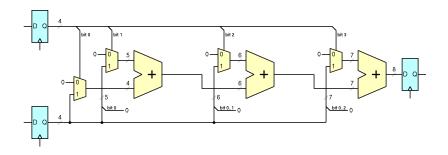
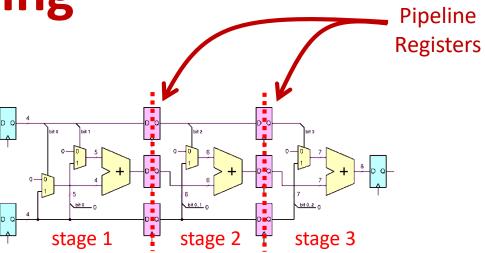
CS-200 Computer Architecture

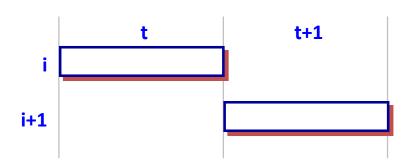
Part 4c. Instruction Level Parallelism Pipelining

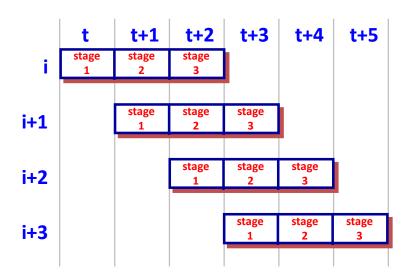
Paolo lenne <paolo.ienne@epfl.ch>

Pipelining

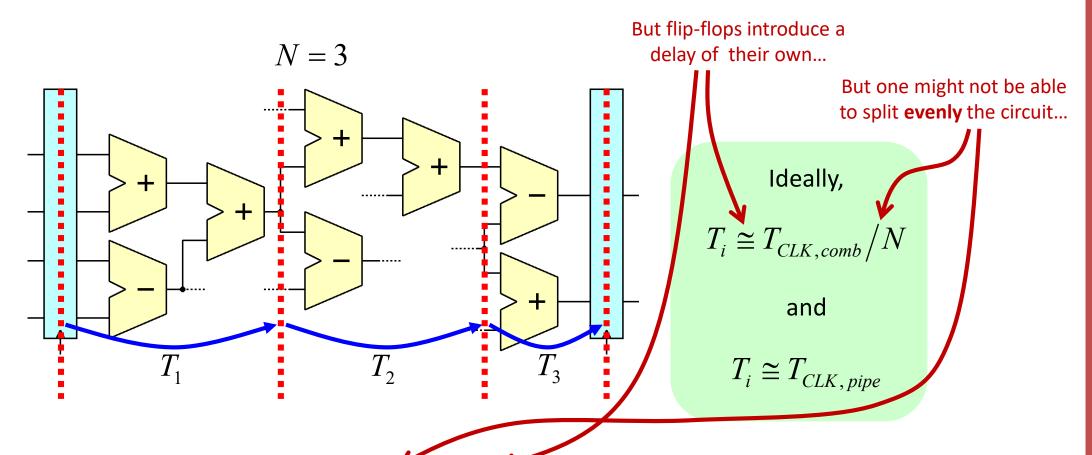








Practical Pipelining



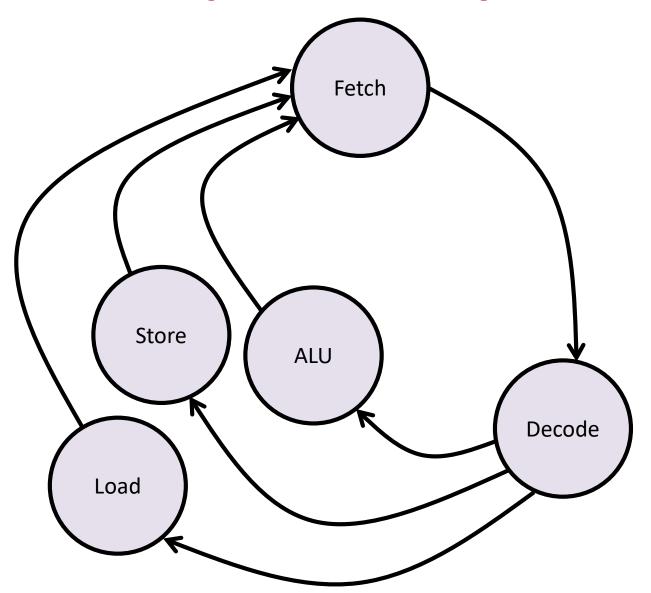
• Latency
$$\lambda_{pipe} = N \cdot \max_{i=0..N-1} (T_i + T_{FF}) = N \cdot T_{CLK, pipe} = N / f_{pipe} > \lambda_{orig}$$

• Throughput
$$\phi_{pipe} = 1/\max_{i=0..N-1} (T_i + T_{FF}) = f_{pipe}$$

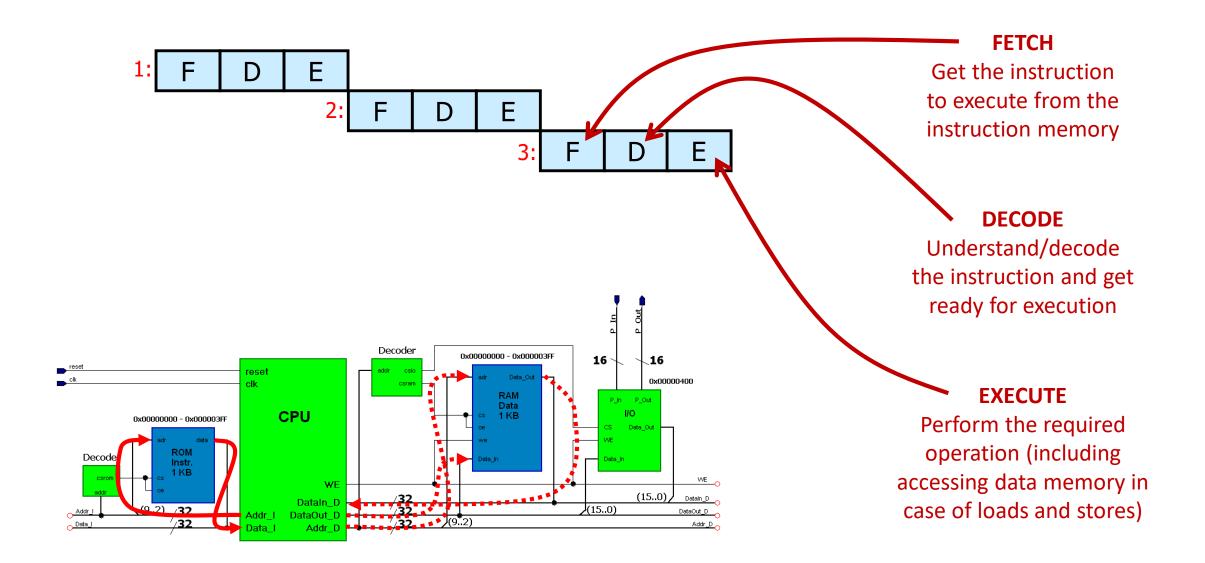
Pipelining for Processors?

- Pipelining useful only if the activity in object needs to be repeated many time
 - We have plenty of instructions to execute!
- Pipelining needs to split the single activity in object into many subactivities
 - We have a good logical split into fetch, decode, execute, etc.

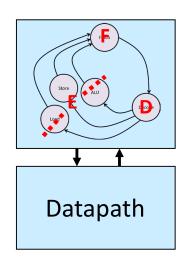
Example: A Simple Multi-Cycle Processor

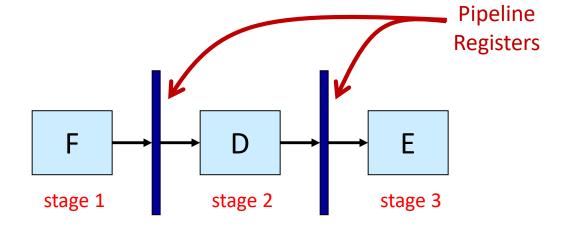


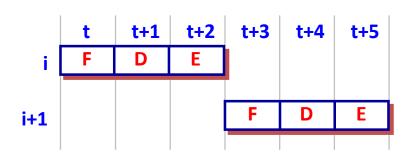
Example: A Simple Schedule

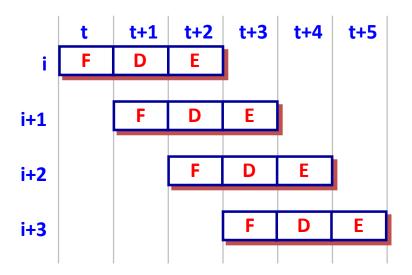


Pipelining the Processor?









No Hardware Reuse across Stages

- In a multicycle processor, some hardware components may be shared across **states**:
 - FETCH typically requires an adder to increment the program counter
 - EXECUTE naturally needs an ALU
 - They are never used at the same time, so the ALU can be used to increment the program counter
- In a pipelined processor, there cannot be sharing across stages, in general:
 - All stages active all the time!
 - Hardware needs to be replicated where appropriate

Two Main Problems

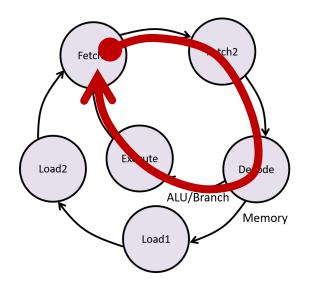
1. CISC vs. RISC

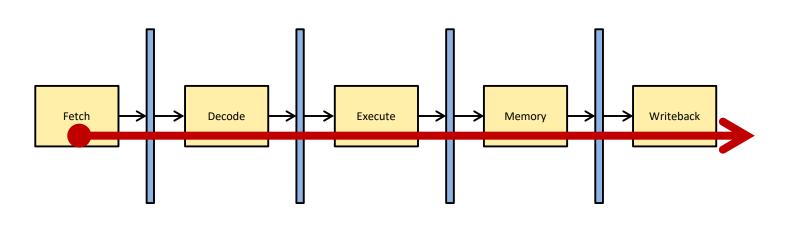
- Can we build equally well a pipeline for a Complex Instruction Set Computer as for a Reduced Instruction Set Computer?
- What does that mean?!

- 2. Instructions are **not** independent
 - Can we execute code correctly?

FSM vs. Pipeline

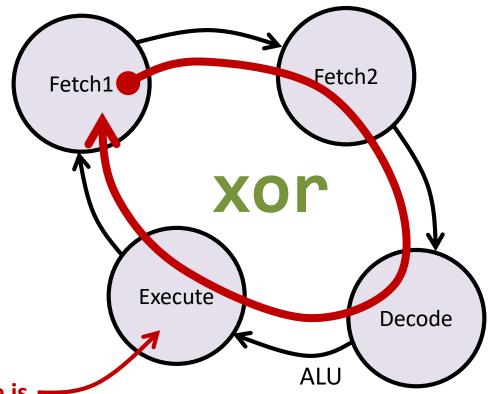
- Any path through our FSM represents the sequence of necessary steps for the execution of an instruction
- The ordered path through the pipeline is the sequence of all possible steps for the execution of any instruction





Adding Instructions to a Multi-Cycle Processor

How do we support add?

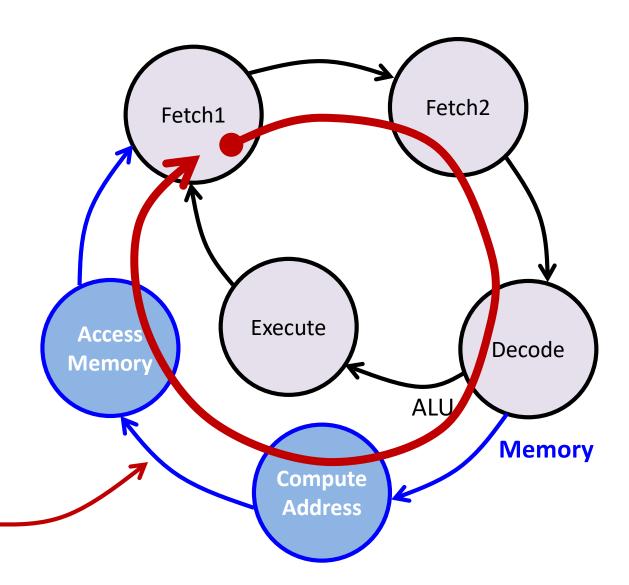


If the sequence of steps to execute is the same (fetch the instruction, read registers, use the ALU, save result in the register), we just need to an ALU that can perform additions

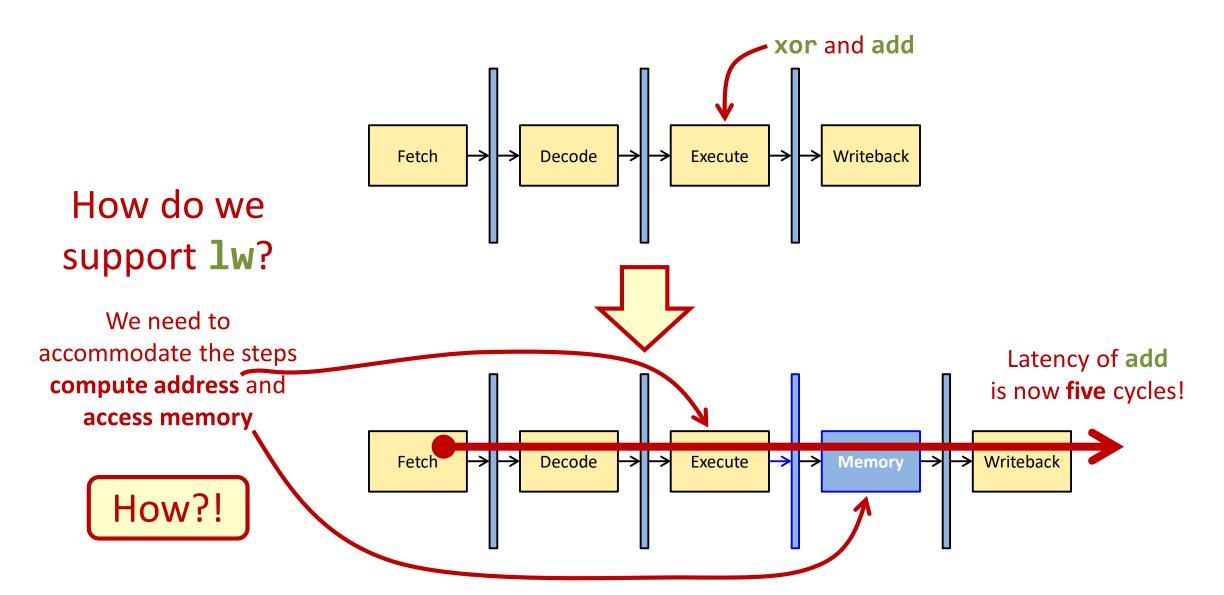
Adding Instructions to a Multi-Cycle Processor

How do we support 1w?

The sequence of steps
to execute is **not** the
same (**compute address and access memory**) so
we add new steps and a
(partially) new path

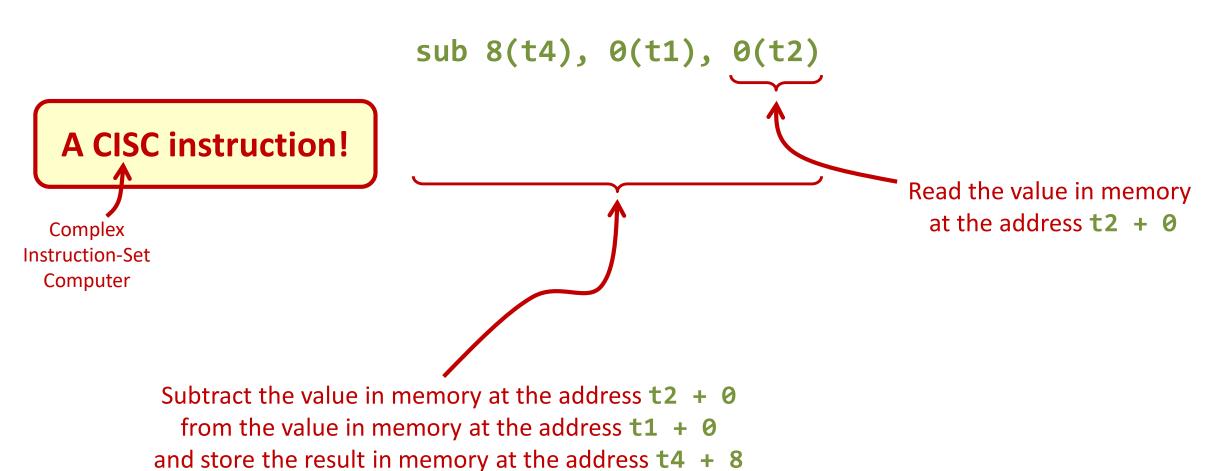


Adding Instructions to a Pipelined Processor



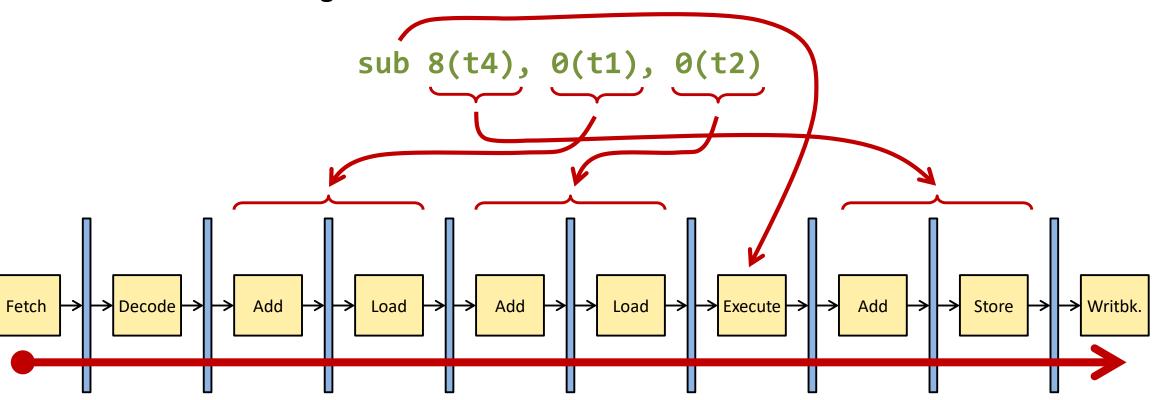
The Importance of the ISA

Imagine that we want to have an instruction



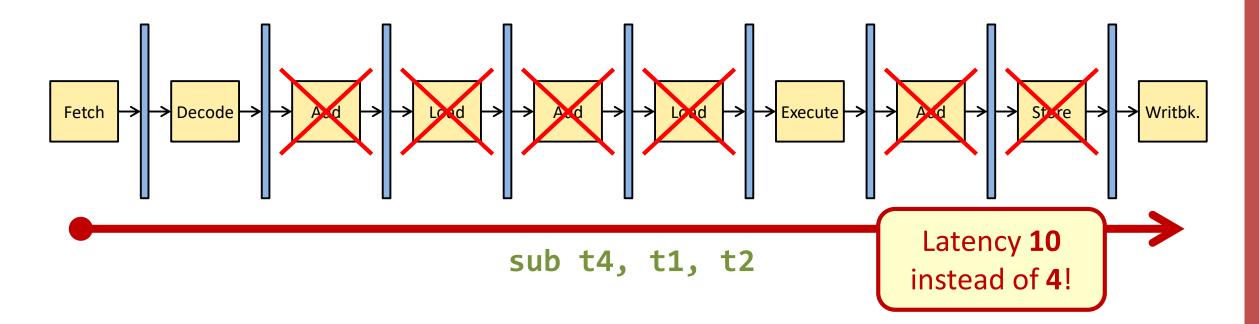
The Importance of the ISA

Imagine that we want to have an instruction



The Importance of the ISA

Imagine that we want to have an instruction



Reduced Instruction-Set Computer

 Instead of imposing a huge penalty to every simple instruction by making complex instructions possible, let's have only similarly simple instructions and build our programs with those:

```
sub 8(t4), 0(t1), 0(t2)
lw t3, 0(t1)
lw t5, 0(t2)
sub t3, t3, t5
sw t3, 8(t4)
```

It turns out that it is not the only way to go, but it is a good one and we will follow it...

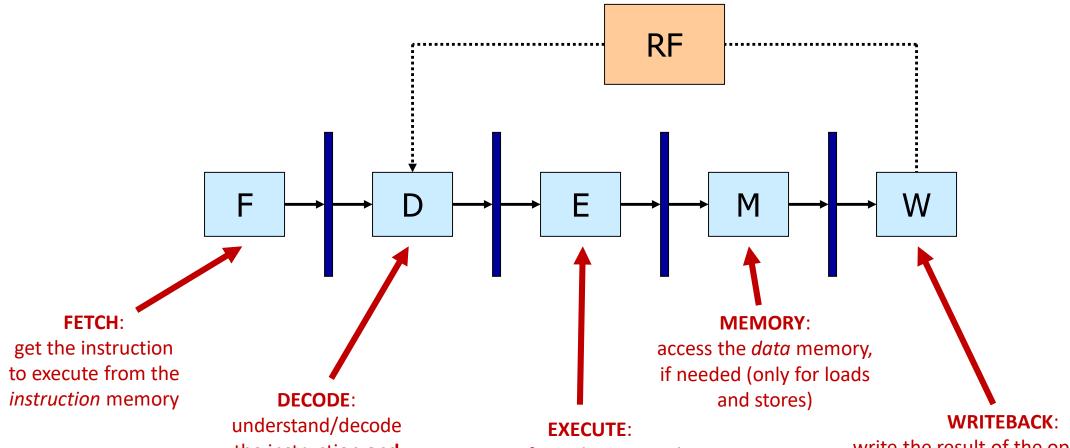
Two Main Problems

1. CISC vs. RISC

- Can we build equally well a pipeline for a Complex Instruction Set Computer as for a Reduced Instruction Set Computer?
- What does that mean?!

- 2. Instructions are **not** independent
 - Can we execute code correctly?

Simple 5-Stage MIPS Pipeline



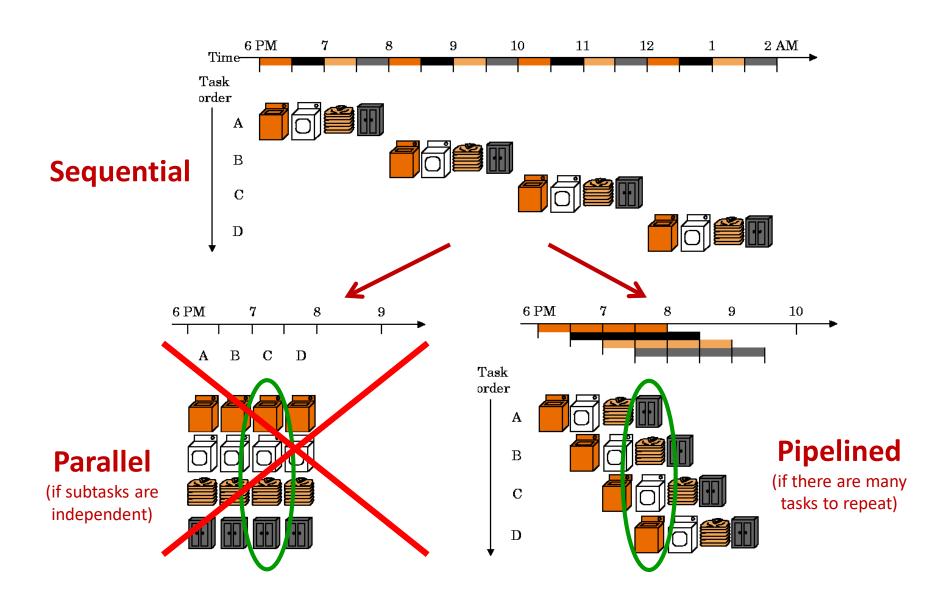
the instruction and

obtain arguments from the register file

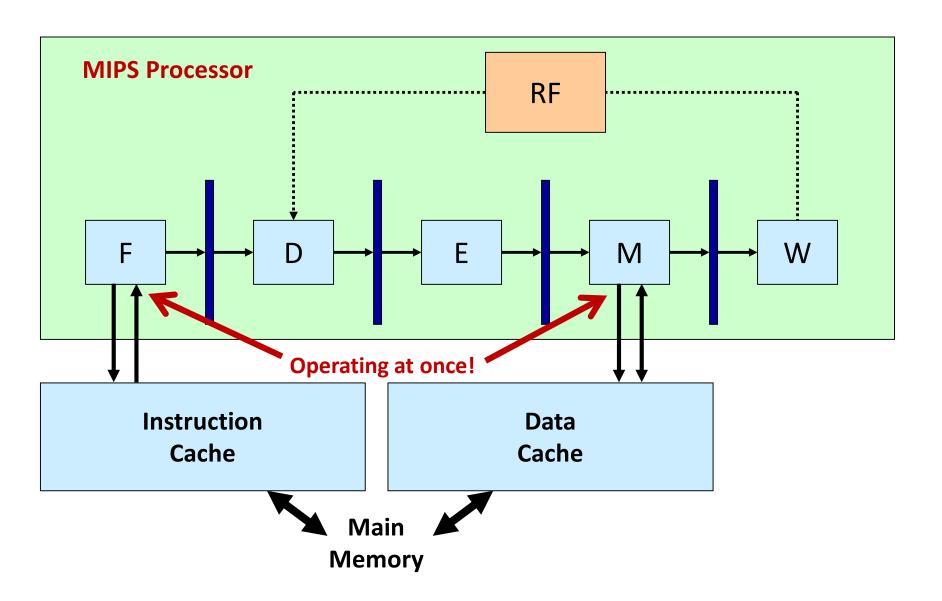
perform the required operation in the ALU (including address calculations for loads and stores)

write the result of the operation, if any, to the register file (either produced from the ALU or received from the data memory)

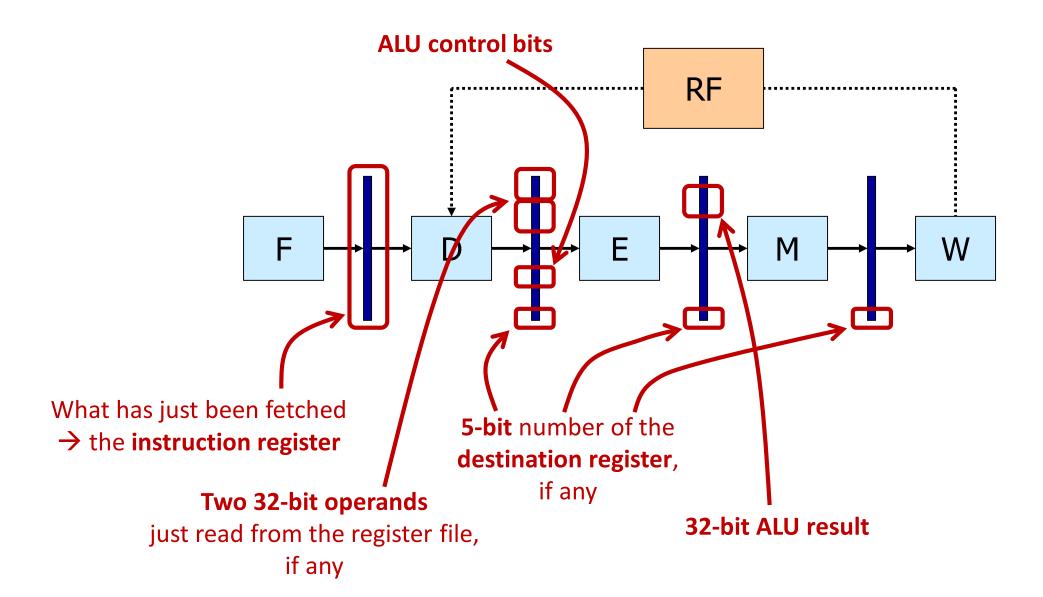
The Laundry Metaphor

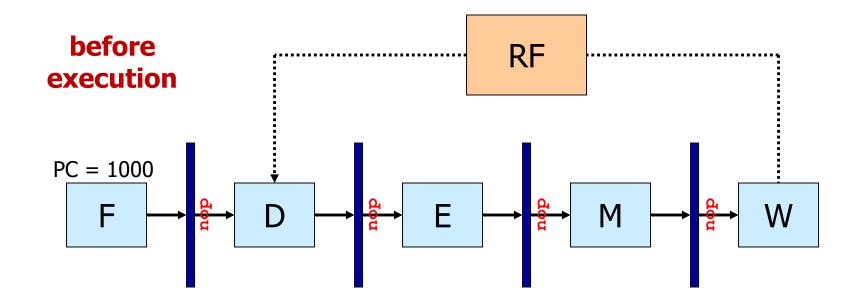


Two Distinct Memory Interfaces



What Is in the Pipeline Registers?





Empty pipeline: all stages at **nop** (no operation)

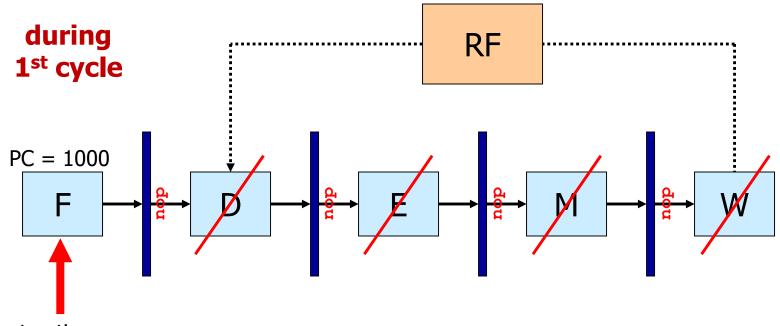
```
1000: add $r2, $r0, $r1

1004: sub $r5, $r3, $r4

1008: sw $r6, 50($r7)

1012: lw $r9, 20($r8)

1016: mul $r12, $r10, $r11
```



From instruction memory:

add \$r2, \$r0, \$r1

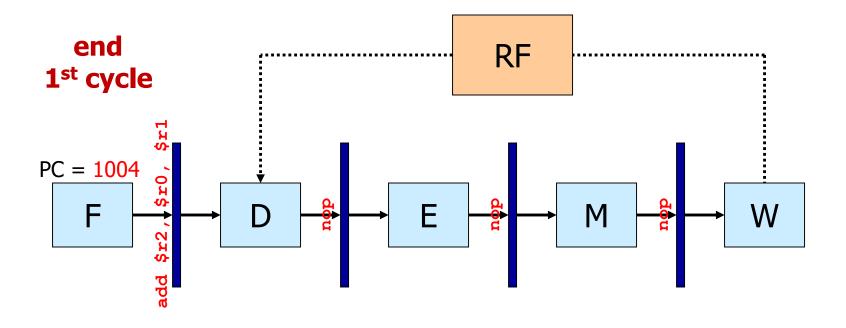
```
1000: add $r2, $r0, $r1

1004: sub $r5, $r3, $r4

1008: sw $r6, 50($r7)

1012: lw $r9, 20($r8)

1016: mul $r12, $r10, $r11
```



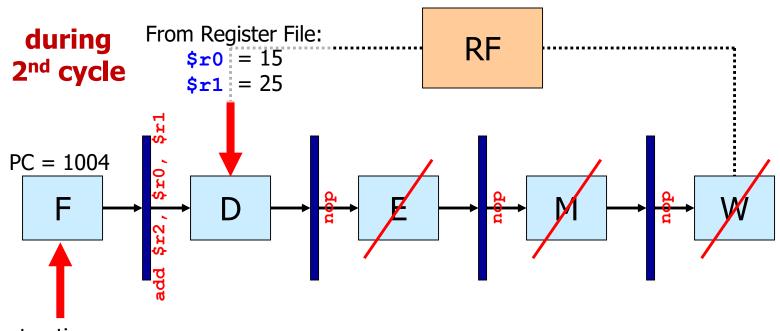
```
1000: add $r2, $r0, $r1

1004: sub $r5, $r3, $r4

1008: sw $r6, 50($r7)

1012: lw $r9, 20($r8)

1016: mul $r12, $r10, $r11
```



From instruction memory:

```
sub $r5, $r3, $r4
```

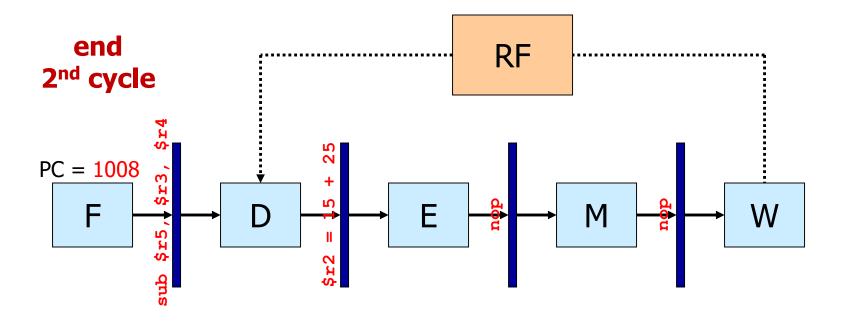
```
1000: add $r2, $r0, $r1

1004: sub $r5, $r3, $r4

1008: sw $r6, 50($r7)

1012: lw $r9, 20($r8)

1016: mul $r12, $r10, $r11
```



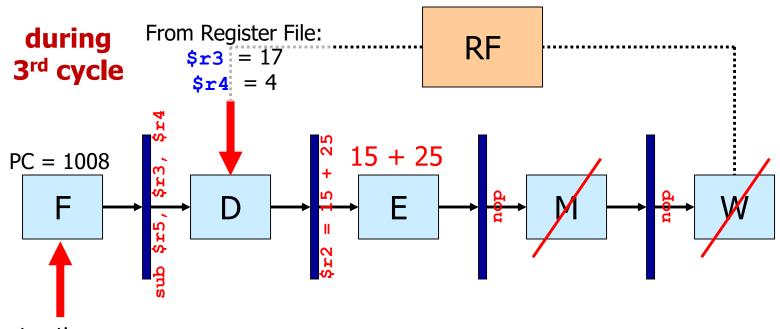
```
1000: add $r2, $r0, $r1

1004: sub $r5, $r3, $r4

1008: sw $r6, 50($r7)

1012: lw $r9, 20($r8)

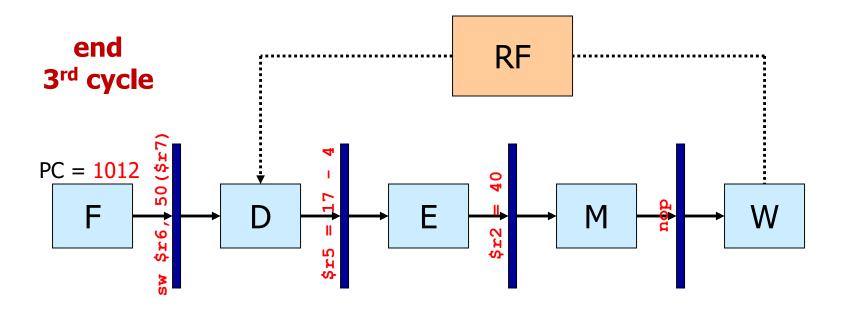
1016: mul $r12, $r10, $r11
```



From instruction memory:

```
sw $r6, 50($r7)
```

```
1000: add
               $r2,
                     $r0,
                           $r1
1004: sub
               $r5,
                     $r3,
                           $r4
1008: sw
               $r6,
                     50 ($r7)
1012: lw
               $r9,
                     20($r8)
1016: mul
               $r12, $r10, $r11
```



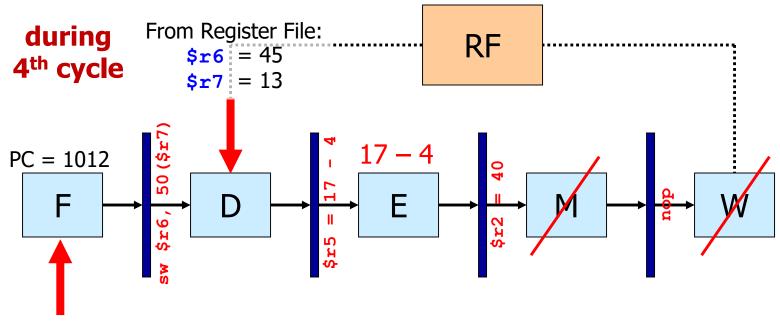
```
1000: add $r2, $r0, $r1

1004: sub $r5, $r3, $r4

1008: sw $r6, 50($r7)

1012: lw $r9, 20($r8)

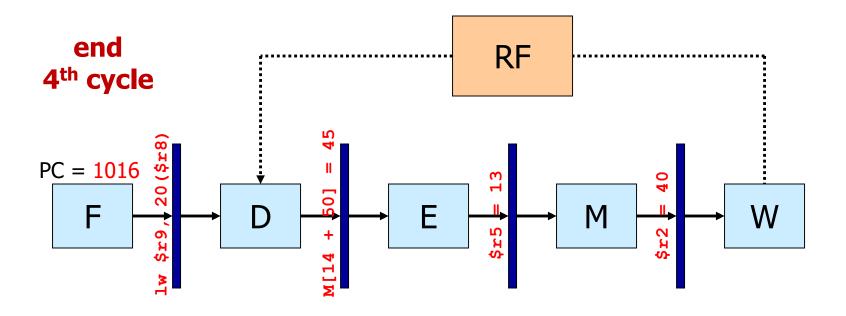
1016: mul $r12, $r10, $r11
```



From instruction memory:

lw \$r9, 20(\$r8)

```
1000: add
               $r2,
                     $r0,
                           $r1
1004: sub
               $r5,
                     $r3,
                           $r4
1008: sw
               $r6,
                     50 ($r7)
1012: lw
               $r9,
                     20($r8)
1016: mul
               $r12, $r10, $r11
```



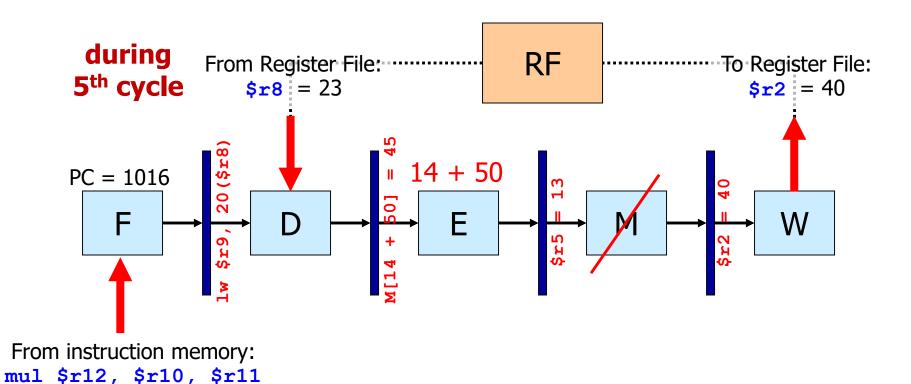
```
1000: add $r2, $r0, $r1

1004: sub $r5, $r3, $r4

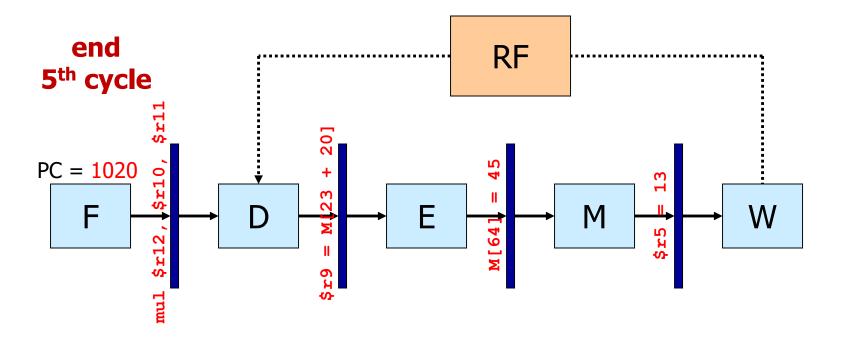
1008: sw $r6, 50($r7)

1012: lw $r9, 20($r8)

1016: mul $r12, $r10, $r11
```



```
1000: add
               $r2,
                     $r0,
                           $r1
1004: sub
               $r5,
                     $r3,
                           $r4
1008: sw
               $r6,
                     50 ($r7)
1012: lw
               $r9,
                     20($r8)
1016: mul
               $r12, $r10, $r11
```



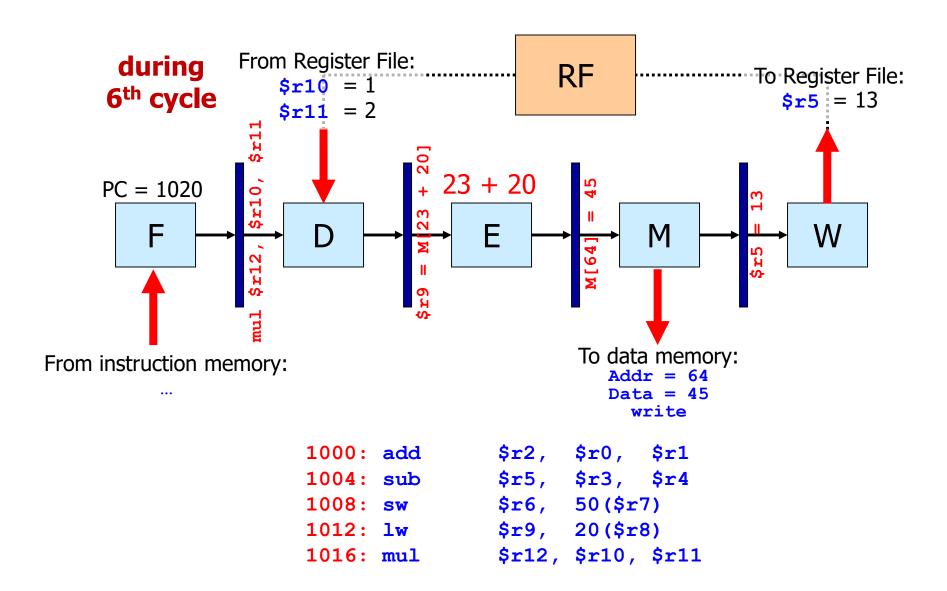
```
1000: add $r2, $r0, $r1

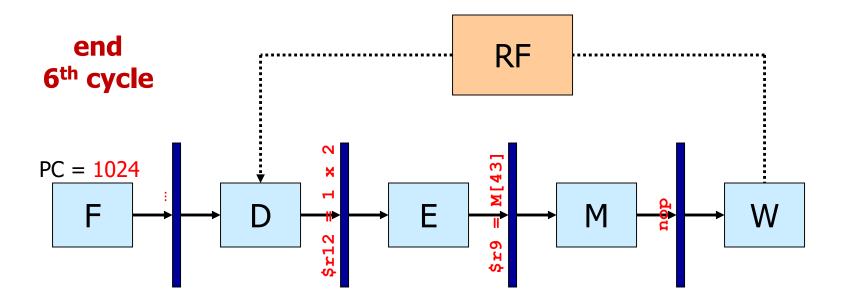
1004: sub $r5, $r3, $r4

1008: sw $r6, 50($r7)

1012: lw $r9, 20($r8)

1016: mul $r12, $r10, $r11
```





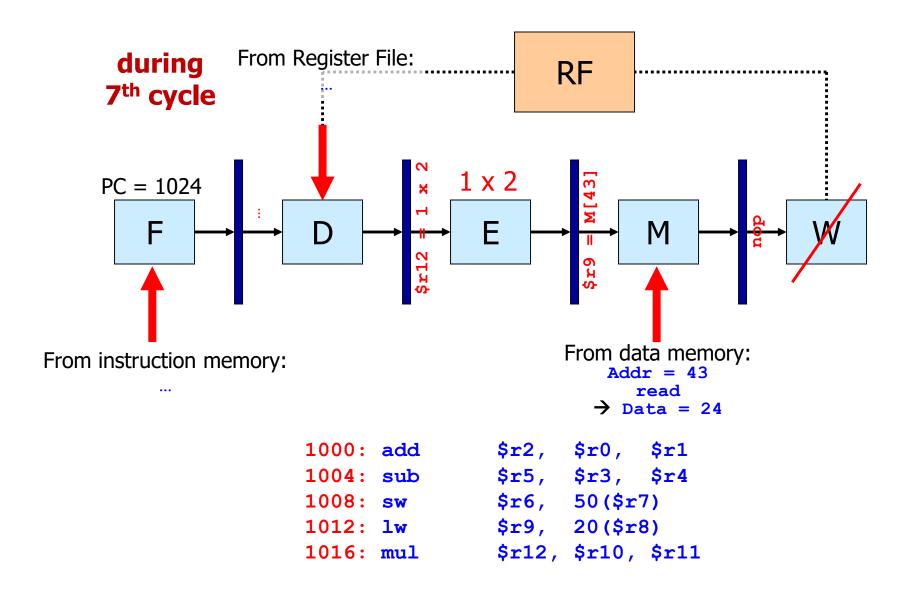
```
1000: add $r2, $r0, $r1

1004: sub $r5, $r3, $r4

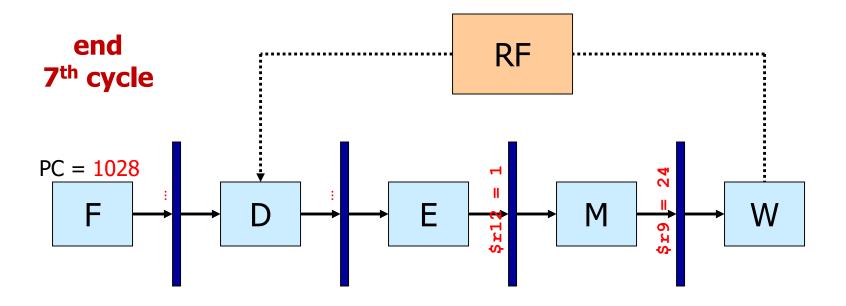
1008: sw $r6, 50($r7)

1012: lw $r9, 20($r8)

1016: mul $r12, $r10, $r11
```



Example of Pipelined Execution



```
1000: add $r2, $r0, $r1

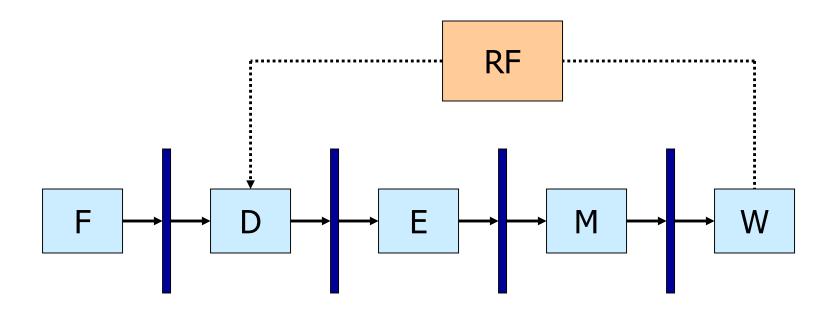
1004: sub $r5, $r3, $r4

1008: sw $r6, 50($r7)

1012: lw $r9, 20($r8)

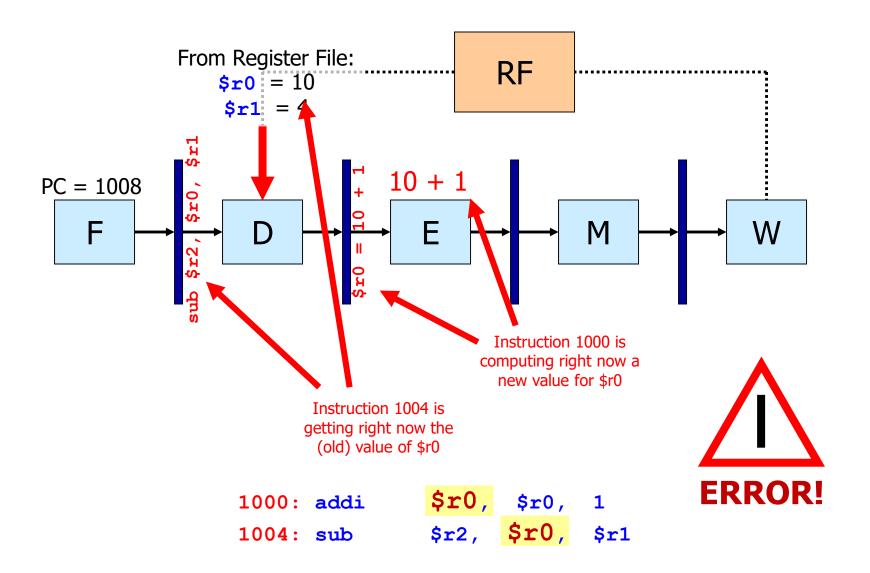
1016: mul $r12, $r10, $r11
```

Pipeline Schedule



1000:	F	D	Ε	М	W				
1004:		F	D	Е	М	W			
1008:			F	D	Е	М	W		
1012:				F	D	Е	М	W	
1016:					F	D	Ε	М	W

Problem!



Two Main Problems

- 1. CISC vs. RISC
 - Can we build equally well a pipeline for a Complex Instruction Set Computer as for a Reduced Instruction Set Computer?
 - What does that mean?!

- 2. Instructions are **not** independent
 - Can we execute code correctly?

RAW, WAR and WAW Dependences

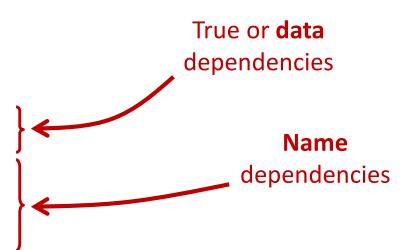


addd has a RAW dependence on divd

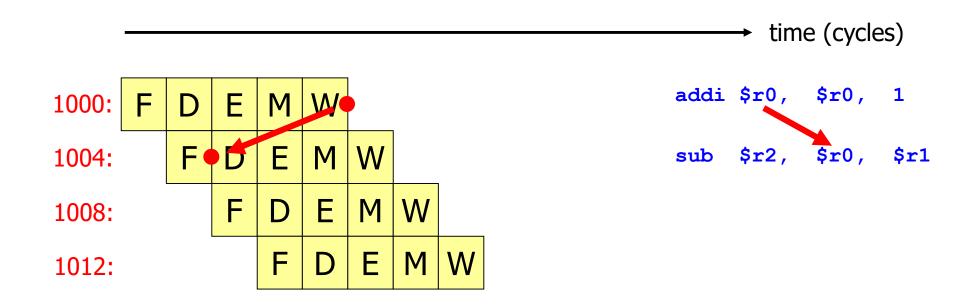
Write

After

- subd has a WAR dependence on addd
- adddi has a WAW dependence on divd



Data Hazards

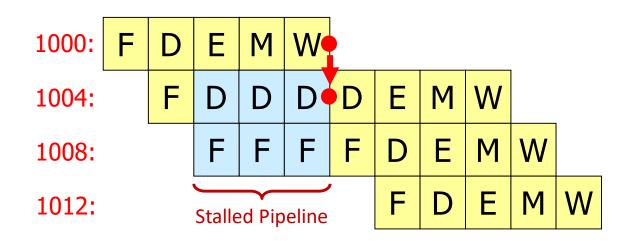


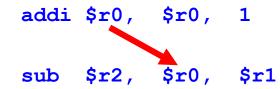


Causality violation!

We try to use a result before it is produced!

Data Hazards Solved by Stalling the Pipeline

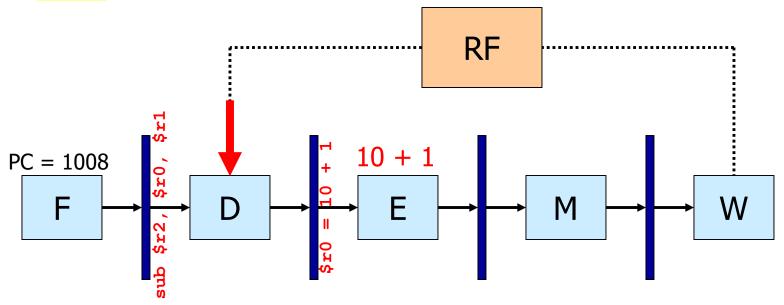




- The natural solution to Data Hazards caused by RAW dependences is to implement some logic in the processor to stop/repeat the decoding until the required value is available
- "Stalling" roughly means introducing nop's in the pipeline
- Due to the rigidity of the pipeline, if one stage is stalled (D in the example), all the
 preceding ones must be stalled too (e.g., F)

1000: addi \$r0, \$r0, 1

1004: sub \$r2, \$r0, \$r1



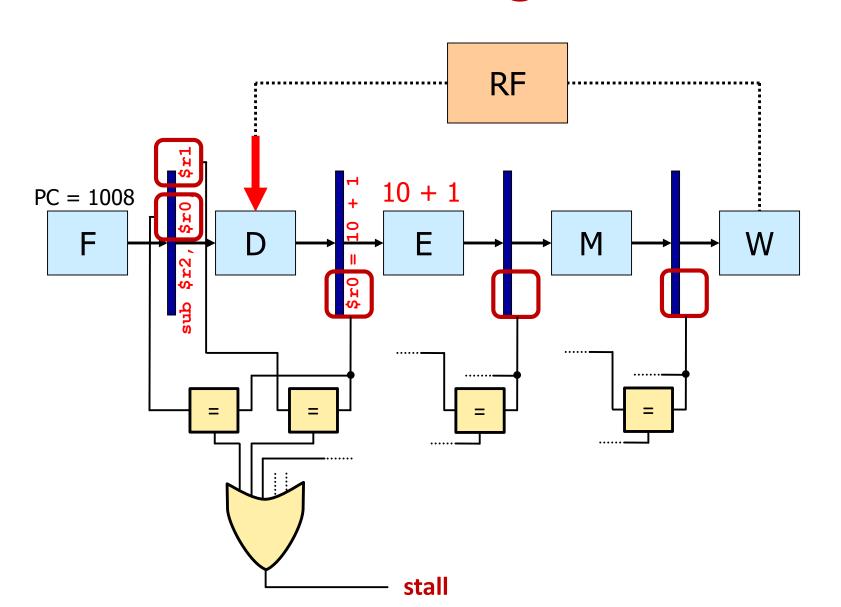
1000: addi

1004: sub

\$r0, \$r0, 1

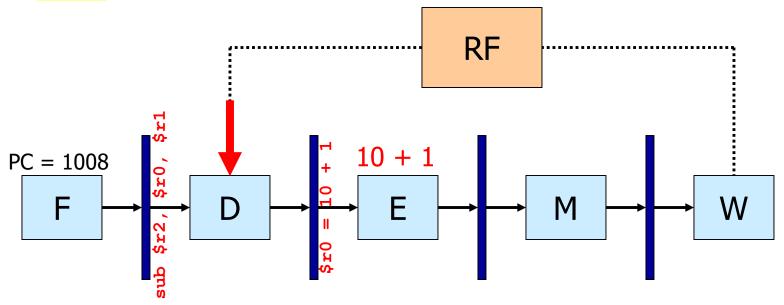
\$r2, **\$r0**, **\$r1**

Detecting



1000: addi \$r0, \$r0, 1

1004: sub \$r2, \$r0, \$r1



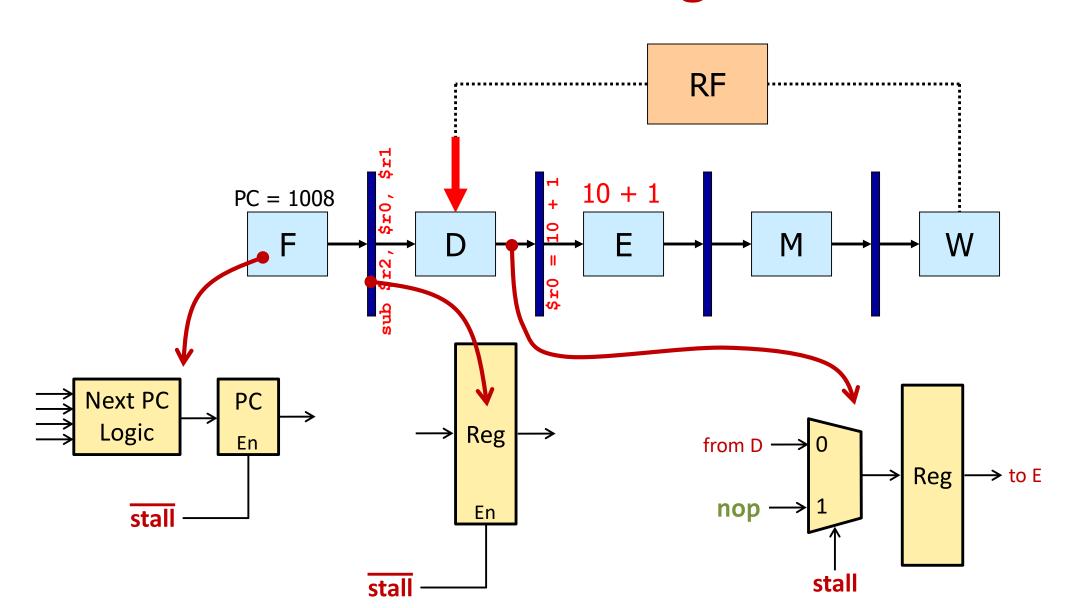
1000: addi

1004: sub

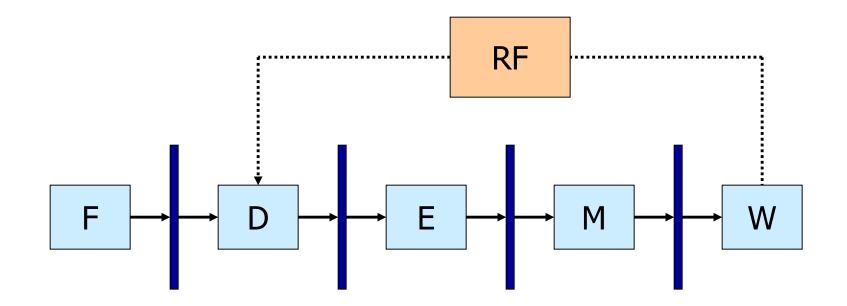
\$r0, \$r0, 1

\$r2, **\$r0**, \$r1

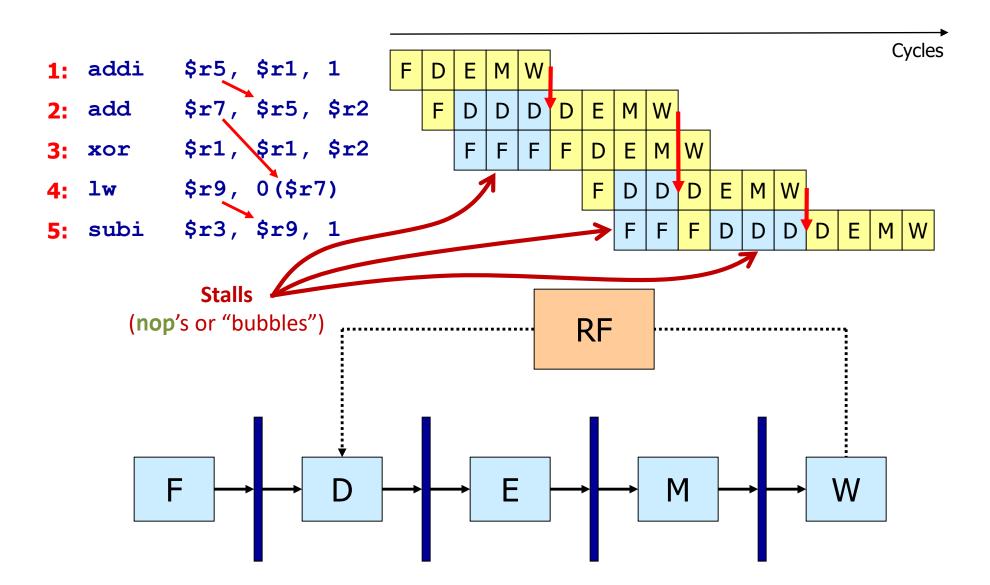
Stalling



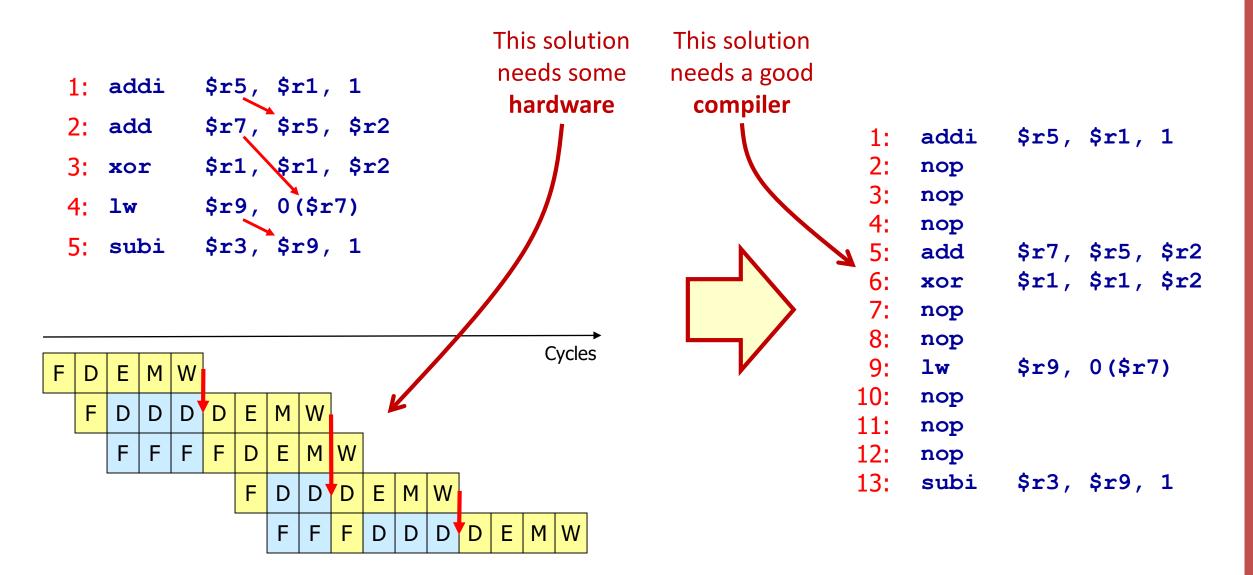
1: addi \$r5, \$r1, 1
2: add \$r7, \$r5, \$r2
3: xor \$r1, \$r1, \$r2
4: lw \$r9, 0(\$r7)
5: subi \$r3, \$r9, 1



Data Hazards



Another Solution?

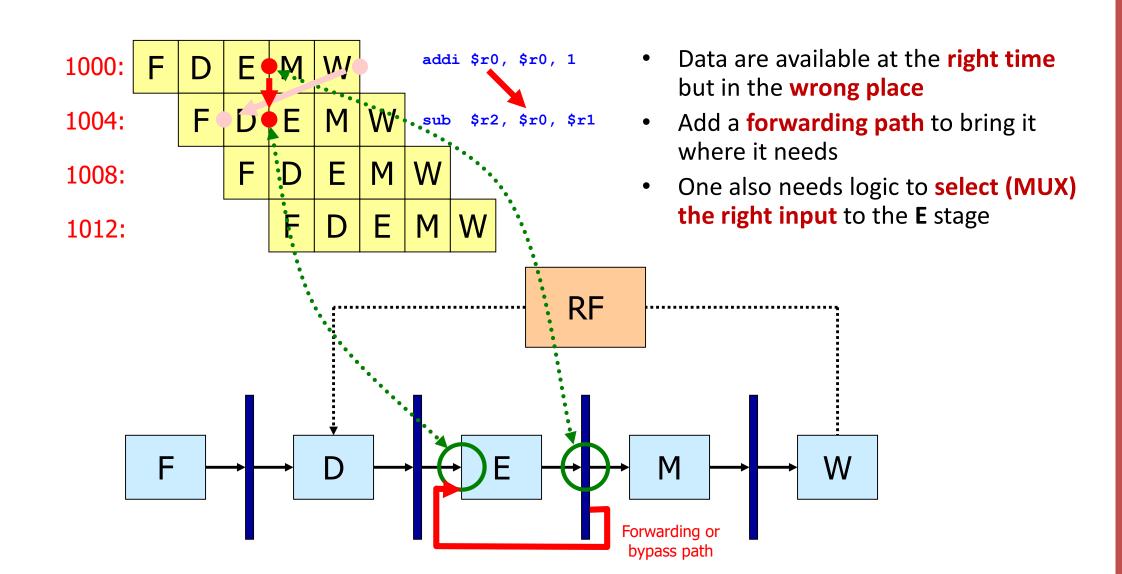


Architecture and Microarchitecture

- Architecture: what is in the ISA contract
 - Instructions, registers, etc.
- Microarchitecture: what is specific of an implementation
 - Multicycle vs. pipelined, FSM or pipeline structure, etc.

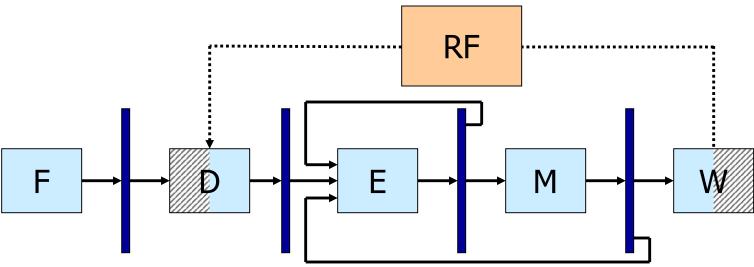
Some of the solutions evoked (reschedule instructions, add **nop**'s, delay slots) **expose typically microarchitectural aspects** (pipeline structure) **in the architecture** \rightarrow the same binary **does not run** on different processors!

Data Hazards Solved by Forwarding Values

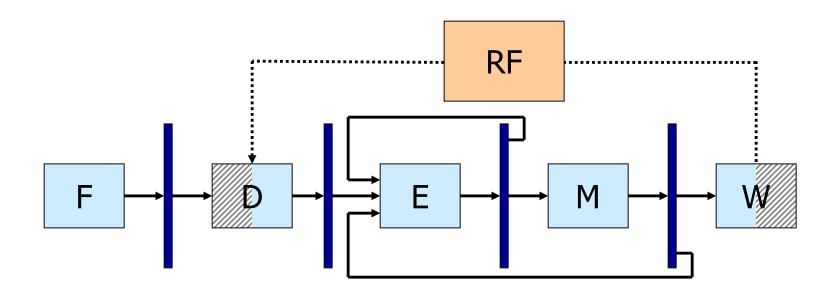


Classic MIPS Pipeline

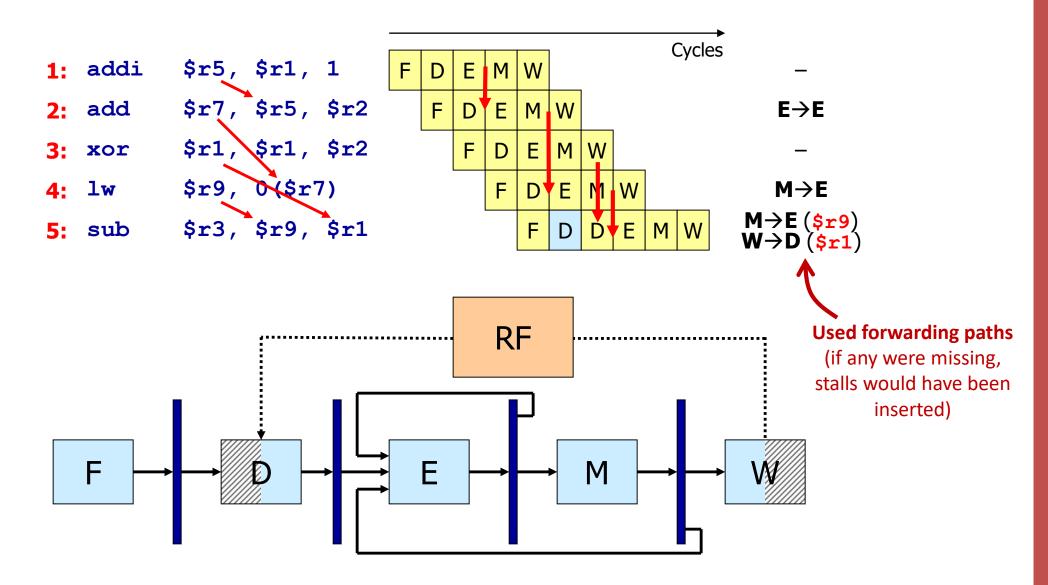
- 5-stage pipeline with all forwarding paths:
 - $E \rightarrow E, M \rightarrow E$
 - $W \rightarrow D$
- The register-file forwarding (W→D) is a special case:
 - During W, registers are written in the first half of the cycle
 - During **D**, registers are read in the second half of the cycle
 - → a register can be correctly written and read in the same cycle
- Not always all forwarding path exist...



- 1: addi \$r5, \$r1, 1
- 2: add \$r7, \$r5, \$r2
- 3: xor \$r1, \$r1, \$r2
- 4: lw \$r9, 0(\$r7)
- 5: subi \$r3, \$r9, 1

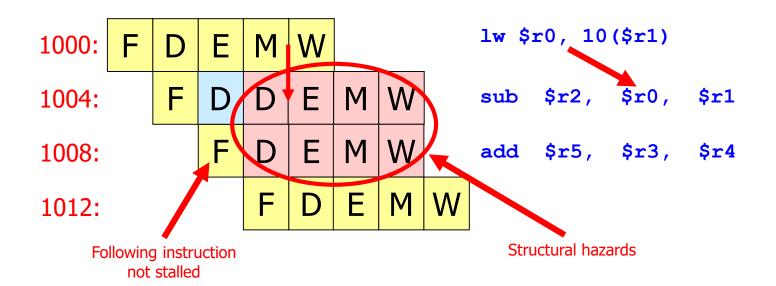


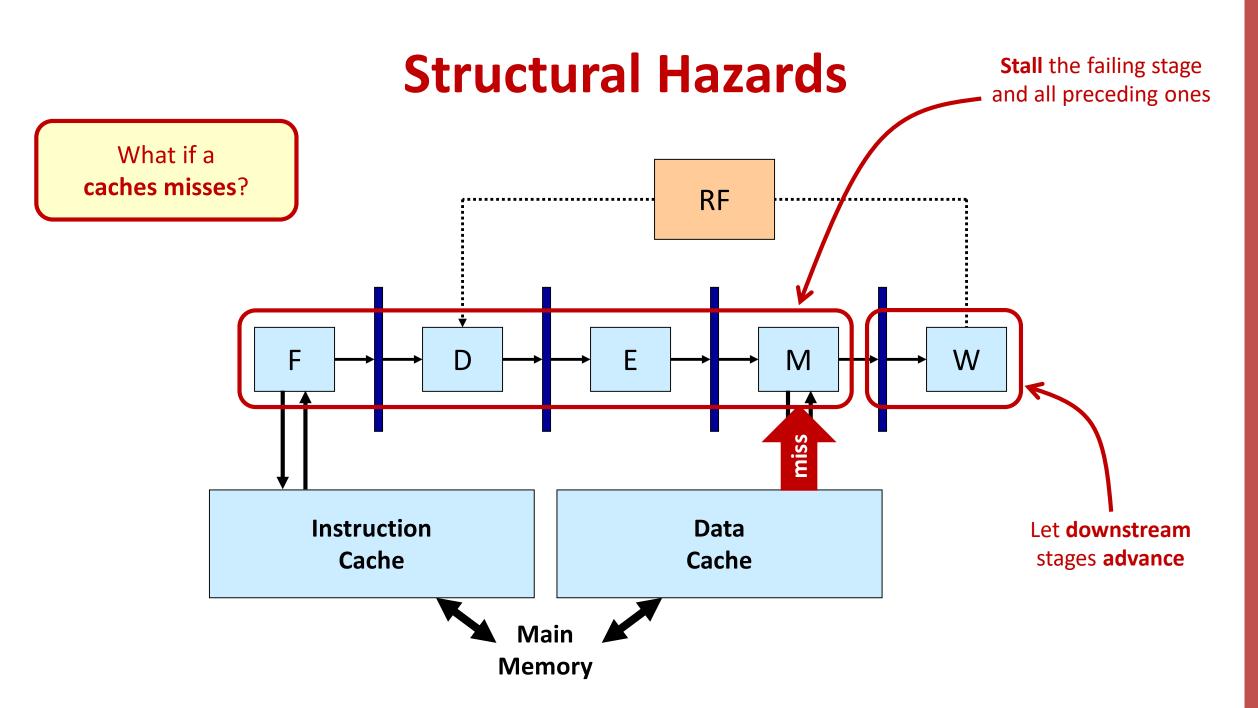
Reduced Data Hazards



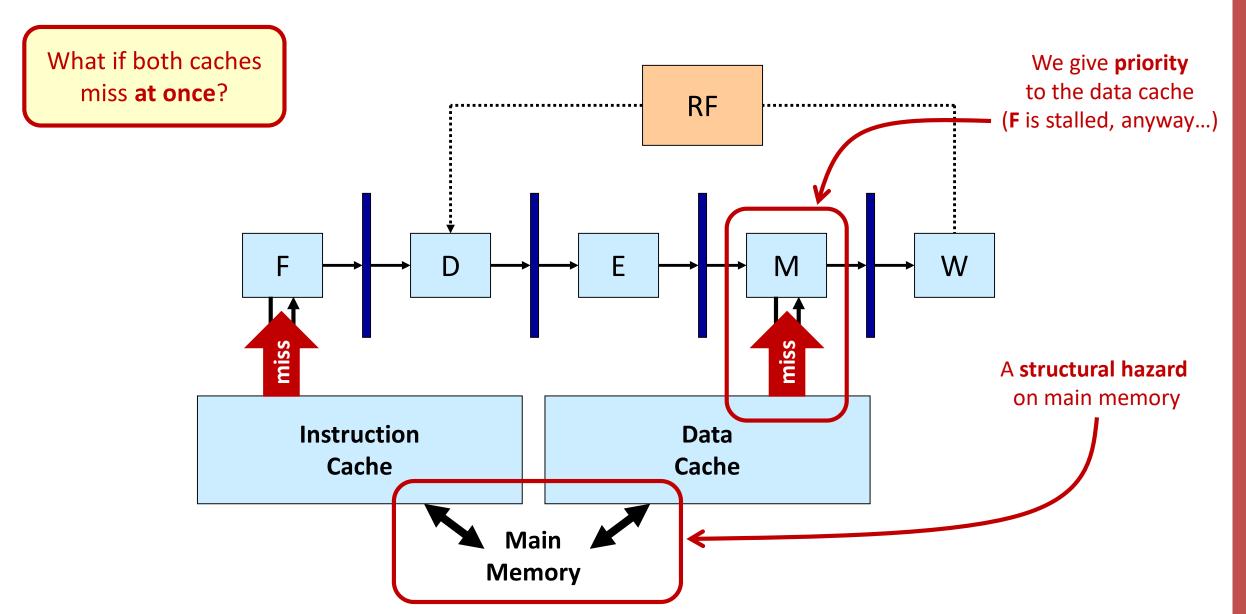
Structural Hazards

- A structural hazard happens when different instructions compete for the same resource (e.g., pipeline stage)
- Structural hazards cannot happen in our pipeline
- If we did not stall also instructions following one missing an operand, we could have structural hazards

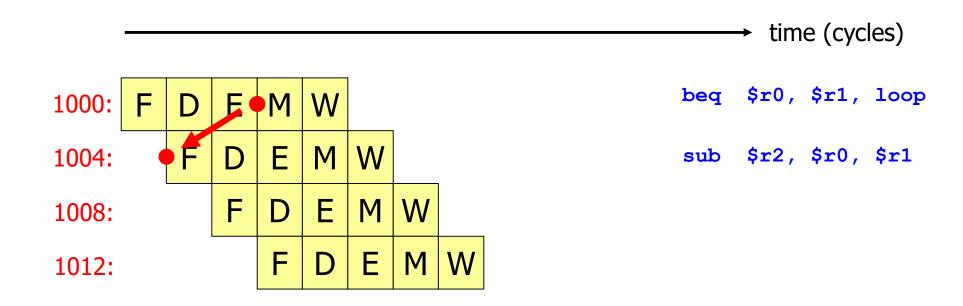




Structural Hazards



Control Hazards

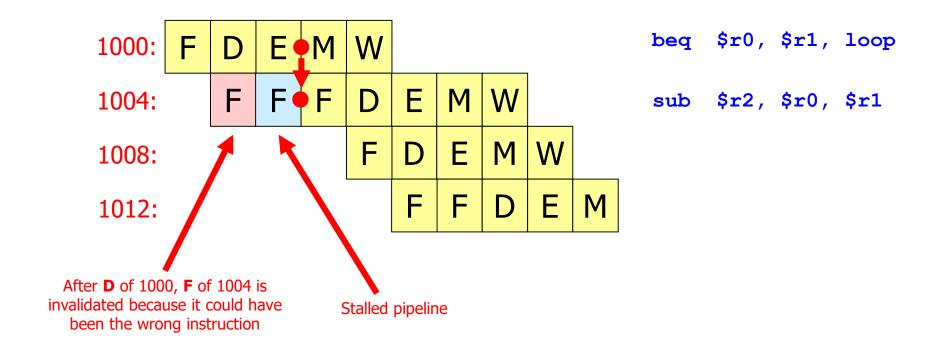




Causality violation!

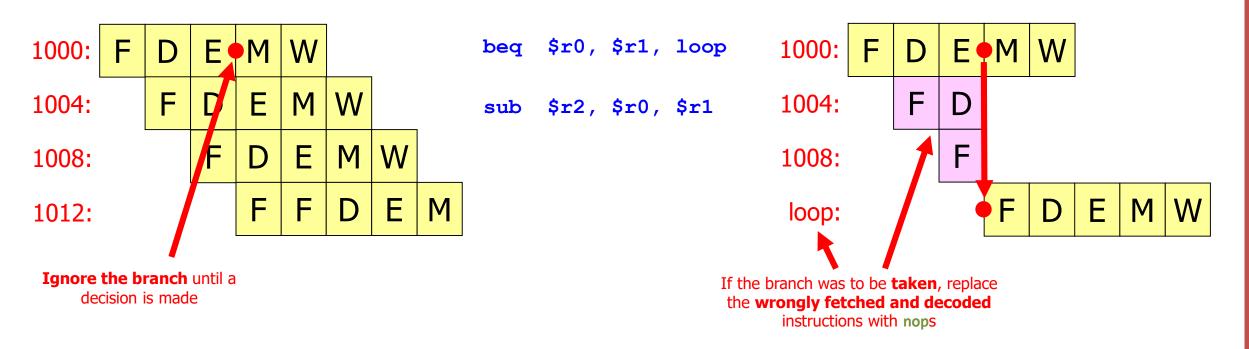
We fetch an instruction before we know which one!

Control Hazards Solved by Stalling the Pipeline



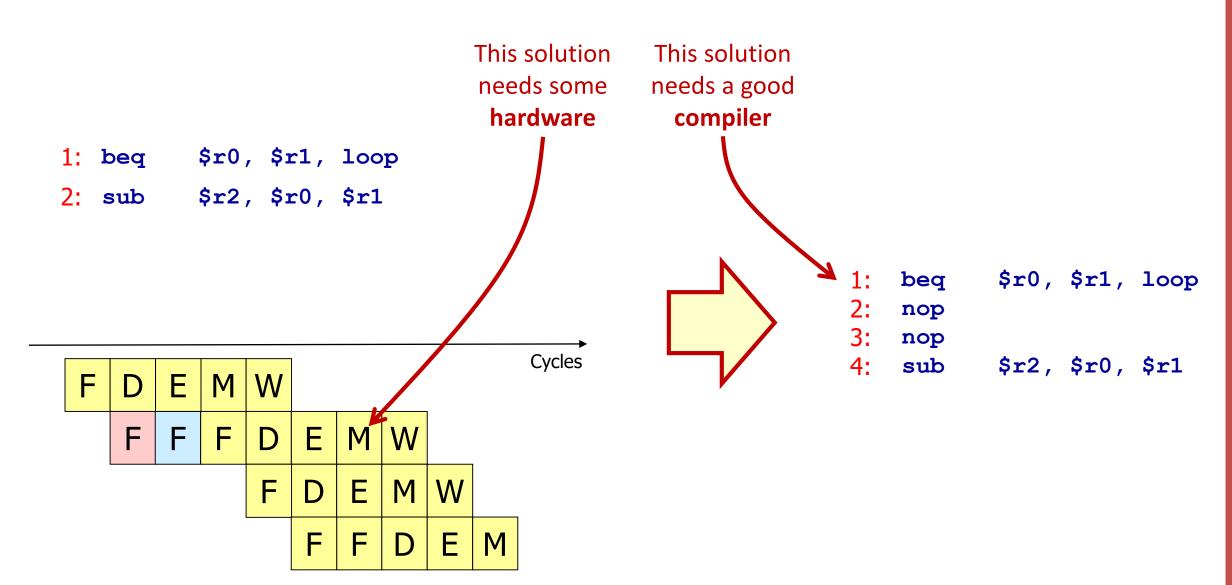
- Similarly to the way we solve data hazards, we can **stall the pipeline** (**F**), once it is discovered, after **D**, that an instruction was a branch, and this **until the branch is resolved**
- If, for instance, the correct address of the next instruction is known at the end of the E stage, 2 cycles
 are lost every branch

Fetching and Decoding Do Not Do Any Damage

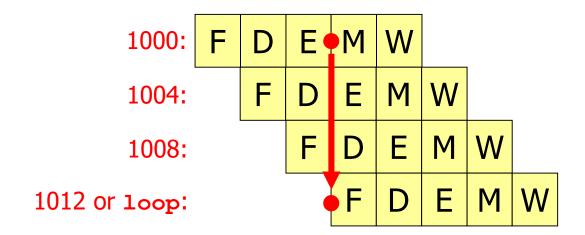


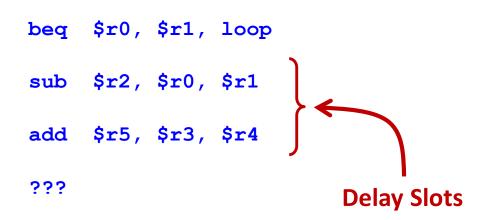
- Fetching or decoding a wrong instruction does not create any problem, provided that the instruction is not also executed
- If the outcome of a branch is known at the end of the **E** stage, we can **wait** until then to **conditionally kill** the following two instructions in the pipeline **if the branch happens to be taken**
- Now 2 cycles are lost only for taken branches and none is lost for nontaken ones

Another Solution?



Control Hazards Solved by Defining Delay Slots





- Alternatively, we can modify the definition of the architecture and decide that the two instructions
 following a branch are executed in any case (branch taken or not) as if they were before
- These instructions after the branches are called **delay slots**
- Note that code becomes counterintuitive!!!
- MIPS did it as some others, but quite rare in current architectures

Use of Delay Slots

- A simple way of using delay slots is to use them for nop's—but then it is not better
 than stalling the pipeline
- A better idea is to put there instructions which precede the branch and on which the branch has no dependence
- Suppose an architecture with two delay slots:

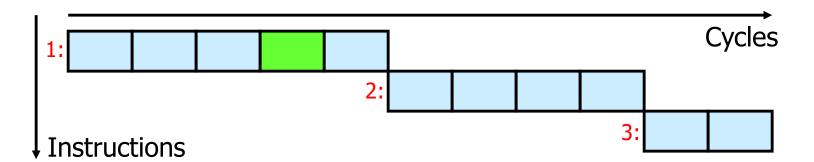
```
1000:
             $r2, $r0, $r7
       sub
1004:
                                        1000:
       mul
            $r1, $r6, $r7
                                               mul
                                                     $r1,
                                                           $r6, $r7
1008:
             $r5, $r3, $r4
       add
1012:
                                        1004:
       beq
             $r0, $r1, loop
                                               beq
                                                     $r0, $r1, loop
1016:
                                        1008:
                                                sub
                                                     $r2, $r0, $r7
       nop
                                                                        Delay
                Delay
                                        1012:
1020:
                 slots
                                                                        slots
                                                add
                                                     $r5, $r3, $r4
       nop
1024:
                                        1016:
             $r8, 12($r9)
       lw
                                                lw
                                                     $r8, 12($r9)
```

Branch Prediction

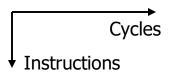
- A better strategy is to guess the branch outcome and fetch the corresponding instruction (either the next instruction or the branch destination, but not necessarily the former)
 - If the guess is correct, no cycle is lost
 - If the guess is wrong, what has been fetched and decoded is thrown away ("squashed")
- Branch predictors of modern processors are extremely sophisticated: dynamic predictors learn from previous executions of a branch...
- Complex predictors can be correct up to 95-99% of the time
- The quality of branch predictors has made architectures with delay slots extremely rare

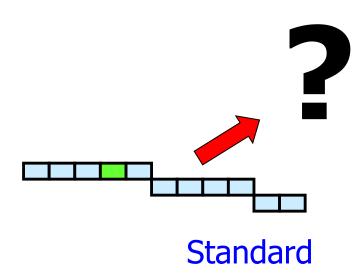
Starting Point (Programmer Model)

Sequential multicycle processor

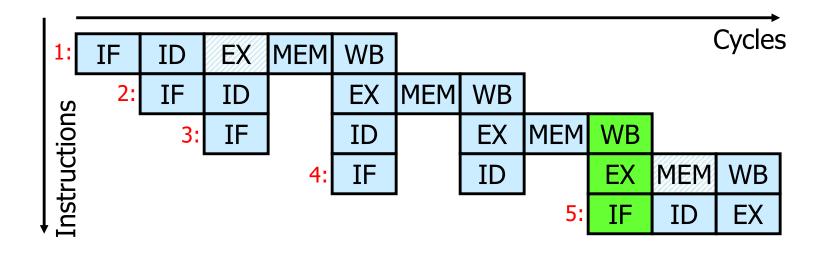


Instruction Level Parallelism?





First Step: Pipelining



 Simplest form of Instruction Level Parallelism (ILP): several instructions are now executed at once

Three Types of Hazards Hinder Pipelining

- Data Hazards (= data dependences)
 - Solutions:
 - Forwarding paths, wherever possible
 - Stalls, in all other cases
- Control Hazards (= jumps and branches)
 - Solutions:
 - Delay slots, if the architecture allows it
 - Branch prediction, to try to do the right thing
 - Stalls, if not
- Structural Hazards (= conflicting need for a resource)
 - Solutions:
 - Rigid pipelines which cannot have structural hazards by construction
 - Stalls, otherwise

References

- Patterson & Hennessy, COD RISC-V Edition
 - Sections 4.6-4.9