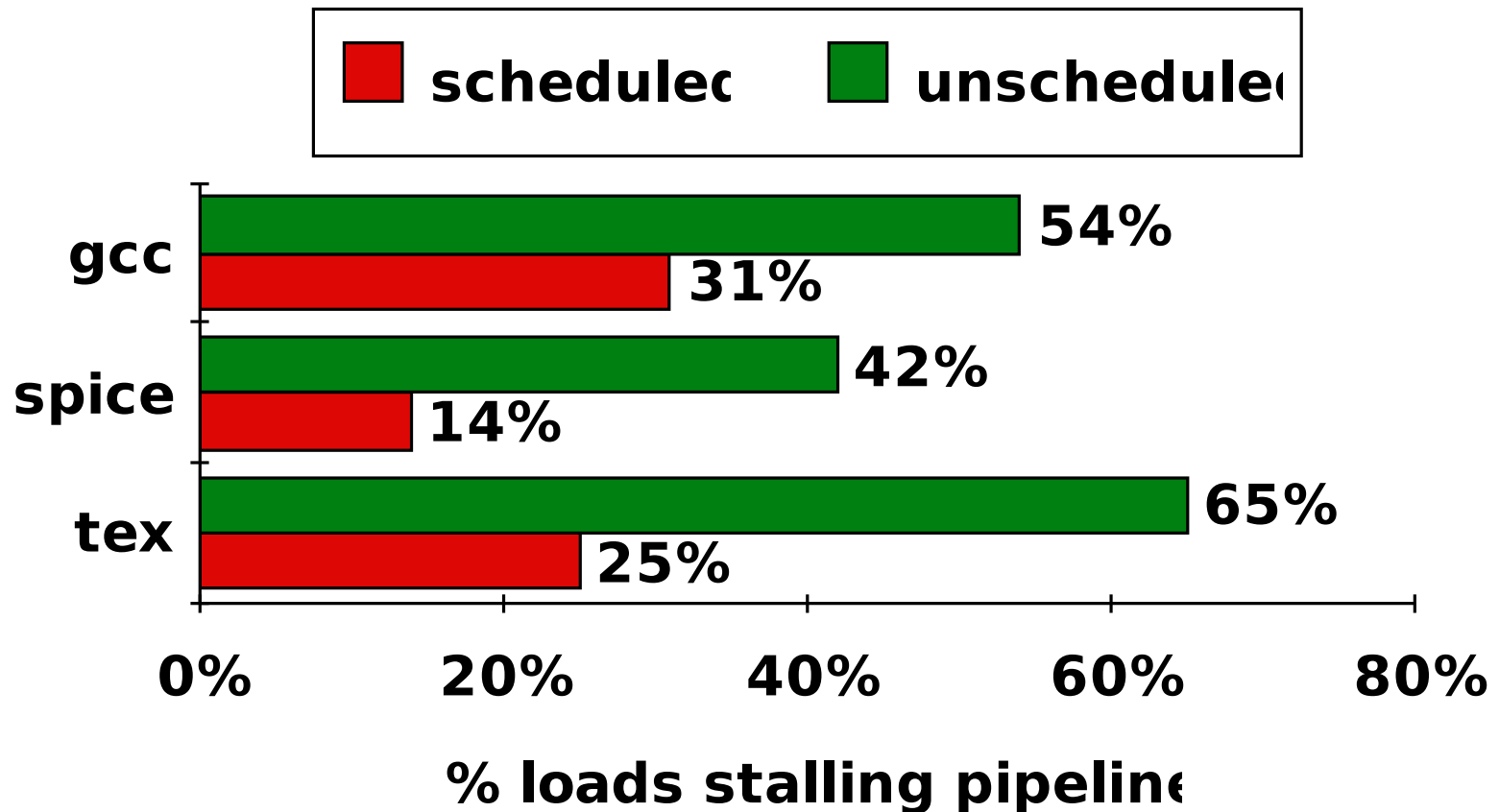
Instruction scheduled by compiler - move instruction in order to reduce stall.

lw Rb, b	-- code sequence for a = b+c before scheduling
lw Rc, c	
Add Ra, Rb, Rc	-- stall
sw a, Ra	
lw Re, e	-- code sequence for d = e+f before scheduling
lw Rf, f	
sub Rd, Re, Rf	-- stall
sw d, Rd	

Arrangement of code after scheduling.

```
lw Rb, b
lw Rc, c
lw Re, e
Add Ra, Rb, Rc
lw Rf, f
sw a, Ra
sub Rd, Re, Rf
sw d, Rd
```

Data Hazards



Control Hazards

A.1 What is Pipelining?

A.2 The Major Hurdle of Pipelining- Structural Hazards

- Structural Hazards
- Data Hazards
- Control Hazards

A.3 How is Pipelining Implemented

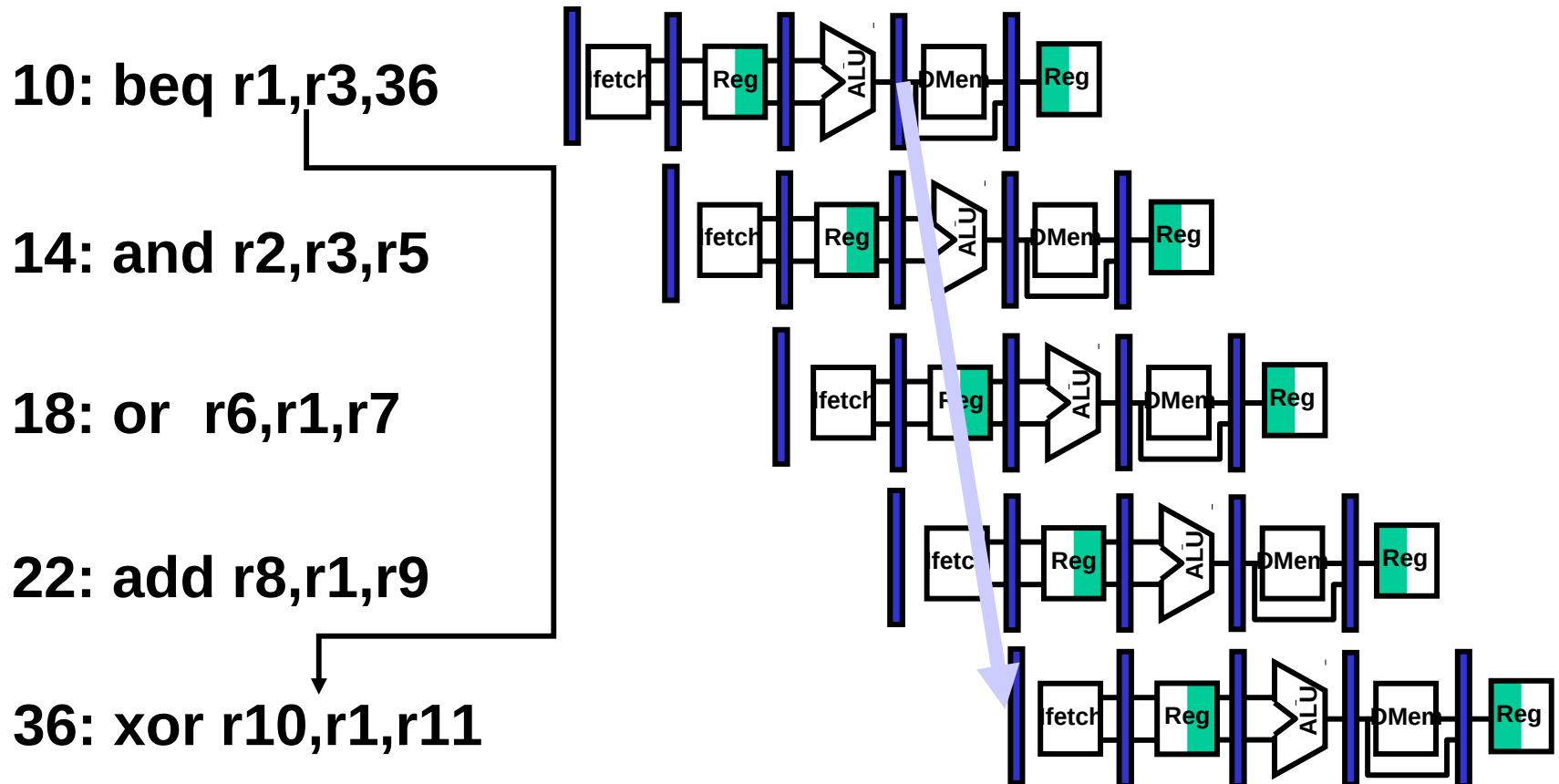
A.4 What Makes Pipelining Hard to Implement?

A.5 Extending the MIPS Pipeline to Handle Multi-cycle Operations

A control hazard is when we need to find the destination of a branch, and can't fetch any new instructions until we know that destination.

Control Hazards

Control Hazard on Branches Three Stage Stall



Control Hazards

Branch Stall Impact

- **If CPI = 1, 30% branch, Stall 3 cycles => new CPI = 1.9!**
(Whoa! How did we get that 1.9???)
- **Two part solution to this dramatic increase:**
 - Determine branch taken or not sooner, AND
 - Compute taken branch address earlier
- **MIPS branch tests if register = 0 or $\neq 0$**
- **MIPS Solution:**
 - Move Zero test to ID/RF stage
 - Adder to calculate new PC in ID/RF stage
 - must be fast
 - can't afford to subtract
 - compares with 0 are simple
 - Greater-than, Less-than test sign-bit, but not-equal must OR all bits
 - more general compares need ALU
 - 1 clock cycle penalty for branch versus 3

In the next chapter, we'll look at ways to avoid the branch all together.

Control Hazards

Five Branch Hazard Alternatives

#1: Stall until branch direction is clear

#2: Predict Branch Not Taken

- Execute successor instructions in sequence
- “Squash” instructions in pipeline if branch actually taken
- Advantage of late pipeline state update
- 47% MIPS branches not taken on average
- PC+4 already calculated, so use it to get next instruction

#3: Predict Branch Taken

- 53% MIPS branches taken on average
- But haven't calculated branch target address in MIPS
 - MIPS still incurs 1 cycle branch penalty
 - Other machines: branch target known before outcome

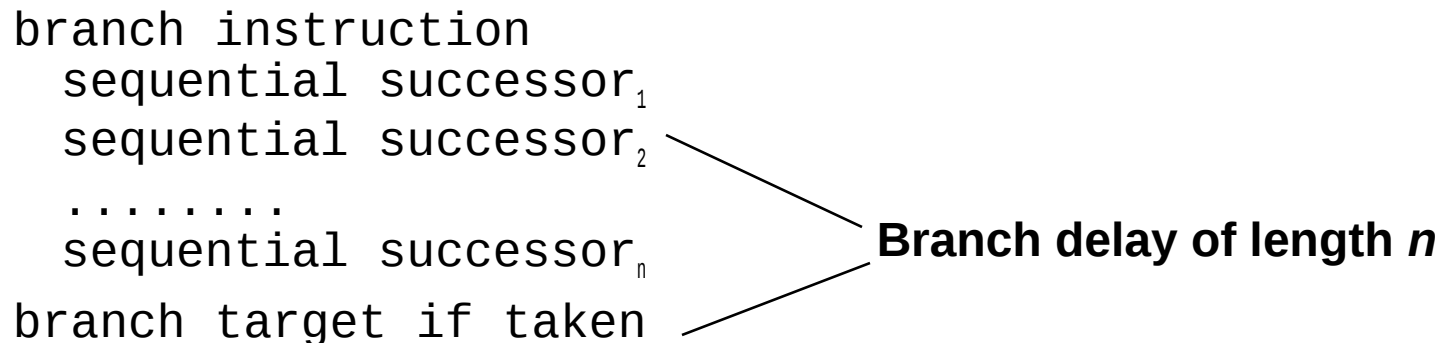
Control Hazards

Five Branch Hazard Alternatives

#4: Execute Both Paths

#5: Delayed Branch

- Define branch to take place **AFTER** a following instruction



- 1 slot delay allows proper decision and branch target address in 5 stage pipeline
- MIPS uses this

Control Hazards

Delayed Branch

- **Where to get instructions to fill branch delay slot?**
 - Before branch instruction
 - From the target address: only valuable when branch taken
 - From fall through: only valuable when branch not taken
 - Cancelling branches allow more slots to be filled
- **Compiler effectiveness for single branch delay slot:**
 - Fills about 60% of branch delay slots
 - About 80% of instructions executed in branch delay slots useful in computation
 - About 50% (60% x 80%) of slots usefully filled
- **Delayed Branch downside: 7-8 stage pipelines, multiple instructions issued per clock (superscalar)**

Control Hazards

Evaluating Branch Alternatives

$$\text{Pipeline speedup} = \frac{\text{Pipeline depth}}{1 + \text{Branch frequency} \times \text{Branch penalty}}$$

<i>Scheduling scheme</i>	<i>Branch penalty</i>	<i>CPI</i>	<i>speedup v. unpipelined</i>	<i>Speedup v. stall</i>
Stall pipeline	3	1.42	3.5	1.0
Predict taken	1	1.14	4.4	1.26
Predict not taken	1	1.09	4.5	1.29
Delayed branch	0.5	1.07	4.6	1.31

Conditional & Unconditional = 14%,

65% change PC

Control Hazards

Pipelining Introduction Summary

- **Just overlap tasks, and easy if tasks are independent**
- **Speed Up \leq Pipeline Depth; if ideal CPI is 1, then:**

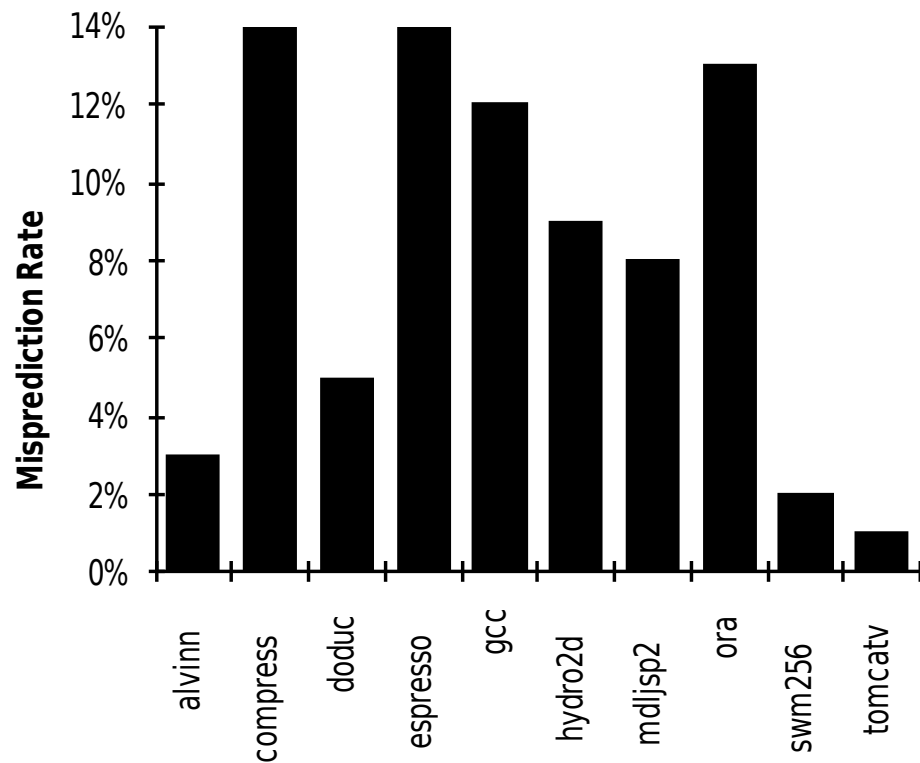
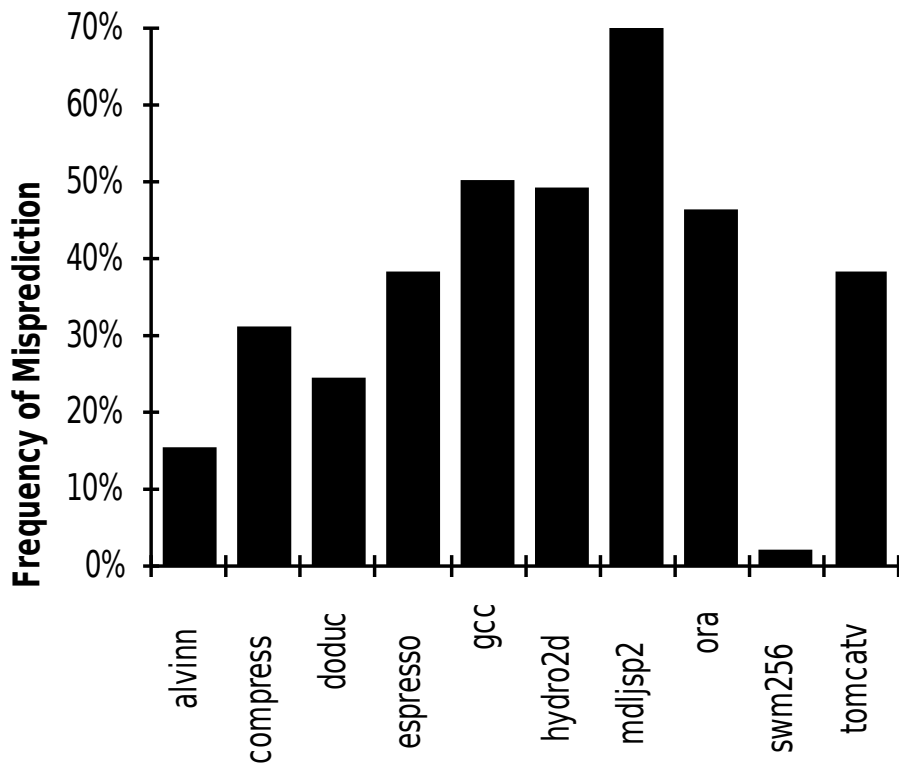
$$\text{Speedup} = \frac{\text{Pipeline Depth}}{1 + \text{Pipeline stall CPI}} \times \frac{\text{Clock Cycle Unpipelined}}{\text{Clock Cycle Pipelined}}$$

- **Hazards limit performance on computers:**
 - Structural: need more HW resources
 - Data (RAW,WAR,WAW): need forwarding, compiler scheduling
 - Control: delayed branch, prediction

Control Hazards

The compiler can program what it thinks the branch direction will be. Here are the results when it does so.

Compiler “Static” Prediction of Taken/Untaken Branches



Always taken

Taken backwards

Not Taken Forwards

Control Hazards

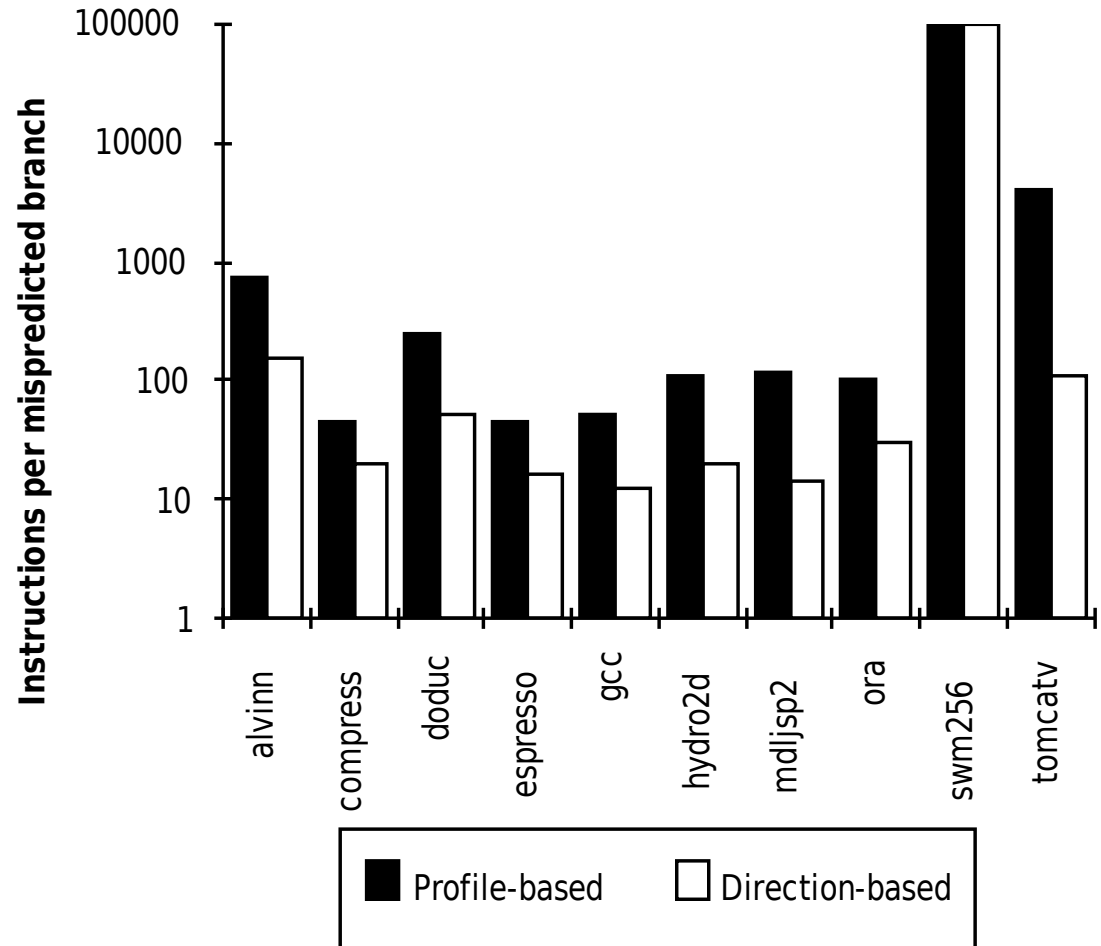
Compiler “Static” Prediction of Taken/Untaken Branches

- **Improves strategy for placing instructions in delay slot**
- **Two strategies**
 - Backward branch predict taken, forward branch not taken
 - Profile-based prediction: record branch behavior, predict branch based on prior run

Control Hazards

Evaluating Static Branch Prediction Strategies

- Misprediction ignores frequency of branch
- “Instructions between mispredicted branches” is a better metric



What Makes Pipelining Hard?

A.1 What is Pipelining?

**A.2 The Major Hurdle of Pipelining-
Structural Hazards**

- Data Hazards
- Control Hazards

A.3 How is Pipelining Implemented

**A.4 What Makes Pipelining Hard to
Implement?**

**A.5 Extending the MIPS Pipeline to
Handle Multi-cycle Operations**

What Makes Pipelining Hard?

Interrupts cause
great havoc!

Examples of interrupts:

- Power failing,
- Arithmetic overflow,
- I/O device request,
- OS call,
- Page fault

Interrupts (also known as: faults, exceptions, traps) often require

- surprise jump (to vectored address)
- linking return address
- saving of PSW (including CCs)
- state change (e.g., to kernel mode)

There are 5 instructions executing in 5 stage pipeline when an interrupt occurs:

- How to stop the pipeline?
- How to restart the pipeline?
- Who caused the interrupt?

What Makes Pipelining Hard?

Interrupts cause
great havoc!

What happens on interrupt while in delay slot ?

- Next instruction is not sequential
- solution #1: save multiple PCs
- Save current and next PC
 - Special return sequence, more complex hardware

solution #2: single PC plus

- Branch delay bit
- PC points to branch instruction

Stage

Problem that causes the interrupt

IF

Page fault on instruction fetch; misaligned memory access; memory-protection violation

ID

Undefined or illegal opcode

EX

Arithmetic interrupt

MEM

Page fault on data fetch; misaligned memory access; memory-protection violation

What Makes Pipelining Hard?

Interrupts cause
great havoc!

- **Simultaneous exceptions in more than one pipeline stage, e.g.,**
 - Load with data page fault in MEM stage
 - Add with instruction page fault in IF stage
 - Add fault will happen BEFORE load fault
- **Solution #1**
 - Interrupt status vector per instruction
 - Defer check until last stage, kill state update if exception
- **Solution #2**
 - Interrupt ASAP
 - Restart everything that is incomplete

Another advantage for state update late in pipeline!

What Makes Pipelining Hard?

Interrupts cause
great havoc!

Here's what happens on a data page fault.

	1	2	3	4	5	6	7	8	9		
i	F	D	X	M	W						
i+1		F	D	X	M	W	<- page fault				
i+2			F	D	X	M	W	<- squash			
i+3				F	D	X	M	W	<- squash		
i+4					F	D	X	M	W	<- squash	
i+5	trap ->					F	D	X	M	W	
i+6	trap handler ->						F	D	X	M	W

What Makes Pipelining Hard?

Complex Instructions

Complex Addressing Modes and Instructions

- **Address modes: Autoincrement causes register change during instruction execution**
 - Interrupts? Need to restore register state
 - Adds WAR and WAW hazards since writes are no longer the last stage.
- **Memory-Memory Move Instructions**
 - Must be able to handle multiple page faults
 - Long-lived instructions: partial state save on interrupt
- **Condition Codes**

Handling Multi-cycle Operations

A.1 What is Pipelining?

**A.2 The Major Hurdle of Pipelining-
Structural Hazards**

- Data Hazards
- Control Hazards

A.3 How is Pipelining Implemented

**A.4 What Makes Pipelining Hard to
Implement?**

**A.5 Extending the MIPS Pipeline to
Handle Multi-cycle Operations**

**Multi-cycle instructions also
lead to pipeline complexity.**

**A very lengthy instruction
causes everything else in
the pipeline to wait for it.**

Multi-Cycle Operations

Floating Point

Floating point gives long execution time.

This causes a stall of the pipeline.

It's possible to pipeline the FP execution unit so it can initiate new instructions without waiting full latency. Can also have multiple FP units.

<i>FP Instruction</i>	<i>Latency</i>	<i>Initiation Rate</i>
Add, Subtract	4	3
Multiply	8	4
Divide	36	35
Square root	112	111
Negate	2	1
Absolute value	2	1
FP compare	3	2

Multi-Cycle Operations

Floating Point

Divide, Square Root take -10X to -30X longer than Add

- Interrupts?
- Adds WAR and WAW hazards since pipelines are no longer same length

	1	2	3	4	5	6	7	8	9	10	11
i	IF	ID	EX	MEM	WB						
i + 1		IF	ID	EX	EX	EX	EX	MEM	WB		
i + 2			IF	ID	EX	MEM	WB				
i + 3				IF	ID	EX	EX	EX	EX	MEM	WB
i + 4					IF	ID	EX	MEM	WB		
i + 5						IF	ID	--	--	EX	EX
i + 6							IF	--	--	ID	EX

Notes:

i + 2: no WAW, but this complicates an interrupt

i + 4: no WB conflict

i + 5: stall forced by structural hazard

i + 6: stall forced by in-order issue

Summary of Pipelining Basics

- **Hazards limit performance**
 - Structural: need more HW resources
 - Data: need forwarding, compiler scheduling
 - Control: early evaluation & PC, delayed branch, prediction
- **Increasing length of pipe increases impact of hazards; pipelining helps instruction bandwidth, not latency**
- **Interrupts, Instruction Set, FP makes pipelining harder**
- **Compilers reduce cost of data and control hazards**
 - Load delay slots
 - Branch delay slots
 - Branch prediction

Credits

I have not written these notes by myself. There's a great deal of fancy artwork here that takes considerable time to prepare.

I have borrowed from:

Wen-mei & Patel: <http://courses.ece.uiuc.edu/ece511/lectures/lecture3.ppt>

Patterson: <http://www.cs.berkeley.edu/~pattsrn/252S98/index.html>

Rabaey: (He used lots of Patterson material):

<http://bwrc.eecs.berkeley.edu/Classes/CS252/index.htm>

Katz: (Again, he borrowed heavily from Patterson):

<http://http.cs.berkeley.edu/~randy/Courses/CS252.F95/CS252.Intro.html>

Mark Hill: (Follows text fairly well): <http://www.cs.wisc.edu/~markhill/cs752/>

Summary

A.1 What is Pipelining?

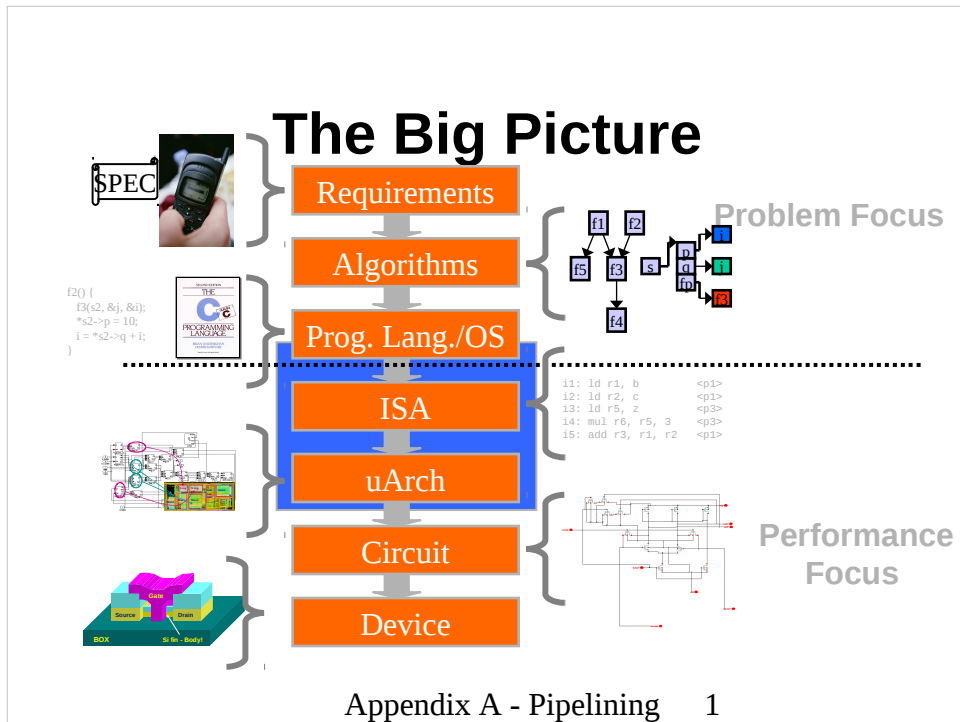
A.2 The Major Hurdle of Pipelining-Structural Hazards

- Data Hazards
- Control Hazards

A.3 How is Pipelining Implemented

A.4 What Makes Pipelining Hard to Implement?

A.5 Extending the MIPS Pipeline to Handle Multi-cycle Operations



The big picture. Talk about the levels of abstraction. Talk about the fact that this is where all programs get ushered into hardware execution.

Circuits are increasingly providing both opportunities (resources, bandwidth) and challenges (noise, power).

Circuits are locally designed; software is globally intertwined

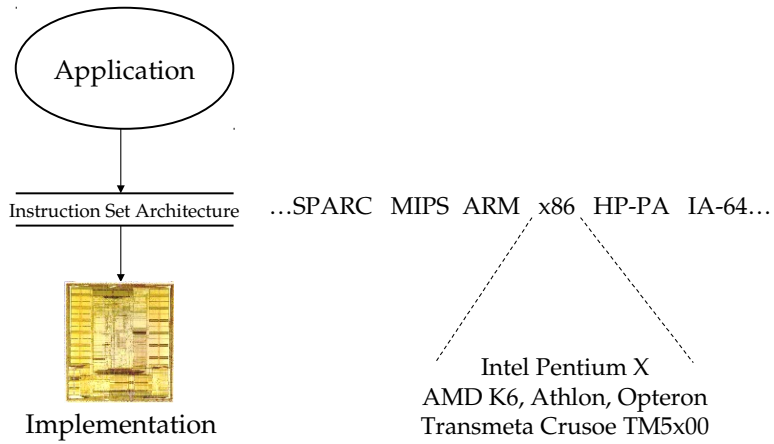
Software is increasingly over designed for portability and productivity.

The path between the two domains is increasingly stressed and inadequate due to this mismatch.

The focus of the thrust is to provide a very strong path from the productivity oriented software domain into the performance oriented hardware domain.

Translate device/circuit level innovations into visible benefit at the application/software level!

Instruction Set Architecture



Instruction Set Architecture

- **Strong influence on cost/performance**
- **New ISAs are rare, but versions are not**
 - 16-bit, 32-bit and 64-bit X86 versions
- **Longevity is a strong function of marketing prowess**

Traditional Issues

- **Strongly constrained by the number of bits available to instruction encoding**
- **Opcodes/operands**
- **Registers/memory**
- **Addressing modes**
- **Orthogonality**
- **0, 1, 2, 3 address machines**
- **Instruction formats**
- **Decoding uniformity**

Introduction

A.1 What is Pipelining?

A.2 The Major Hurdle of Pipelining-Structural Hazards

- Data Hazards
- Control Hazards

A.3 How is Pipelining Implemented

A.4 What Makes Pipelining Hard to Implement?

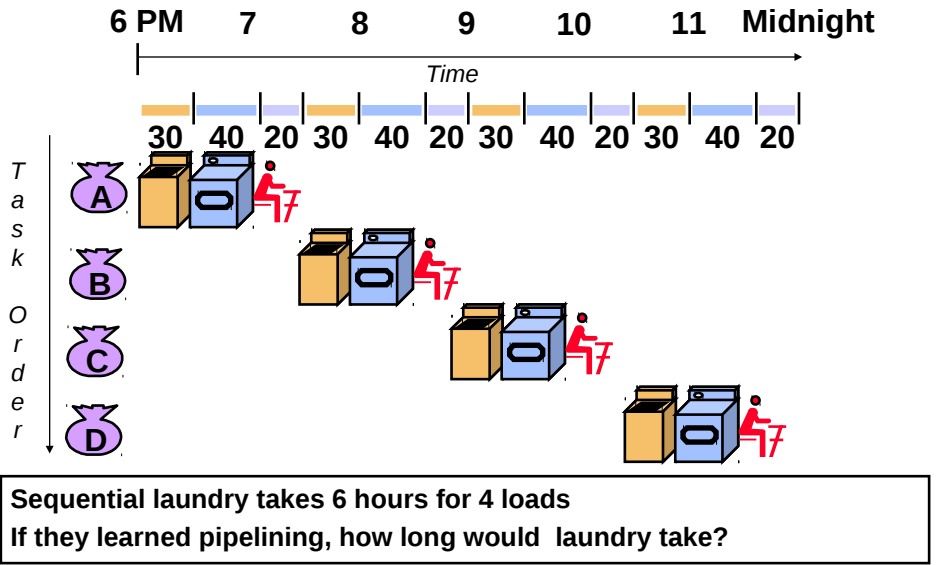
A.5 Extending the MIPS Pipeline to Handle Multi-cycle Operations

What Is Pipelining

- Laundry Example
- Ann, Brian, Cathy, Dave each have one load of clothes to wash, dry, and fold
- Washer takes 30 minutes
- Dryer takes 40 minutes
- “Folder” takes 20 minutes

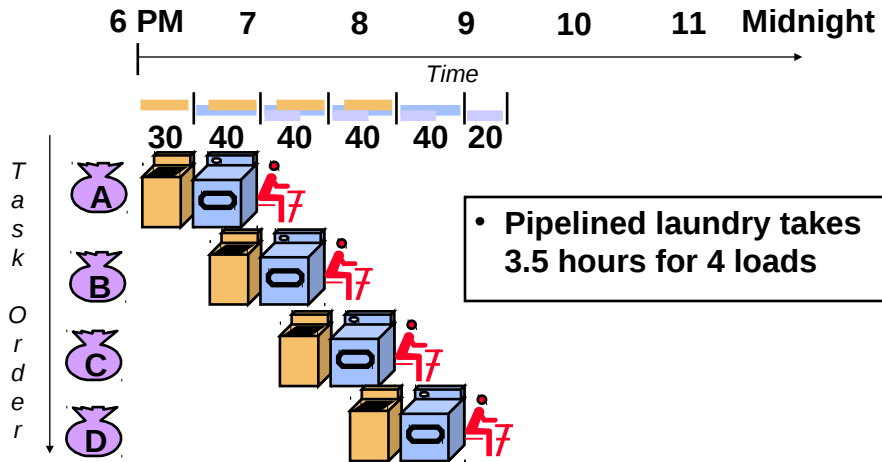


What Is Pipelining



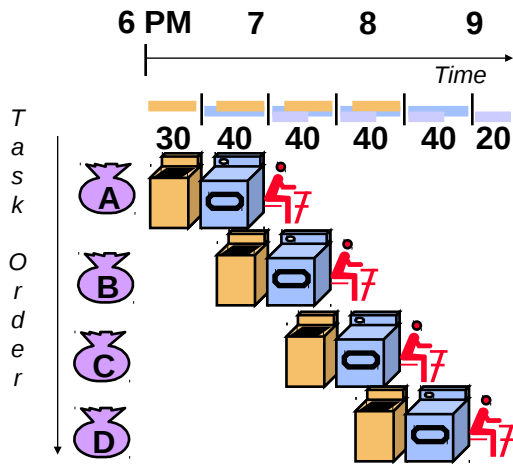
What Is Pipelining

Start work ASAP



Appendix A - Pipelining 8

What Is Pipelining

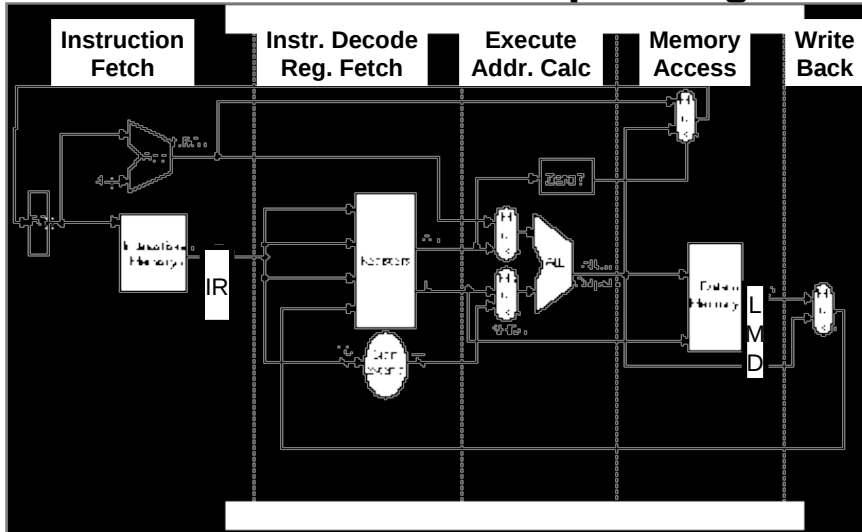


Pipelining Lessons

- 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
- Potential speedup = **Number pipe stages**
- Unbalanced lengths of pipe stages reduces speedup
- Time to "**fill**" pipeline and time to "**drain**" it reduces speedup

What Is Pipelining

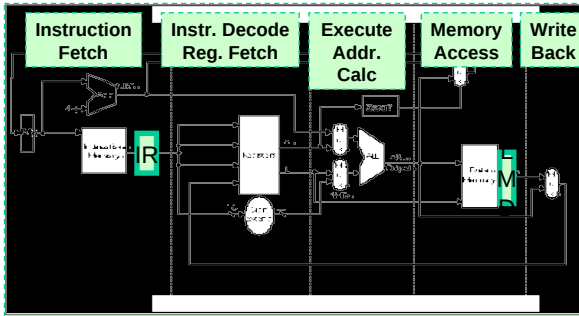
MIPS Without Pipelining



Appendix A - Pipelining 10

What Is Pipelining

MIPS Functions



Passed To Next Stage

$IR \leftarrow Mem[PC]$
 $NPC \leftarrow PC + 4$

Instruction Fetch (IF):

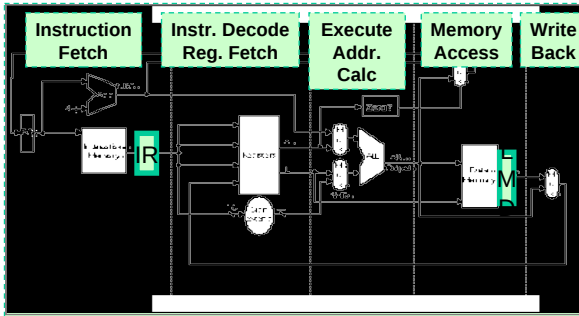
Send out the PC and fetch the instruction from memory into the instruction register (IR); increment the PC by 4 to address the next sequential instruction.

IR holds the instruction that will be used in the next stage.

NPC holds the value of the next PC.

What Is Pipelining

MIPS Functions



Passed To Next Stage

```
A <- Regs[IR6..IR10];  
B <- Regs[IR10..IR15];  
Imm <- ((IR16) ##IR16-31
```

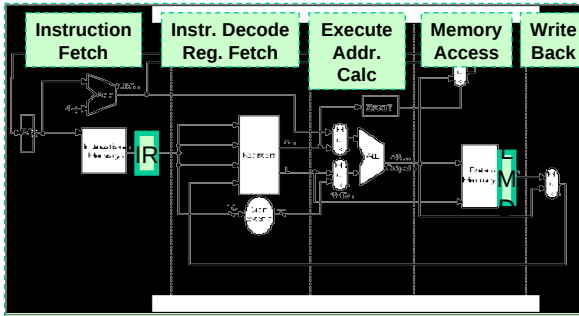
Instruction Decode/Register Fetch Cycle (ID):

Decode the instruction and access the register file to read the registers. The outputs of the general purpose registers are read into two temporary registers (A & B) for use in later clock cycles.

We extend the sign of the lower 16 bits of the Instruction Register.

What Is Pipelining

MIPS Functions



Passed To Next Stage
A <- A func. B
cond = 0;

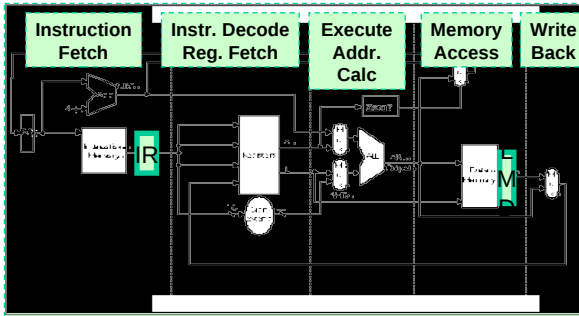
Execute Address Calculation (EX):

We perform an operation (for an ALU) or an address calculation (if it's a load or a Branch).

If an ALU, actually do the operation. If an address calculation, figure out how to obtain the address and stash away the location of that address for the next cycle.

What Is Pipelining

MIPS Functions



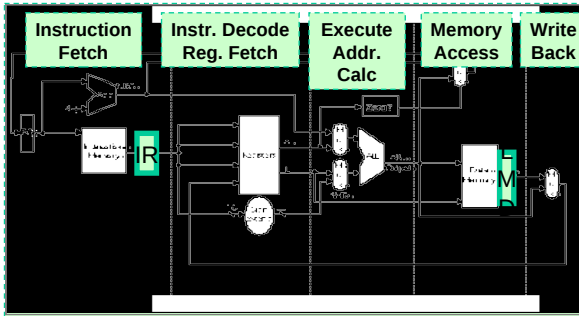
Passed To Next Stage

$A = \text{Mem}[\text{prev. } B]$
or
 $\text{Mem}[\text{prev. } B] = A$

MEMORY ACCESS (MEM):
If this is an ALU, do nothing.
If a load or store, then access memory.

What Is Pipelining

MIPS Functions



Passed To Next Stage
Regs <- A, B;

WRITE BACK (WB):

Update the registers from either the ALU or from the data loaded.

The Basic Pipeline For MIPS

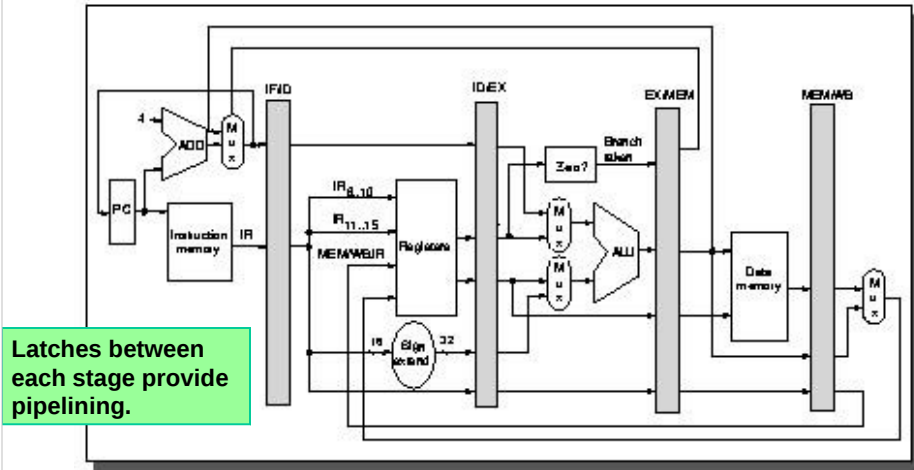


FIGURE 3.4 The datapath is pipelined by adding a set of registers, one between each pair of pipe stages.

The Basic Pipeline For MIPS

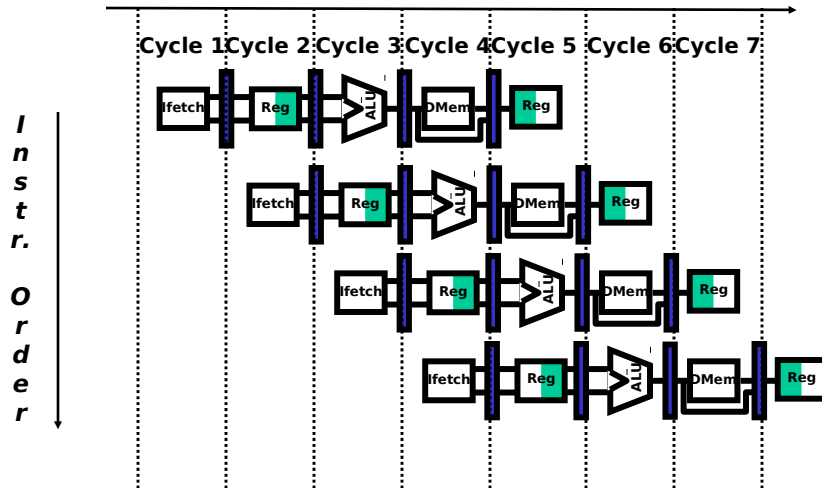


Figure 3.3

Pipeline Hurdles

A.1 What is Pipelining?

A.2 The Major Hurdle of Pipelining- Structural Hazards

- Structural Hazards
- Data Hazards
- Control Hazards

A.3 How is Pipelining Implemented

A.4 What Makes Pipelining Hard to Implement?

A.5 Extending the MIPS Pipeline to Handle Multi-cycle Operations

Limits to pipelining: Hazards prevent next instruction from executing during its designated clock cycle

- **Structural hazards:** HW cannot support this combination of instructions (single person to fold and put clothes away)
- **Data hazards:** Instruction depends on result of prior instruction still in the pipeline (missing sock)
- **Control hazards:** Pipelining of branches & other instructions that change the PC
- Common solution is to **stall** the pipeline until the hazard is resolved, inserting one or more "**bubbles**" in the pipeline

Pipeline Hurdles

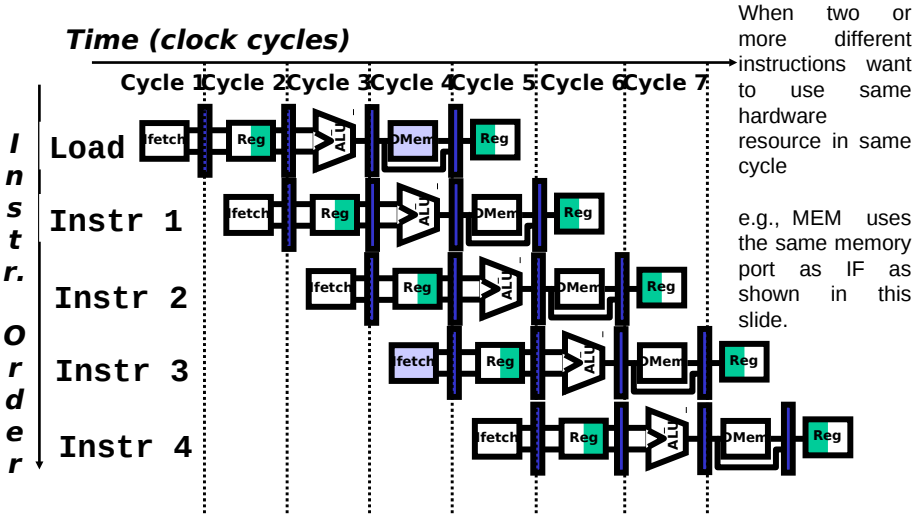
Definition

- conditions that lead to incorrect behavior if not fixed
- Structural hazard
 - two different instructions use same h/w in same cycle
- Data hazard
 - two different instructions use same storage
 - must appear as if the instructions execute in correct order
- Control hazard
 - one instruction affects which instruction is next

Resolution

- Pipeline interlock logic detects hazards and fixes them
- simple solution: stall -
- increases CPI, decreases performance
- better solution: partial stall -
- some instruction stall, others proceed better to stall early than late

Structural Hazards



When two or more different instructions want to use same hardware resource in same cycle

e.g., MEM uses the same memory port as IF as shown in this slide.

Figure 3.6

Structural Hazards

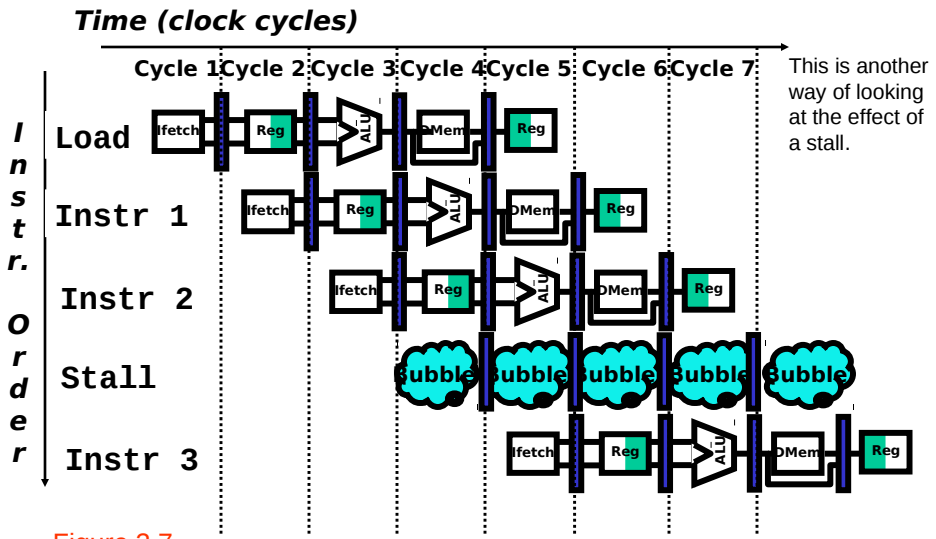


Figure 3.7

Structural Hazards

Instruction	Clock cycle number									
	1	2	3	4	5	6	7	8	9	10
Load instruction	IF	ID	EX	MEM	WB					
Instruction $j + 1$		IF	ID	EX	MEM	WB				
Instruction $j + 2$			IF	ID	EX	MEM	WB			
Instruction $j + 3$				stall	IF	ID	EX	MEM	WB	
Instruction $j + 4$						IF	ID	EX	MEM	WB
Instruction $j + 5$							IF	ID	EX	MEM
Instruction $j + 6$								IF	ID	EX

This is another way to represent the stall we saw on the last few pages.

Structural Hazards

Dealing with Structural Hazards

Stall

- low cost, simple
- Increases CPI
- use for rare case since stalling has performance effect

Pipeline hardware resource

- useful for multi-cycle resources
- good performance
- sometimes complex e.g., RAM

Replicate resource

- good performance
- increases cost (+ maybe interconnect delay)
- useful for cheap or divisible resources

Structural Hazards

Structural hazards are reduced with these rules:

- Each instruction uses a resource at most once
- Always use the resource in the same pipeline stage
- Use the resource for one cycle only

Many RISC ISA's are designed with this in mind

Sometimes very complex to do this. For example, memory of necessity is used in the IF and MEM stages.

Some common Structural Hazards:

- Memory - we've already mentioned this one.
- Floating point - Since many floating point instructions require many cycles, it's easy for them to interfere with each other.
- Starting up more of one type of instruction than there are resources. For instance, the PA-8600 can support two ALU + two load/store instructions per cycle - that's how much hardware it has available.

Structural Hazards

This is the example on Page 144.

We want to compare the performance of two machines. Which machine is faster?

- Machine A: Dual ported memory - so there are no memory stalls
- Machine B: Single ported memory, but its pipelined implementation has a 1.05 times faster clock rate

Assume:

- Ideal CPI = 1 for both
- Loads are 40% of instructions executed

$$\begin{aligned}\text{SpeedUp}_A &= \text{Pipeline Depth} / (1 + 0) \times (\text{clock}_{\text{unpipe}} / \text{clock}_{\text{pipe}}) \\ &= \text{Pipeline Depth}\end{aligned}$$

$$\begin{aligned}\text{SpeedUp}_B &= \text{Pipeline Depth} / (1 + 0.4 \times 1) \\ &\quad \times (\text{clock}_{\text{unpipe}} / (\text{clock}_{\text{unpipe}} / 1.05)) \\ &= (\text{Pipeline Depth} / 1.4) \times 1.05 \\ &= 0.75 \times \text{Pipeline Depth}\end{aligned}$$

$$\text{SpeedUp}_A / \text{SpeedUp}_B = \text{Pipeline Depth} / (0.75 \times \text{Pipeline Depth}) = 1.33$$

- Machine A is 1.33 times faster

Data Hazards

- A.1 What is Pipelining?**
- A.2 The Major Hurdle of Pipelining- Structural Hazards**
 - Structural Hazards
 - Data Hazards
 - Control Hazards
- A.3 How is Pipelining Implemented**
- A.4 What Makes Pipelining Hard to Implement?**
- A.5 Extending the MIPS Pipeline to Handle Multi-cycle Operations**

These occur when at any time, there are instructions active that need to access the same data (memory or register) locations.

Where there's real trouble is when we have:

**instruction A
instruction B**

and B manipulates (reads or writes) data before A does. This violates the order of the instructions, since the architecture implies that A completes entirely before B is executed.

Data Hazards

Execution Order is:
Instr_i
Instr_j

Read After Write (RAW)
Instr_j tries to read operand before Instr_i writes it

I: add r1, r2, r3
J: sub r4, r1, r3

- Caused by a “**Dependence**” (in compiler nomenclature). This hazard results from an actual need for communication.


Data Hazards

Execution Order is:
Instr_i
Instr_j

Write After Read (WAR)

Instr_j tries to write operand *before* Instr_i reads it

- Gets wrong operand

 **I: sub r4, r1, r3**
J: add r1, r2, r3
K: mul r6, r1, r7

- Called an “[anti-dependence](#)” by compiler writers. This results from reuse of the name “r1”.

- **Can't happen in MIPS 5 stage pipeline because:**
 - All instructions take 5 stages, and
 - Reads are always in stage 2, and
 - Writes are always in stage 5


Data Hazards

Execution Order is:
Instr_i
Instr_j

Write After Write (WAW)

Instr_j tries to write operand *before* Instr_i writes it

- Leaves wrong result (Instr_i not Instr_j)

 I: sub r1, r4, r3
J: add r1, r2, r3
K: mul r6, r1, r7

- Called an “**output dependence**” by compiler writers
This also results from the reuse of name “r1”.
- Can't happen in MIPS 5 stage pipeline because:
 - All instructions take 5 stages, and
 - Writes are always in stage 5
- Will see WAR and WAW in later more complicated pipes

Data Hazards

Simple Solution to RAW

- Hardware detects RAW and stalls
- Assumes register written then read each cycle
 - + low cost to implement, simple
 - reduces IPC
- Try to minimize stalls

Minimizing RAW stalls

- Bypass/forward/short-circuit (We will use the word "forward")
- Use data before it is in the register
 - + reduces/avoids stalls
 - complex
- Crucial for common RAW hazards

Data Hazards

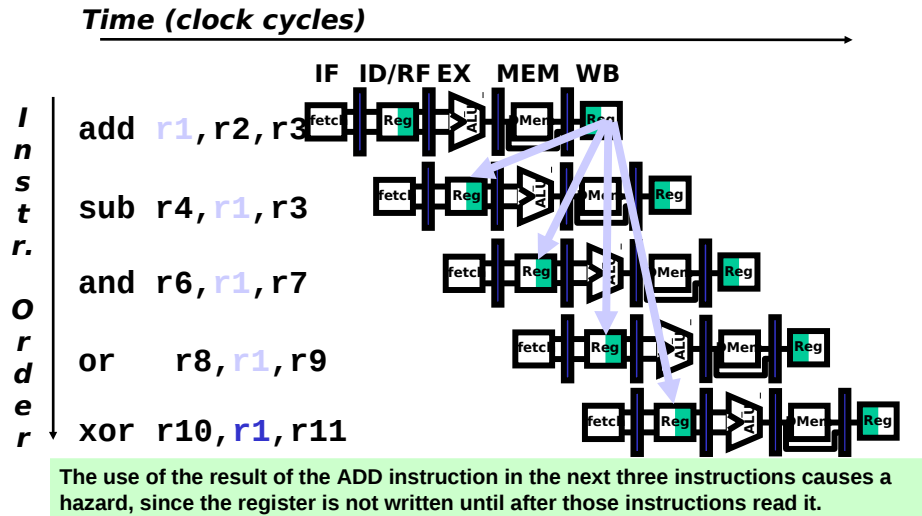


Figure 3.9

Data Hazards

Forwarding To Avoid Data Hazard

Forwarding is the concept of making data available to the input of the ALU for subsequent instructions, even though the generating instruction hasn't gotten to WB in order to write the memory or registers.

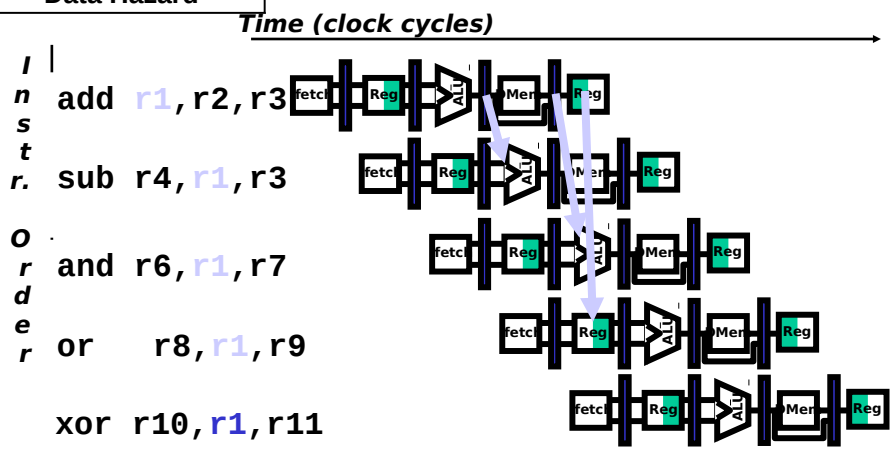
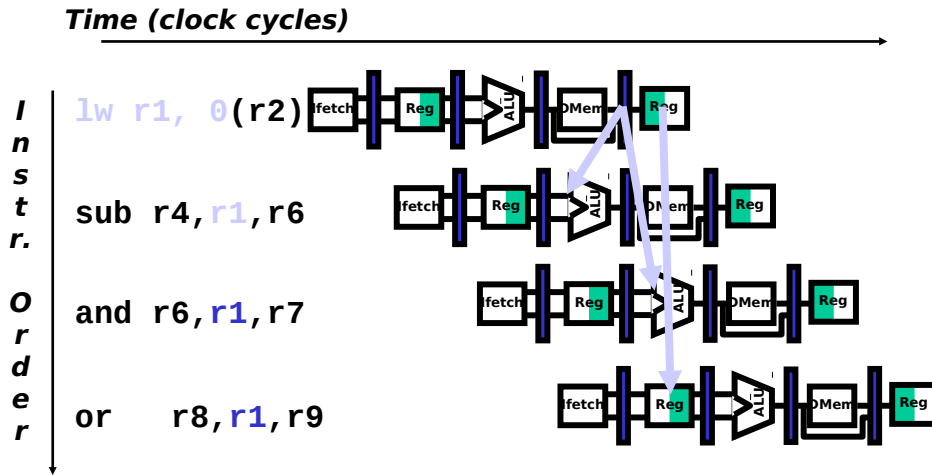


Figure 3.10

Data Hazards

The data isn't loaded until after the MEM stage.

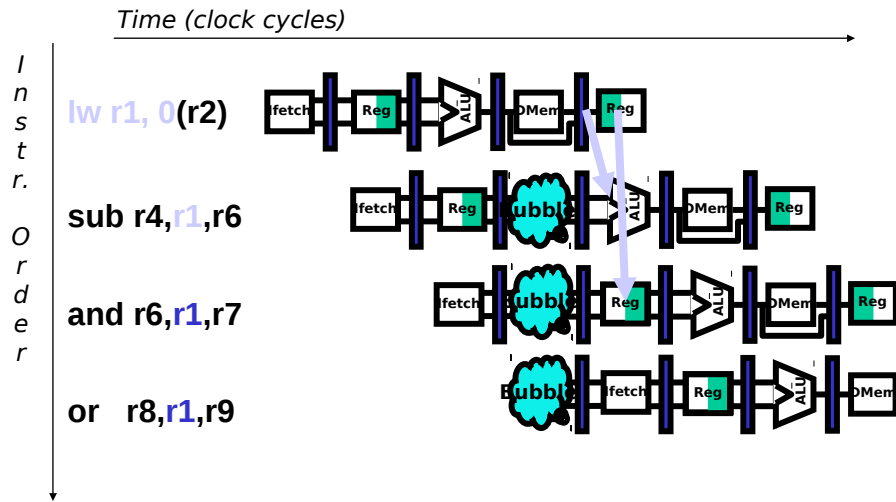


There are some instances where hazards occur, even with forwarding.

Figure 3.12

Data Hazards

The stall is necessary as shown here.



There are some instances where hazards occur, even with forwarding.

Figure 3.13

Data Hazards

This is another representation of the stall.

LW	R1, 0(R2)	IF	ID	EX	MEM	WB				
SUB	R4, R1, R5		IF	ID	EX	MEM	WB			
AND	R6, R1, R7			IF	ID	EX	MEM	WB		
OR	R8, R1, R9				IF	ID	EX	MEM	WB	

LW	R1, 0(R2)	IF	ID	EX	MEM	WB				
SUB	R4, R1, R5		IF	ID	stall	EX	MEM	WB		
AND	R6, R1, R7			IF	stall	ID	EX	MEM	WB	
OR	R8, R1, R9				stall	IF	ID	EX	MEM	WB

Data Hazards

Pipeline Scheduling

Instruction scheduled by compiler - move instruction in order to reduce stall.

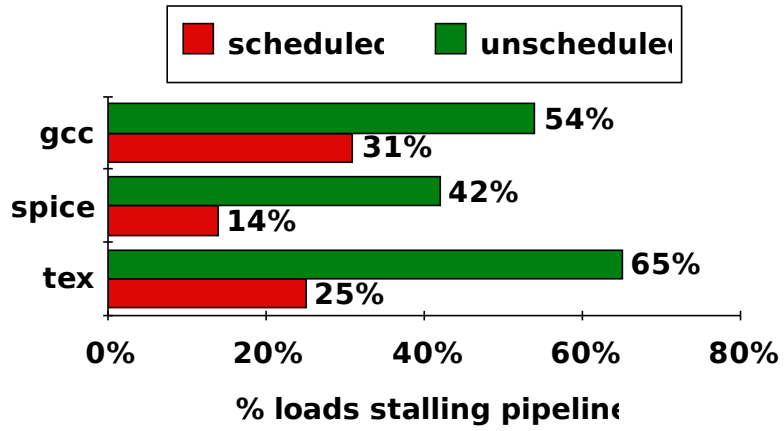
```
lw Rb, b          -- code sequence for a = b+c before scheduling
lw Rc, c
Add Ra, Rb, Rc    -- stall
sw a, Ra
lw Re, e          -- code sequence for d = e+f before scheduling
lw Rf, f
sub Rd, Re, Rf    -- stall
sw d, Rd
```

Arrangement of code after scheduling.

```
lw Rb, b
lw Rc, c
lw Re, e
Add Ra, Rb, Rc
lw Rf, f
sw a, Ra
sub Rd, Re, Rf
sw d, Rd
```


Data Hazards

Pipeline Scheduling



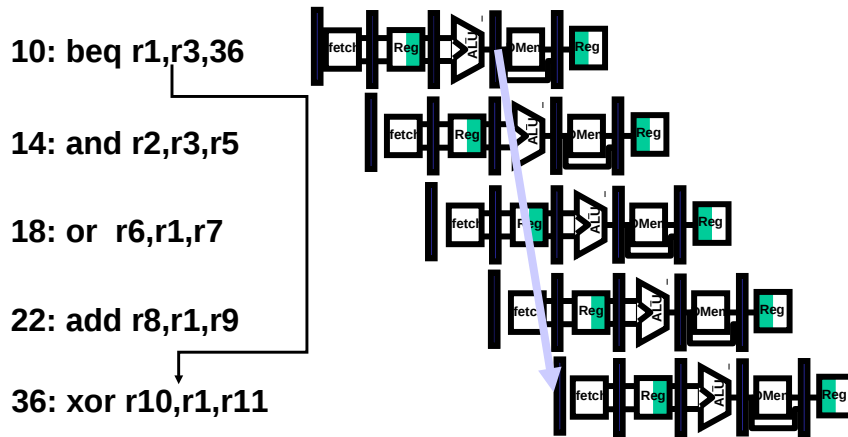
Control Hazards

A.1 What is Pipelining?
A.2 The Major Hurdle of Pipelining- Structural Hazards
-- Structural Hazards
-- Data Hazards
-- Control Hazards
A.3 How is Pipelining Implemented
A.4 What Makes Pipelining Hard to Implement?
A.5 Extending the MIPS Pipeline to Handle Multi-cycle Operations

A control hazard is when we need to find the destination of a branch, and can't fetch any new instructions until we know that destination.

Control Hazards

Control Hazard on Branches
Three Stage Stall



Control Hazards

Branch Stall Impact

- **If CPI = 1, 30% branch, Stall 3 cycles => new CPI = 1.9!**
(Whoa! How did we get that 1.9???)
- **Two part solution to this dramatic increase:**
 - Determine branch taken or not sooner, AND
 - Compute taken branch address earlier
- **MIPS branch tests if register = 0 or $\neq 0$**
- **MIPS Solution:**
 - Move Zero test to ID/RF stage
 - Adder to calculate new PC in ID/RF stage
 - must be fast
 - can't afford to subtract
 - compares with 0 are simple
 - Greater-than, Less-than test sign-bit, but not-equal must OR all bits
 - more general compares need ALU
 - 1 clock cycle penalty for branch versus 3

In the next chapter, we'll look at ways to avoid the branch all together.

Control Hazards

Five Branch Hazard Alternatives

#1: Stall until branch direction is clear

#2: Predict Branch Not Taken

- Execute successor instructions in sequence
- “Squash” instructions in pipeline if branch actually taken
- Advantage of late pipeline state update
- 47% MIPS branches not taken on average
- PC+4 already calculated, so use it to get next instruction

#3: Predict Branch Taken

- 53% MIPS branches taken on average
- **But haven't calculated branch target address in MIPS**
 - MIPS still incurs 1 cycle branch penalty
 - Other machines: branch target known before outcome

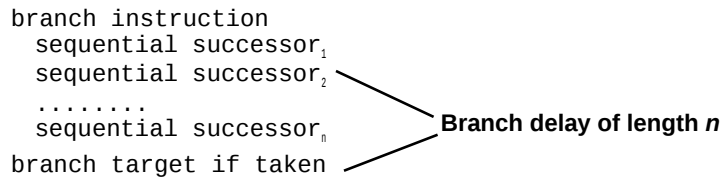
Control Hazards

Five Branch Hazard Alternatives

#4: Execute Both Paths

#5: Delayed Branch

- Define branch to take place **AFTER** a following instruction



- 1 slot delay allows proper decision and branch target address in 5 stage pipeline
- MIPS uses this

Control Hazards

Delayed Branch

- **Where to get instructions to fill branch delay slot?**
 - Before branch instruction
 - From the target address: only valuable when branch taken
 - From fall through: only valuable when branch not taken
 - Cancelling branches allow more slots to be filled
- **Compiler effectiveness for single branch delay slot:**
 - Fills about 60% of branch delay slots
 - About 80% of instructions executed in branch delay slots useful in computation
 - About 50% (60% x 80%) of slots usefully filled
- **Delayed Branch downside: 7-8 stage pipelines, multiple instructions issued per clock (superscalar)**

Control Hazards

Evaluating Branch Alternatives

$$\text{Pipeline speedup} = \frac{\text{Pipeline depth}}{1 + \text{Branch frequency} \times \text{Branch penalty}}$$

<i>Scheduling scheme</i>	<i>Branch penalty</i>	<i>CPI</i>	<i>speedup v. unpipelined</i>	<i>Speedup v. stall</i>
Stall pipeline	3	1.42	3.5	1.0
Predict taken	1	1.14	4.4	1.26
Predict not taken	1	1.09	4.5	1.29
Delayed branch	0.5	1.07	4.6	1.31

Conditional & Unconditional = 14%,

65% change PC

Control Hazards

Pipelining Introduction Summary

- Just overlap tasks, and easy if tasks are independent
- Speed Up \propto Pipeline Depth; if ideal CPI is 1, then:

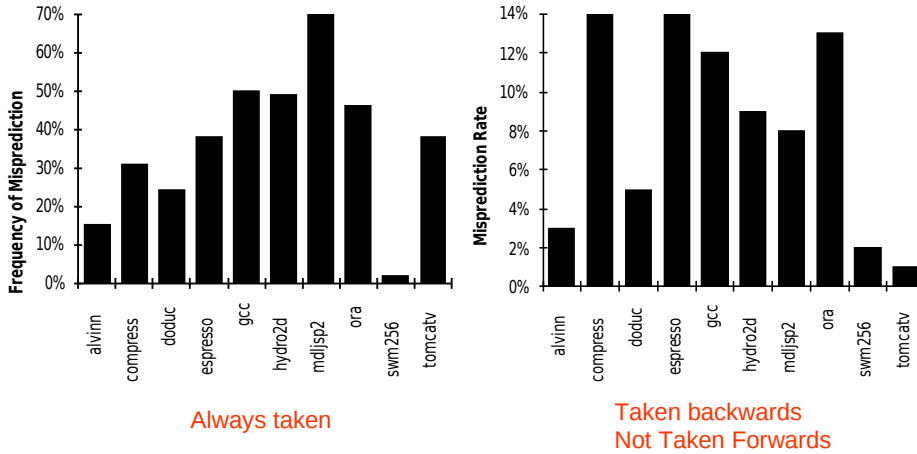
$$\text{Speedup} = \frac{\text{Pipeline Depth}}{1 + \text{Pipeline stall CPI}} \times \frac{\text{Clock Cycle Unpipelined}}{\text{Clock Cycle Pipelined}}$$

- Hazards limit performance on computers:
 - Structural: need more HW resources
 - Data (RAW,WAR,WAW): need forwarding, compiler scheduling
 - Control: delayed branch, prediction

Control Hazards

The compiler can program what it thinks the branch direction will be. Here are the results when it does so.

Compiler "Static" Prediction of Taken/Untaken Branches



Control Hazards

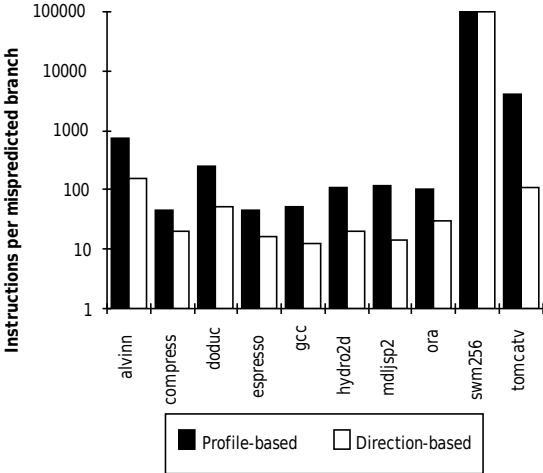
Compiler “Static” Prediction of Taken/Untaken Branches

- **Improves strategy for placing instructions in delay slot**
- **Two strategies**
 - Backward branch predict taken, forward branch not taken
 - Profile-based prediction: record branch behavior, predict branch based on prior run

Control Hazards

Evaluating Static Branch Prediction Strategies

- Misprediction ignores frequency of branch
- “Instructions between mispredicted branches” is a better metric



What Makes Pipelining Hard?

A.1 What is Pipelining?

A.2 The Major Hurdle of Pipelining- Structural Hazards

- Data Hazards
- Control Hazards

A.3 How is Pipelining Implemented

A.4 What Makes Pipelining Hard to Implement?

A.5 Extending the MIPS Pipeline to Handle Multi-cycle Operations

What Makes Pipelining Hard?

Interrupts cause
great havoc!

Examples of interrupts:

- Power failing,
- Arithmetic overflow,
- I/O device request,
- OS call,
- Page fault

There are 5 instructions executing in 5 stage pipeline when an interrupt occurs:

- How to stop the pipeline?
- How to restart the pipeline?
- Who caused the interrupt?

Interrupts (also known as: faults, exceptions, traps) often require

- surprise jump (to vectored address)
- linking return address
- saving of PSW (including CCs)
- state change (e.g., to kernel mode)

What Makes Pipelining Hard?

Interrupts cause
great havoc!

What happens on interrupt while in delay slot ?

- Next instruction is not sequential
- solution #1: save multiple PCs
- Save current and next PC
 - Special return sequence, more complex hardware
- solution #2: single PC plus
- Branch delay bit
 - PC points to branch instruction

<i>Stage</i>	<i>Problem that causes the interrupt</i>
IF	Page fault on instruction fetch; misaligned memory access; memory-protection violation
ID	Undefined or illegal opcode
EX	Arithmetic interrupt
MEM	Page fault on data fetch; misaligned memory access; memory-protection violation

What Makes Pipelining Hard?

Interrupts cause
great havoc!

- **Simultaneous exceptions in more than one pipeline stage, e.g.,**
 - Load with data page fault in MEM stage
 - Add with instruction page fault in IF stage
 - Add fault will happen BEFORE load fault
- **Solution #1**
 - Interrupt status vector per instruction
 - Defer check until last stage, kill state update if exception
- **Solution #2**
 - Interrupt ASAP
 - Restart everything that is incomplete

Another advantage for state update late in pipeline!

What Makes Pipelining Hard?

Interrupts cause
great havoc!

Here's what happens on a data page fault.

	1	2	3	4	5	6	7	8	9
i	F	D	X	M	W				
i+1		F	D	X	M	W	<- page fault		
i+2			F	D	X	M	W	<- squash	
i+3				F	D	X	M	W	<- squash
i+4					F	D	X	M	W <- squash
i+5	trap ->					F	D	X	M W
i+6	trap handler ->						F	D	X M W

What Makes Pipelining Hard?

Complex Instructions

Complex Addressing Modes and Instructions

- **Address modes: Autoincrement causes register change during instruction execution**
 - Interrupts? Need to restore register state
 - Adds WAR and WAW hazards since writes are no longer the last stage.
- **Memory-Memory Move Instructions**
 - Must be able to handle multiple page faults
 - Long-lived instructions: partial state save on interrupt
- **Condition Codes**

Handling Multi-cycle Operations

A.1 What is Pipelining?

A.2 The Major Hurdle of Pipelining- Structural Hazards

- Data Hazards
- Control Hazards

A.3 How is Pipelining Implemented

A.4 What Makes Pipelining Hard to Implement?

A.5 Extending the MIPS Pipeline to Handle Multi-cycle Operations

**Multi-cycle instructions also
lead to pipeline complexity.**

**A very lengthy instruction
causes everything else in
the pipeline to wait for it.**

Multi-Cycle Operations

Floating Point

Floating point gives long execution time.

This causes a stall of the pipeline.

It's possible to pipeline the FP execution unit so it can initiate new instructions without waiting full latency. Can also have multiple FP units.

<i>FP Instruction</i>	<i>Latency</i>	<i>Initiation Rate</i>
Add, Subtract	4	3
Multiply	8	4
Divide	36	35
Square root	112	111
Negate	2	1
Absolute value	2	1
FP compare	3	2

Multi-Cycle Operations

Floating Point

Divide, Square Root take -10X to -30X longer than Add

- Interrupts?
- Adds WAR and WAW hazards since pipelines are no longer same length

	1	2	3	4	5	6	7	8	9	10	11
i	IF	ID	EX	MEM	WB						
I+1		IF	ID	EX	EX	EX	EX	MEM	WB		
I+2			IF	ID	EX	MEM	WB				
I+3				IF	ID	EX	EX	EX	EX	MEM	WB
I+4					IF	ID	EX	MEM	WB		
I+5						IF	ID	--	--	EX	EX
I+6							IF	--	--	ID	EX

Notes:
 I + 2: no WAW, but this complicates an interrupt
 I + 4: no WB conflict
 I + 5: stall forced by structural hazard
 I + 6: stall forced by in-order issue

Summary of Pipelining Basics

- **Hazards limit performance**
 - Structural: need more HW resources
 - Data: need forwarding, compiler scheduling
 - Control: early evaluation & PC, delayed branch, prediction
- **Increasing length of pipe increases impact of hazards; pipelining helps instruction bandwidth, not latency**
- **Interrupts, Instruction Set, FP makes pipelining harder**
- **Compilers reduce cost of data and control hazards**
 - Load delay slots
 - Branch delay slots
 - Branch prediction

Credits

I have not written these notes by myself. There's a great deal of fancy artwork here that takes considerable time to prepare.

I have borrowed from:

Wen-mei & Patel: <http://courses.ece.uiuc.edu/ece511/lectures/lecture3.ppt>

Patterson: <http://www.cs.berkeley.edu/~pattsrn/252S98/index.html>

Rabaey: (He used lots of Patterson material):

<http://bwrc.eecs.berkeley.edu/Classes/CS252/index.htm>

Katz: (Again, he borrowed heavily from Patterson):

<http://http.cs.berkeley.edu/~randy/Courses/CS252.F95/CS252.Intro.html>

Mark Hill: (Follows text fairly well): <http://www.cs.wisc.edu/~markhill/cs752/>

Summary

A.1 What is Pipelining?

A.2 The Major Hurdle of Pipelining-Structural Hazards

- Data Hazards
- Control Hazards

A.3 How is Pipelining Implemented

A.4 What Makes Pipelining Hard to Implement?

A.5 Extending the MIPS Pipeline to Handle Multi-cycle Operations