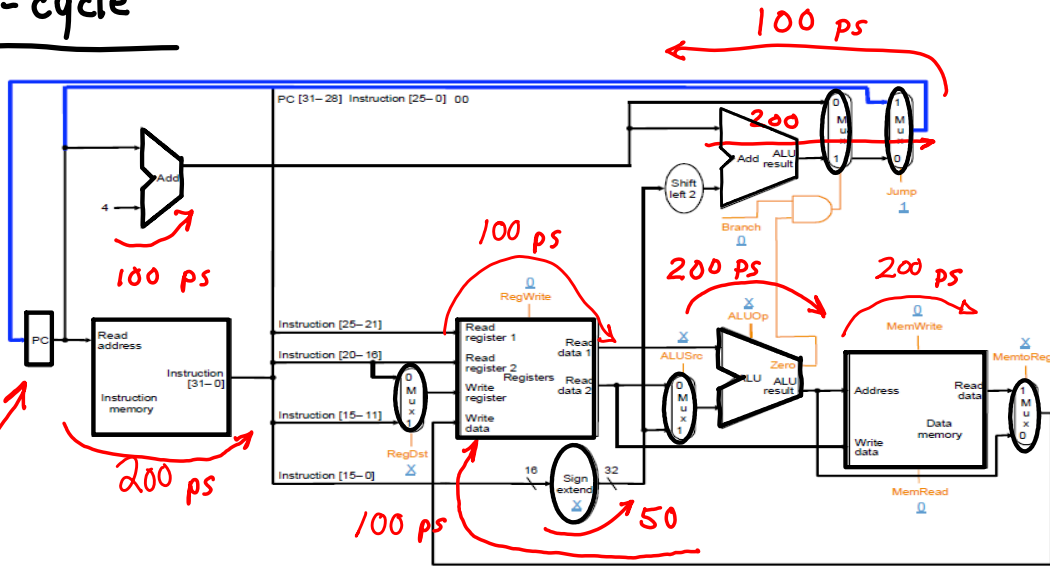


# Performance

How fast can we clock?

1. 1 cycle MIPS performance
3. general pipelining
8. MIPS pipe, LW
12. MIPS performance, piped vs non-piped
14. arrays, piped
15. hazards

1-cycle



clock changes PC output.

How long before we can clock PC, Regfile? critical path?

## Single Cycle Processor Performance

- Functional unit delay
  - Memory: 200ps
  - ALU and adders: 200ps
  - Register file: 100 ps

$$ps = 10^{-12} \text{ sec}$$

Instruction Class	Instruction memory	Register read	ALU operation	Data memory	Register write	Total
R-type	200	100	200		100	600
load	200	100	200	200	100	800
store	200	100	200	200		700
branch	200	100	200			500
jump	200					200

$$\text{max delay} = T_{\text{clock}}$$

$$\frac{1}{T_{\text{clock}}} = \frac{1}{0.8 \text{ ns}} = 1.25 \text{ GHz}$$

- CPU clock cycle = 800 ps = 0.8ns (1.25GHz)

What if we make control more complex?  
Set clock by opcode?

• Instruction Mix

- 45% ALU
- 25% loads
- 10% stores
- 15% branches
- 5% jumps

Instruction Class	Instruction memory	Register read	ALU operation	Data memory	Register write	Total
R-type	200	100	200		100	600
load	200	100	200	200	100	800
store	200	100	200	200		700
branch	200	100	200			500
jump	200					200

ps → 0.6 ns  
0.8  
0.7  
0.5  
0.2

• CPU clock cycle =  $0.6 \times 45\% + 0.8 \times 25\% + 0.7 \times 10\% + 0.5 \times 15\% + 0.2 \times 5\%$   
= 0.625 ns (1.6GHz)

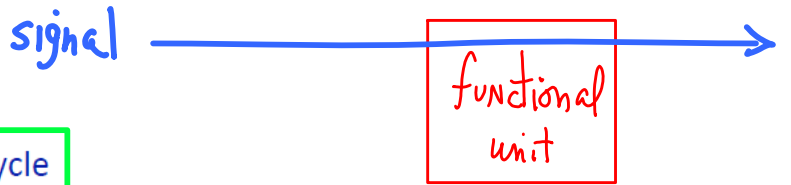
⇒ 1.6 GHz

what is speedup?  $S_{new-old} = \frac{n_{CPI} (1/CR)}{n_{CPI} (1/CR)} = \frac{CR_{new}}{CR_{old}} = \frac{1.25}{1.6}$  worth it? skew?

what's the

• Problem:

- Each functional unit used **once per cycle**
- Most of the time it is sitting waiting for its turn
  - Well it is calculating all the time, but it is **waiting for valid data**
- There is no parallelism in this arrangement



wait do

• Making instructions take **more cycles** can make machine **faster**!?

- Each instruction takes roughly the same time
  - While the CPI is much worse, the **clock freq is much higher**
- **Overlap execution** of multiple instructions at the same time
  - Different instructions will be active at the same time
- This is called "Pipelining"
- We will look at a 5 stage pipeline
  - Modern machines (Core 2) have order **15 cycles/instruction**

$CPI \uparrow CR \uparrow ?$

$Perf = \frac{n}{T} = \frac{n}{n_{CPI} (1/CR)}$   
 $= \uparrow CR / CPI \uparrow$

Can we win?

If  $CR \uparrow$  by slicing path into  $k$  small pieces,  
 delay → delay/ $k$  →  $T_{clock}$  →  $T_{clock}/k$  →  $CR \rightarrow k CR$   
 #cycles →  $k$  →  $CPI \rightarrow k CPI$

$Perf \rightarrow \frac{n}{n (k CPI) (1/k CR)}$  no change?

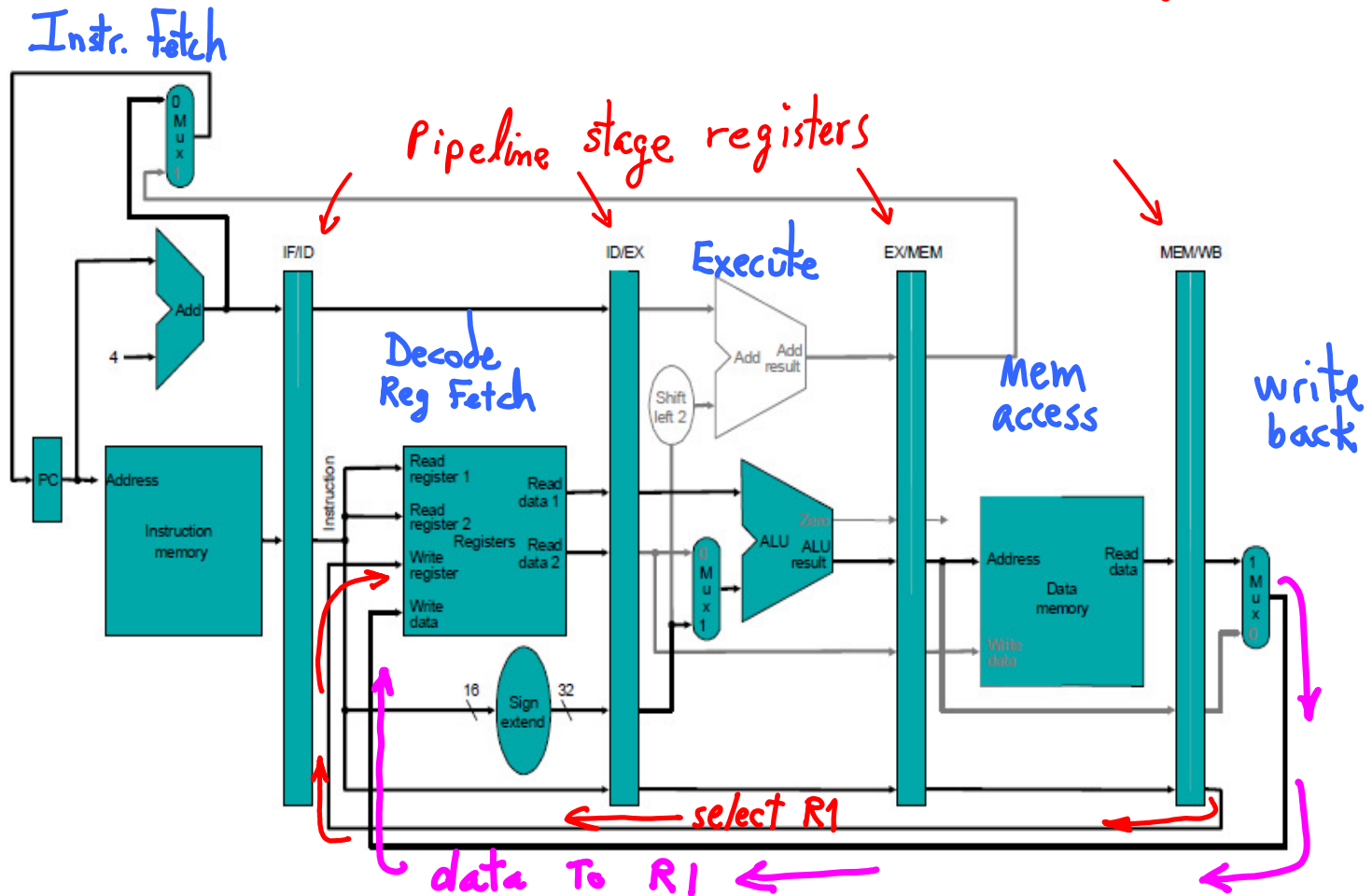
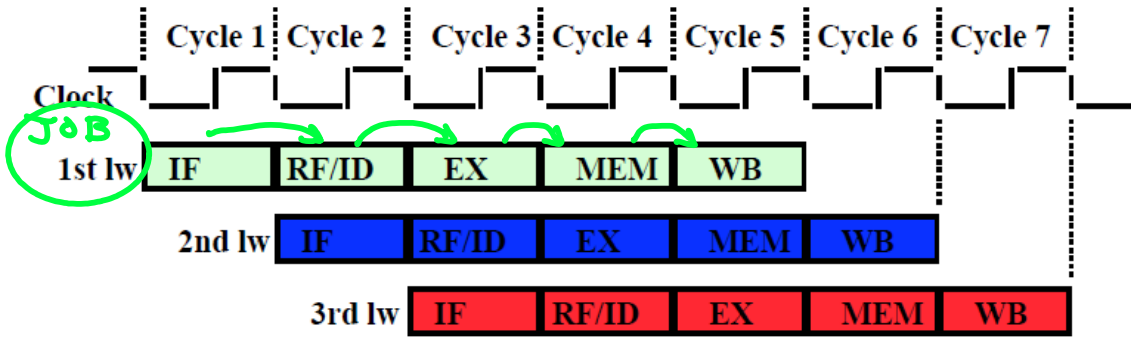
# Pipelining Load *lw*

- Load instruction takes 5 stages
  - Five independent functional units work on each stage
    - Each functional unit used only once *for a single instr.*
  - Another load can start as soon as 1<sup>st</sup> finishes IF stage
  - Each load still takes 5 cycles to complete
  - The *throughput*, however, is much higher

*all stages busy*  
*1 instr. exits per cycle*

$$CPI = \frac{1 \text{ instr.}}{1 \text{ cycle}}$$

$$T = 5 \text{ cycles latency}$$



**SUB R4, R5, R1**

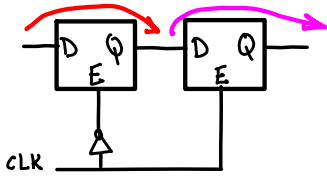
Required delay before using written data.

**ADD R1, R2, R3**

Insert NOPs (BRnzp #0)? Compiler does this, or HW?

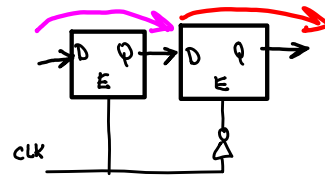
# fix delay? → negedge FF for Regfile

Positive edge-triggered FF:  
output changes on rising clock



CLK  
Sample input      change output

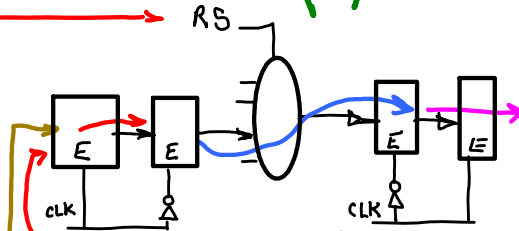
Negative edge-triggered FF:  
output changes on rising clock



Sample input      change output

Reg file      ID/EX stage reg

① posedge change  
instruction decode,  
Reg fetch: RS  
WB data

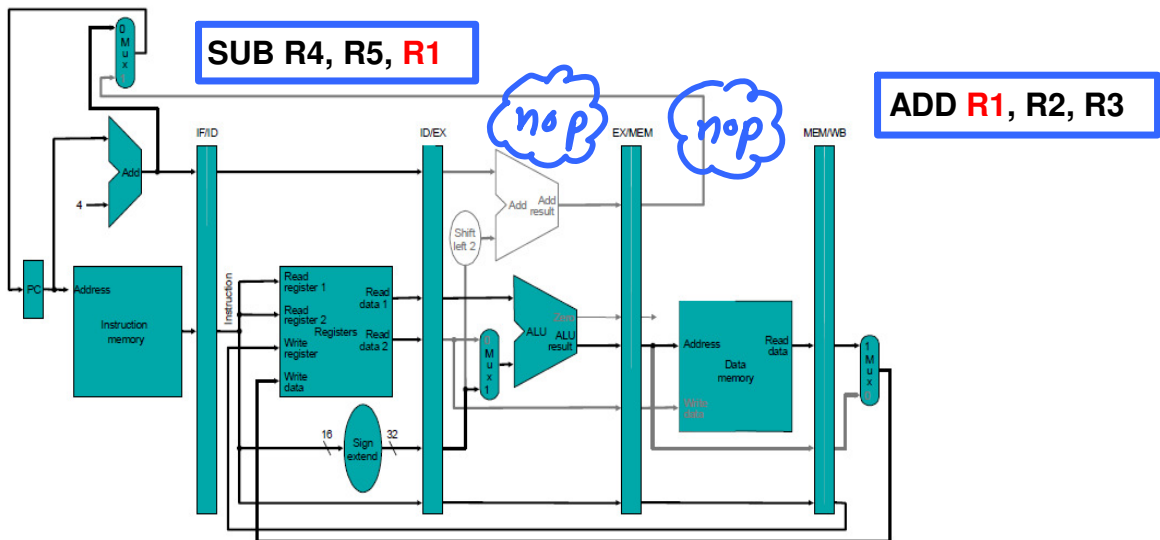


Write Data from WB stage

③ next posedge change  
data written to Regfile  
read for EX stage

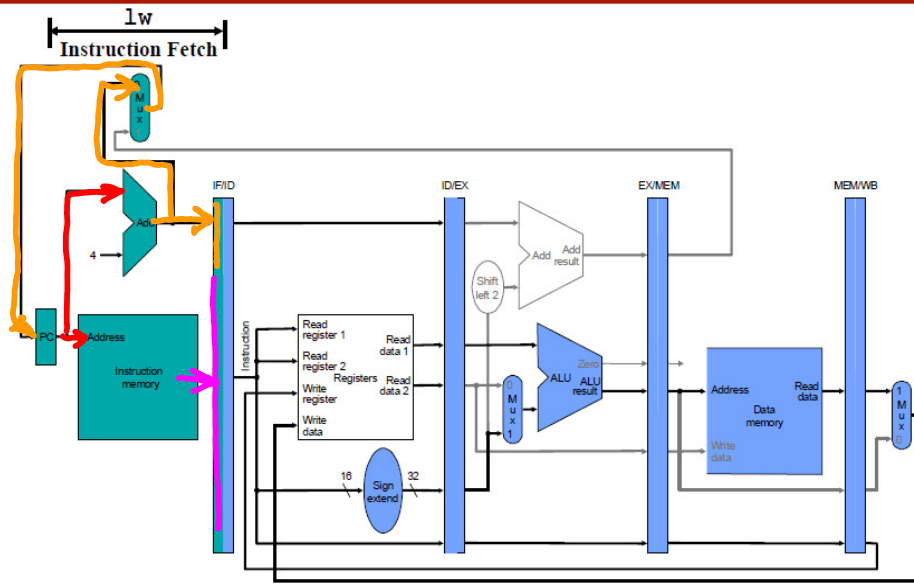
② negedge change  
Regfile output changes,  
WB data on output

Written data available for read in same cycle.



MIPS: LW \$2, 15(\$1)  
[ op | rs | rt | off ]

LC4: LDR R2, R1, #15  
[ op | SR1 | DR | off ]

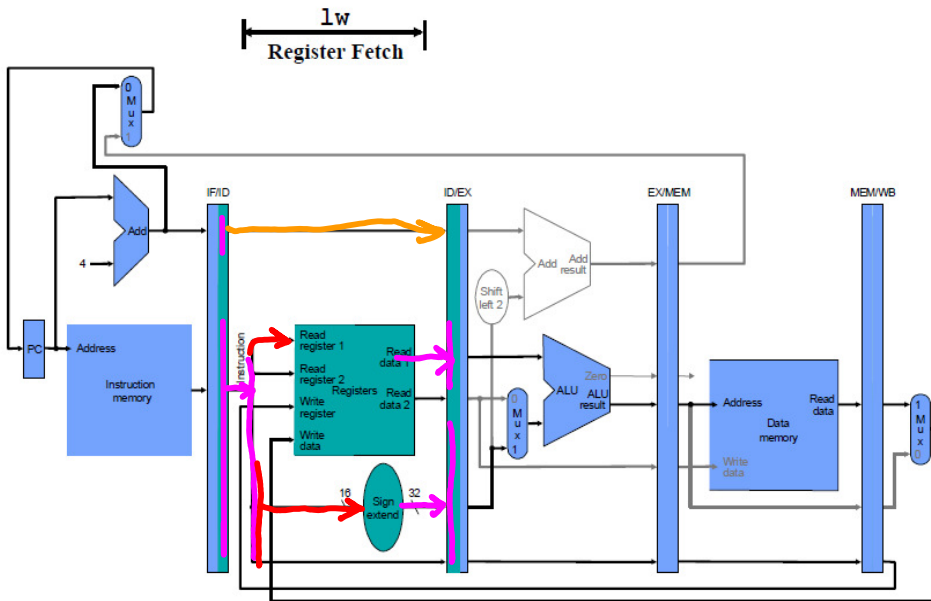


Instruction fetch:

PC++ ==> PC

PC++ ==> IF/ID PC  
Instruction ==> Instr

**Register Fetch**



Decode/reg fetch:

IF/ID ID/EX

PC ==> PC

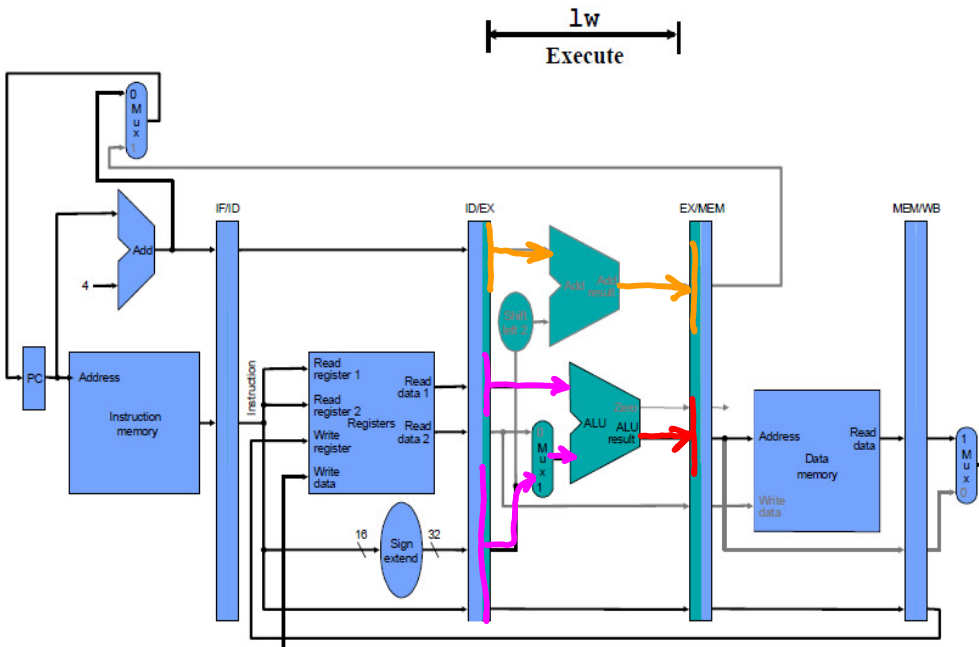
SR1 ==> SR1  
SR1out ==> SR1out

offset ==> SEXT  
==> offset

Instr.SR2 ==> DR

Instr.OP ==> decode  
==> CTL

**Execute**



Execute:

ID/EX EX/MEM

PC ==> ADD ==> PC

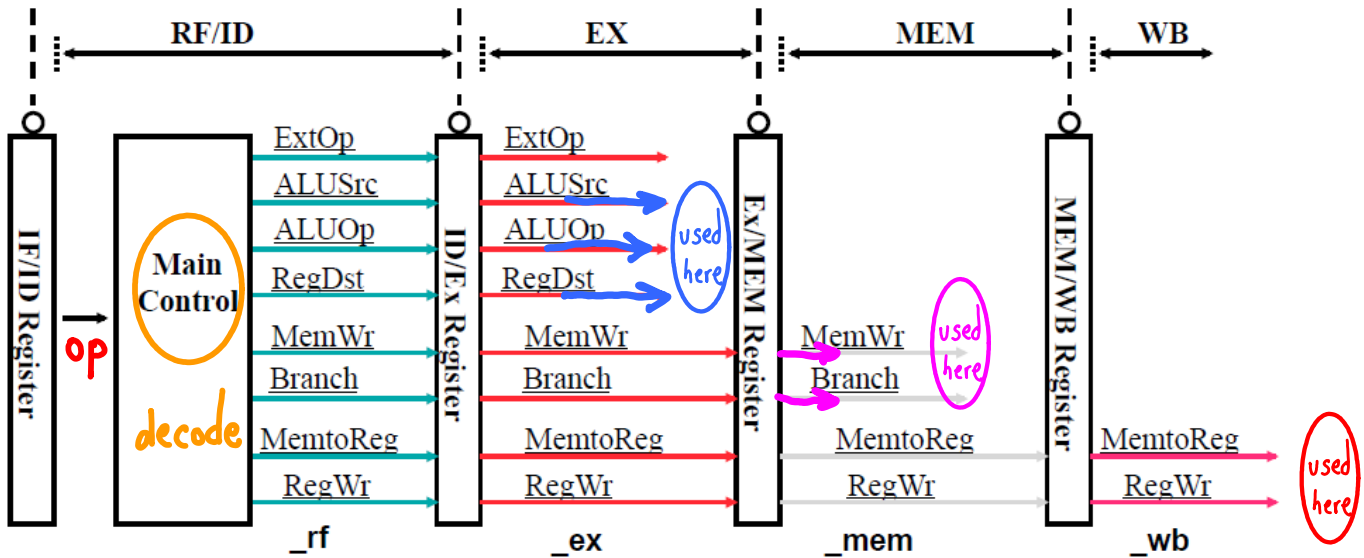
SR1out ==> ALU  
offset ==> ALU  
==> Res

DR ==> DR

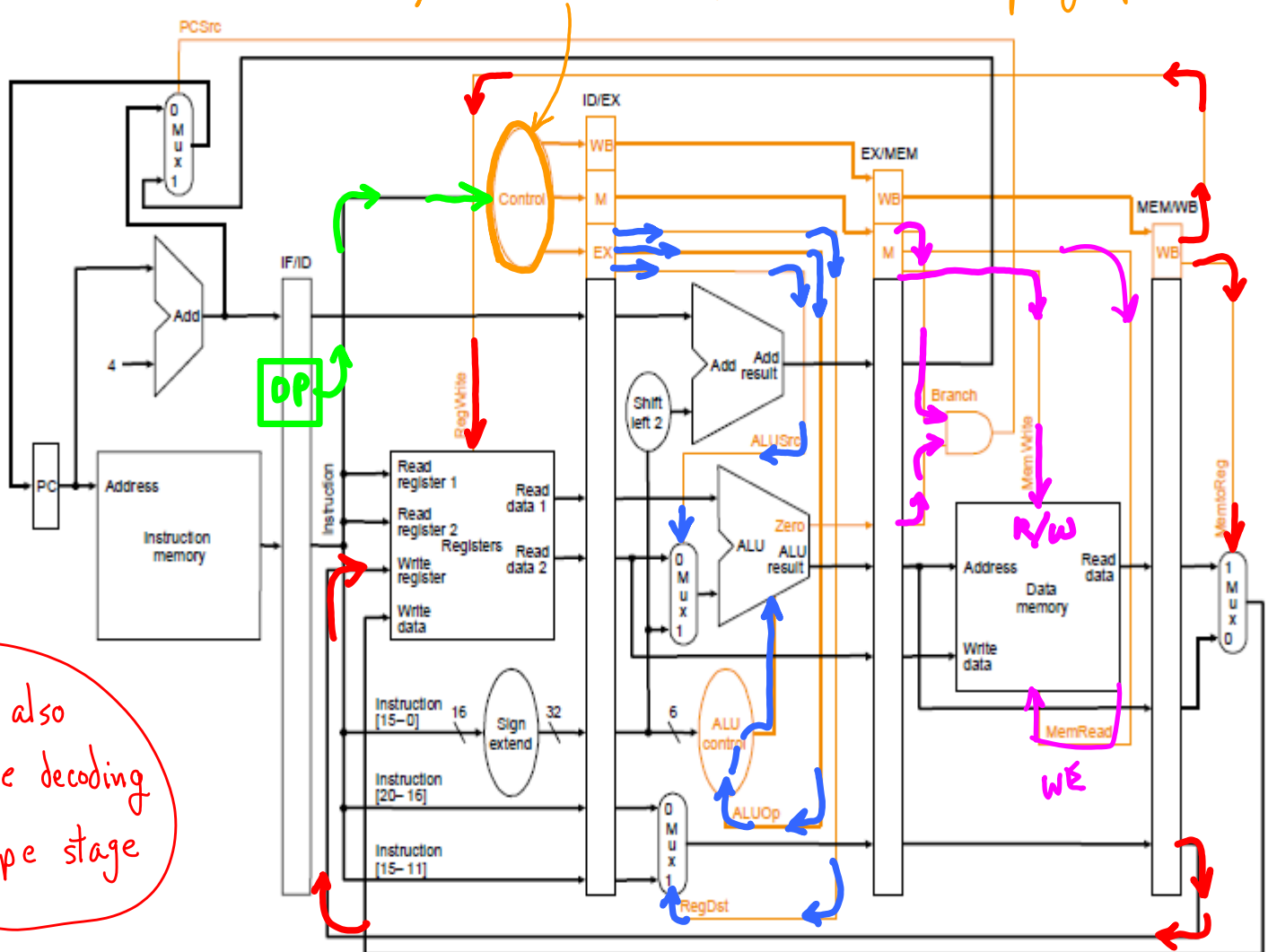
CTL ==> CTL



# Implementing Control



Decoding: combinational/ $\mu$ Code control signals (table lookup by opcode)

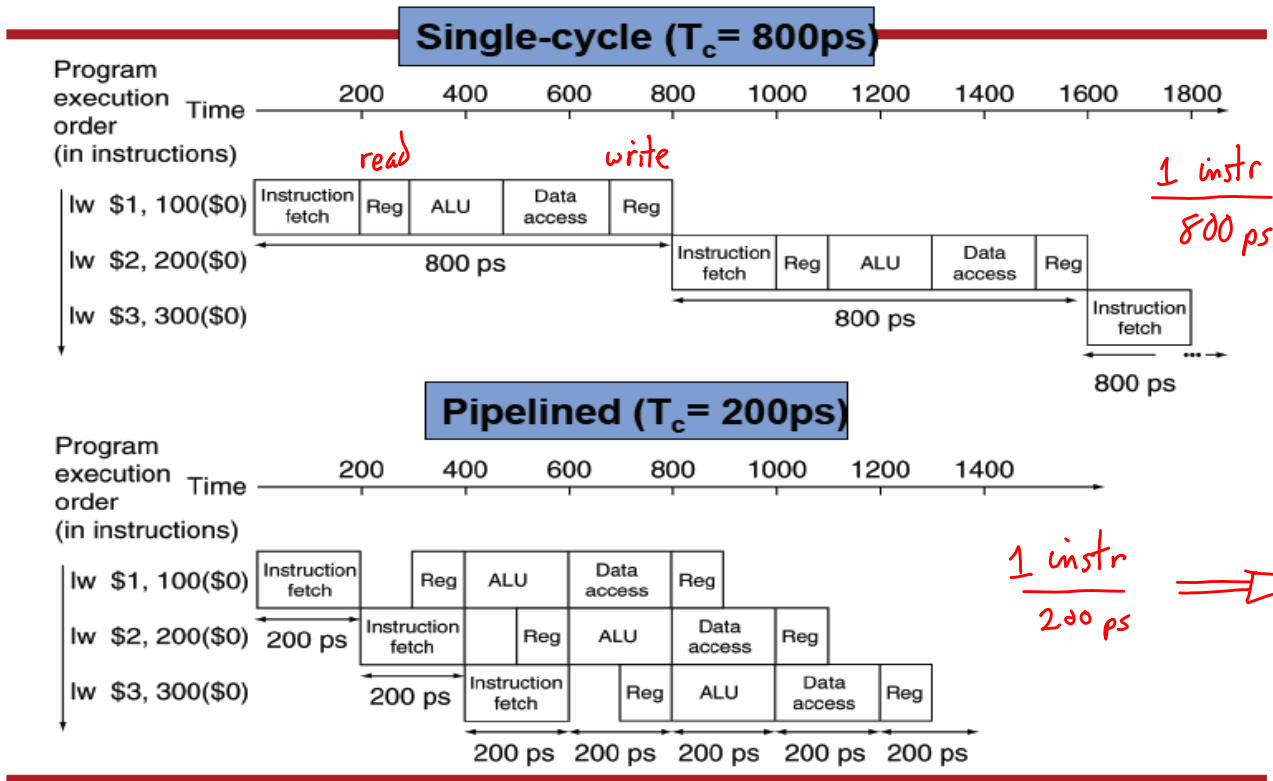


Could also pipeline decoding by pipe stage



- Assume time for stages is
  - 100ps for register read or write
  - 200ps for other stages

## Pipeline Performance



- MIPS SA designed for pipelining
  - All instructions are 32-bits
    - Easier to fetch and decode in one cycle
    - c.f. x86: 1- to 17-byte instructions
  - Few and regular instruction formats
    - Can decode and read registers in one step
  - Load/store addressing
    - Can calculate address in 3<sup>rd</sup> stage, access memory in 4<sup>th</sup> stage
  - Alignment of memory operands
    - Memory access takes only one cycle

orthogonality

vs. 2 for misaligned  $\Rightarrow$  stall



## But Something Is Fishy Here

- If dividing it into 5 parts made the clock faster  $800\text{ ps} \rightarrow 200\text{ ps}$ 
  - And the effective CPI is still one (if?)
- Then dividing it into 10 parts would make the clock even faster  $200\text{ ps} \rightarrow 100\text{ ps} ?$ 
  - And wouldn't the CPI still be one?
- Then why not go to twenty cycles?
- Really two issues
  - Some things really have to complete in a cycle *cannot divide every operation*
    - Find next PC from current PC
  - CPI is not really one
    - Sometimes you need the results a previous instruction that is not done

*An instruction does not complete in 1 cycle, need result before latency is done?*

## Can Pipelining Lead to an Arbitrary Short Clock Cycle?

- Min clock cycle = longest combinatorial delay + FF setup + clock skew
- Pipelining reduces the combinatorial delay
  - Less work per pipeline stage
  - Ideally, N stages reduce delay to 1/N
  - Best you can achieve is Clock cycle  $\rightarrow$  FF setup + clock skew
    - Diminishing returns from ever longer pipelines...
- Imbalance between stages also reduces benefits from subdividing
- Even if you could continuously improve clock frequency
  - ↑ – Power consumption  $\propto$  Frequency ↑

# Dependencies and Hazards

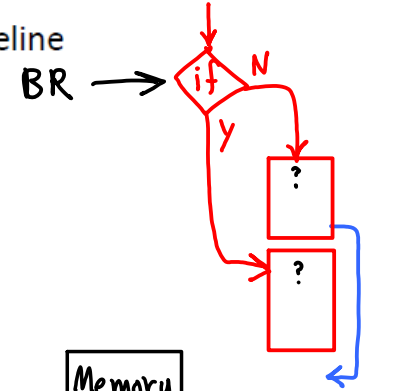
- Hazards: situations that prevent starting the next instruction in the next cycle
  - Wasted cycles,  $CPI > 1$
- Hazards are due to dependencies between instructions
  - Two instructions share resources or data
  - Pipelining may lead to overlapping their execution

## Types of hazards

- Structural Hazard (resource conflict)
  - Two instructions need to use the same piece of hardware
- Data Hazard
  - Instruction depends on result of instruction still in the pipeline
- Control Hazard
  - Instruction fetch depends on the result of instruction in pipeline

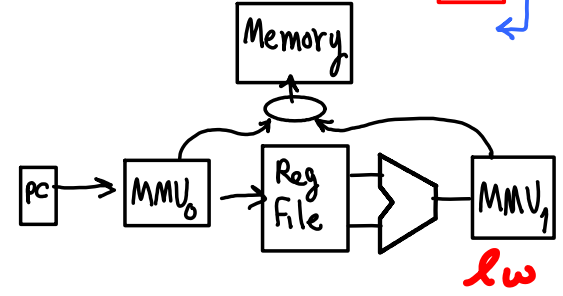
LC3: LDI  
2 refs to data memory

lw  
add

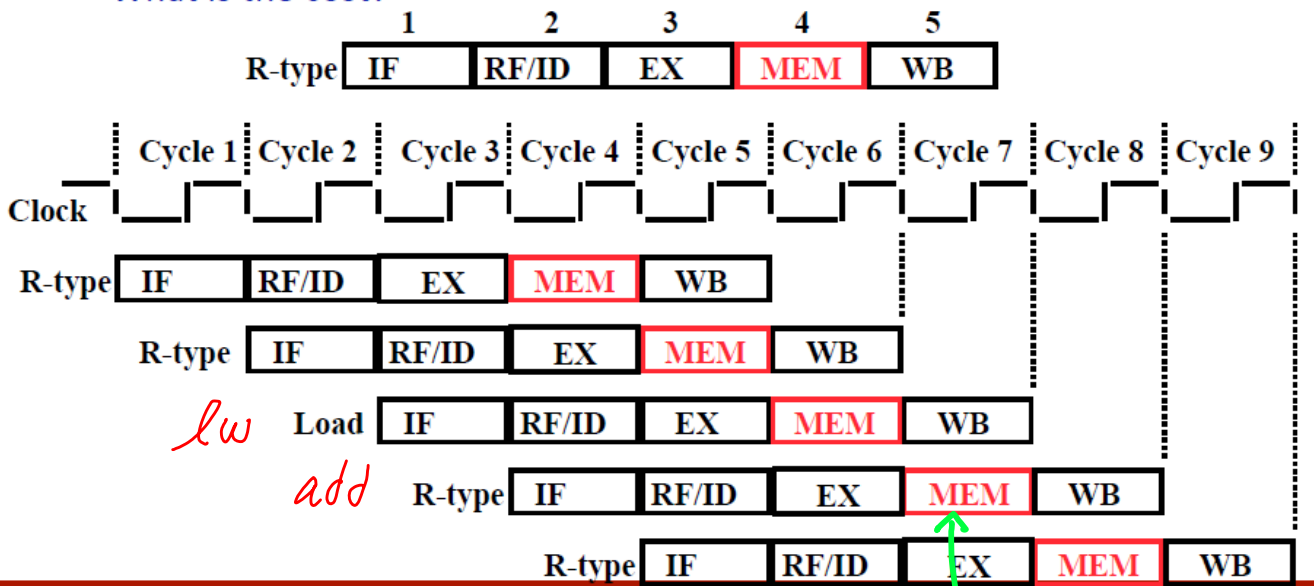


## STRUCTURAL HAZARD

- Simple example: MIPS pipeline with a single unified memory
  - No separate instruction & data memories
- Load/store requires data access
- Instruction fetch would have to stall for that cycle
  - Would cause a pipeline "bubble"



- Delay R-type register write by one cycle → *don't short-circuit*
- Does this increase the CPI of instruction? → *what was CPI above?*
- What is the cost?



## *sequential consistency* Data Dependencies

- Data dependencies for instruction  $j$  following instruction  $i$

Read after Write (RAW) (true dependence)

- Instruction  $j$  tries to read before instruction  $i$  tries to write it

Write after Write (WAW) (output dependence)

- Instruction  $j$  tries to write an operand before  $i$  writes its value

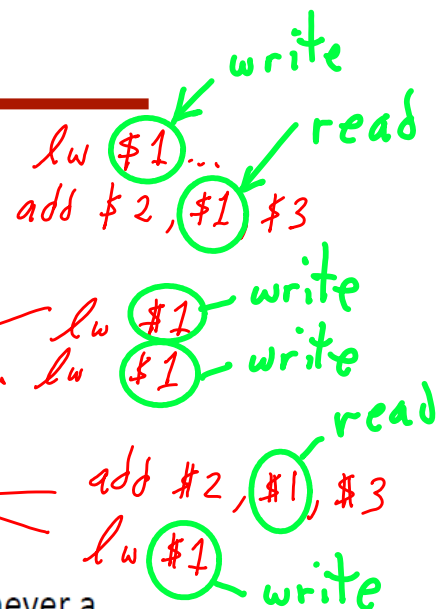
Write after Read (WAR) (anti dependence)

- Instruction  $j$  tries to write a destination before it is read by  $i$

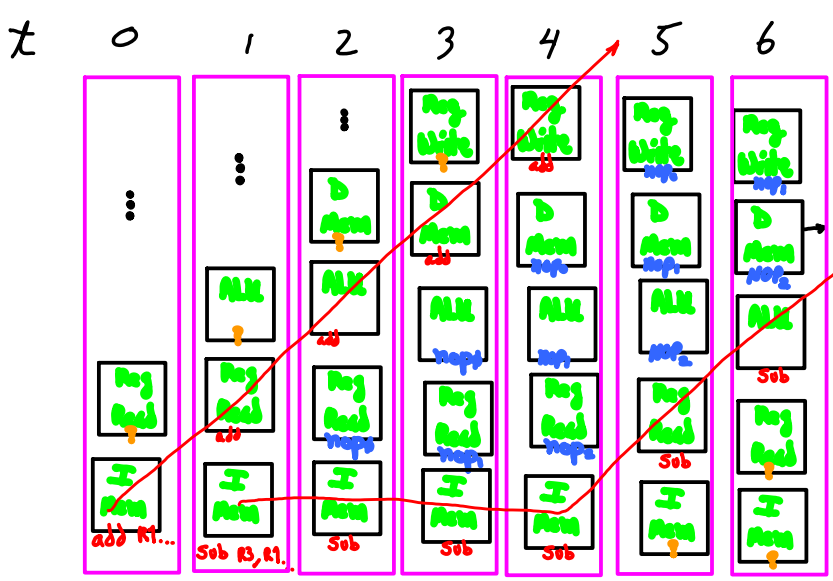
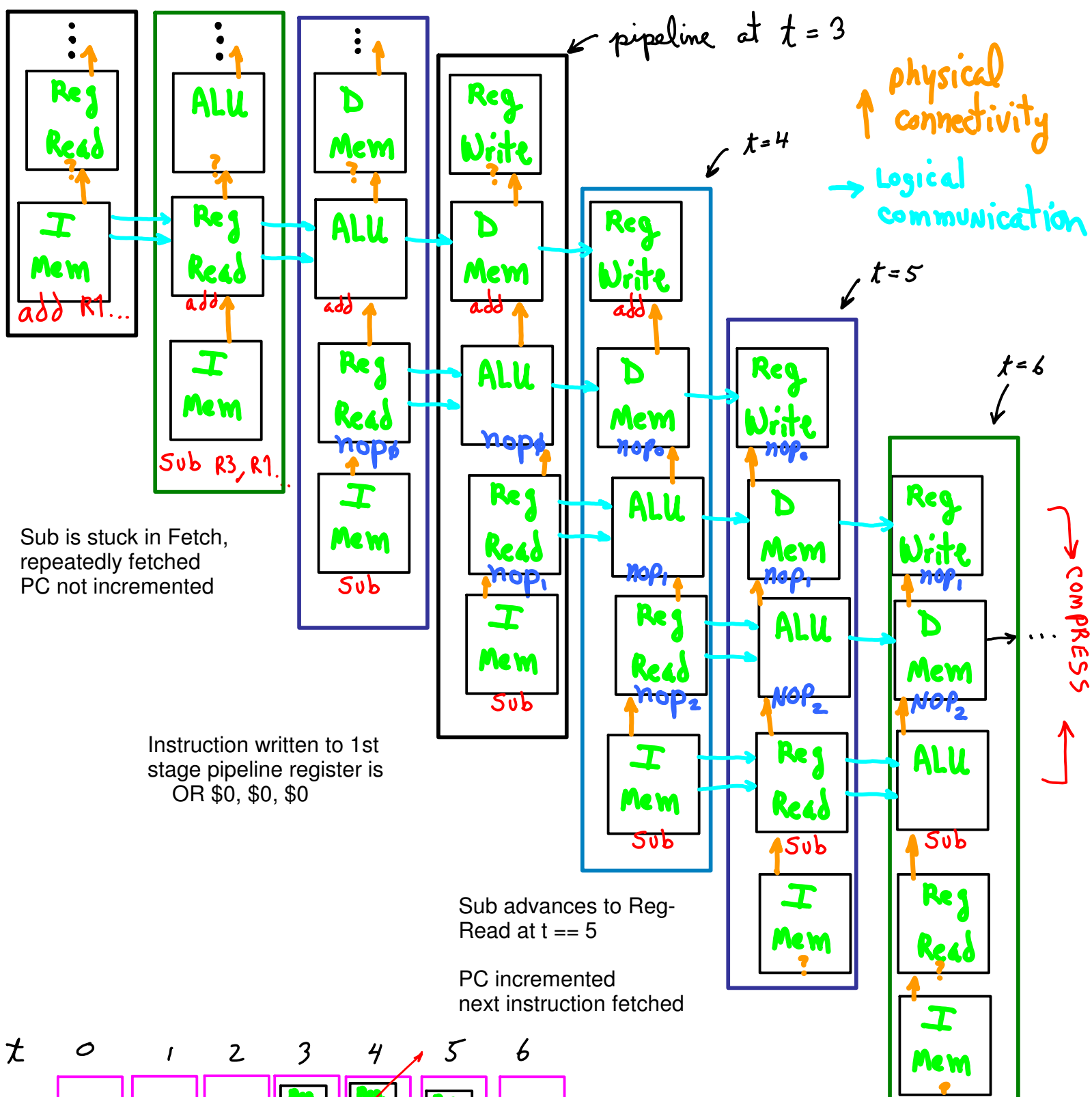
- No such thing as a Read after Read (RAR) hazard since there is never a problem reading twice

*Cannot re-order effects of operations!*

- Dependencies are a property of your program (always there)
- Dependencies may lead to hazards on a specific pipeline







if we line up each physical pipeline in parallel, instructions flow up through the stages. Logical communication is along diagonals.

# STALLS + Performance

Suppose 40% of instructions cause 3-bubble stalls

$$T_{w/stalls} = (n \text{ instructions}) \left[ (60\%)(1 \text{ cycle}) + (40\%)(1 \text{ cycle} + 3 \text{ bubbles}) \right] (1/CR)$$

$$T_{w/stalls} = (n \text{ instructions}) \left[ (100\%)(1 \text{ cycle}) \right] (1/CR)$$

$$S_{w/-w/o} = \frac{[1 \cdot 1]}{[0.6 + 0.4(4)]} = \frac{1}{2.2} < \frac{1}{2}$$

added logic

## How to Stall the Pipeline

### OR How to Insert a NOP or Bubble

- You discover the need to stall when 2<sup>nd</sup> instruction is in ID stage
  - Idea: repeat its ID stage until hazard resolved; let all instructions ahead of it move forward; stall all instructions behind it

1. Force control values in ID/EX register a NOP instruction

- As if you fetched or \$0, \$0, \$0
- When it propagates to EX, MEM and WB on following cycles, nothing will happen (nop = no-operation)

NOP  
OR R1, R1, R1  
AND R1, R1, R1

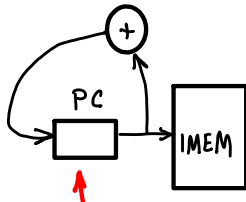
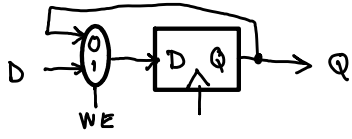
1 or 1 = 1  
0 or 0 = 0  
1 and 1 = 1  
0 and 0 = 0

OR, set WE=0

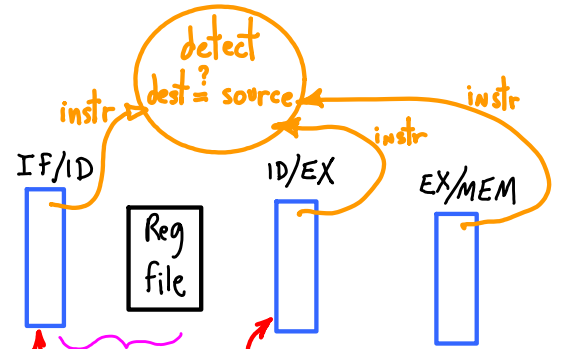
2. Prevent update of PC and IF/ID register

- Using instruction is decoded again
- Following instruction is fetched again

0 → WE



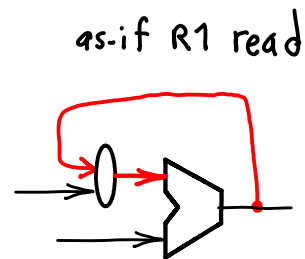
we=0  
following instr stalled



we=0  
dependent instr. stalled  
NOPs enter pipeline

# BETTER THAN STALLING

send data directly to ALU input? Do write later?  
(new feedback path)

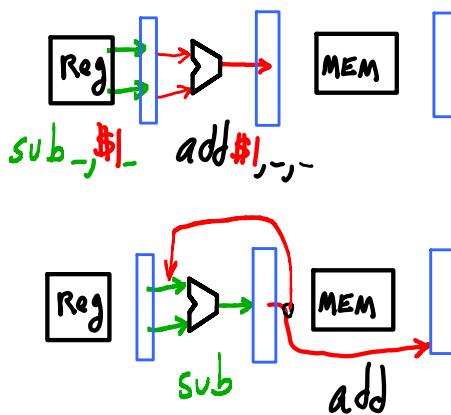


RAW reg-mode?

add \$1, \_, \_  
sub \_, \$1, \_

Data available next tick.

Forwarding (feedback) works.



LW?

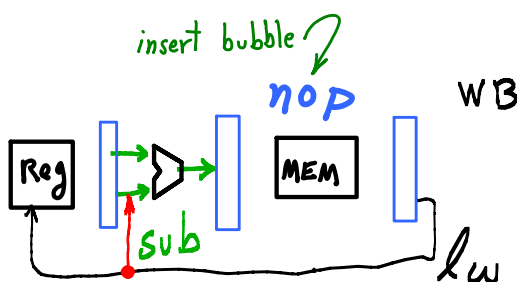
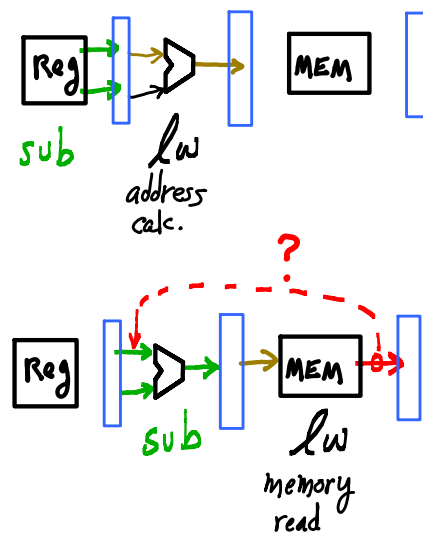
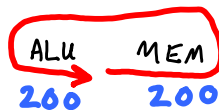
lw \$1, \_, \_  
sub \_, \$1, \_

WHY NOT forward Dmem.out?

DELAY = 200ps (memory) + 200ps (ALU)

SLOW down clock to 400ps?!?

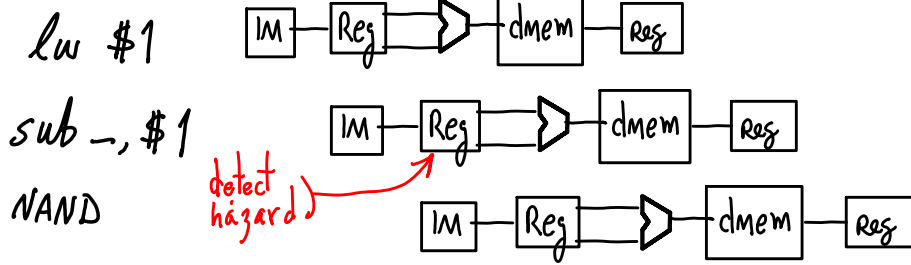
Forward from WB instead, insert NOP



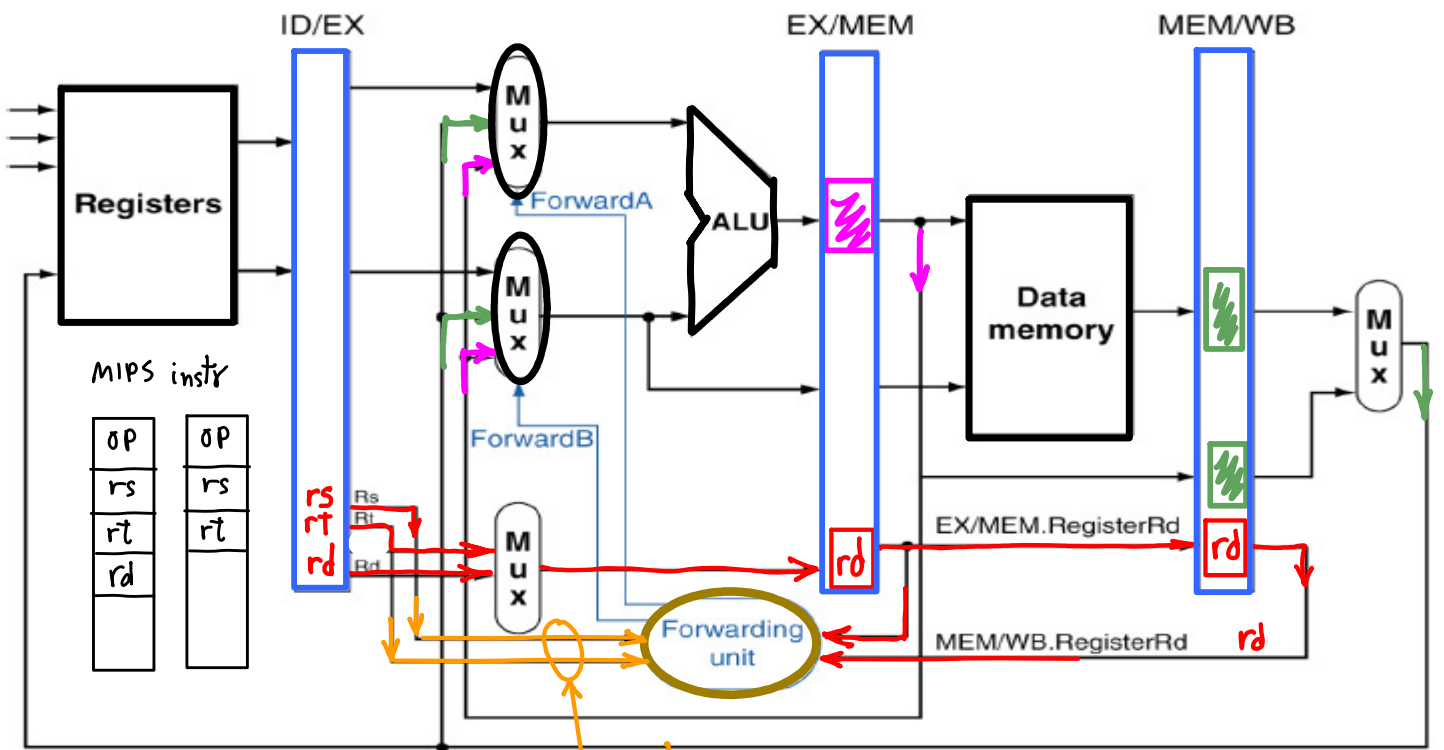
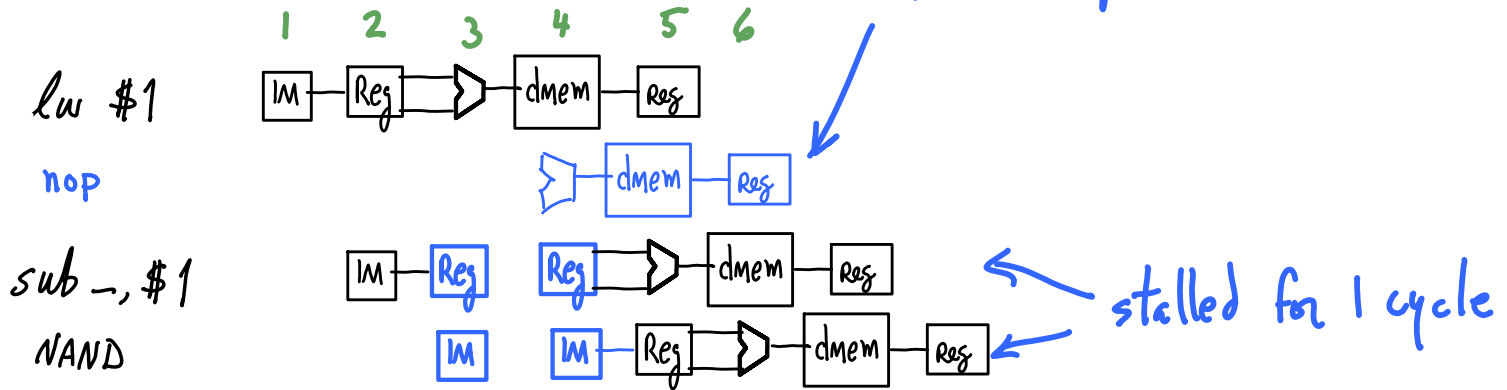


# LW stall

cycle:



as if nop was fetched



b. With forwarding

Forwarding Control:  $ID/EX.(rs, rt) = ? (EX/MEM.rd \text{ or } MEM/WB.rd)$   
 $\implies$  Set MUXes

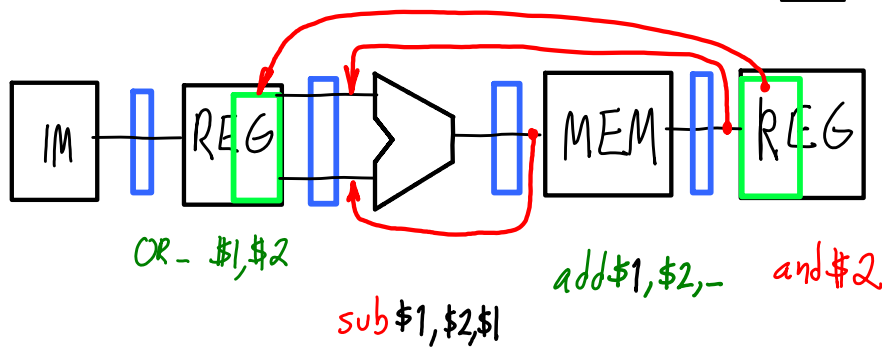
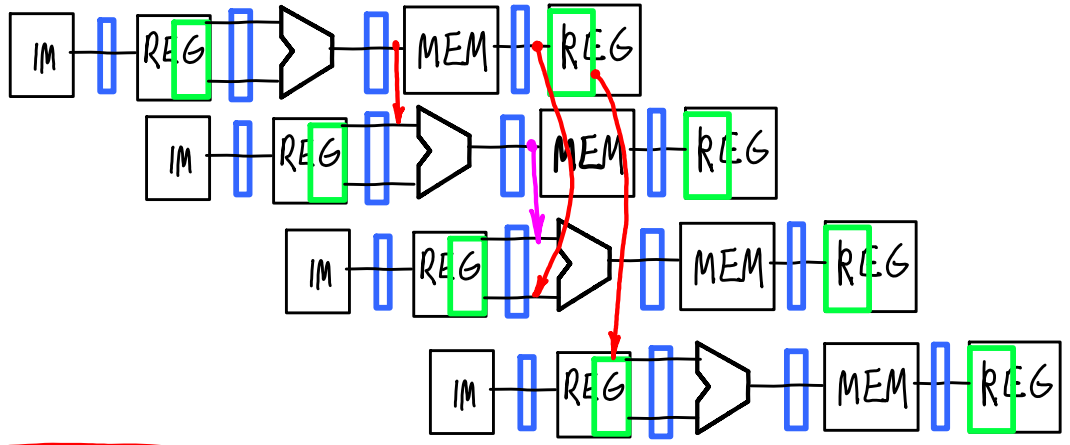
# multiple feedback at once?

AND \$2, \_, \_

ADD \$1, \$2, \_

SUB \_, \$2, \$1

OR \_, \$1, \$2



Can we demonstrate that there aren't any structural hazards for forwarding paths for operate instructions?

Feedback paths to ALU go to both inputs.  
 Hazard detection sets MUXes: Opcode needed in pipe stage registers for detection.

