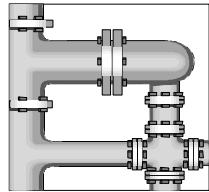


How Computers Work

Lecture 12

Introduction to Pipelining



How Computers Work Lecture 12 Page 1

A Common Chore of College Life



How Computers Work Lecture 12 Page 2

Propagation Times



$T_{pd_wash} =$ _____

$T_{pd_dry} =$ _____

How Computers Work Lecture 12 Page 3

Doing 1 Load

Step 1:



Step 2:







Total Time = _____

= _____

How Computers Work Lecture 12 Page 4

Doing 2 Loads




Combinational (Harvard) Method

Step 1:		
Step 2:		Total Time
Step 3:		= _____
Step 4:		= _____

How Computers Work Lecture 12 Page 5

Doing 2 Loads

Pipelined (MIT) Method

Step 1:		
Step 2:		Total Time
Step 3:		= _____
		= _____

How Computers Work Lecture 12 Page 6

Doing N Loads

- Harvard Method: _____
- MIT Method: _____

How Computers Work Lecture 12 Page 7

A Few Definitions

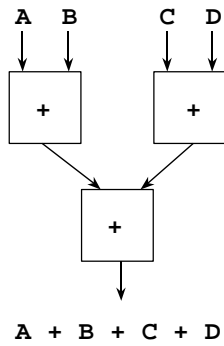
Latency: Time for 1 object to pass
through entire system.
(= _____ for Harvard laundry)
(= _____ for MIT laundry)

Throughput: Rate of objects going
through.
(= _____ for Harvard laundry)
(= _____ for MIT laundry)

How Computers Work Lecture 12 Page 8

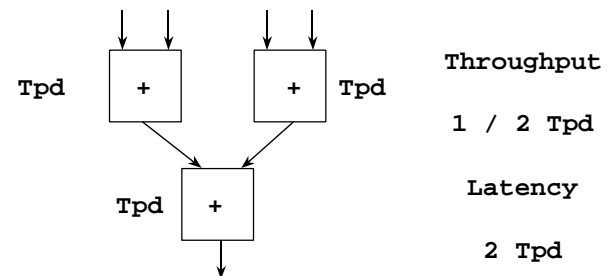
A Computational Problem

Add 4 Numbers:



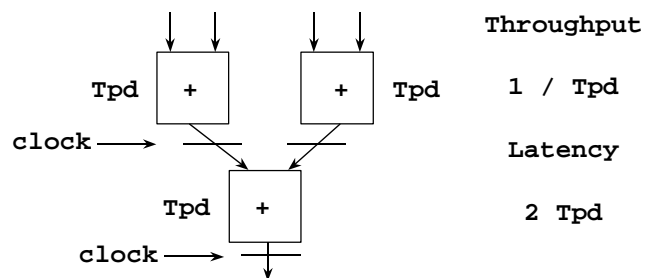
How Computers Work Lecture 12 Page 9

As a Combinational Circuit



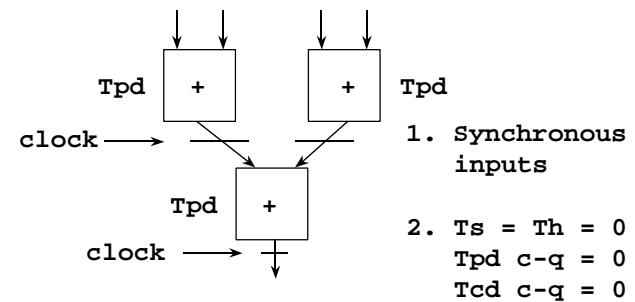
How Computers Work Lecture 12 Page 10

As a Pipelined Circuit



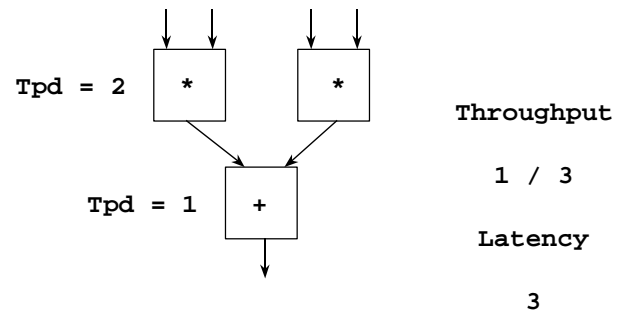
How Computers Work Lecture 12 Page 11

Simplifying Assumptions



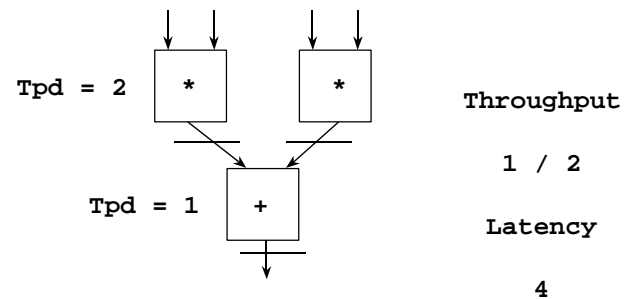
How Computers Work Lecture 12 Page 12

An Inhomogeneous Case (Combinational)



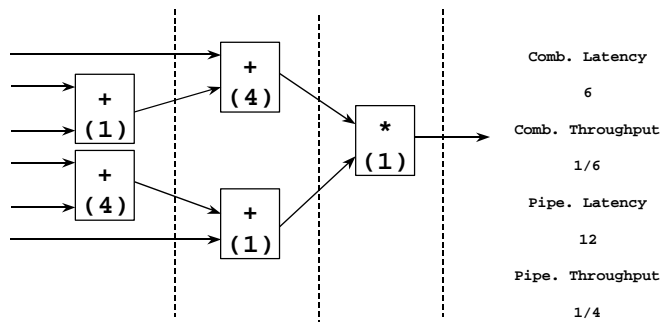
How Computers Work Lecture 12 Page 13

An Inhomogeneous Case (Pipelined)



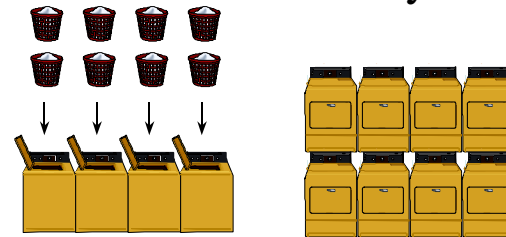
How Computers Work Lecture 12 Page 14

How about this one?



How Computers Work Lecture 12 Page 15

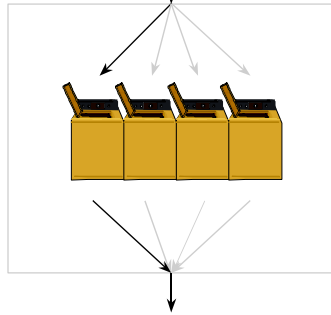
How MIT Students
REALLY do Laundry



Steady State Throughput = _____
Steady State Latency = _____

How Computers Work Lecture 12 Page 16

Interleaving (an alternative to Pipelining)



For N Units
of delay T_{pd} ,
steady state

Throughput

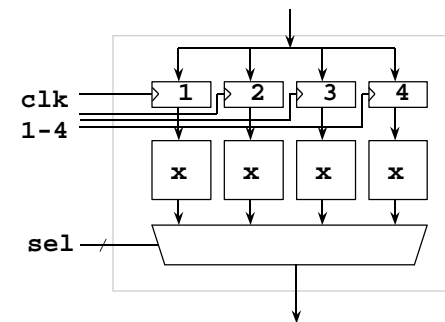
N / T_{pd}

Latency

T_{pd}

How Computers Work Lecture 12 Page 17

Interleaving Parallel Circuits



How Computers Work Lecture 12 Page 18

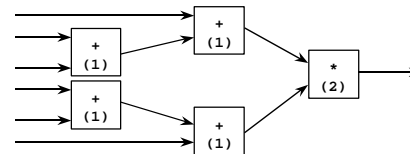
Definition of a Well-Formed Pipeline

- Same number of registers along path from any input to every computational unit
 - Insures that every computational unit sees inputs IN PHASE
- Is true (non-obvious) whenever the # of registers between all inputs and all outputs is the same.

How Computers Work Lecture 12 Page 19

Method for Forming Well-Formed Pipelines

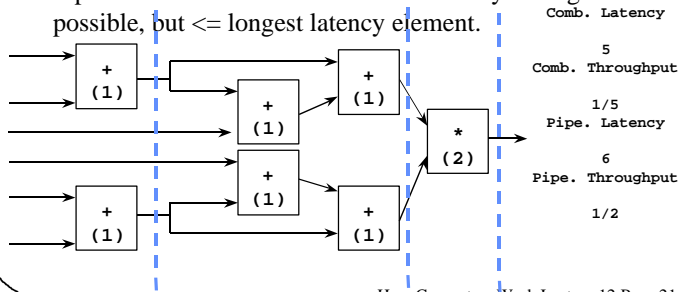
- Add registers to system output at will
- Propagate registers from intermediate outputs to intermediate inputs, cloning registers as necessary.



How Computers Work Lecture 12 Page 20

Method for Maximizing Throughput

- Pipeline around longest latency element
- Pipeline around other sections with latency as large as possible, but \leq longest latency element.



How Computers Work Lecture 12 Page 21

A Few Questions

- Assuming a circuit is pipelined for optimum throughput with 0 delay registers, is the pipelined throughput always greater than or equal to the combinational throughput?
 - A: Yes
- Is the pipelined latency ever less than combinational latency?
 - A: No
- When is the pipelined latency equal to combinational latency?
 - A: If contents of all pipeline stages have equal combinational latency

How Computers Work Lecture 12 Page 22

CPU Performance

MIPS = Millions of Instructions Per Second

Freq = Clock Frequency, MHz

CPI = Clocks per Instruction

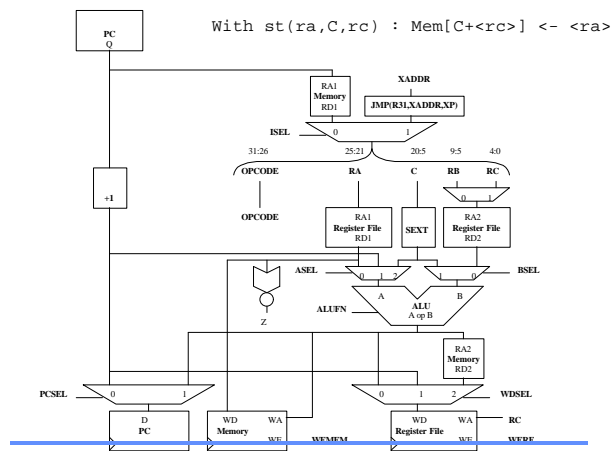
$$\text{MIPS} = \frac{\text{Freq}}{\text{CPI}}$$

To Increase MIPS:

1. DECREASE CPI.
 - RISC reduces CPI to 1.0.
 - CPI < 0? Tough... we'll see multiple instruction issue machines at end of term.
2. INCREASE Freq.
 - Freq limited by delay along longest combinational path; hence
 - **PIPELINING** is the key to improved performance through fast clocks.

How Computers Work Lecture 12 Page 23

Review: A Top-Down View of the Beta Architecture



How Computers Work Lecture 12 Page 24

Pipeline Stages

GOAL: Maintain (nearly) 1.0 CPI, but increase clock speed.

APPROACH: structure processor as 4-stage pipeline:

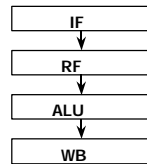
Instruction Fetch stage: Maintains PC, fetches one instruction per cycle and passes it to

Register File stage: Reads source operands from register file, passes them to

ALU stage: Performs indicated operation, passes result to

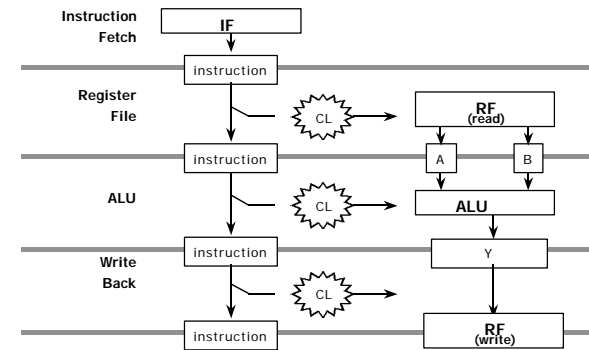
Write-Back stage: writes result back into register file.

WHAT OTHER information do we have to pass down the pipeline?

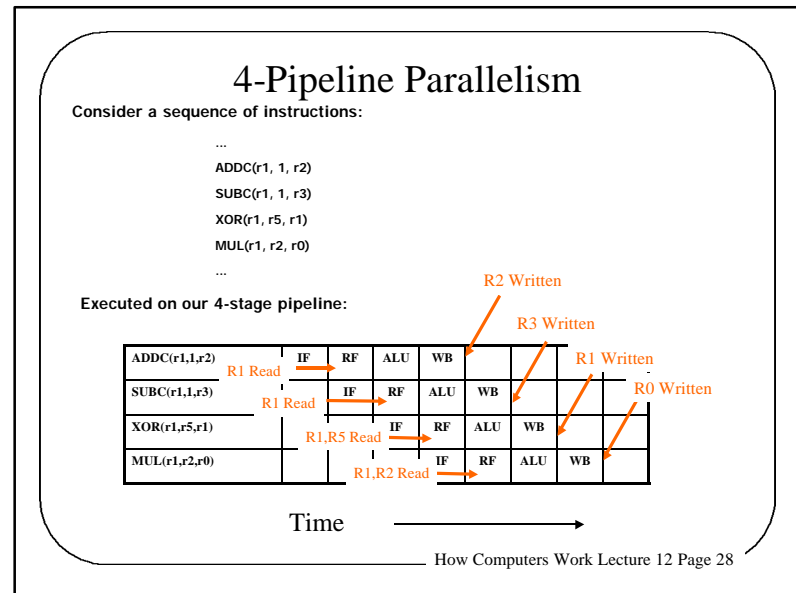
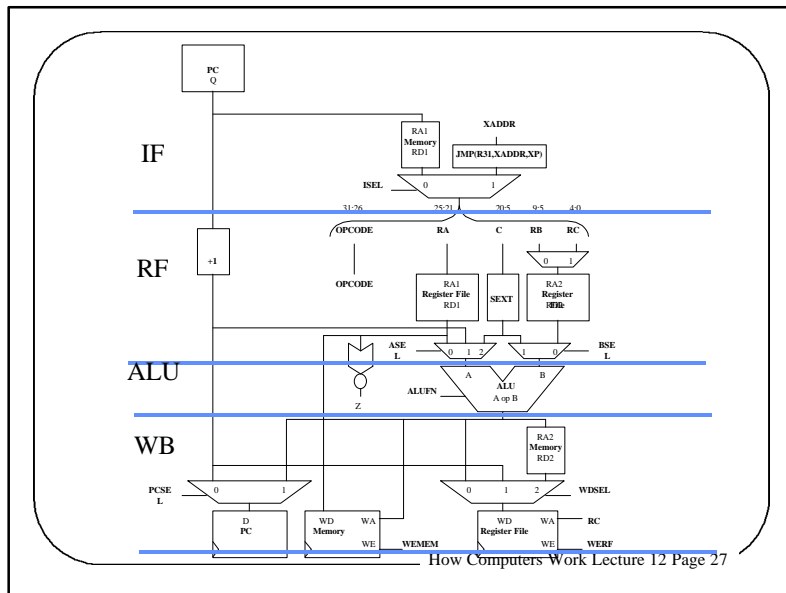


How Computers Work Lecture 12 Page 25

Sketch of 4-Stage Pipeline



How Computers Work Lecture 12 Page 26



Pipeline Problems

BUT, consider instead:

```

LOOP:  ADD(r1, r2, r3)
       CMPLC(r3, 100, r0)
       BT(r0, LOOP)
       XOR(r31, r31, r3)
       MUL(r1, r2, r2)
       ...
    
```

ADD(r1,r2,r3)	IF	RF	ALU	WB				
CMPLC(r3,100,r0)		IF	RF	ALU	WB			
BT(r0,LOOP)			IF	RF	ALU	WB		
XOR(r31,r31,r3)				IF	RF	ALU	WB	
MUL(r1,r2,r2)					IF	RF	ALU	WB

Time →

How Computers Work Lecture 12 Page 29

Pipeline Hazards

PROBLEM:

Contents of a register WRITTEN by instruction k is READ by instruction k+1... before its stored in RF! EG:

```

ADD(r1, r2, r3)
CMPLC(r3, 100, r0)
MULC(r1, 100, r4)
SUB(r1, r2, r5)
    
```

fails since CMPLC sees "stale" <r3>.

ADD(r1,r2,r3)	IF	RF	ALU	WB				
CMPLC(r3,100,r0)		IF	RF	ALU	WB			
BT(r0,LOOP)			IF	RF	ALU	WB		
XOR(r31,r31,r3)				IF	RF	ALU	WB	
MUL(r1,r2,r2)					IF	RF	ALU	WB

Time →

How Computers Work Lecture 12 Page 30

ADD(r1,r2,r3)	IF	RF	ALU	WB					
CMPLE(r3,100,r0)		IF	RF	ALU	WB				
			IF	RF	ALU	WB			
CMPLE(r3,100,r0)				IF	RF	ALU	WB		
					IF	RF	ALU	WB	

R3 Written

R3 Read

SOLUTIONS:

1. "Program around it".
 - ... document weirdo semantics, declare it a software problem.
 - Breaks sequential semantics!
 - Costs code efficiency.

EXAMPLE: Rewrite

```

ADD(r1, r2, r3)          ADD(r1, r2, r3)
CMPLE(r3, 100, r0)      MULC(r1, 100, r4)
MULC(r1, 100, r4)       as  SUB(r1, r2, r5)
SUB(r1, r2, r5)         CMPLE(r3, 100, r0)

```

HOW OFTEN can we do this?

SOLUTIONS:

2. Stall the pipeline.
 - Freeze IF, RF stages for 2 cycles, inserting NOPs into ALU IR...

ADD(r1,r2,r3)	IF	RF	ALU	WB					
NOP		IF	RF	ALU	WB				
NOP			IF	RF	ALU	WB			
CMPLE(r3,100,r0)				IF	RF	ALU	WB		
BT(r0,LOOP)					IF	RF	ALU	WB	
XOR(r31,r31,r3)						IF	RF	ALU	WB
MUL(r1,r2,r2)							IF	RF	ALU
								IF	RF

R3 Written

R3 Read

DRAWBACK: SLOW

SOLUTIONS:

3. Bypass Paths.

Add extra data paths & control logic to re-route data in problem cases.

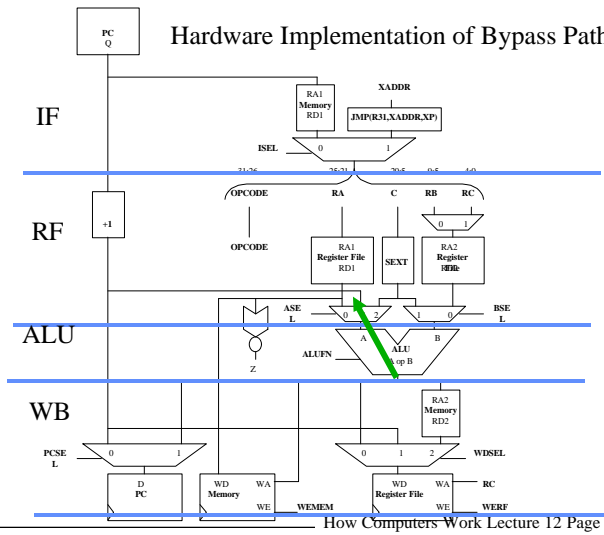
ADD(r1,r2,r3)	IF	RF	ALU	WB				
CMPLC(r3,100,r0)		IF	RF	ALU	WB			
BT(r0,LOOP)			IF	RF	ALU	WB		
XOR(r31,r31,r3)				IF	RF	ALU	WB	
MUL(r1,r2,r2)				IF	RF	ALU	WB	

<R1>+<R2> Produced

<R1>+<R2> Used

How Computers Work Lecture 12 Page 33

Hardware Implementation of Bypass Paths



How Computers Work Lecture 12 Page 34

Next Time:

- Detailed Design of
 - Bypass Paths + Control Logic
- What to do when Bypass Paths Don't Work
 - Branch Delays / Tradeoffs
 - Load/Store Delays / Tradeoffs
 - Multi-Stage Memory Pipeline

How Computers Work Lecture 12 Page 35