# Computer Architecture

## Pipelining and Instruction Level Parallelism–An Introduction

## Outline of This Lecture

**Introduction to the Concept of Pipelined Processor**
- Pipelined Datapath and Pipelined Control
- Pipeline Example: Instructions Interaction

**Pipeline Hazards**
- Forwarding
- Stalls

**Introduction to Instruction Level Parallelism**
- Superscalar, VLIW
- Out-of-order execution
- Branch Prediction
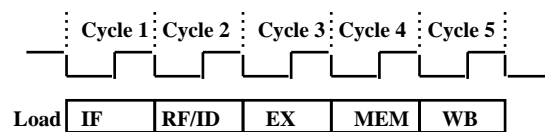- Future

# The Five Stages of Load

**IF: Instruction Fetch**
– **Fetch the instruction from the Instruction Memory**

**RF/ID: Registers Fetch and Instruction Decode**

**EX: Calculate the memory address**

**MEM: Read the data from the Data Memory**

**WB: Write the data back to the register file**

| | Cycle 1 | Cycle 2 | Cycle 3 | Cycle 4 | Cycle 5 |
|------|---------|---------|---------|---------|---------|
| Load | IF | RF/ID | EX | MEM | WB |

# Key Ideas Behind Pipelining

**Analogy–Grading the mid term exams:**
– **6 problems, six people grading the exam**
– **Each person grades ONE problem**
– **Pass exam to next person as soon as one finishes her part**
– **Assume each problem takes 0.15 hour to grade**
  • **Each individual exam still takes 0.9 hours to grade**
  • **But with 6 people, all exams can be graded much quicker:**
    – **100 exams: 90 hours, vs. 90 hrs x 6 = 540 hours**

**The load instruction has 5 stages:**
– **Five independent functional units to work on each stage**
  • **Each functional unit is used only once**
– **Another load can start as soon as 1st finishes its IF stage**
– **Each load still takes five cycles to complete**
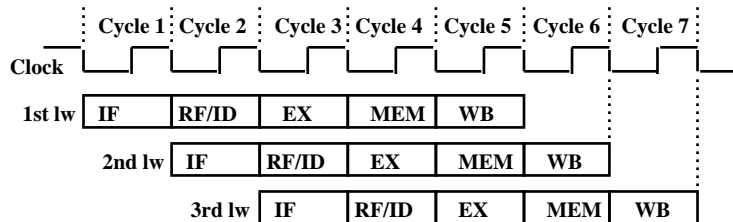– **The throughput, however, is much higher**

# Pipelining the Load Instruction

**Five independent functional units in pipeline are:**
- **Instruction Memory for the IF stage**
- **Register file's read ports for the RF/ID stage**
- **ALU for the EX stage**
- **Data Memory for the MEM stage**
- **Register File's Write port (bus W) for the WB stage**

**1 instruction enters the pipeline every cycle**
- **1 instruction comes out of pipeline (completes) every cycle**
- **"Effective" Cycles per Instruction (CPI) is 1**

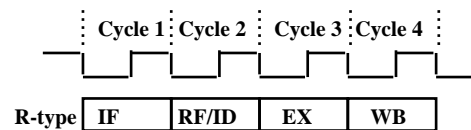| | Cycle 1 | Cycle 2 | Cycle 3 | Cycle 4 | Cycle 5 | Cycle 6 | Cycle 7 |
|---|---|---|---|---|---|---|---|
| **1st lw** | IF | RF/ID | EX | MEM | WB | | |
| **2nd lw** | | IF | RF/ID | EX | MEM | WB | |
| **3rd lw** | | | IF | RF/ID | EX | MEM | WB |

# Four Stages of R-type

**IF: Instruction Fetch**
- **Fetch the instruction from the Instruction Memory**

**RF/ID: Registers Fetch  and Instruction Decode**

**EX: ALU operates on the two register operands**

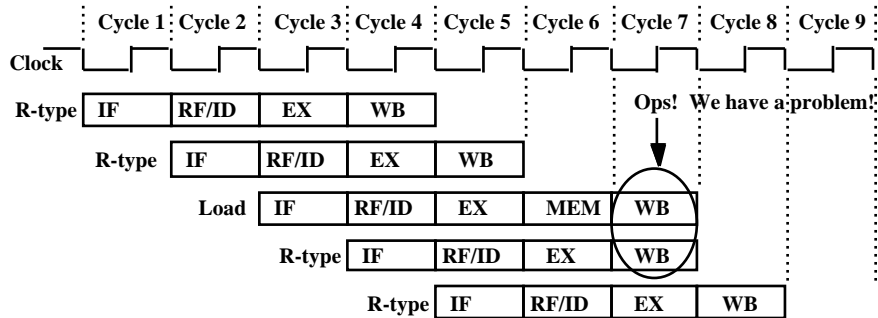**WB: Write the ALU output back to the register file**

| | Cycle 1 | Cycle 2 | Cycle 3 | Cycle 4 |
|---|---|---|---|---|
| **R-type** | IF | RF/ID | EX | WB |

# Pipelining R-type + Load

**We have a problem:**

– **Two instructions try to write to register file at same time!**

| | Cycle 1 | Cycle 2 | Cycle 3 | Cycle 4 | Cycle 5 | Cycle 6 | Cycle 7 | Cycle 8 | Cycle 9 |
|---|---|---|---|---|---|---|---|---|---|

Clock

Ops! We have a problem!

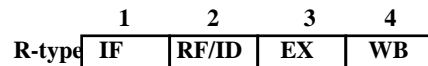| R-type | IF | RF/ID | EX | WB | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| R-type | | IF | RF/ID | EX | WB | | | | |
| Load | | | IF | RF/ID | EX | MEM | WB | | |
| R-type | | | | IF | RF/ID | EX | WB | | |
| R-type | | | | | IF | RF/ID | EX | WB | |

# Important Observation

**A functional unit can be used *once* per instruction**

**Each functional unit must be used at same stage for all instructions:**
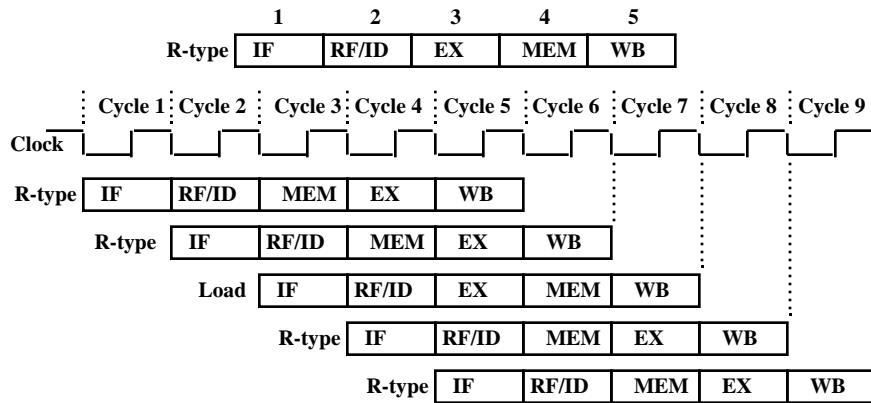
– **Load uses Register File's Write Port during its 5th stage**

–

    •

–

–

| | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| Load | IF | RF/ID | EX | MEM | WB |

–

– **R-type uses Register File's Write Port during its 4th stage**

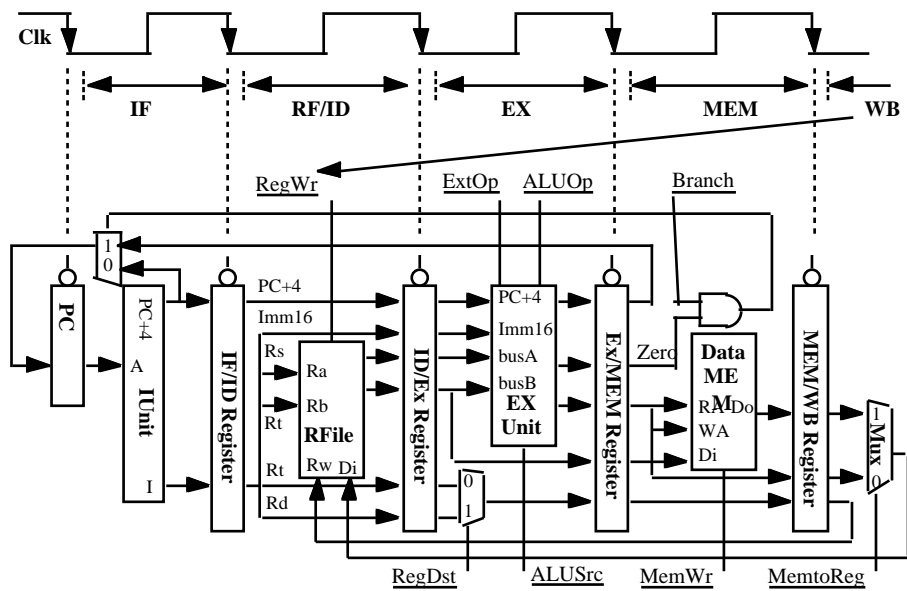| | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| R-type | IF | RF/ID | EX | WB |

# Solution: Delay R-type WB a Cycle

**Delay R-type's register write by one cycle:**
- **R-type instructions also use Reg File's write port at Stage 5**
- **MEM stage is a NOOP stage: nothing is being done**

|  | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| **R-type** | IF | RF/ID | EX | MEM | WB |

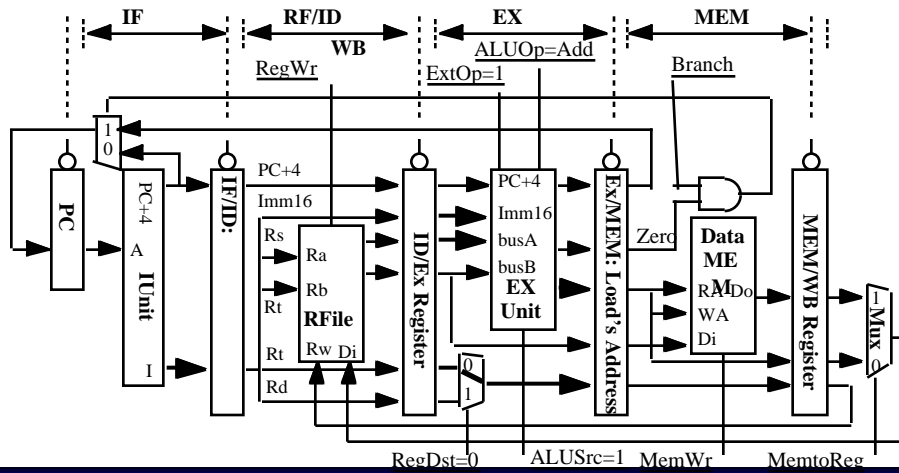| | Cycle 1 | Cycle 2 | Cycle 3 | Cycle 4 | Cycle 5 | Cycle 6 | Cycle 7 | Cycle 8 | Cycle 9 |
|---|---|---|---|---|---|---|---|---|---|
| **Clock** | | | | | | | | | |
| **R-type** | IF | RF/ID | MEM | EX | WB | | | | |
| **R-type** | | IF | RF/ID | MEM | EX | WB | | | |
| **Load** | | | IF | RF/ID | EX | MEM | WB | | |
| **R-type** | | | | IF | RF/ID | MEM | EX | WB | |
| **R-type** | | | | | IF | RF/ID | MEM | EX | WB |

# A Pipelined Datapath

# How About Control Signals?

**Control Signals at Stage N  = Func (Instr. at Stage N)**

– **N  =  EX, MEM, or WB**

**Example: Controls Signals at EX Stage**
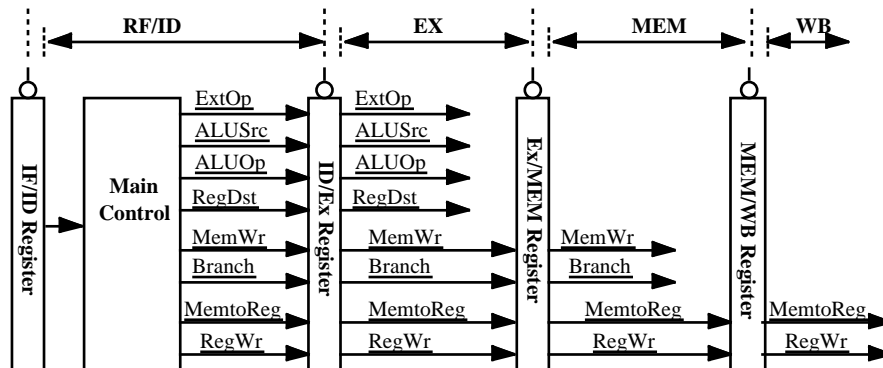
– **Func(Load's EX)**

# Pipeline Control

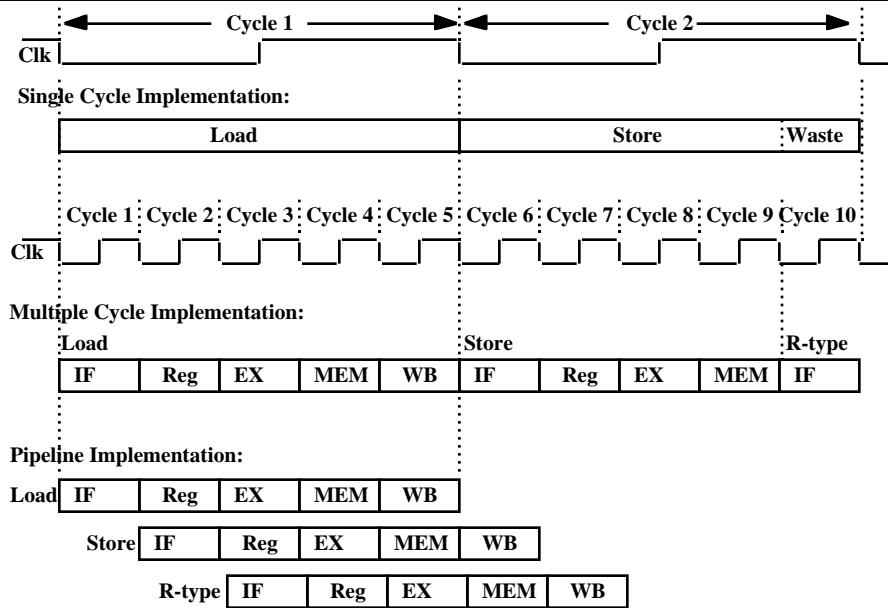**The Main Control generates the control signals during RF/ID**

– **Control signals for EX (ExtOp, ALUSrc, ...) used 1 cycle later**

– **Control signals for MEM (MemWr, Branch) used 2 cycles later**

– **Control signals for WB (MemtoReg MemWr) used 3 cycles later**
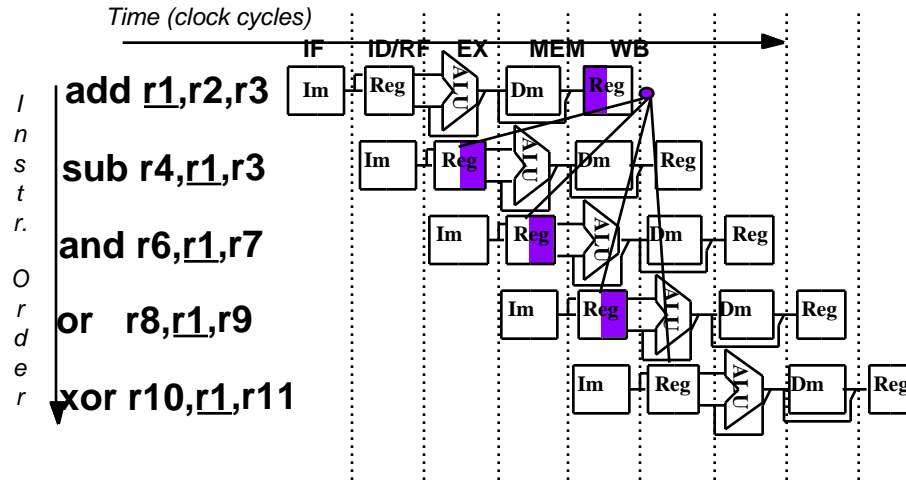
# Single Cycle, Multi-Cycle, Pipelined

# Hazards–Challenge to Pipelining

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

- **structural hazards: HW cannot support this combination of instructions**
  - **earlier case of load and R-typ like a structural hazard, but normally cannot fix by retiming instruction.**
- **data hazards: instruction depends on result of prior instruction still in the pipeline**
- **control hazards: pipelining of branches & other instructionsCommon solution is to stall the later part of the pipeline until the hazard pipeline**
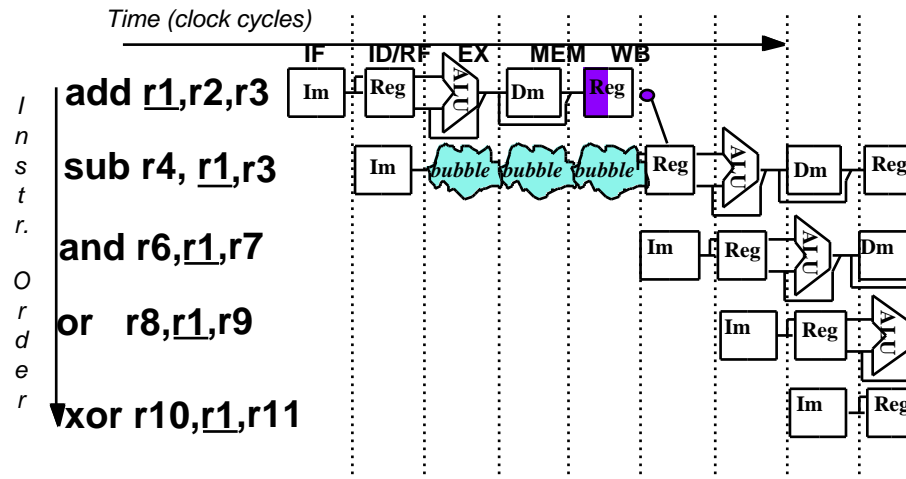
# Data Hazard on r1

**Dependencies backwards in time are hazards**

*Time (clock cycles)*

# HW Stalls to Resolve Hazard

**Dependencies backwards in time are hazards**
 – **eliminate "reverse time" by a stall**

*Time (clock cycles)*
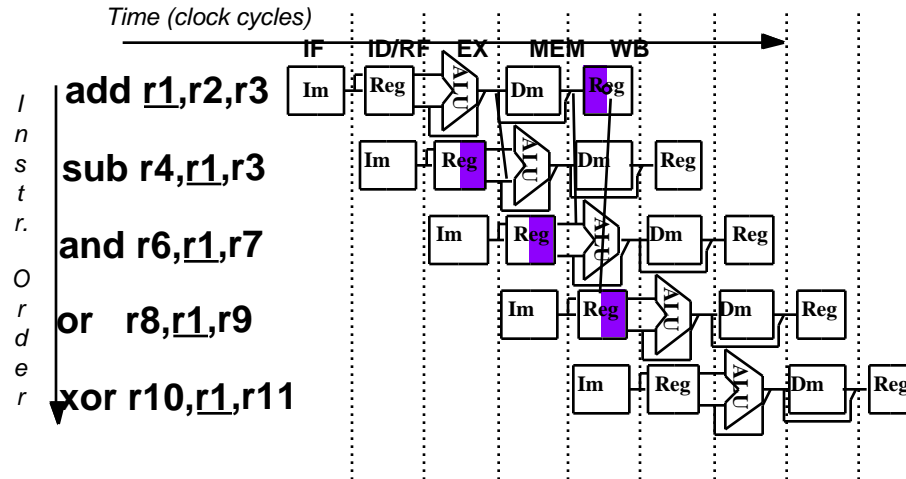
# Insight: Data is available!

**Pipeline registers already contain needed data**
- "Forward" the data to the appropriate unit

*Time (clock cycles)*

IF   ID/RF   EX   MEM   WB

*Instr. Order*

add <u>r1</u>,r2,r3

sub r4,<u>r1</u>,r3

and r6,<u>r1</u>,r7

or   r8,<u>r1</u>,r9

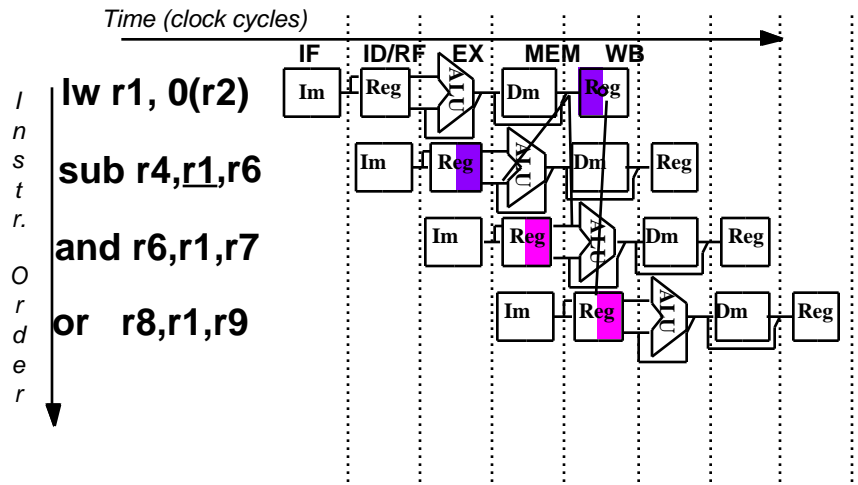xor r10,<u>r1</u>,r11

# HW for "Forwarding" (Bypassing)

**Increase multiplexors to add paths from registers**
- **Assumes register read during write gets new value (otherwise more results to be forwarded)**

ID/EX                    EX/MEM                    MEM/WB

Zero ?

Mux

ALU

Mux

Data Memory

# Forwarding Cannot Hide All Hazards

Time (clock cycles)

|  | IF | ID/RF | EX | MEM | WB |

lw r1, 0(r2)

sub r4,r1,r6

and r6,r1,r7

or  r8,r1,r9

*Instr. Order*

# Option: HW Stalls to Resolve Hazard

**"Interlock": checks for hazard & stalls**

Time (clock cycles)

|  | IF | ID/RF | EX | MEM | WB |

lw r1, 0(r2)

stall

*bubble bubble bubble bubble*

sub r4,r1,r3

and r6,r1,r7

or  r8,r1,r9

*Instr. Order*

# Option: SW resolves hazard

**SW inserts independent instuctions**

– **Worst case: performance no better/worse**

*Time (clock cycles)*

```
           IF   ID/RF  EX    MEM   WB
```

*Instr. Order*

**lw r1, 0(r2)**

**unrelated instruction**

**sub r4,r1,r3**

**and r6,r1,r7**

**or  r8,r1,r9**

# Control Hazard on Branches

Time (in Clock Cycles)

CC 1   CC 2   CC 3   CC 4   CC 5   CC 6   CC 7   CC 8   CC 9

Program Execution Order (in instructions)

40 beq $1,$3,36

44 and $12,$2,$5

48 or $13,$6,$2

52 add $14,$2,$2

80 ld $4,$7,100

# Hazards on Branches

*Time (clock cycles)*

| | IF | ID/RF | EX | MEM | WB |
|---|---|---|---|---|---|

**beq r1,r2,L**

**sub r4,r1,r3**

**and r6,r2,r7**

**or   r8,r7,r9**

**L: add r1,r2,r1**

**Stall for two cycles on every branch!**

# Branch Stall Impact

**CPI Impact:**
- If CPI = 1, 30% branch, Stall 2 cycles => new CPI = 1.6!

**Reducing the branch penalty**
- MIPS branch already more aggressive than most
- limited eq/neq allows us to determine branch condition early (after EX), rather than later (e.g., after MEM)
- doing better
  - use separate comparator rather than ALU and move branch decision to RF (hard!!!)
  - reduces penalty to 1 cycle

**Going further**
- Variety of techniques:
  - separating branch and destination
  - separating branch condition and branch decision
  - hardware prediction of branche

# When is pipelining hard?

Interrupts: 5 instructions executing in 5 stage pipeline
- How to stop the pipeline?
- Restrart?
- Who caused the interrupt?

**Stage**          **Problem interrupts occurring**

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

Load with data page fault, Add with instruction page fault?

Solution 1: interrupt vector/instruction, restart everything incomplete

# First Generation RISC Pipelines

All instructions: 1 pipeline order ("static schedule").

Register write in last stage + reads performed in first stage after issue.
- Simpliy/eliminate hazards

Memory access in stage 4
- Avoid all memory hazards

Control hazards use delayed branch (with fast path)

RAW hazards use bypass, except on load results
- Load resolved by delayed load or stall

Good pipeline performance at little cost/complexity.

# Summary of Pipelining Basics

**Speed Up = Pipeline Depth**

**Hazards limit performance on computers:**
- **structural: need more HW resources**
- **data: need forwarding, compiler scheduling**
- **control: early evaluation & PC, delayed branch, prediction**

**Increasing length of pipe increases  hazards**
- **since pipelining helps instruction bandwidth, not latency**

**Compilers can reduce cost of data & control hazards**
- **load delay slots**
- **branch delay slots**

**Exceptions (also FP, ISA) make pipelining harder**

# Advanced Pipelining

**Pipelining exploits parallelism among instructions by overlapping them**
- **Called Instruction Level Parallelism (ILP)**
- **Limited by a variety of things:**
  - **parallelism in the program**
  - **compiler technology in exposing parallelism**
  - **functional unit capability: how many ovrlapping instructions**
  - **ability of hardware to find instructions to run in parallel**
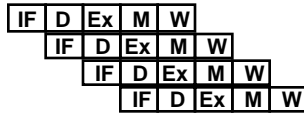
**Exploiting ILP is "hot topic" in processor design:**
- **Lots of different approaches**
  - **Multiple instuctions/cycle**
    - **compiler vs. HW for scheduling instructions**
  - **Both architecture approaches and compiler approaches**
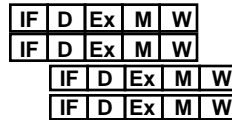
# Exploiting Available ILP

**Technique**

**Pipelining**

| IF | D | Ex | M | W |
| IF | D | Ex | M | W |
| IF | D | Ex | M | W |
| IF | D | Ex | M | W |

**Super-scalar**

**Issue multiple scalar**

**instructions per cycle**

**VLIW**

**Each instruction specifies**

**multiple scalar operations**

**HW Limitation**

**Issue rate,**
**FU stalls,**
**FU depth**

**Hazard resolution**

**Packing**

# Easy Superscalar

I-Cache

Int Reg | Inst Issue and Bypass | FP Reg

Int Unit | Load / Store Unit | FP Add | FP Mul

D-Cache

**Issue integer and FP operations in parallel!**

– **potential hazards?**

– **expected speedup?**

– **what combinations of instructions make sense?**

# Issuing Multiple Instruction/ Cycle

**Superscalar: 2 instructions, 1 FP & 1 anything else**
- Fetch 64-bits/clock cycle; Int on left, FP on right
- Can only issue 2nd instruction if 1st instruction issues
- More ports for FP registers to do FP load & FP op in a pair

| *Type* | *Pipe* | *Stages* | | | | |
|---|---|---|---|---|---|---|
| Int. instruction IF | ID | EX | MEM | W | | |
| FP instruction  IF | ID | EX | MEM | WB | | |
| Int. instruction | IF | ID | EX | MEM | WB | |
| FP instruction | IF | ID | EX | MEM | WB | |
| Int. instruction | | IF | ID | EX | MEM | WB |
| FP instruction | | IF | ID | EX | MEM | WB |

**1 cycle load delay expands to 3 instruction in SS**
- instruction in right half can't use result, nor can either instruction in next slot

# Dynamic Branch Prediction

**Predict direction of branches on past behavior**
- keep a cache of branch behavior, look up prediction

**Performance = f(accuracy, cost of misprediction)**

**Branch prediction buffer:**
- lower bits of PC address index table of 1-bit values
- says whether or not branch taken last time
- evaluate actual banch condition, if prediction incorrect:
  - recover by flushing pipeline, restarting fetch
  - reset prediction

# Speculative Superscalar Execution

**Get all available parallelism**

- **across branches**
- **in face cache misses**
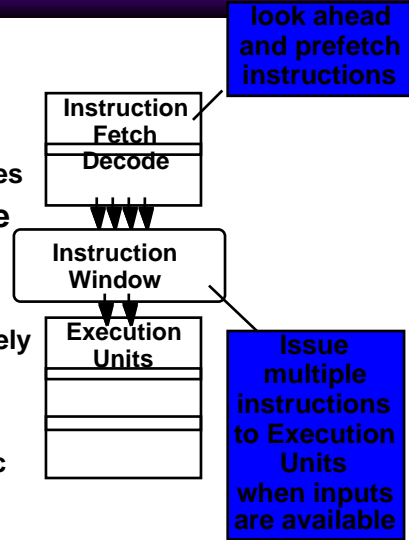- **limited only by data dependences**

**Goal: resources and available bandwidth are only HW limit**

**Branch prediction**

- **execute instructions speculatively**

**Hazard detection and aggressive resolution**

- **out-of-order execution (dynamic scheduling)**
- **in-order completion**
  - **Exception handling easier**
  - **handles incorrect speculation**

**look ahead and prefetch instructions**

```
Instruction
Fetch
Decode

Instruction
Window

Execution
Units
```

**Issue multiple instructions to Execution Units when inputs are available**

# Variety of Modern Microprocessor

| Processor | Instruction Completion Rate | Scheduling of pipeline | Branch prediction |
|-----------|------------------------------|------------------------|-------------------|
| PowerPC 604 | 4 | Dynamic, nonspeculative | HW |
| MIPS R10000 | 4 | Dynamic, speculative | HW |
| Pentium II | 4 | Dynamic, nonspeculative | HW |
| UltraSPARC | 4 | Static | HW |
| Merced | ? | Static? | Static? |

# Limits to Multi-Issue Machines

**Inherent limitations of ILP**

- 1 branch in 5 => 5-way VLIW busy?
- Latencies of units => many operations must be scheduled
- Need about Pipeline Depth x No. Functional Units of independentDifficulties in building HW
- Duplicate FUs to get parallel execution
- Increase ports to Register File (3 x integer/FP rate)
- Increase ports to memory
- Decoding challenge and impact on clock rate, pipeline depth

**Limitations specific to either SS or VLIW implementation**

- Decode issue in SS
- VLIW code size: unroll loops + wasted fields in VLIW
- VLIW lock step => 1 hazard & all instructions stall
- VLIW & binary compatibility

# Summary

**Instruction Level Parallelism in SW or HW**

**Loop level parallelism is easiest to see**

**SW dependencies/Compiler sophistication determine if compiler can unroll loops**

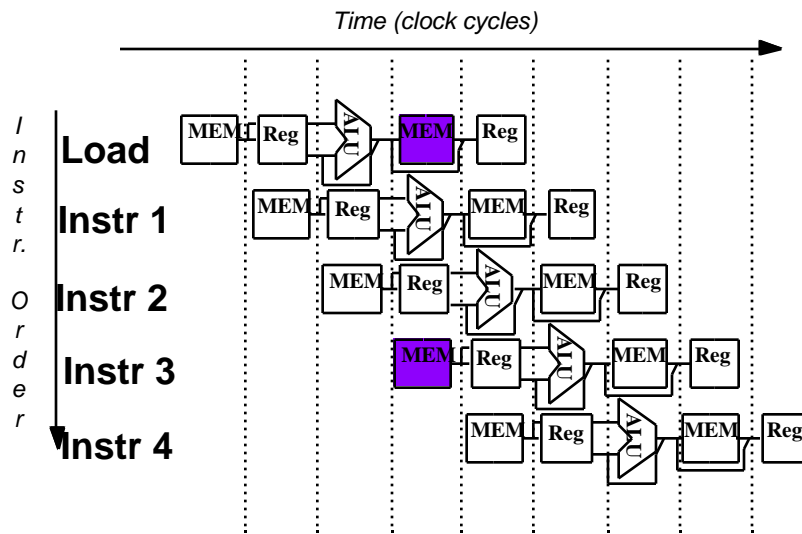**SW Scheduling**

**HW scheduling**

**Branch Prediction**
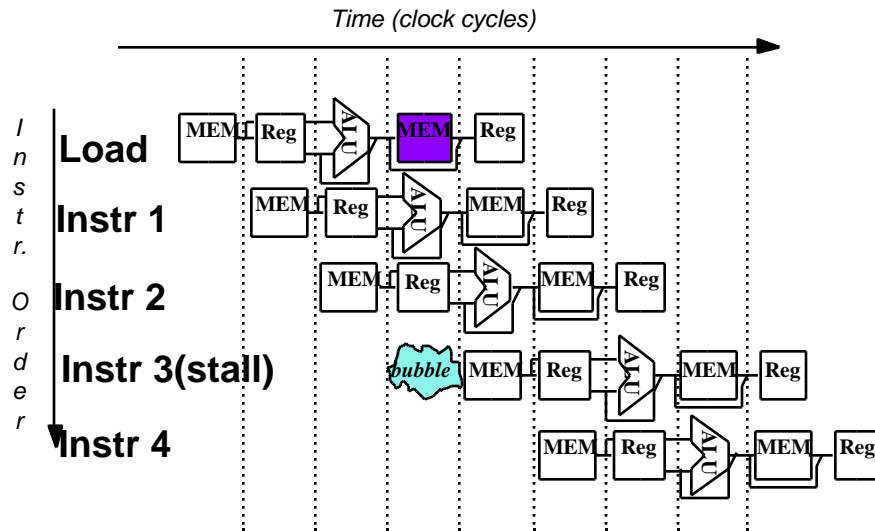
**SuperScalar and VLIW**

- CPI < 1
- Dynamic issue vs. Static issue
- More instructions issue/clock, larger penalty of hazards

**Future? Stay tuned…**

# Single Memory=>Structural Hazard

*Time (clock cycles)*

Instr. Order

Load — MEM Reg ALU **MEM** Reg

Instr 1 — MEM Reg ALU MEM Reg

Instr 2 — MEM Reg ALU MEM Reg

Instr 3 — **MEM** Reg ALU MEM Reg

Instr 4 — MEM Reg ALU MEM Reg

# Stall to resolve Structural Hazard

*Time (clock cycles)*

Instr. Order

Load — MEM Reg ALU **MEM** Reg

Instr 1 — MEM Reg ALU MEM Reg

Instr 2 — MEM Reg ALU MEM Reg

Instr 3(stall) — *bubble* MEM Reg ALU MEM Reg

Instr 4 — MEM Reg ALU MEM Reg

# Duplicate to Resolve Hazard

**Separate Instruction Cache (Im) & Data Cache (Dm)**