

# Pipelining

Computer Architecture

ECE 6913

Brandon Reagen



**NYU**

**TANDON SCHOOL  
OF ENGINEERING**

# What we'll cover today

- 1) Review: ISAs and single-cycle MIPS machine
- 2) Pipelining: improving performance by overlapping execution
- 3) Hazards and optimizations



NYU

TANDON SCHOOL  
OF ENGINEERING

# We have TAs!!

As promised.. For weeks 😊

Yichen Wang: [yw4604@nyu.edu](mailto:yw4604@nyu.edu)

Jiazhen Han: [jh6419@nyu.edu](mailto:jh6419@nyu.edu)

Lab and HW experts. Post questions on Piazza and they'll get back to you.

HW 1 solutions also posted.



**NYU**

TANDON SCHOOL  
OF ENGINEERING

# Did anything exciting happen this week?



NVIDIA = GPU Giant. Why?

ARM = Mobile Giant.

What does that even mean?

FB: >2B MAUs. >90% ad rev from mobile.

“A53 represents more than 48%”

“A7 represents more than 15%”

2 old ARM cores (e.g., ARM gives A53s to universities) account for >63% of the mobile market\*

Paper: Machine Learning at Facebook:  
Understanding Inference at the Edge



NYU

TANDON SCHOOL  
OF ENGINEERING

# Trends in Computer Architecture (From GaTech~2008)

50s to 60s: Computer Architecture ~ Computer Arithmetic

70s to mid 80s: Instruction Set Design, especially ISA appropriate for compilers

90s: Speculation: Predict this, predict that; memory system; I/O system;  
Multiprocessors; Networks

2000s: Power efficiency , Communication, On-die Interconnection Network,  
Multi-this, Multi-that. (**We are here**)

2015 and beyond: Thousand-core processors, Self adapting systems? Self  
organizing structures? DNA Systems/Quantum Computing?



NYU

TANDON SCHOOL  
OF ENGINEERING

# Today

2015 and beyond: Thousand-core processors, Self adapting systems? Self organizing structures? DNA Systems/Quantum Computing? **Not so much**

- 1) Hardware for ML (e.g., nvidia, TPU, cerebras, etc.)
- 2) HW-SW co-design: understand how code is running, find bottleneck, fix it (in both hw and sw)
- 3) SoC design and Accelerator-centric computing!
- 4) New ISAs? New killer-workloads making us go back in time
- 5) Some quantum/dna stuff,

Hard job. Need to be expert in multiple fields



NYU

TANDON SCHOOL  
OF ENGINEERING

# On getting active

Just go for it..

Read papers: Top venues: ISCA, MICRO, HPCA, ASPLOS

Come up with ideas and ask NYU profs what they think



NYU

TANDON SCHOOL  
OF ENGINEERING

# Instruction Set Architecture (ISA)

Contract between the programmer (SW) and machine (HW)

- ISA defines instructions and their functionality
  - Visible state of the machine, what the programmer/compiler can see
  - Changes to state as instructions are processed
- Underlying machine implementation faithfully executes instructions
- Programmers/compiler express functions as series of instructions

Specification: Instructions and State

- Exactly how each instruction changes the machine's state
- Instruction and memory representation

What does an ISA NOT do?

- How instructions are implemented
- Instruction performance or energy usage
- These are left to the microarchitecture





# MIPS Summary

## 32b RISC architecture

- Fixed instruction length
- 32b registers and memory addresses
- Byte addressable memory.. How much?

## Load-store architecture

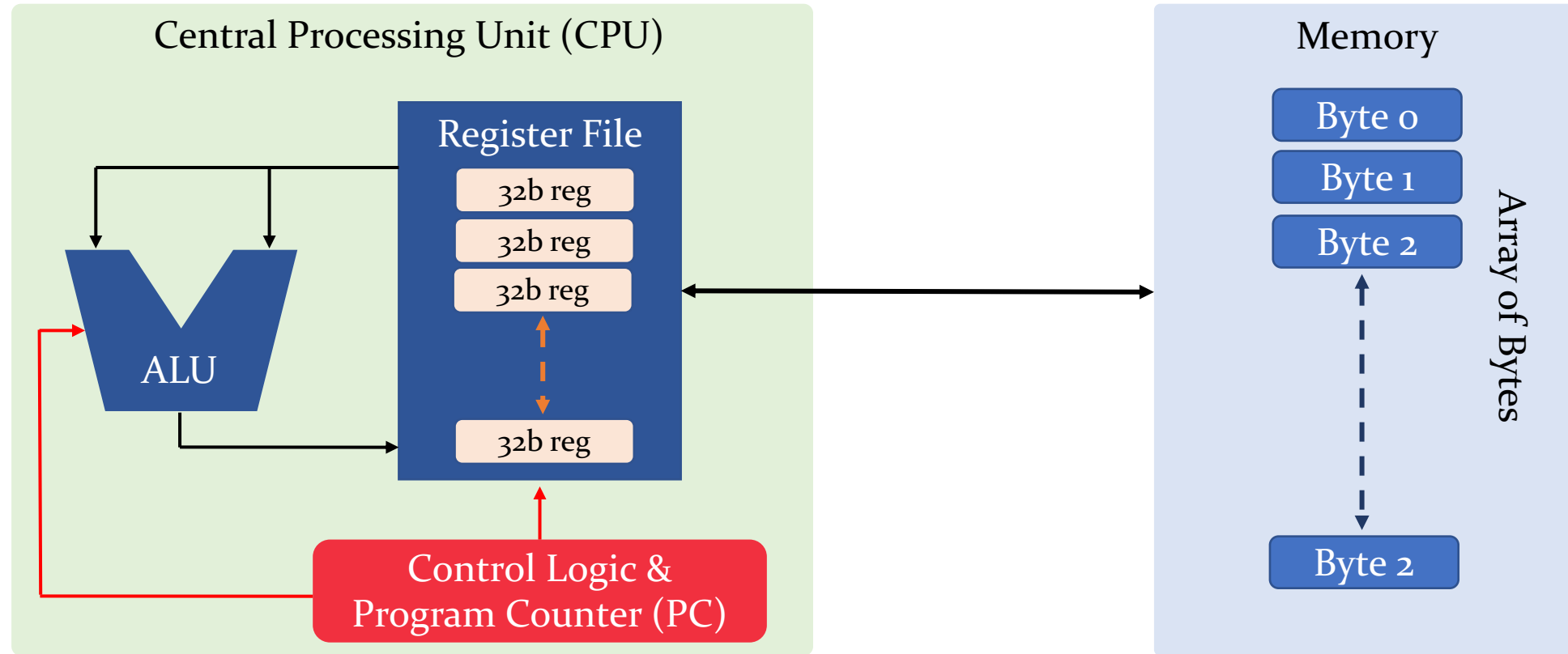
- All compute use values from registers
- Only special instructions access memory

## 32 general purpose registers

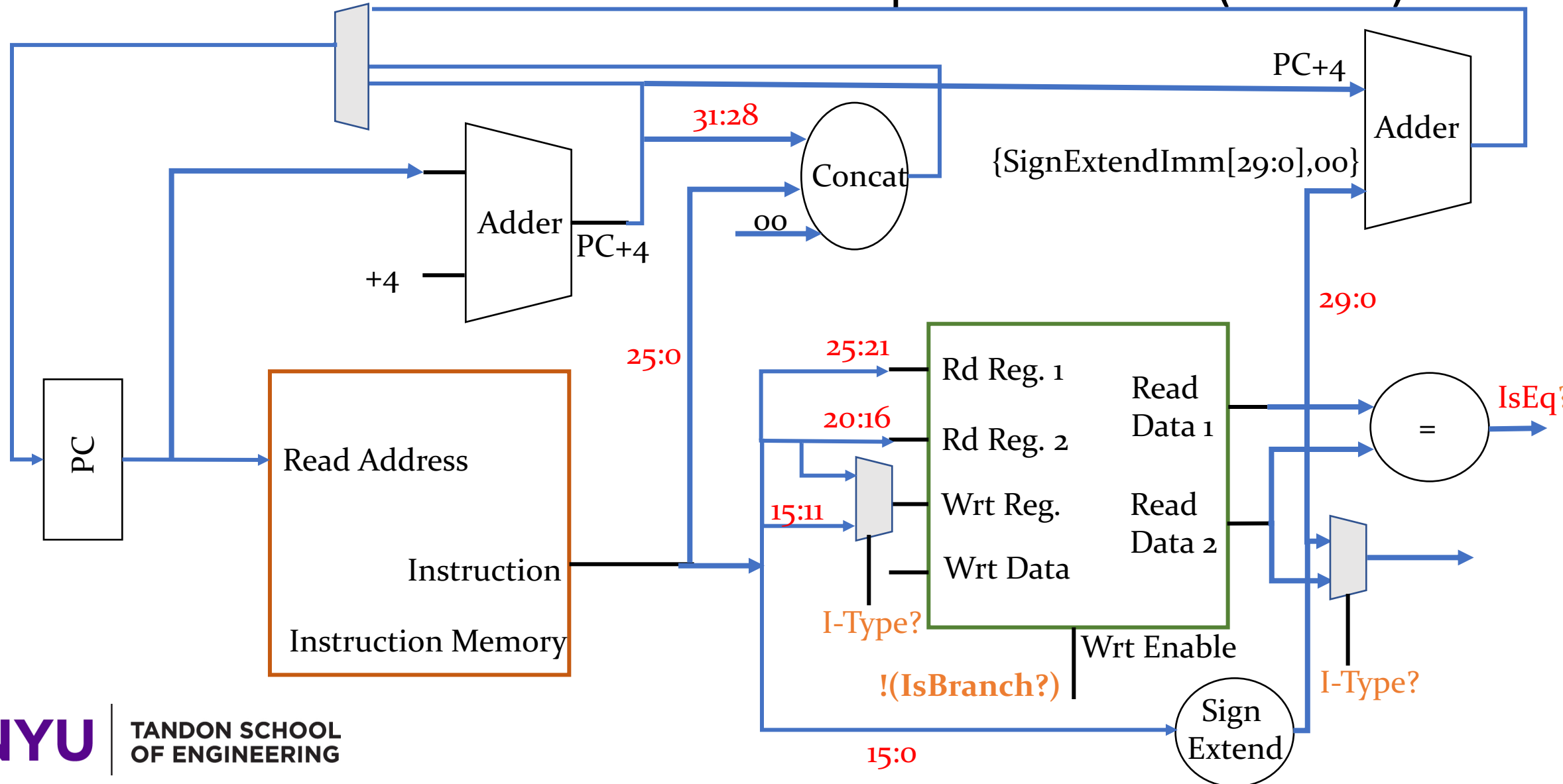
- We will assume the simplified calling convention
- Program counter for finding next instruction
- Hi/lo registers for multiplication



# Architectural (ISA) view of a processor



# Microarchitectural view of a processor (front)

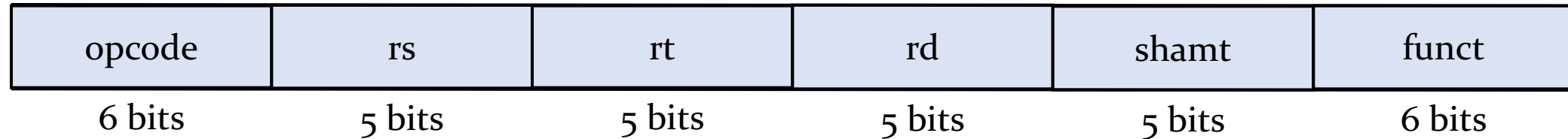


The diagram illustrates the internal structure of a processor, showing the flow of data and control signals. Key components and their interactions are as follows:

- Instruction Register:** Provides the instruction bits (25:21, 20:16, 15:11, 15:0) to the Register File and other control logic.
- Register File:** Contains Read Registers (Rd Reg. 1, Rd Reg. 2) and Write Registers (Wrt Reg., Wrt Data). It outputs Read Data 1 and Read Data 2 to the ALU.
- ALU (Arithmetic Logic Unit):** Performs the ADD operation on Read Data 1 and Read Data 2. It outputs the ALU result to the Memory Address and the Sign Extend block.
- Sign Extend:** Takes the 15:0 bits of the instruction and the ALU result to produce the final 32-bit address for the Memory.
- Memory:** Consists of Address, Write Data, Read Data, and Data Memory. It outputs Read Mem and Write Mem signals.
- Control Signals:**
  - isLoad?:** Indicates if the instruction is a load instruction.
  - isStore?:** Indicates if the instruction is a store instruction.
  - Wrt Enable:** Controls the Write Enable signal for the Register File.

The diagram is labeled with "ADD" and "isLoad?" and "isStore?". The bottom left corner features the NYU TANDON SCHOOL OF ENGINEERING logo.

# Question from last time: ISA design: Shamt

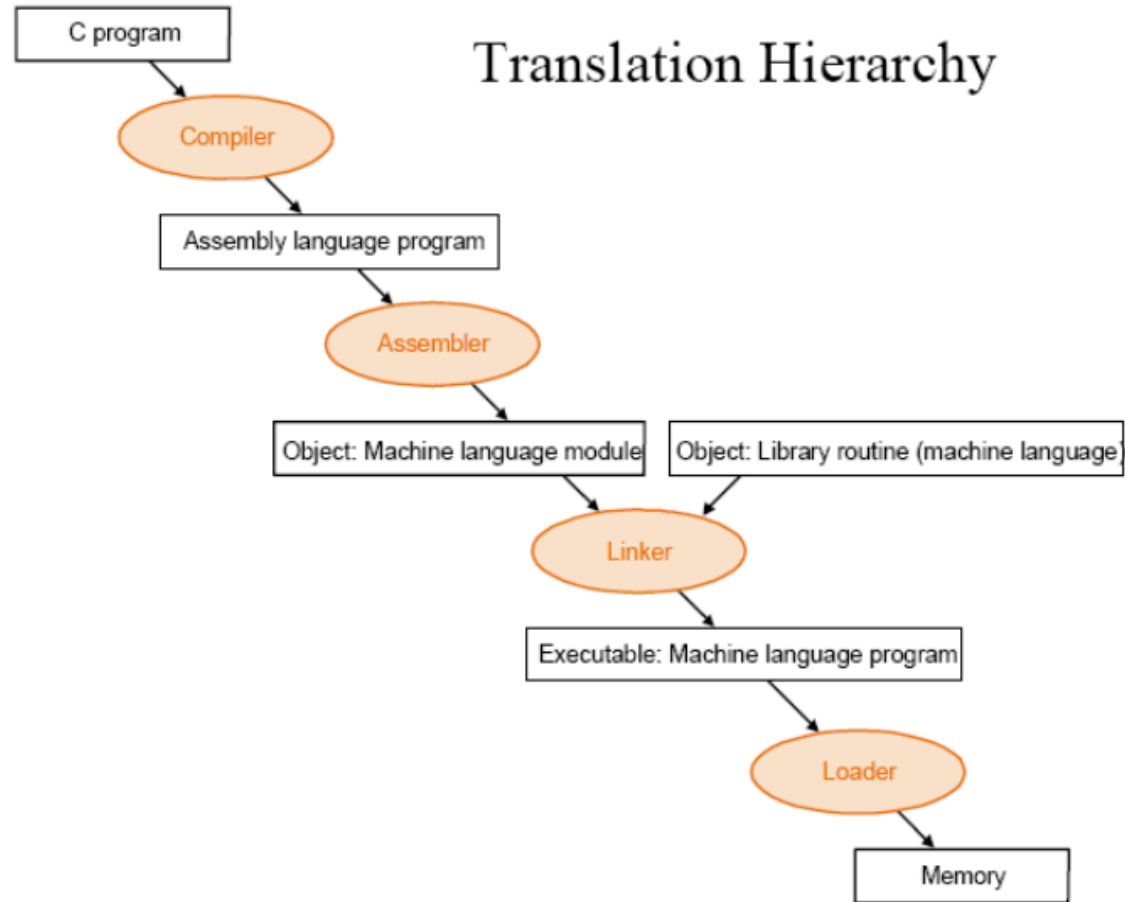


**Shamt** says how many bits to shift (left or right)  
In MIPS, it's part of the ISA.

Why is this a clever design decision?  
How else could you do it?  
What are the performance tradeoffs?



# Question from last time: Pseudoinstruction set



**Move \$r1, \$r2** is actually **or \$r1, \$r2, \$r0**  
What's \$r0?

Why do we have these?



**NYU**

**TANDON SCHOOL  
OF ENGINEERING**

# What makes a good ISA?

## Simplicity favors regularity

- Instruction size, instruction format, data format
- Makes hardware implementation cleaner

## Smaller is faster

- Fewer bit to move, write, and read per instruction
- Register file is faster than memory

## Make the common case fast

- Constants tend to be small, immediate field optimized for this

## Good design demands compromise

- Special formats for important exception
- E.g., jumping far away (as we saw)



# Synchronous Implementation

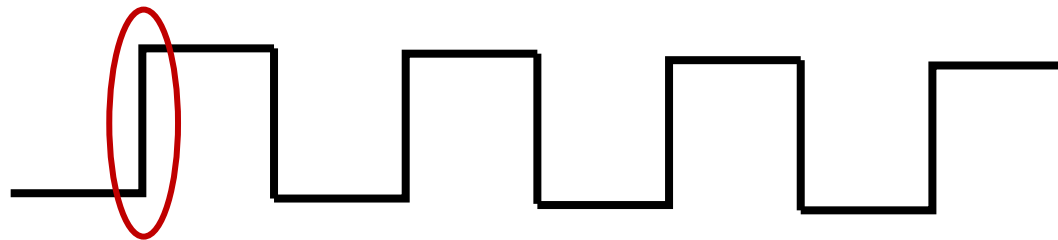
Most digital circuits make use of a periodic “clock”

Synchronous elements enable us to store state

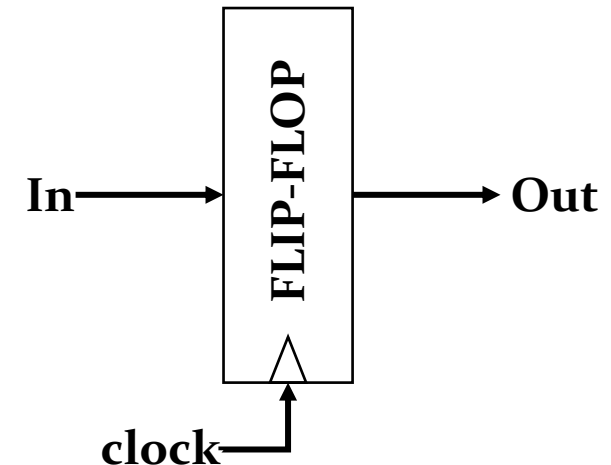
Flip-flop: basic component used in synchronous designs

- Looks at In at every positive clock edge and sets Out = In
- Holds the value of Out steady till next positive edge

Positive Edge



←→  
Clock Period (seconds)

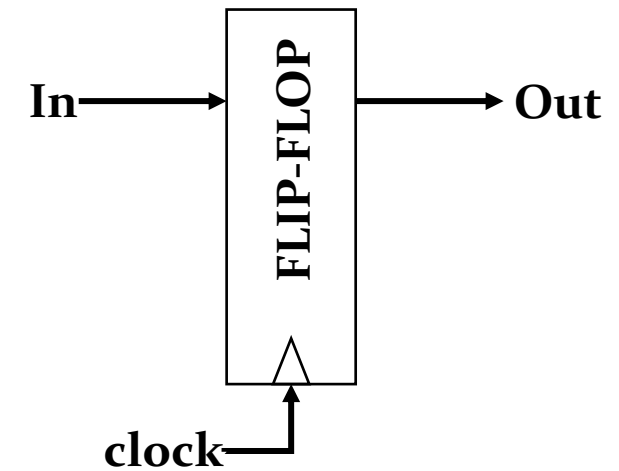
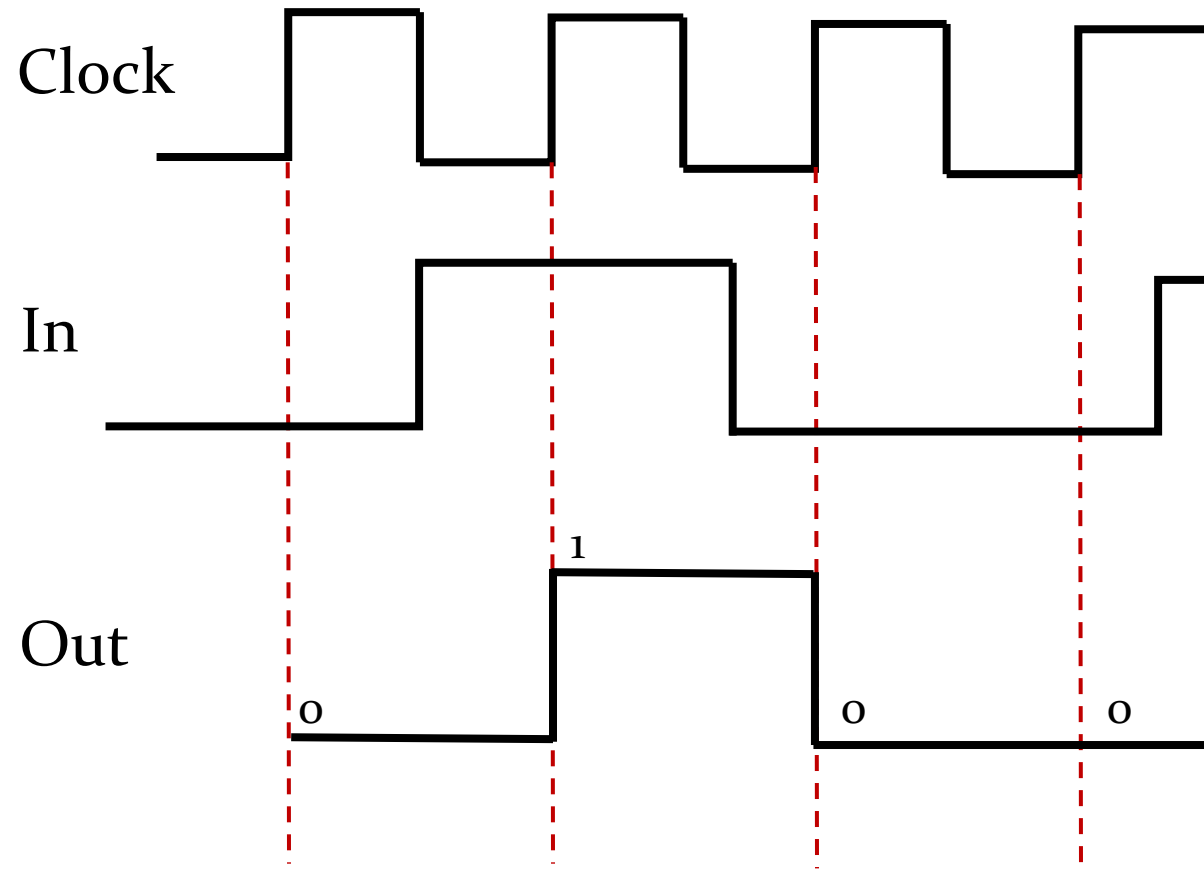


NYU

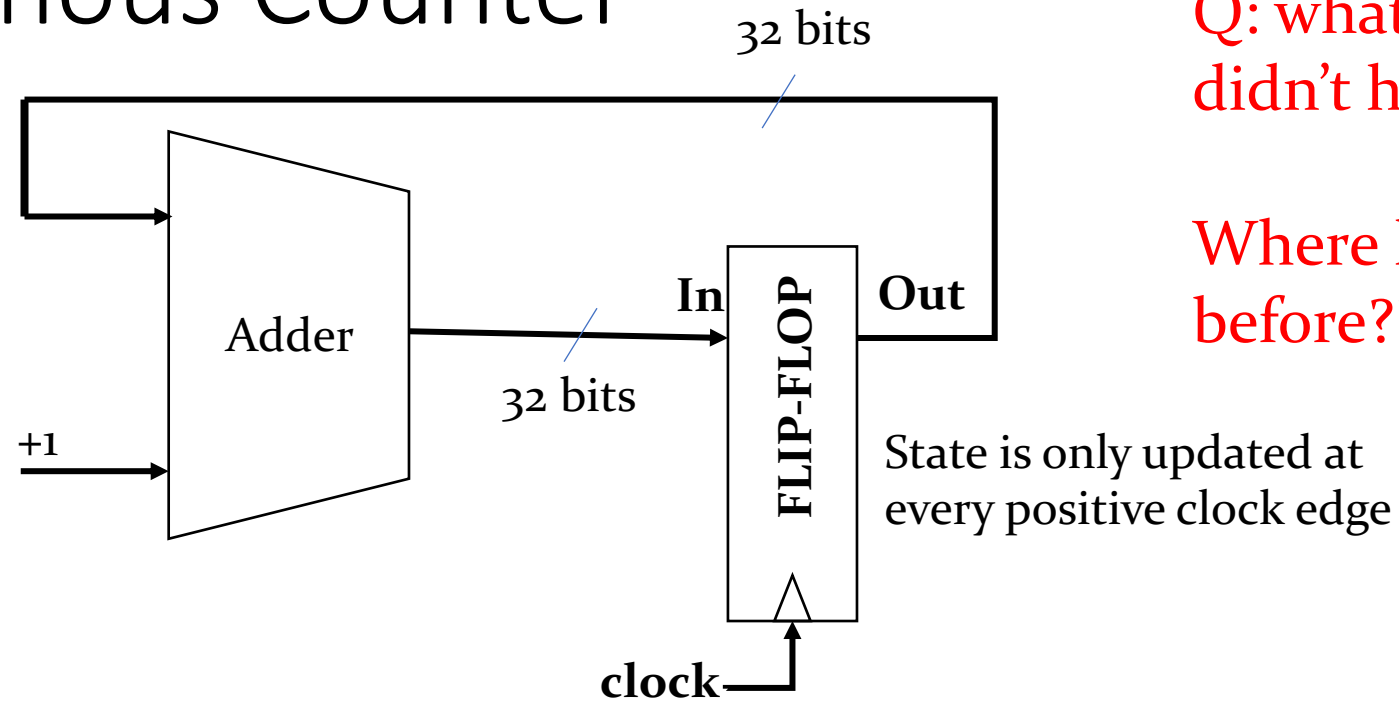
TANDON SCHOOL  
OF ENGINEERING



# Flip-flop Operation

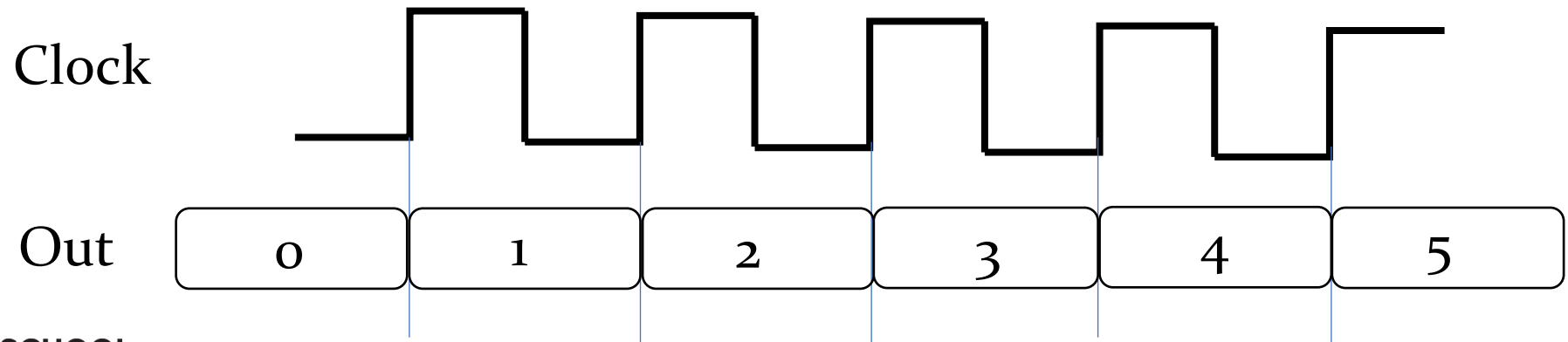


# Synchronous Counter



Q: what happens if we didn't have a FF?

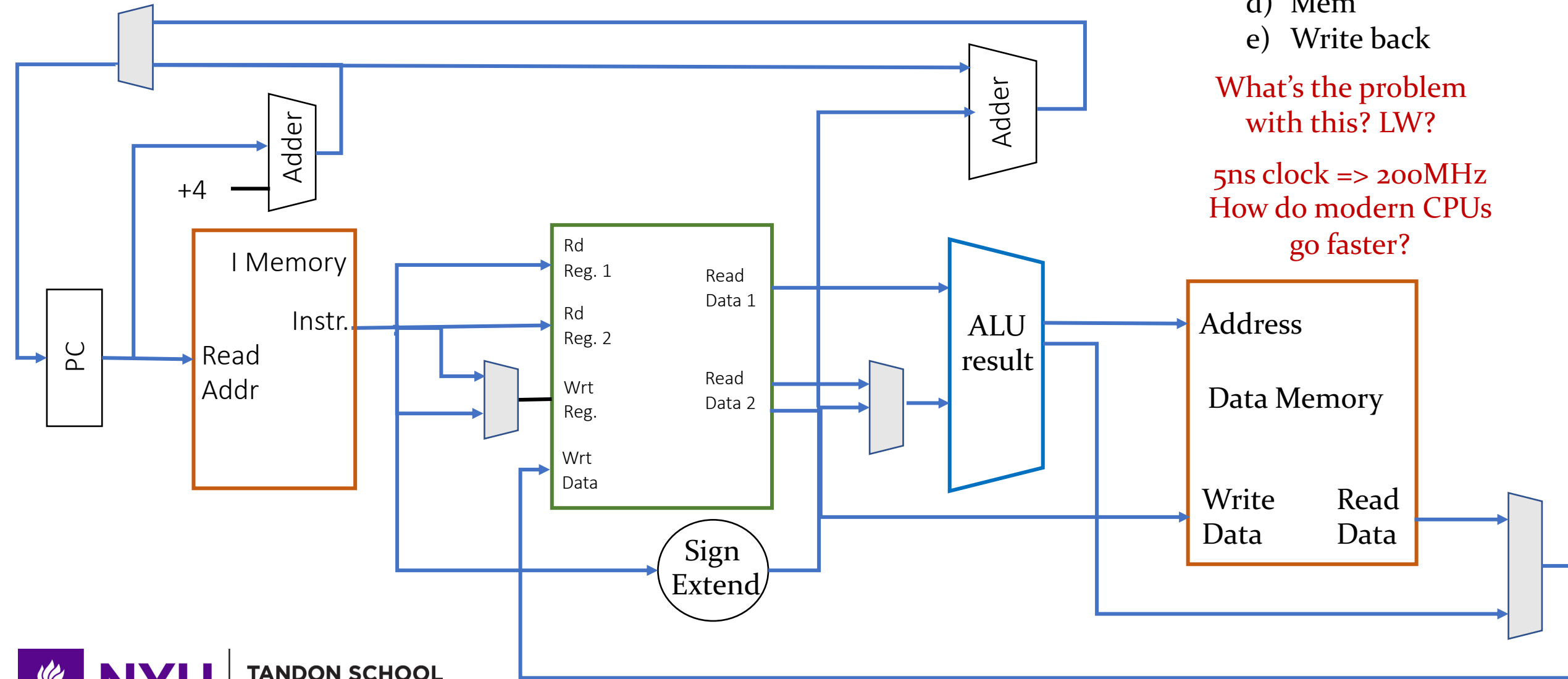
Where have we seen this before??



NYU

TANDON SCHOOL  
OF ENGINEERING

# From last time (simplified)



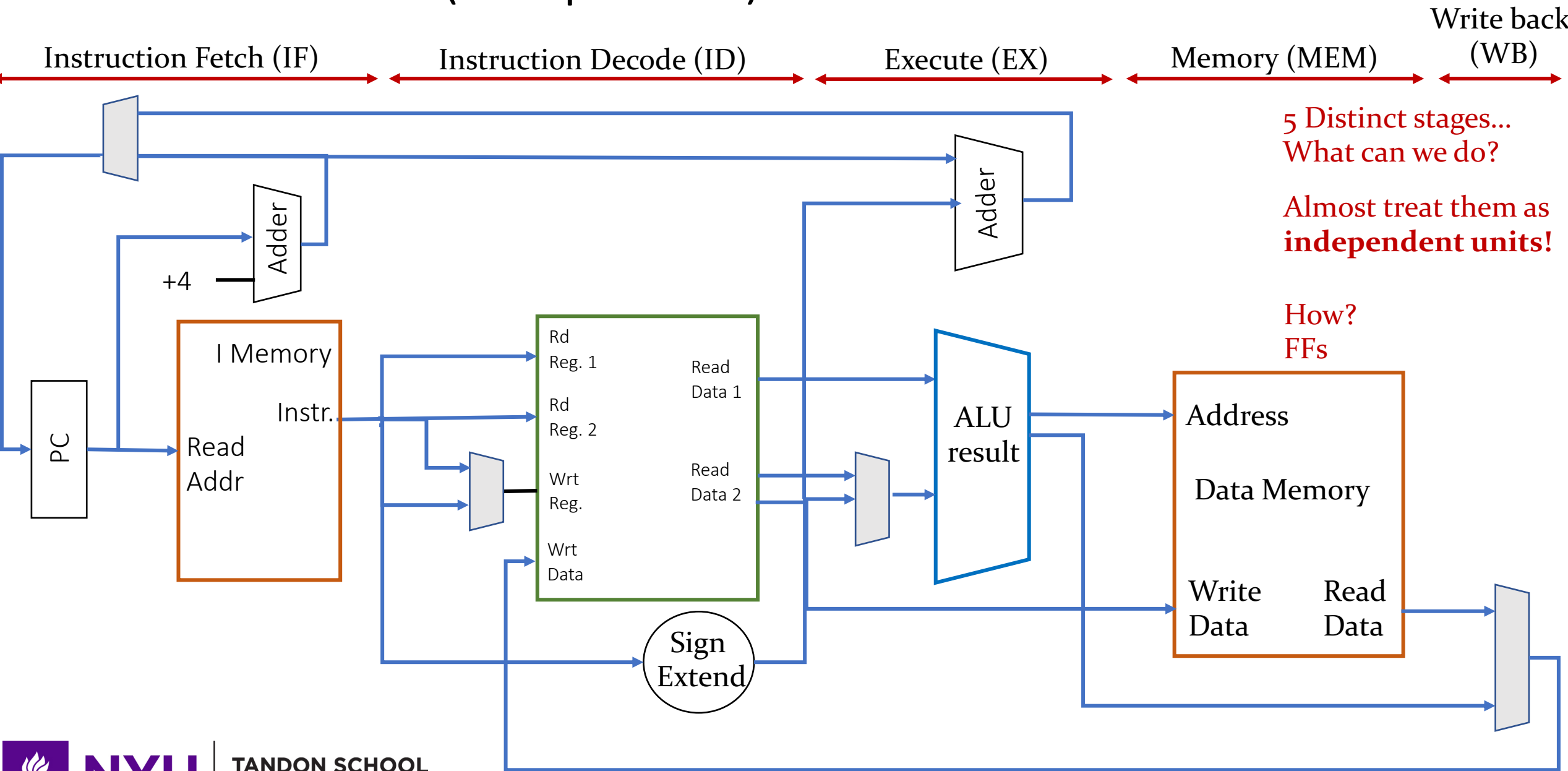
Each takes 1ns..

- a) Imem
- b) RF
- c) ALU
- d) Mem
- e) Write back

## What's the problem with this? LW?

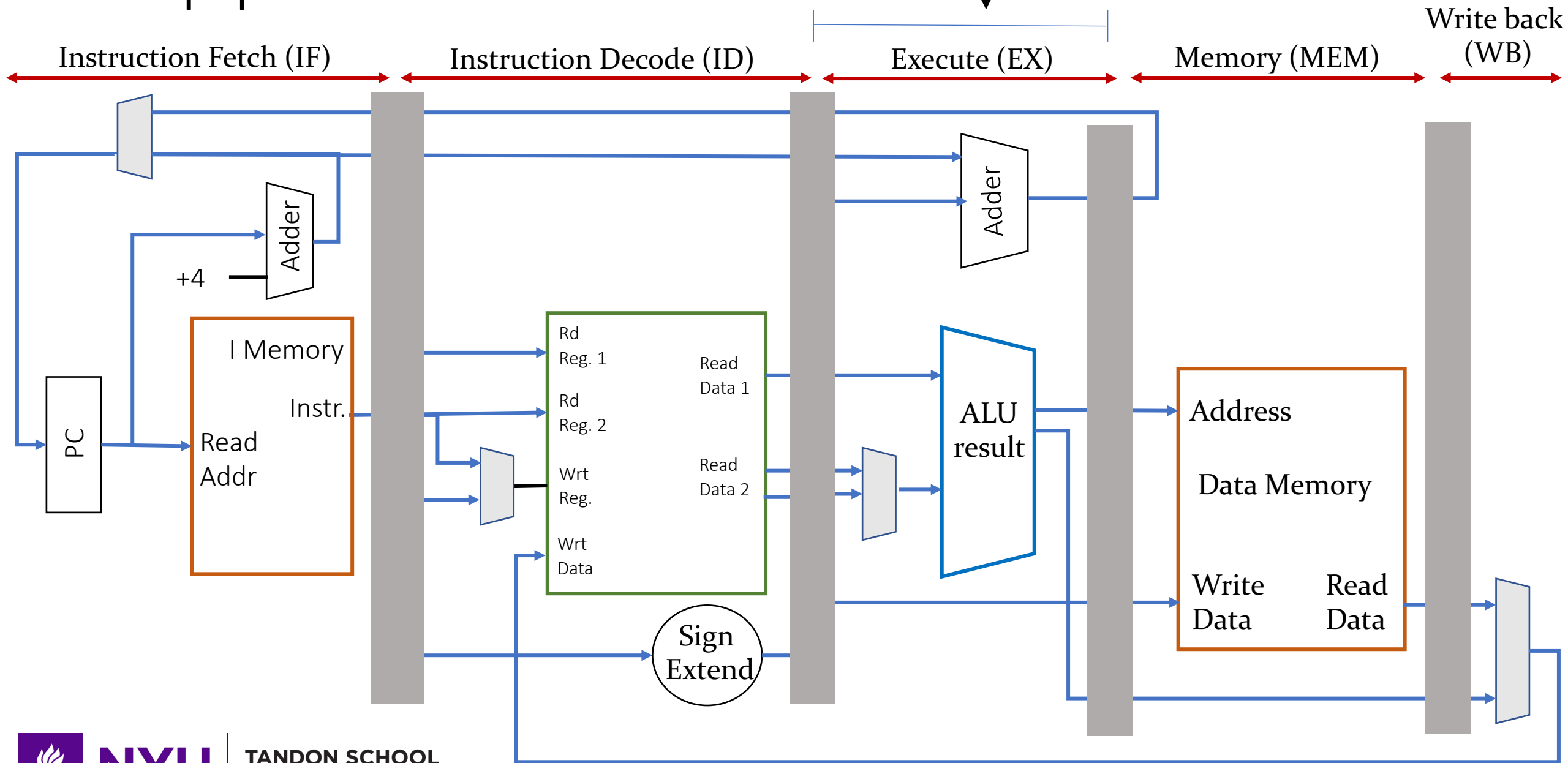
5ns clock => 200MHz  
How do modern CPUs  
go faster?

# From last time (simplified)



# Basic pipeline

1 ns clock!! 5x speedup?



# Why does pipelining help performance?

The previous diagram actually hurts perf.. Why?

single cycle machine = 5ns/inst

Pipeline = (1ns logic delay + FF delay) \* 5 stages = >5ns/inst

So we improved our clock speed, but things got slower?

What'd we miss?

Improving performance is related to **hardware utilization**.



NYU

TANDON SCHOOL  
OF ENGINEERING

# Core concept: Laundry example

Ann, Brian, Cathy, Dave  
need to do laundry



Wash = 30m



Dry = 40m



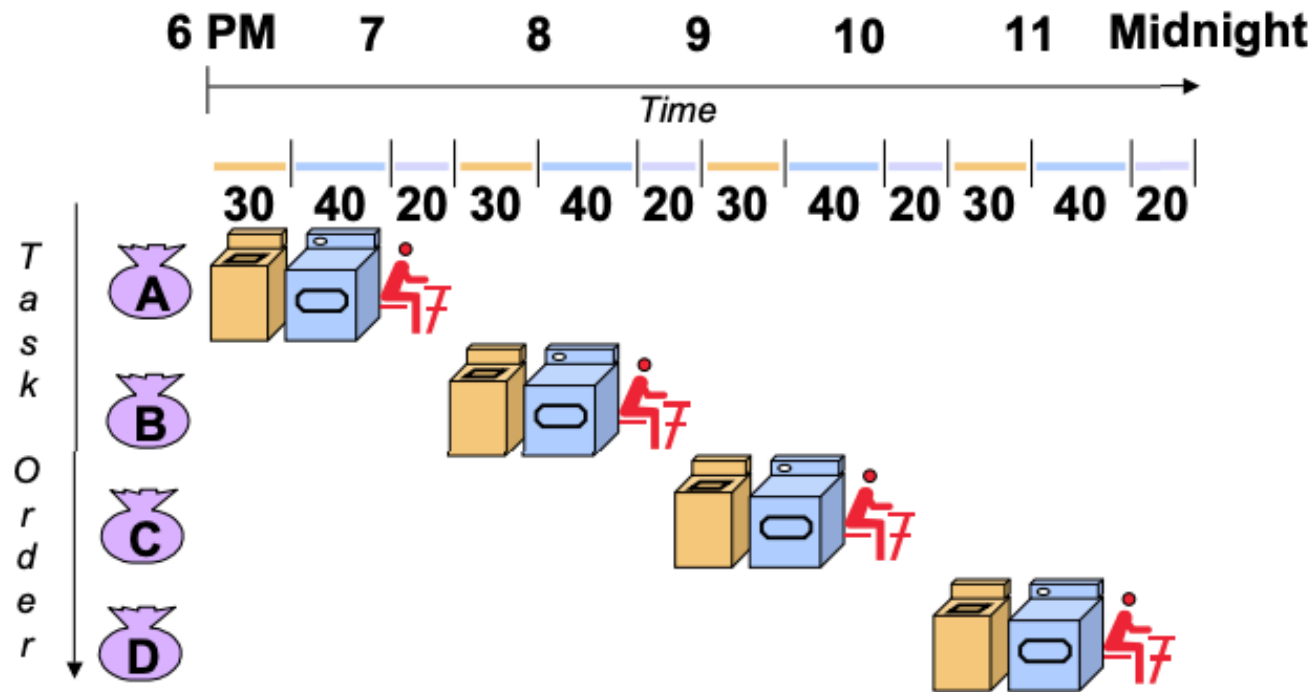
Fold = 20m



NYU

TANDON SCHOOL  
OF ENGINEERING

# “Naïve” Solution: like single-cycle MIPS



Takes 6 hours.

What's wrong with this?

You're not utilizing available resources!

Wasting time.

How can we do better?

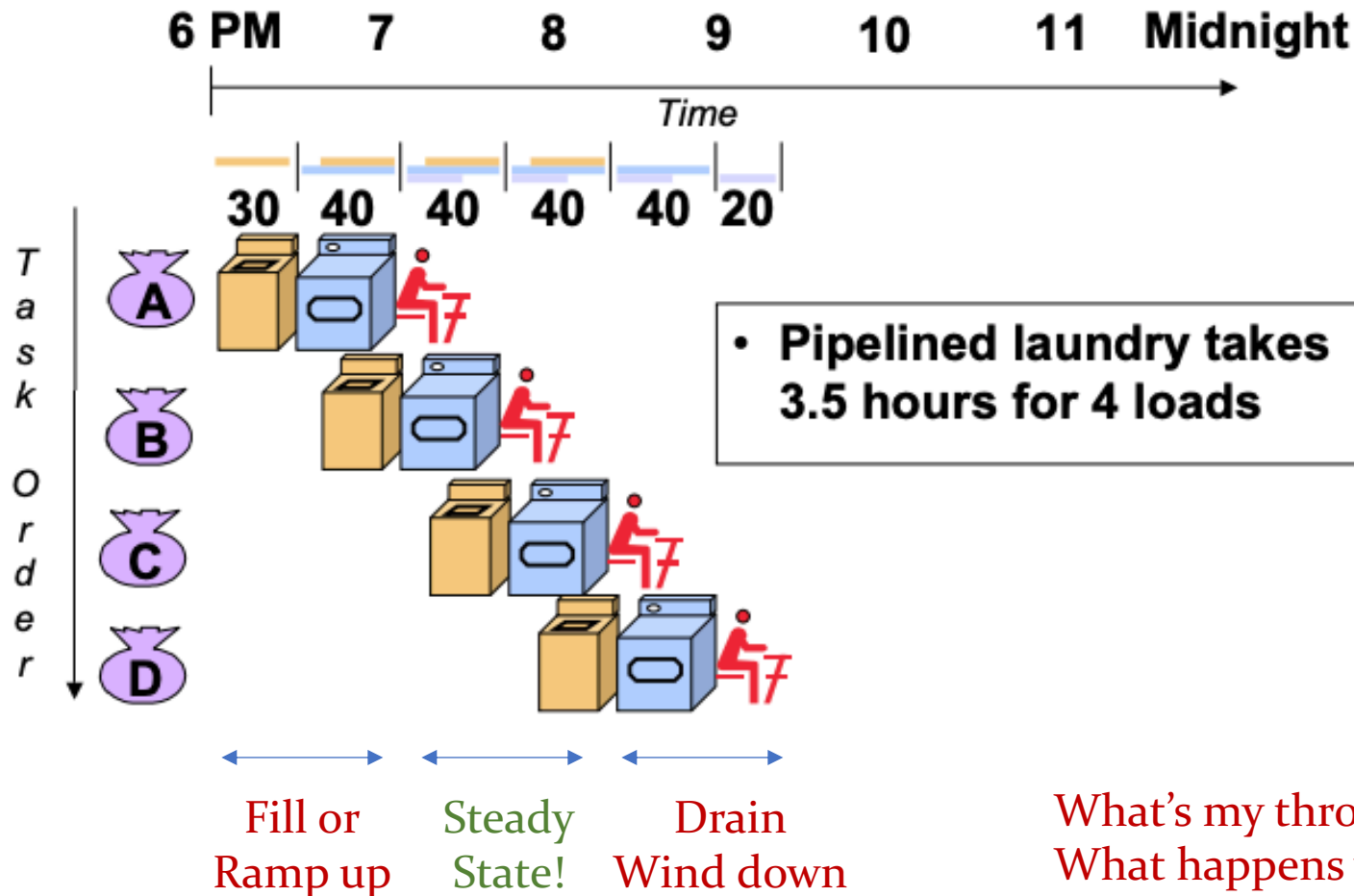


NYU

TANDON SCHOOL  
OF ENGINEERING



# What happens if we overlap work?



Pipelining takeaways:

- Doesn't improve single task latency, it improves throughput
- Pipeline "rate" governed by slowest pipeline stage
- Speedup  $\sim$  number of pipeline stages
- Unbalanced stages reduces speedup
- Time to fill and drain pipe also reduces speedup

What's my throughput?

What happens when I have more loads to do?

What if I own a laundry mat, is this good or bad?

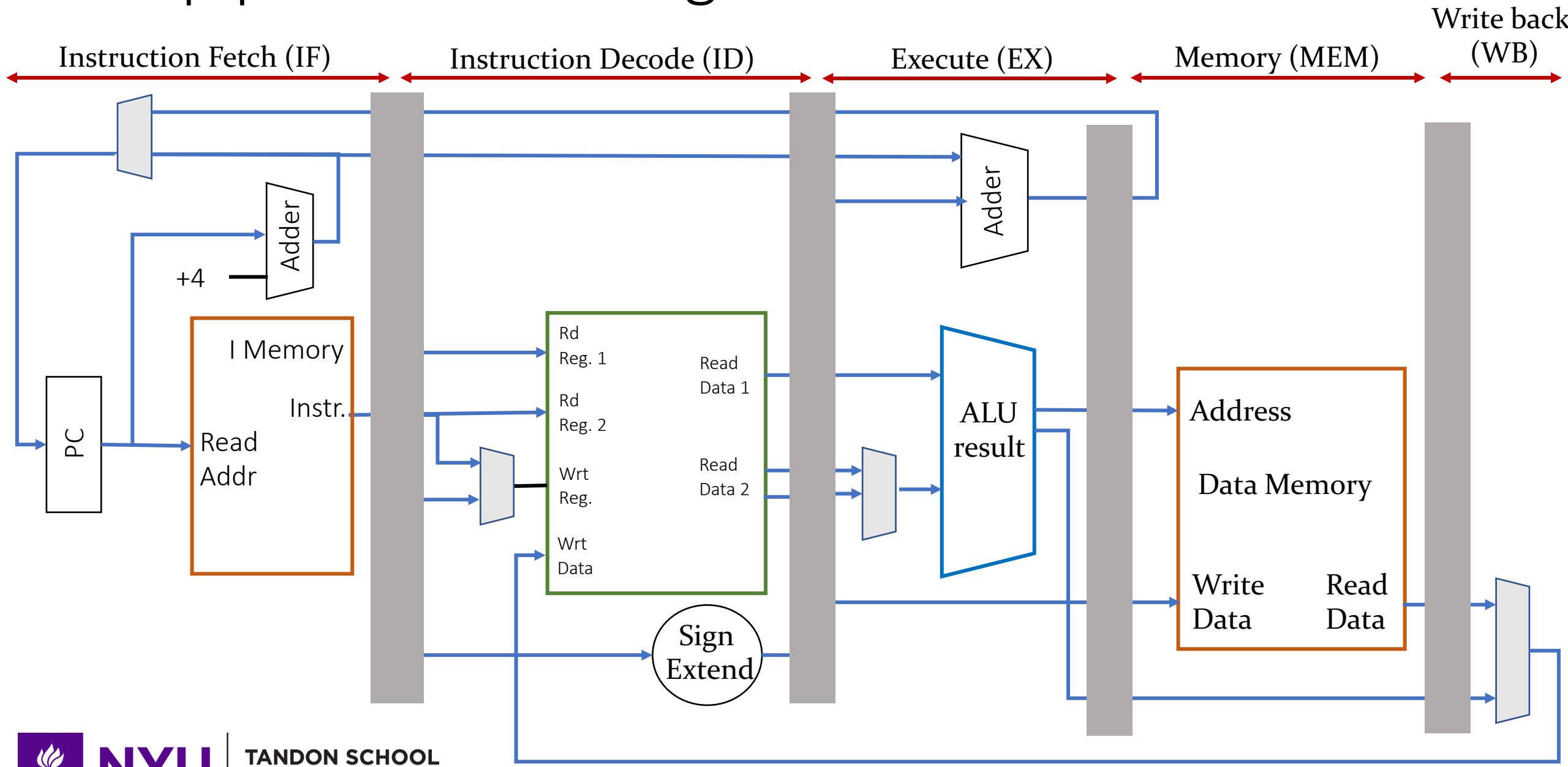
Power bill?



NYU

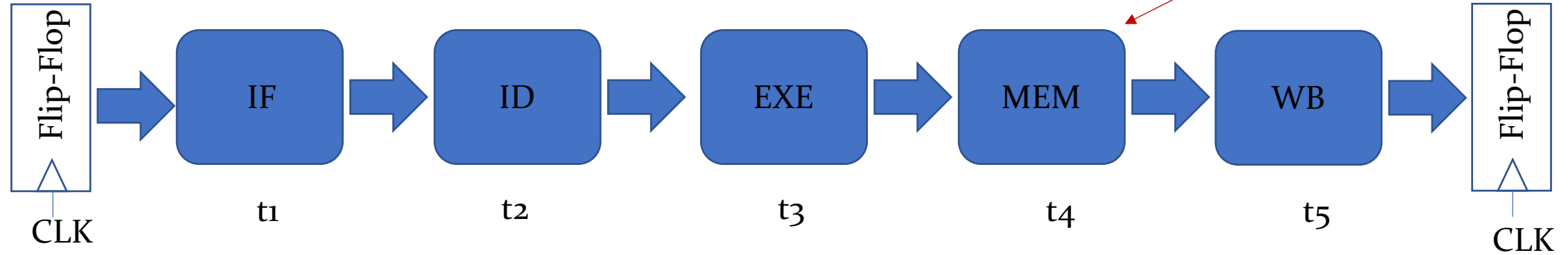
TANDON SCHOOL  
OF ENGINEERING

# Basic pipeline: each stage its own machine!

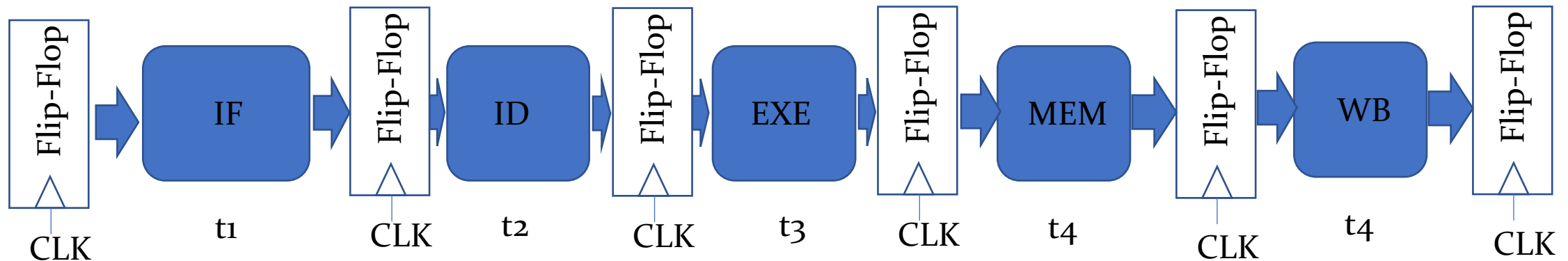


# Improving Throughput: Pipelining

Think about HW utilization when accessing memory...



$$\text{Throughput} = 1/(t_1+t_2+t_3+t_4+t_5)$$



$$\text{Throughput} = ?$$

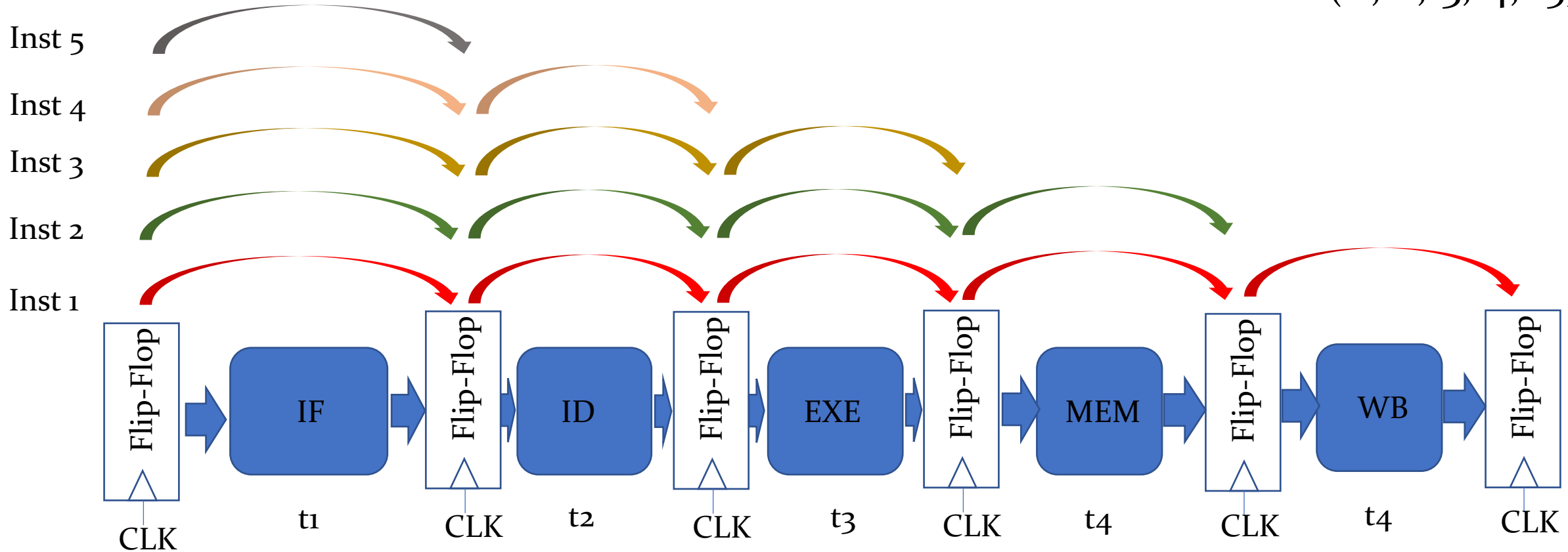


NYU

TANDON SCHOOL  
OF ENGINEERING

# Pipeline Operation

$$\text{Clock Period} = \max(t_1, t_2, t_3, t_4, t_5)$$



In steady state, each module processes data in parallel

$$\text{Throughput} = 1 / \max(t_1, t_2, t_3, t_4, t_5)$$

Q: If each stage takes 1ns, what's our speedup?

Did we miss anything?



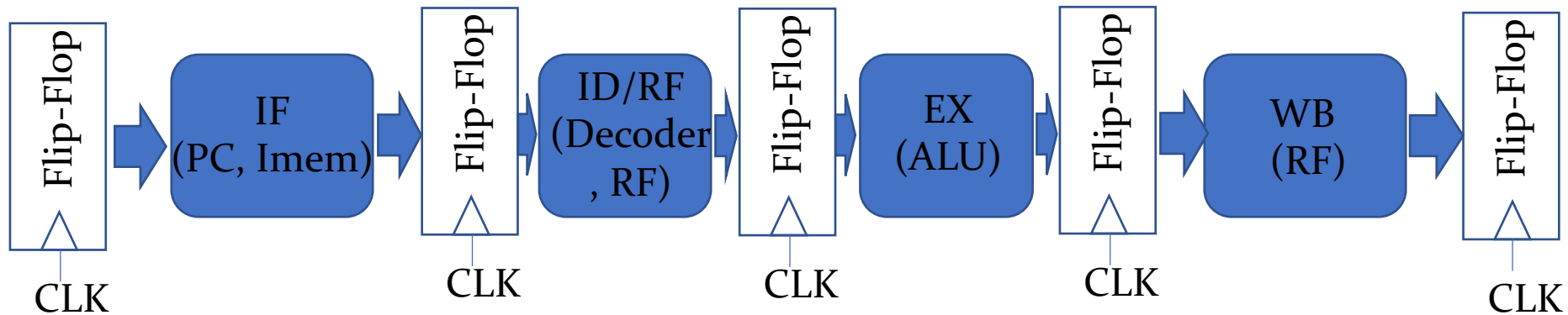
NYU

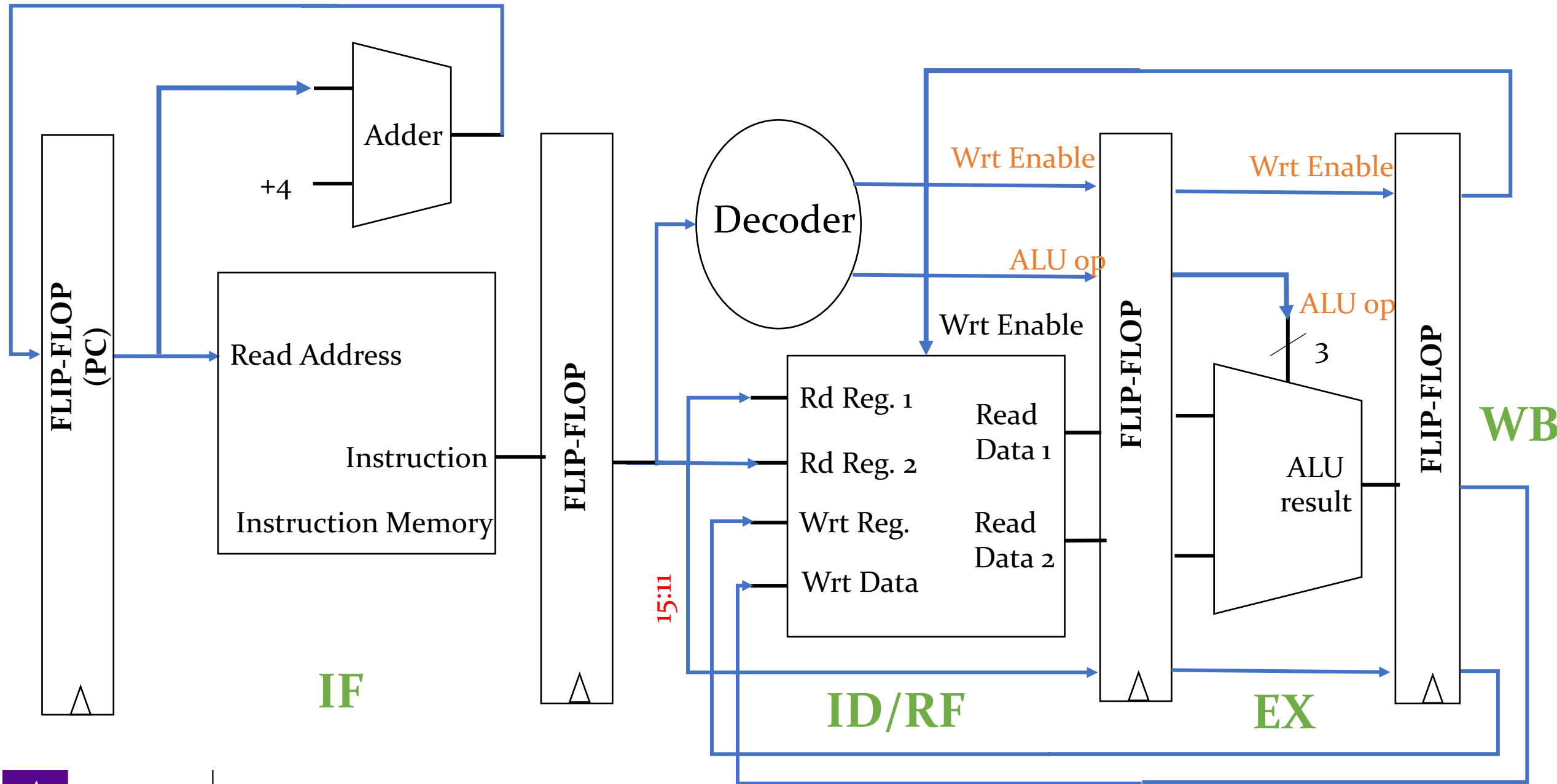
TANDON SCHOOL  
OF ENGINEERING

# Pipelining the Single-Cycle **R-Type** MIPS machine

## Basic steps in executing an R-type instruction

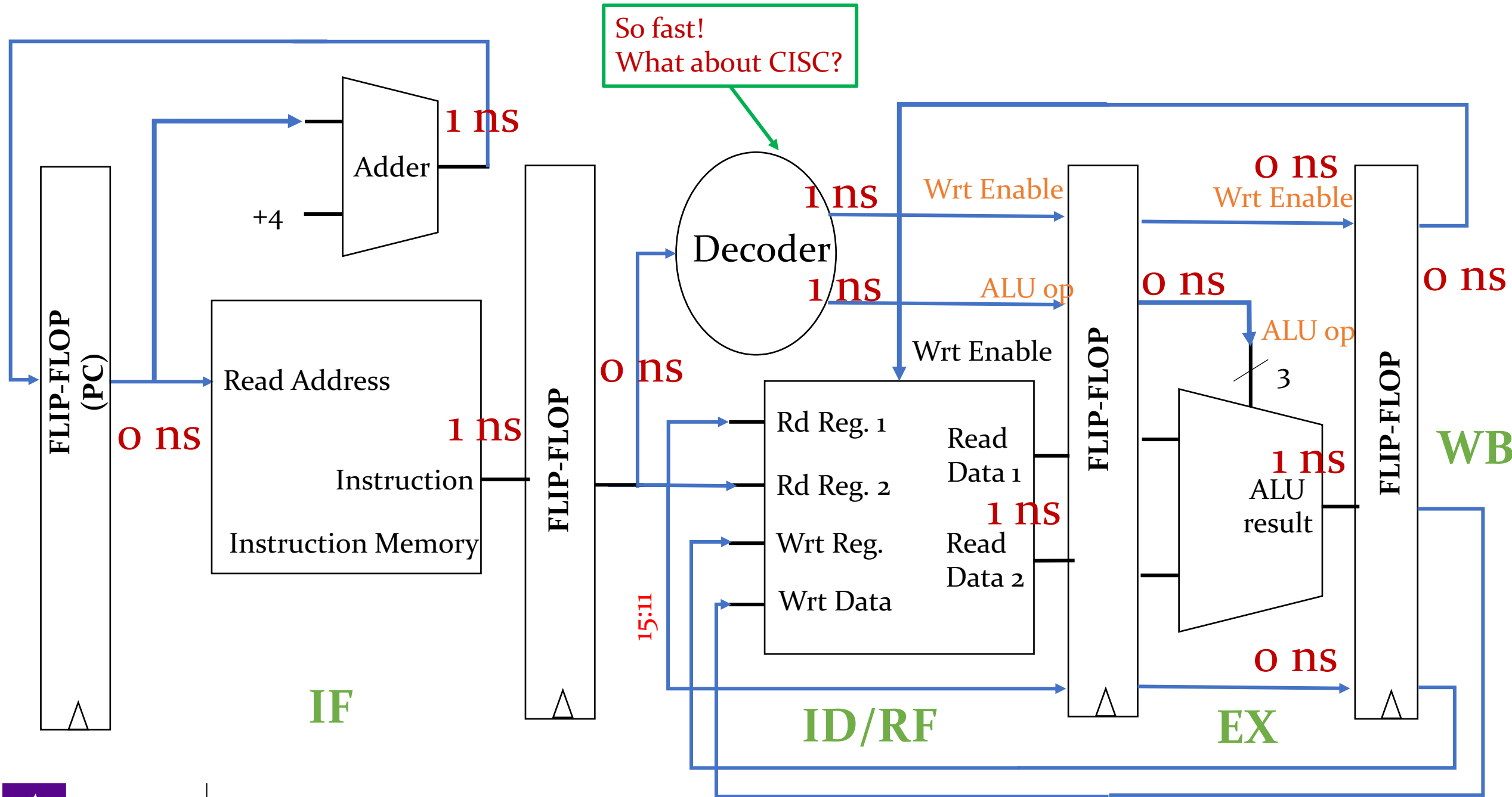
- Instruction fetch (IF)
- Instruction Decode/Register File Read (ID/RF)
- Execute in ALU (EX)
- Write-back Result to RF (WB)



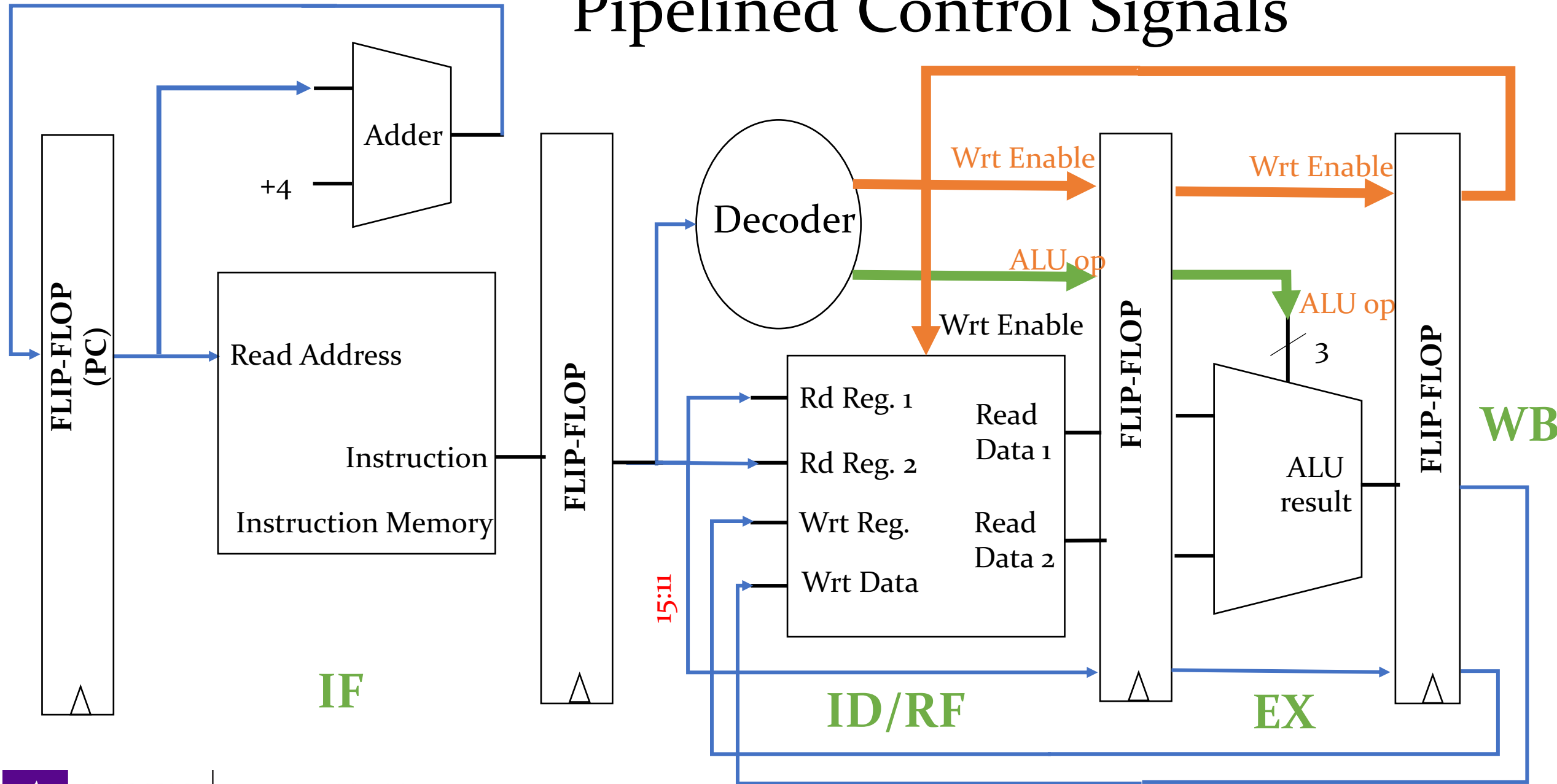


NYU

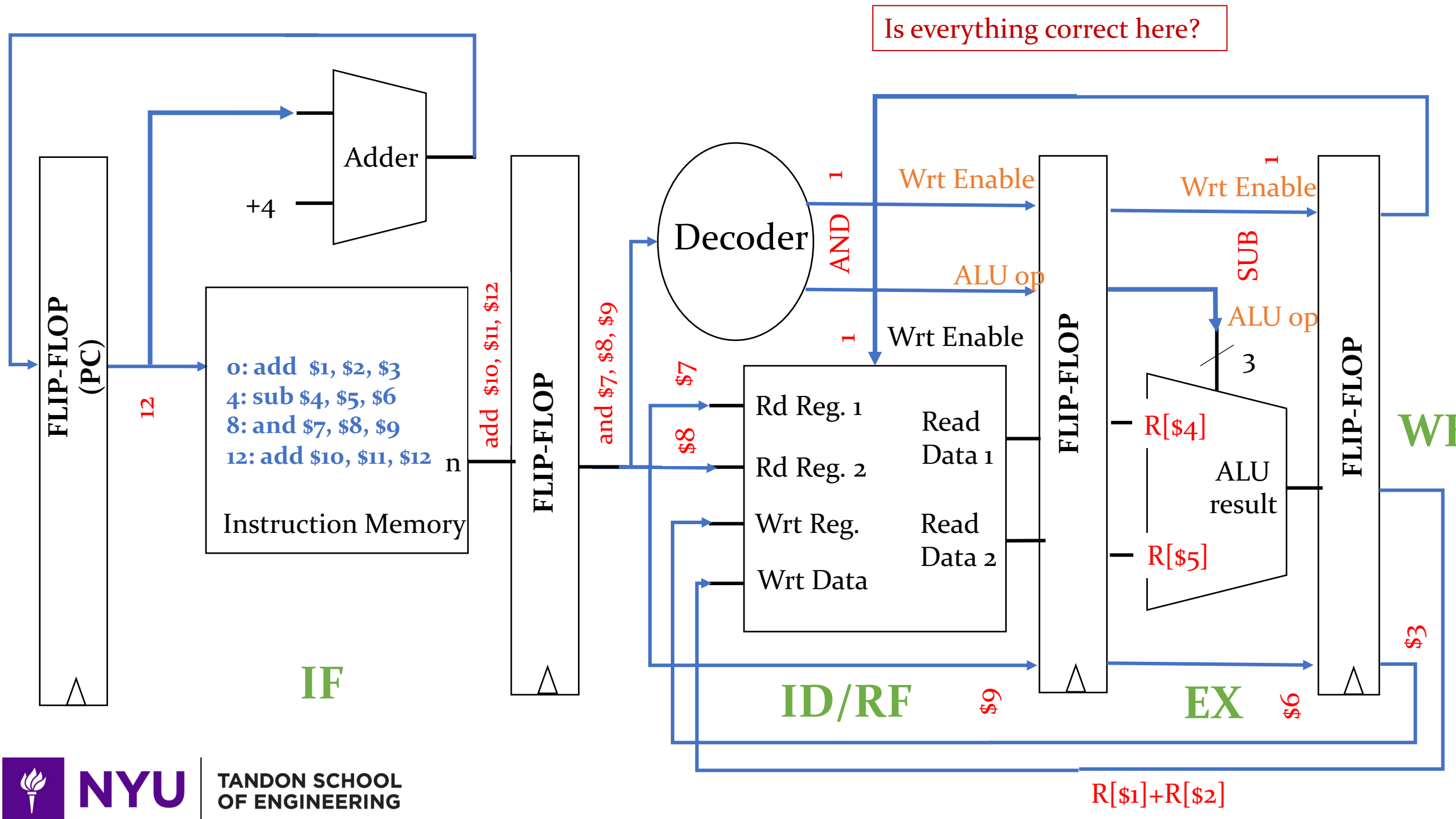
TANDON SCHOOL  
OF ENGINEERING

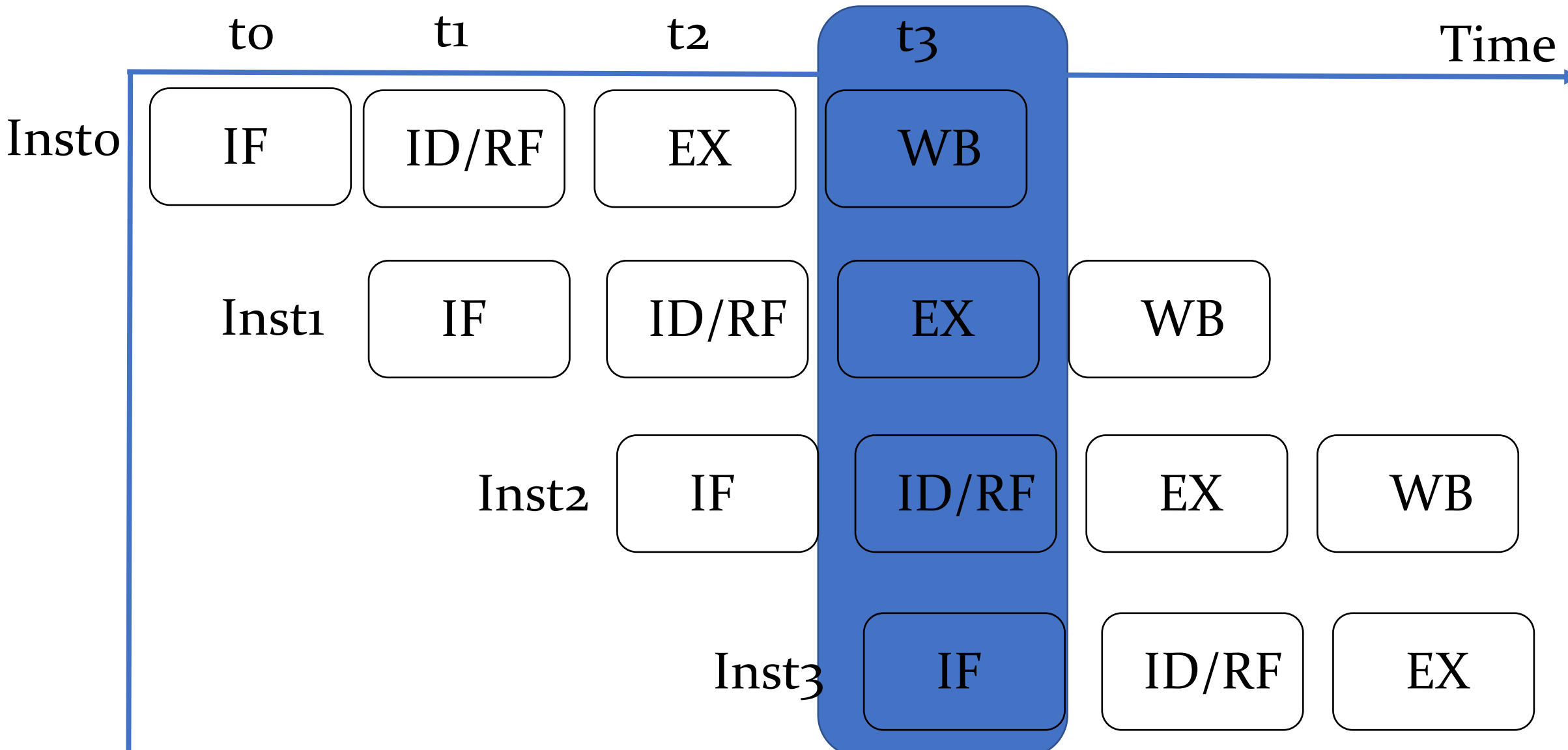


# Pipelined Control Signals









Instructions

This is great! Utilizing all our HW, maximized performance  
But “pipeline” machines are different than laundry..

Why?

# Hazards!

Hazard (formal) := “when the next instruction cannot execute in the following clock cycle”

Hazard (me) := “you need something you don’t have”

There are three of them..

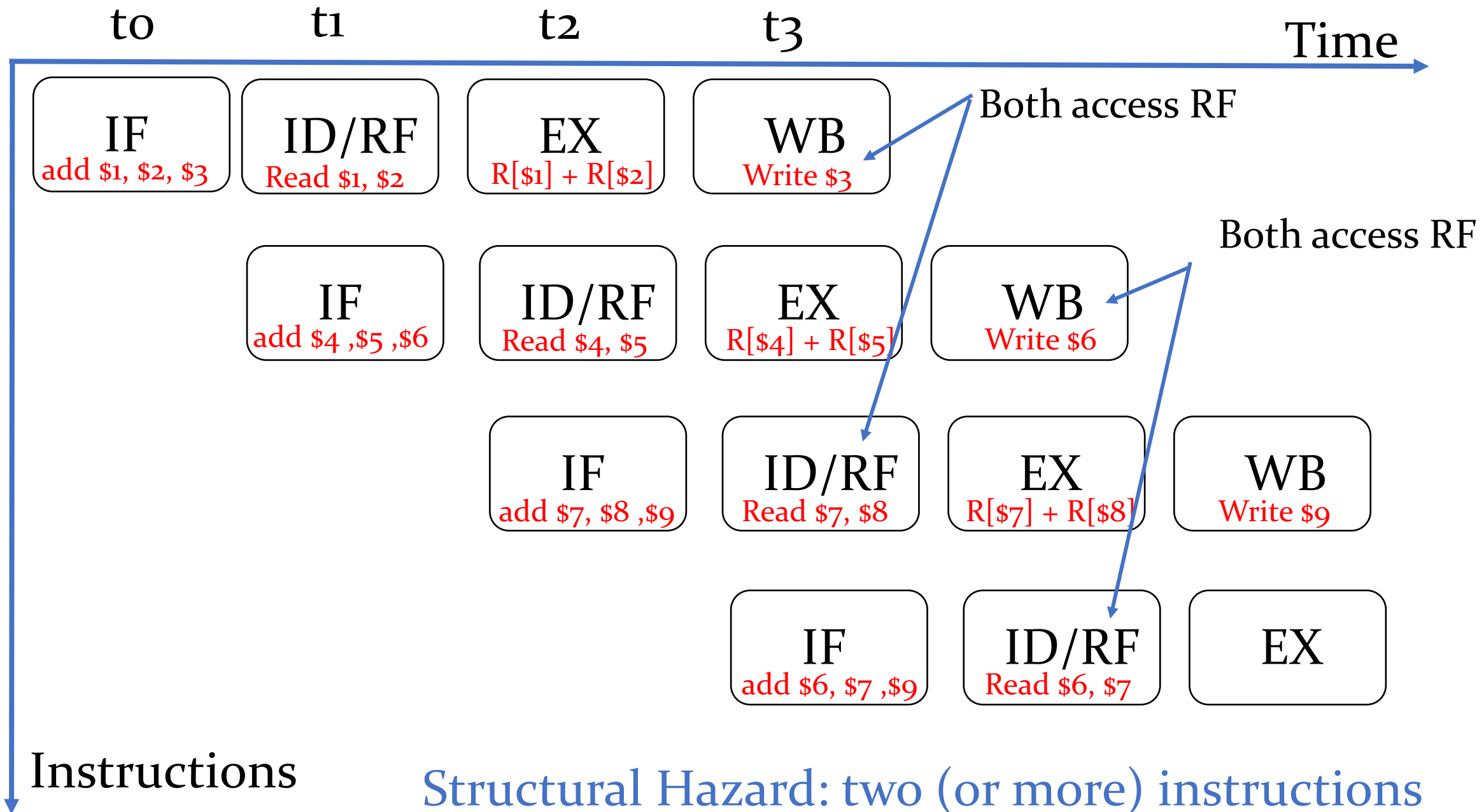
Think about what goes into a program and how they execute

- 1) Structural : hardware resources
- 2) Data : need result from prior computation
- 3) Control : need to know where to go next



NYU

TANDON SCHOOL  
OF ENGINEERING



Structural Hazard: two (or more) instructions accessing the same hardware module



NYU

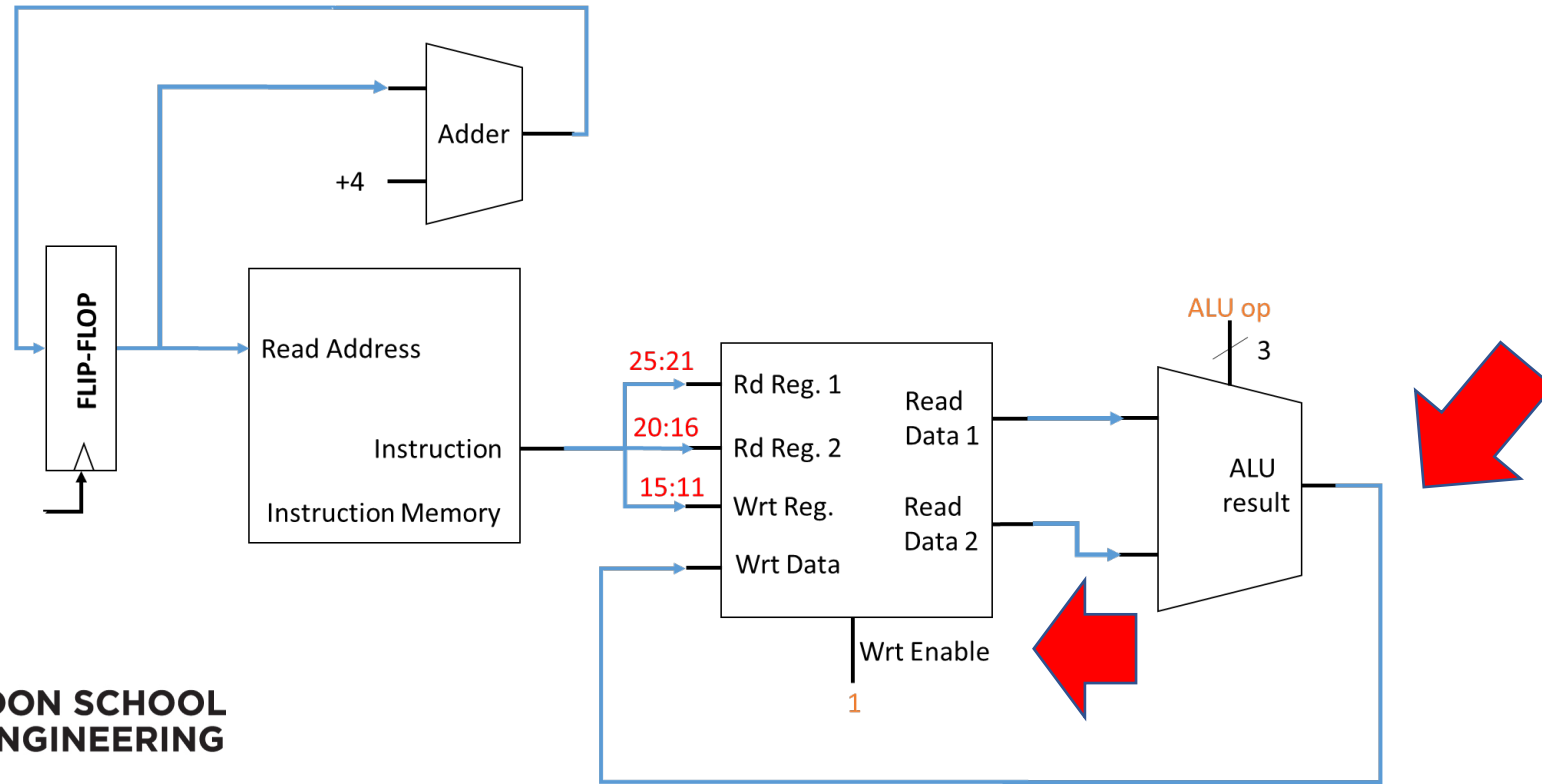
TANDON SCHOOL  
OF ENGINEERING

# Execute/Write-Back Stage Operation

All updates to **architectural state** should occur in first half of period

- PC (already done!)
- **Register file**
- Data memory (not relevant for R-type instructions)

If we achieve this, we can write and read in same cycle!



NYU

TANDON SCHOOL  
OF ENGINEERING

Test understanding of FFs and hazards with example..



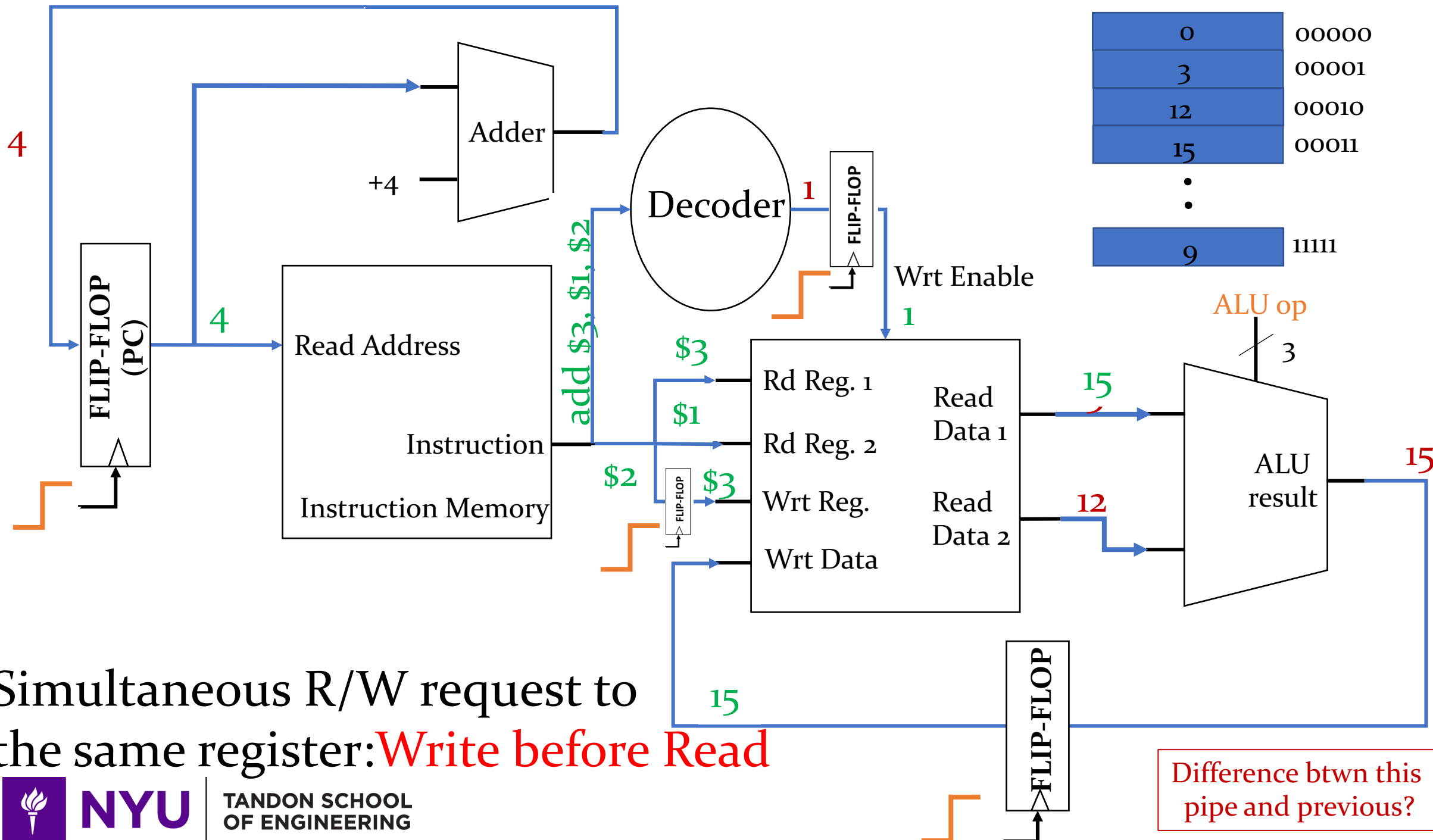
**NYU**

**TANDON SCHOOL  
OF ENGINEERING**









Simultaneous R/W request to the same register: **Write before Read**



NYU

TANDON SCHOOL OF ENGINEERING

# Hazards!

Hazard (formal) := “when the next instruction cannot execute in the following clock cycle”

Hazard (me) := “you need something you don’t have”

There are three of them..

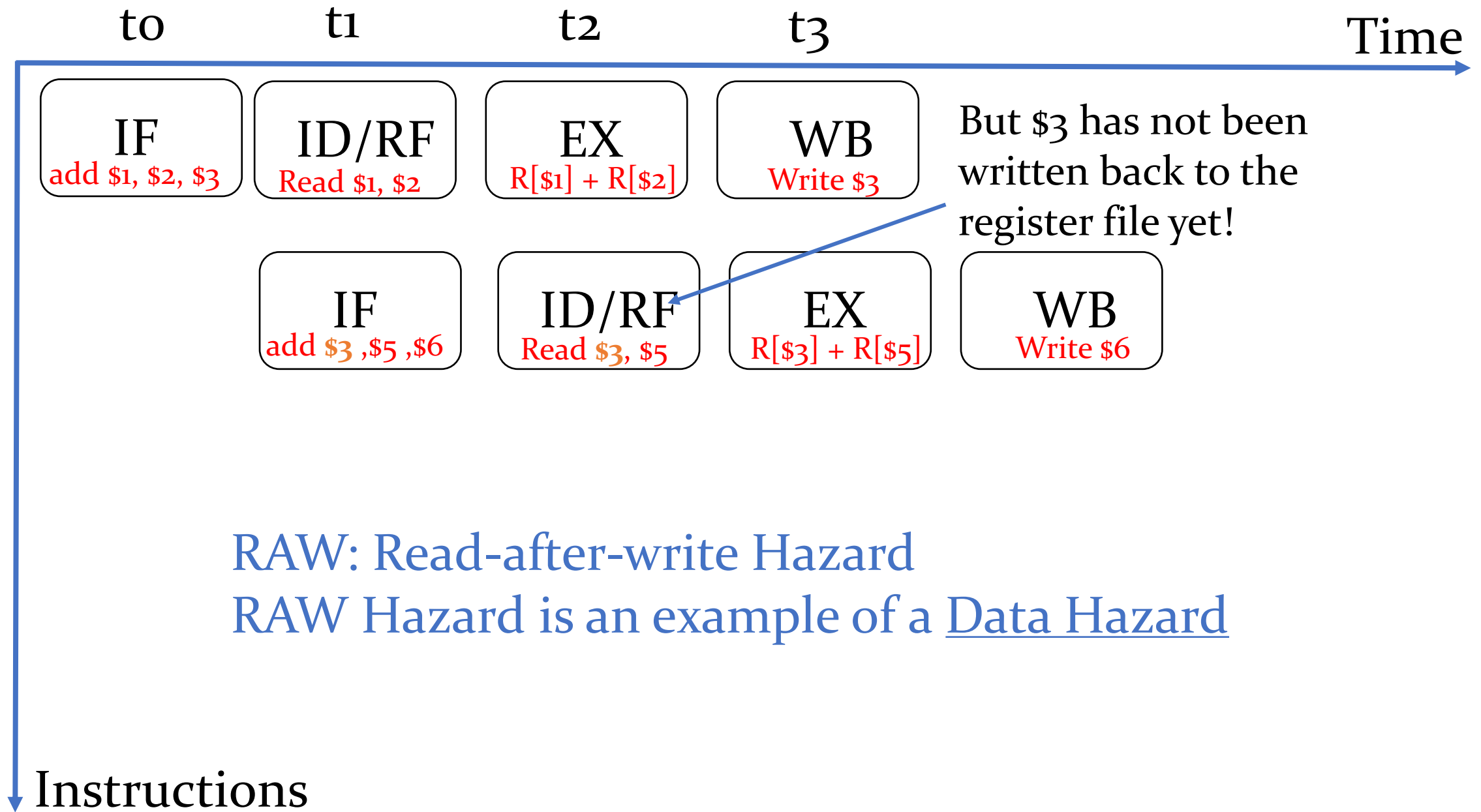
Think about what goes into a program and how they execute

- 1) Structural : hardware resources : **Solve with write before read**
- 2) Data : need result from prior computation
- 3) Control : need to know where to go next



**NYU**

TANDON SCHOOL  
OF ENGINEERING



RAW: Read-after-write Hazard

RAW Hazard is an example of a Data Hazard



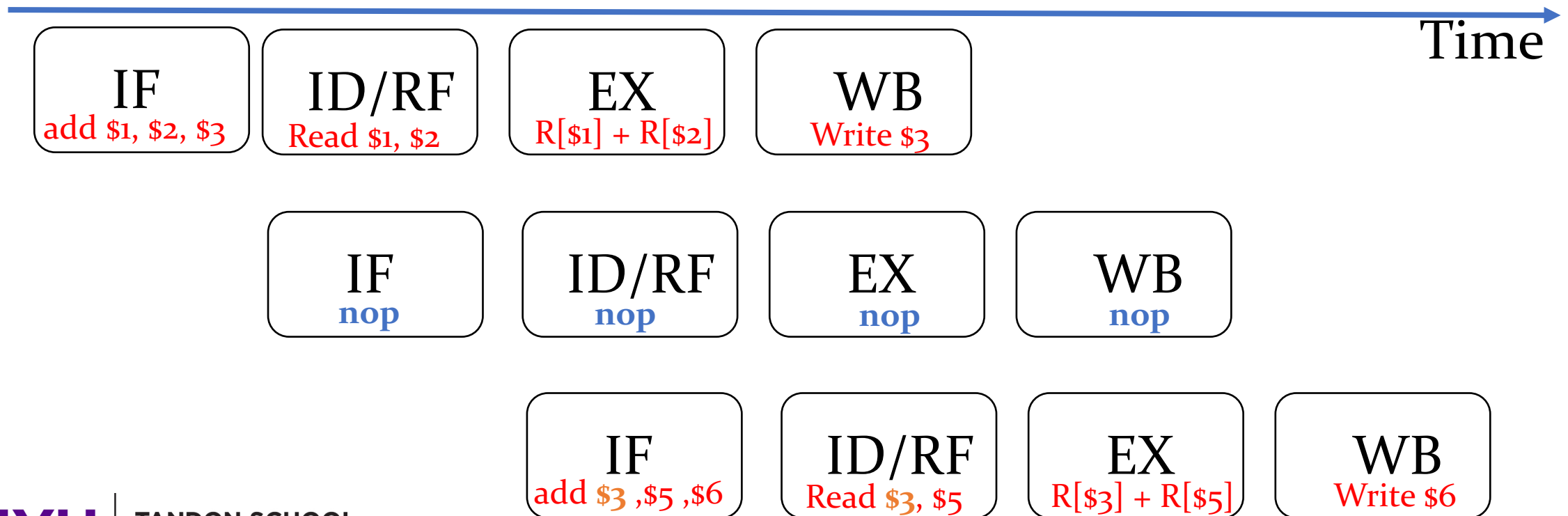
NYU

TANDON SCHOOL  
OF ENGINEERING

# Handling RAW Hazards

## Solution 1

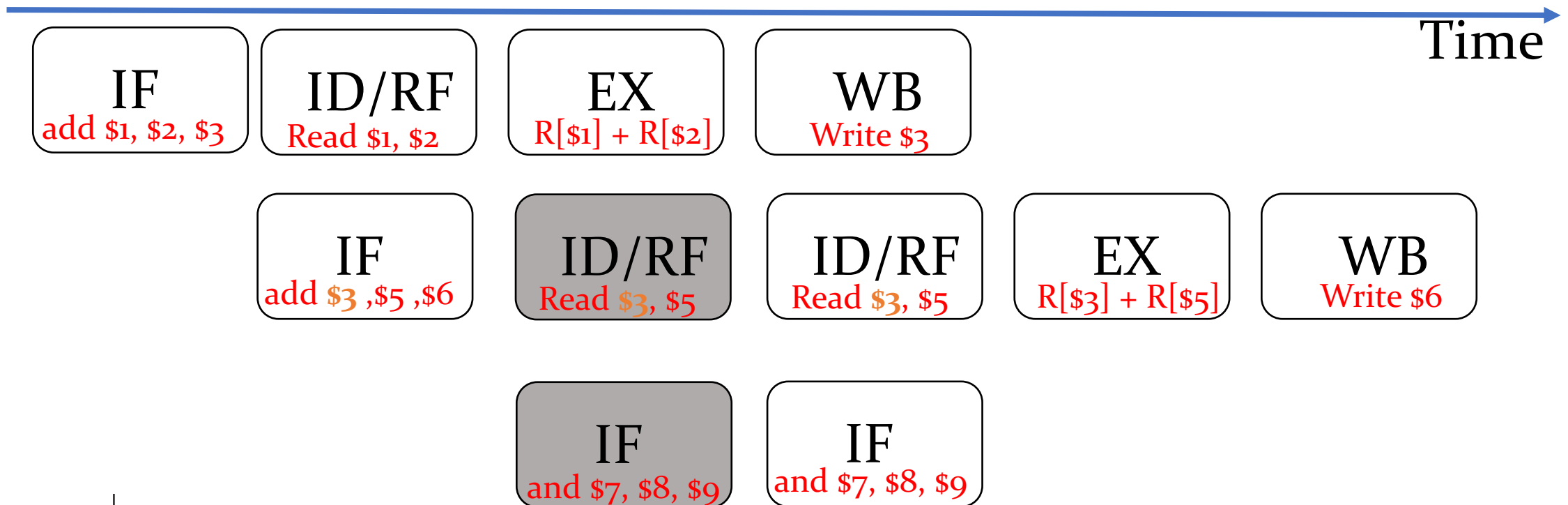
- Detect and avoid in software: compiler inserts **nops** (example: sll \$0, \$0, 0)
- Pros? Cons?



# Handling RAW Hazards

## Solution 2

- Detect and avoid in hardware: **stall** the pipeline
- A stall does not update the pipeline register's values, keeps them the same
- **Pros? Cons?**



# Handling RAW Hazards

## Solution 2

- Detect and avoid in hardware: **stall** the pipeline
- A stall does not update the pipeline register's values, keeps them the same
- **Pros? Cons?**



# Implementing Stalls: Details

Detect Stall in IF/ID stage by checking

- Destination register of previous inst. == source register of current inst

When Stall is detected

- PC flip-flop is “stalled”
  - Pipeline register outputs the same value as in previous clock cycle
  - In other words: no update, ignore the positive edge of the clock
- ID/RF flip-flop is stalled
- Nop inserted in EX flip-flop
  - Else we'll keep detecting the hazard

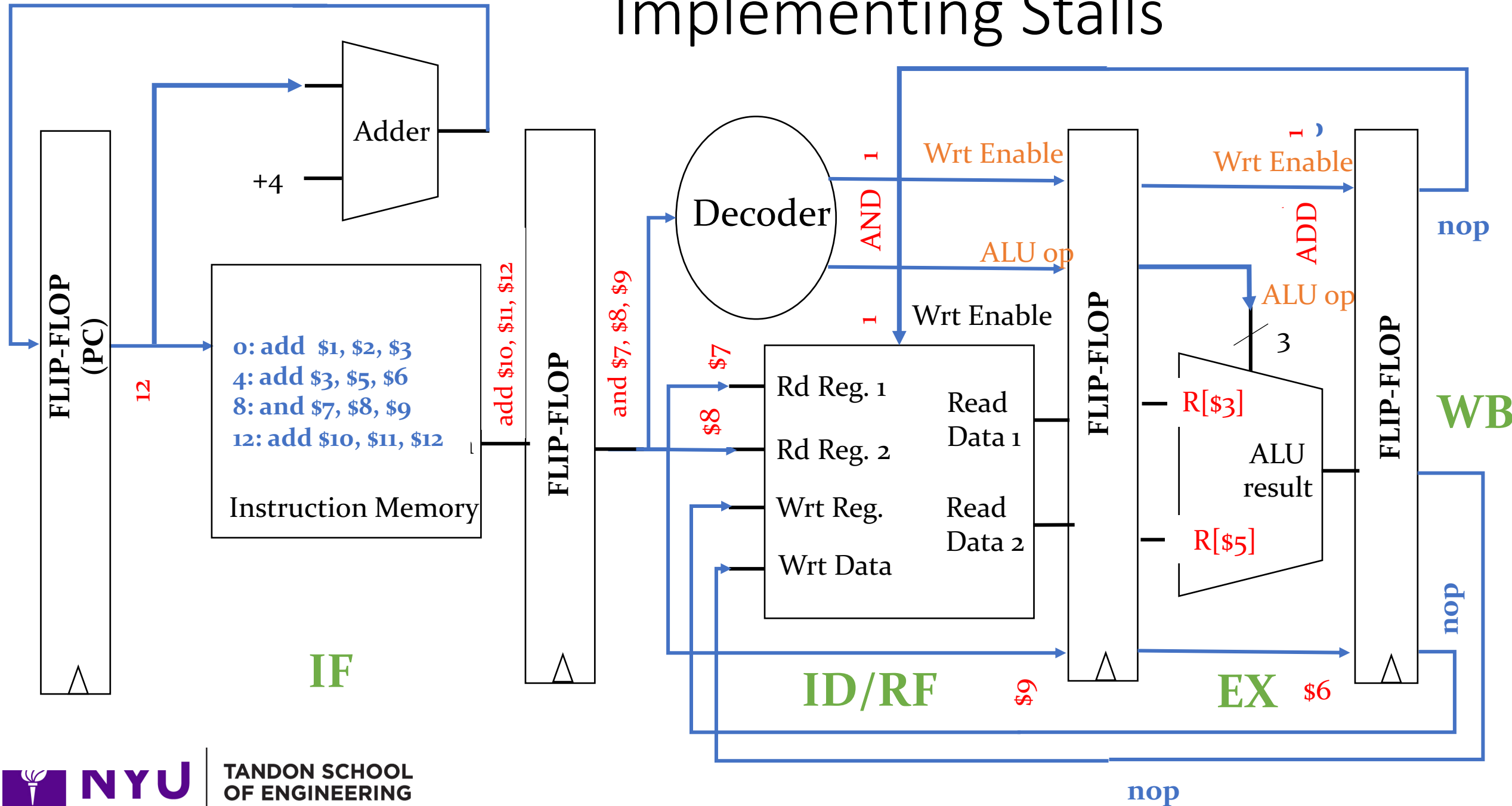
Stall clears in next clock cycle



NYU

TANDON SCHOOL  
OF ENGINEERING

# Implementing Stalls

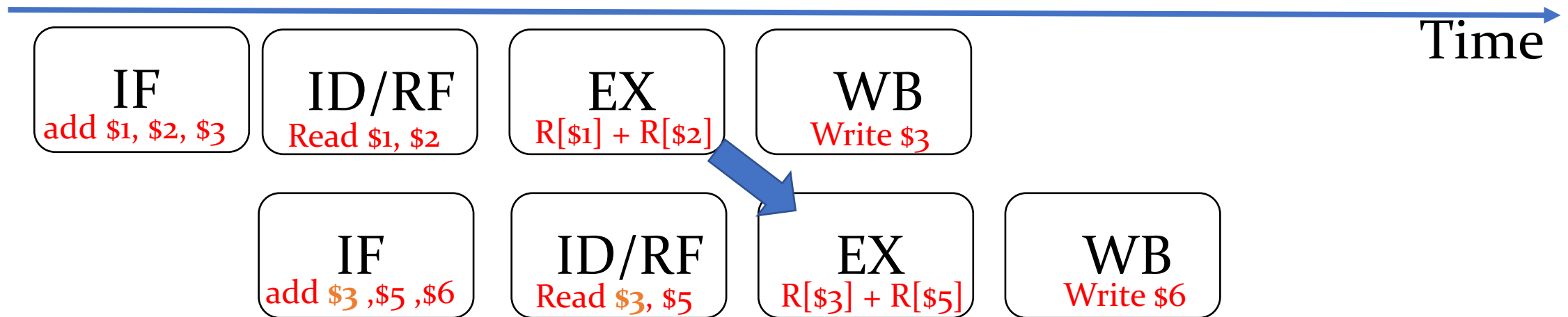




# Handling RAW Hazards

## Solution 3

- Forward data from the EX stage
- Pros? Cons?



How would we implement this in our pipeline?

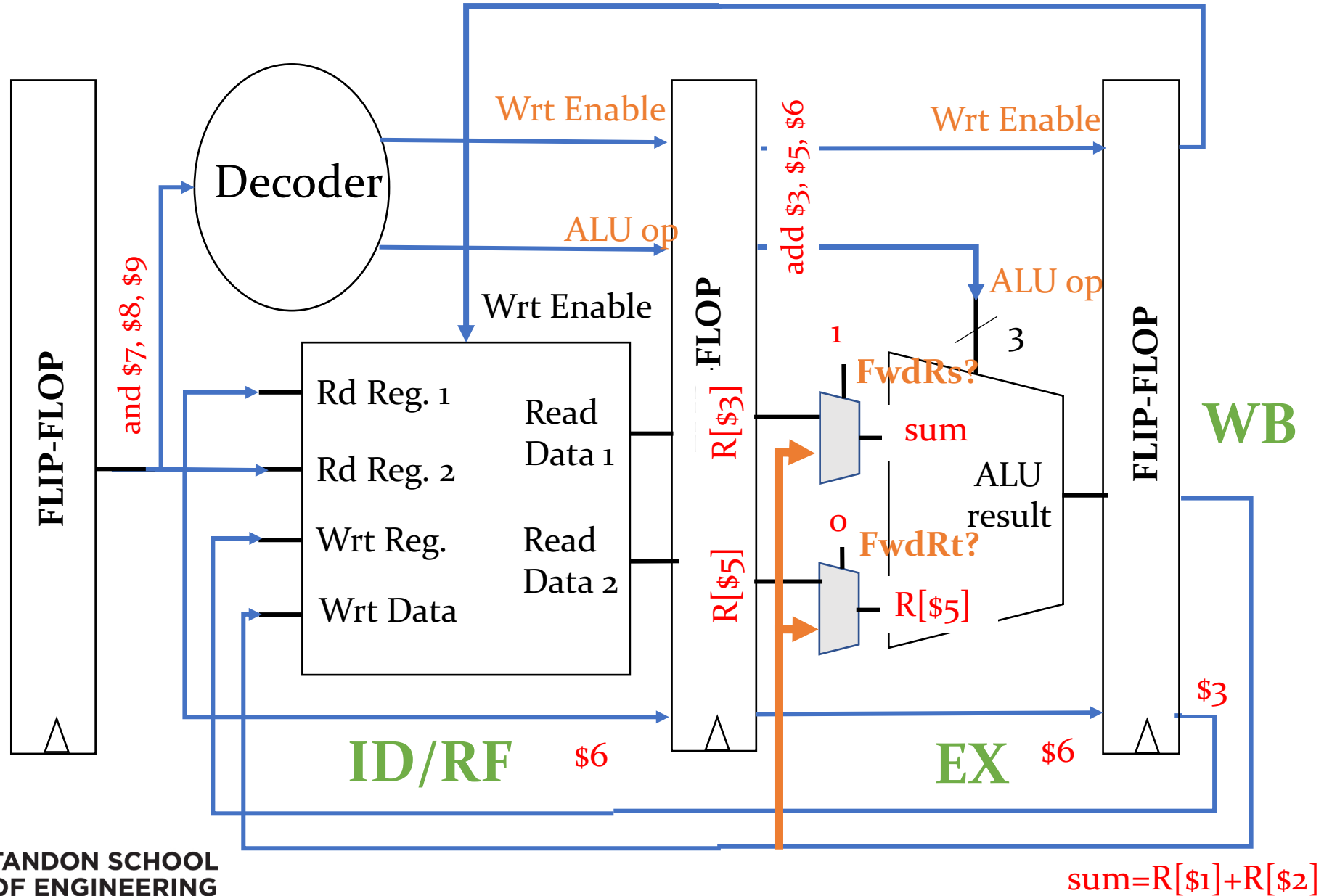


NYU

TANDON SCHOOL  
OF ENGINEERING

# Implementing Forwarding

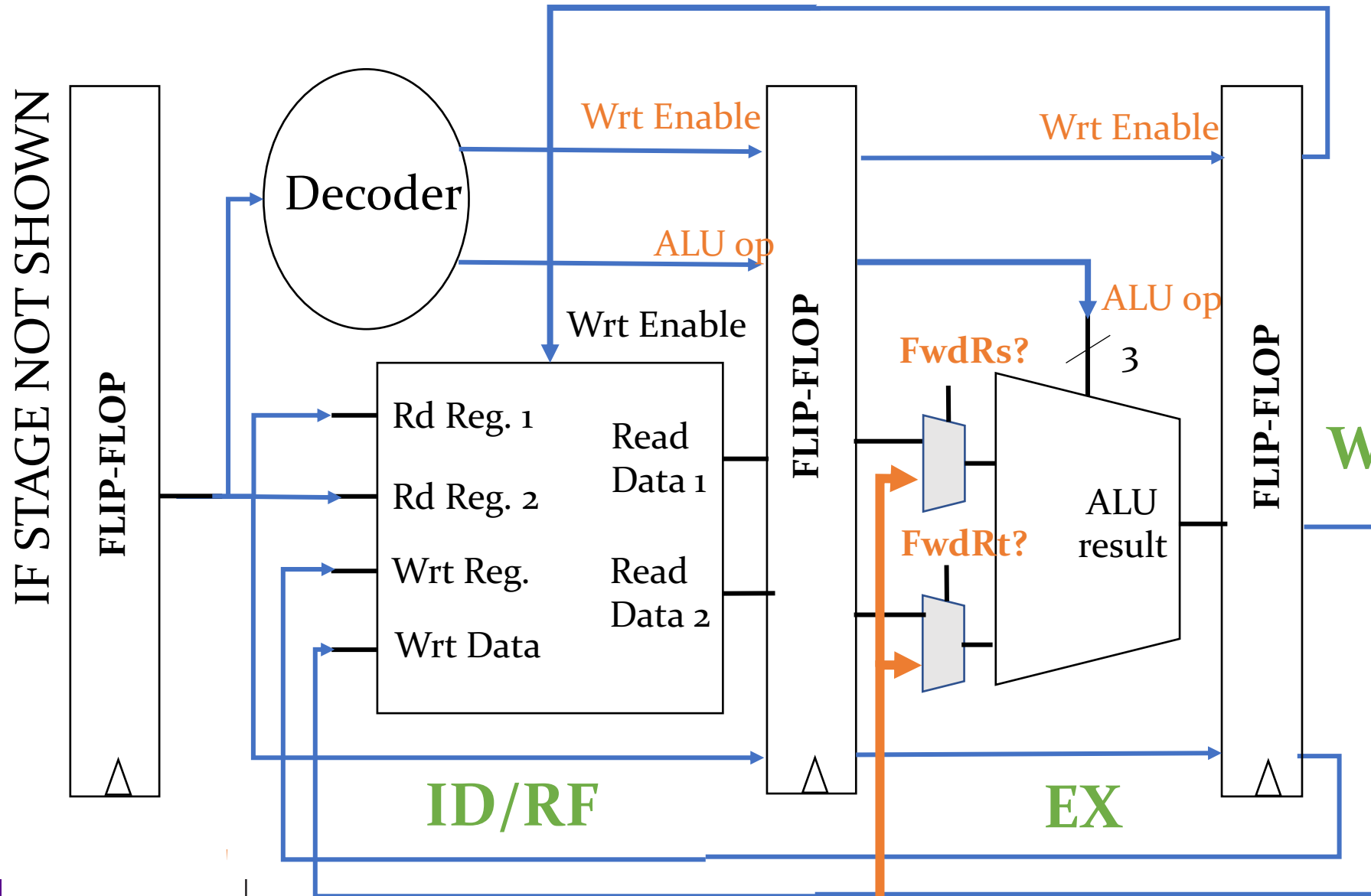
IF STAGE NOT SHOWN



NYU

TANDON SCHOOL  
OF ENGINEERING

# Forwarding Control

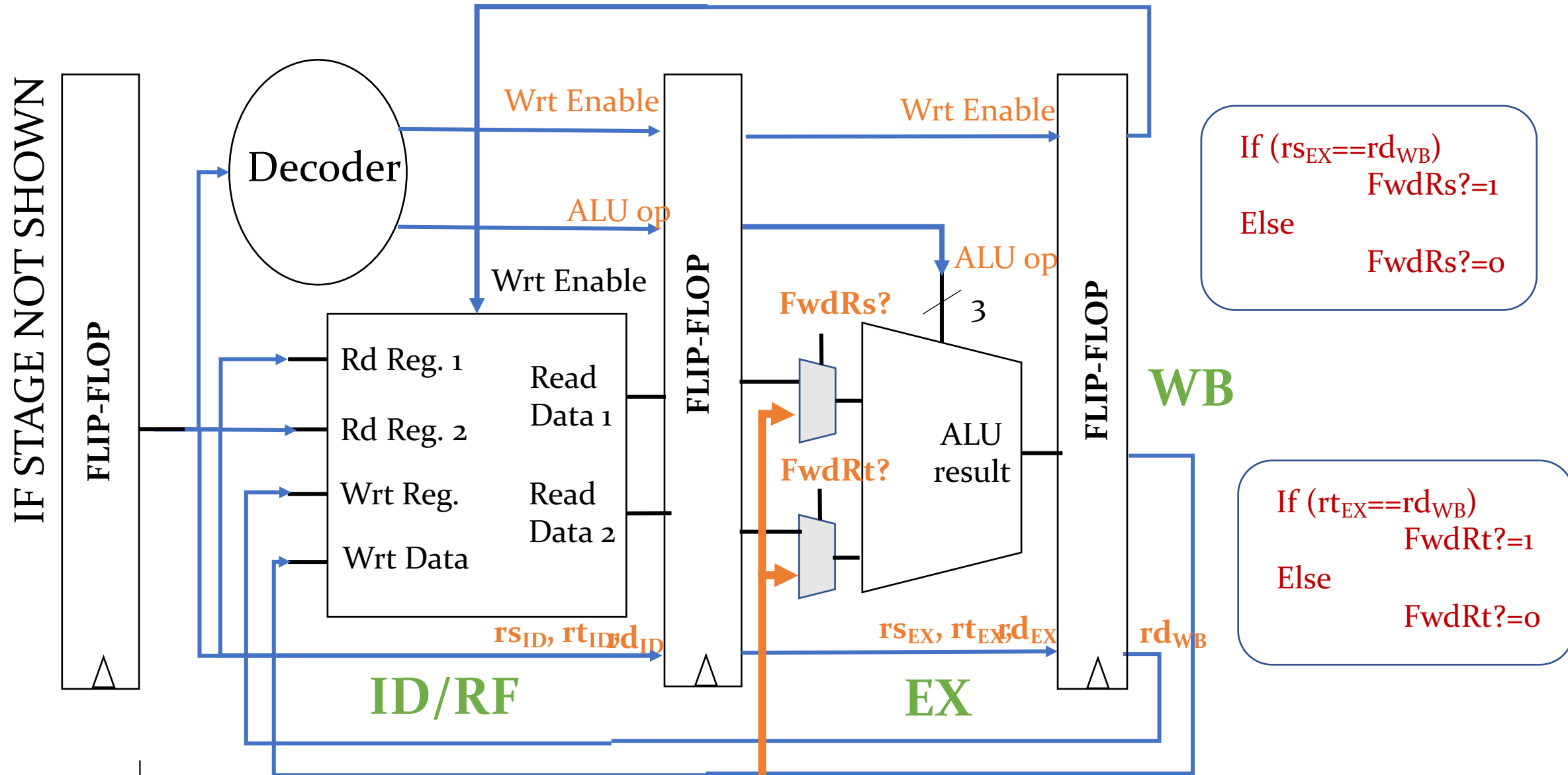


FwdRs? = 1 if destination register (rd) for inst in WB stage is same as source register (rs) for instruction in EX stage

**WB** But we don't know rs and rt in EX stage!

FwdRt? = 1 if destination register (rd) for inst in WB stage is same as source register (rt) for instruction in EX stage

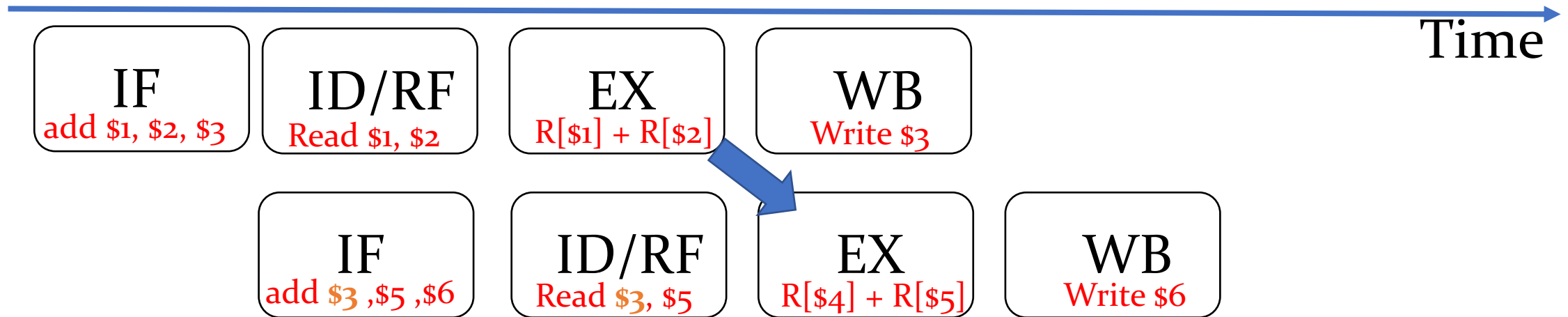
# Forwarding Control



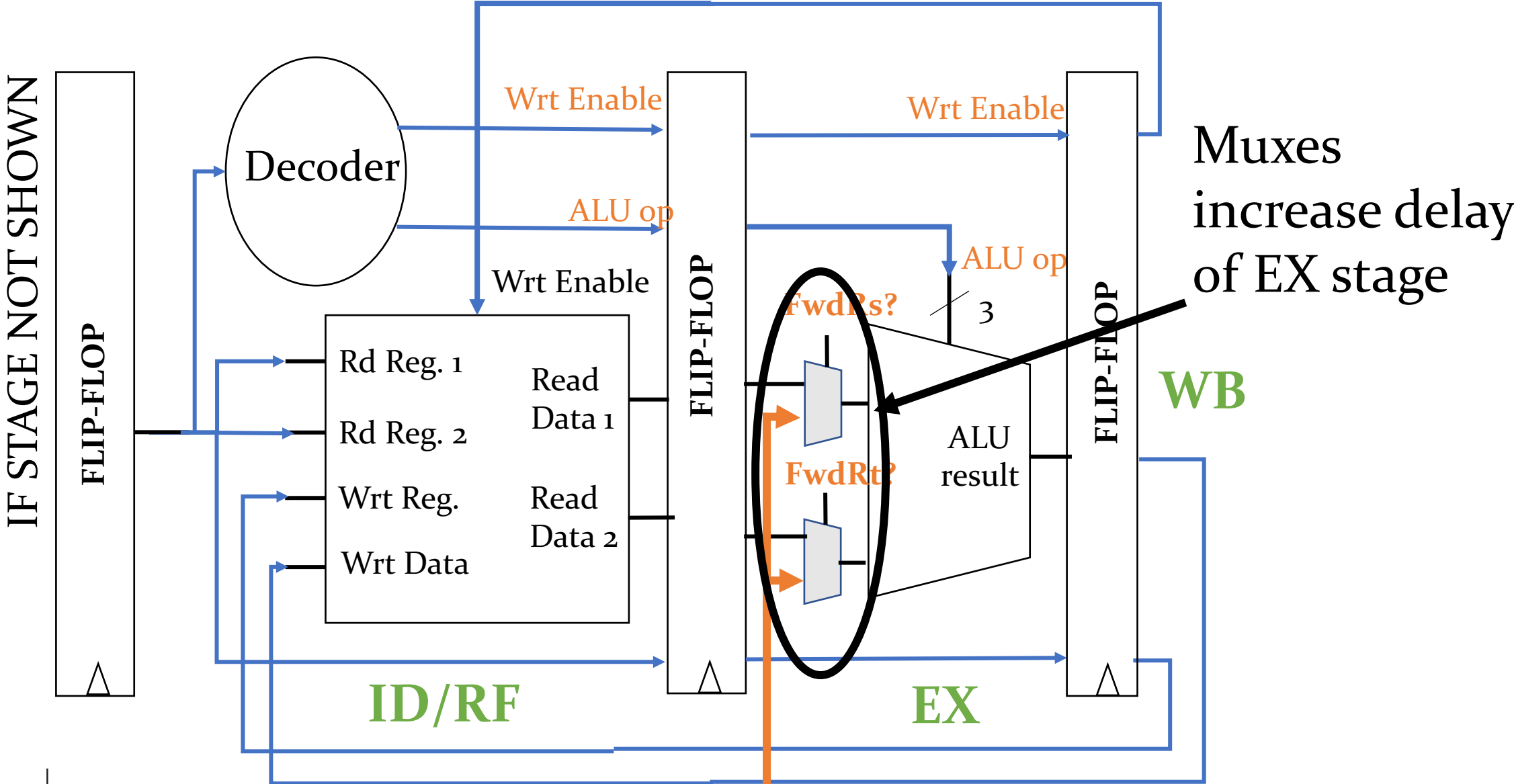
# Handling RAW Hazards

## Solution 3

- Forward data from the EX stage
- Impact on IPC (Instructions/cycle)?
- Impact on clock period?



# Impact on Cycle Time



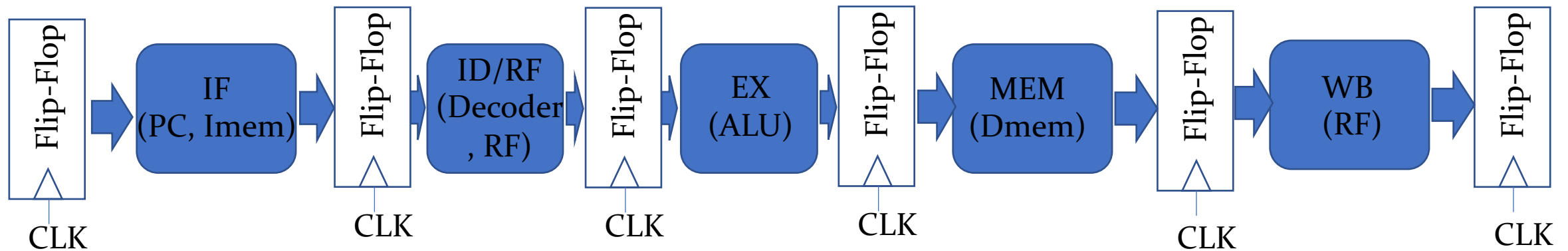
NYU

TANDON SCHOOL  
OF ENGINEERING

# 5 Stage Pipeline (with support for load/store)

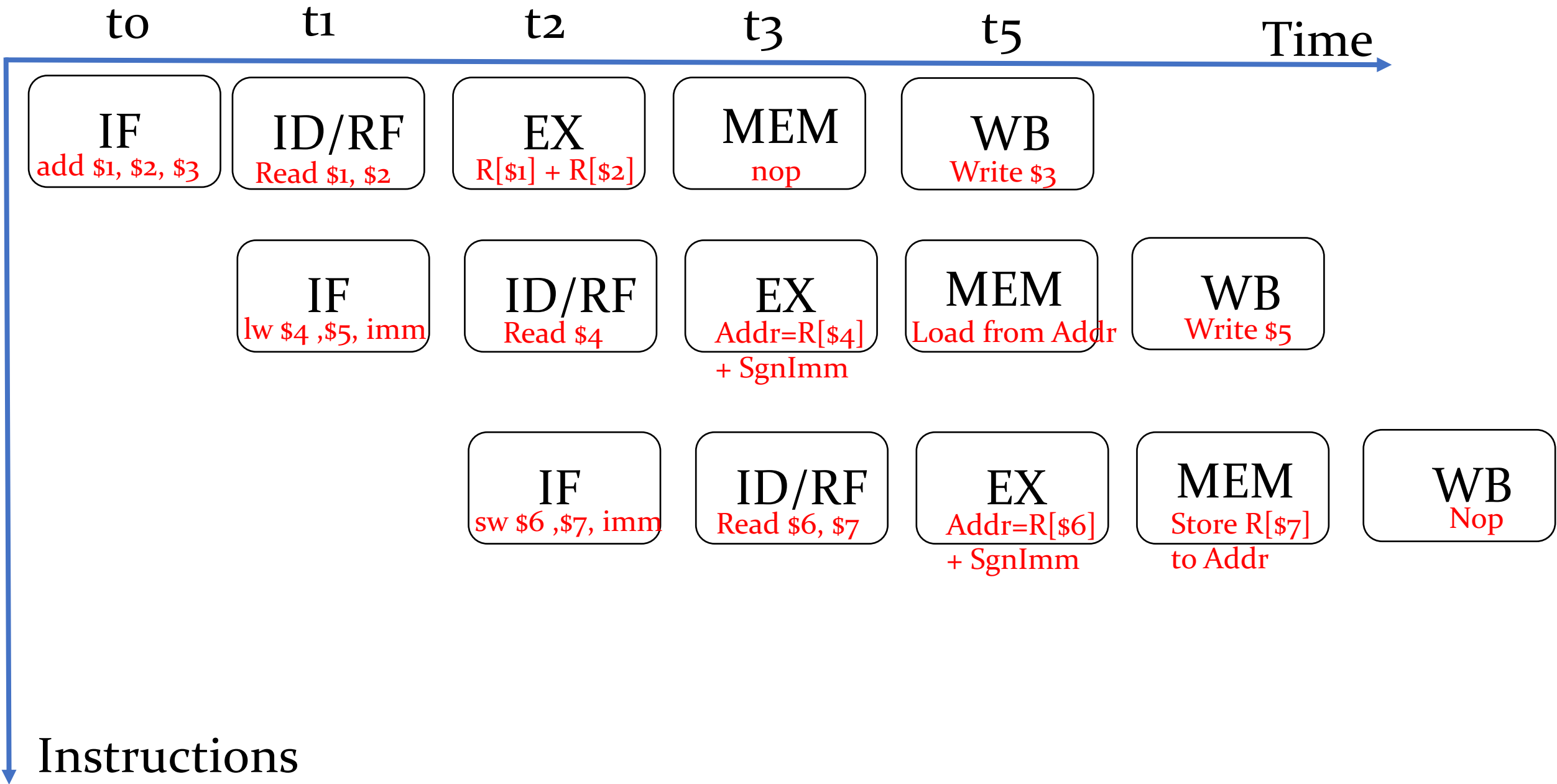
Basic steps in executing an R-type instruction

- Instruction fetch (IF)
- Instruction Decode/Register File Read (ID/RF)
- Execute in ALU (EX)
- **Load result from or store result to data memory (MEM)**
- Write-back Result to RF (WB)



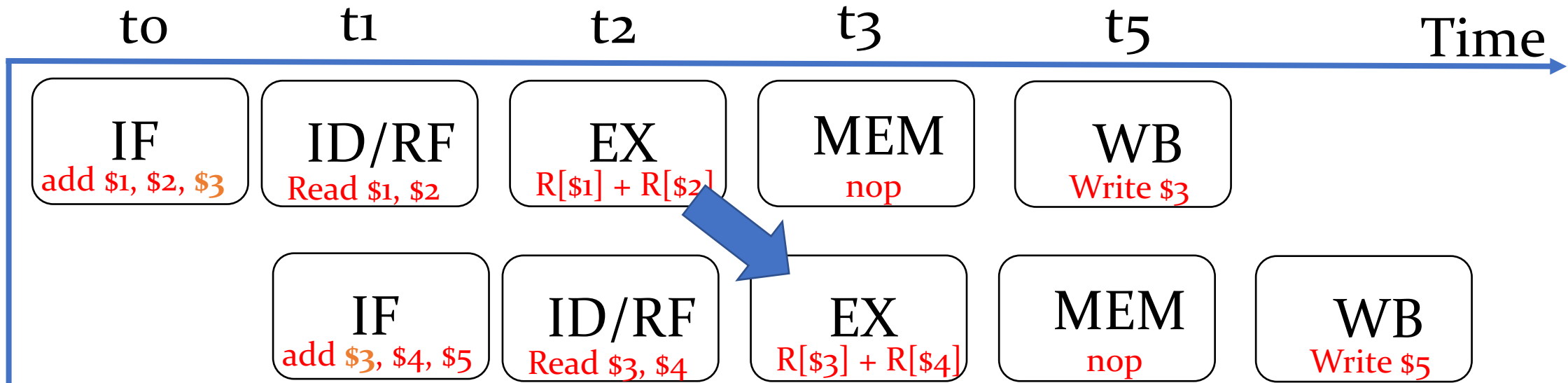
NYU

TANDON SCHOOL  
OF ENGINEERING





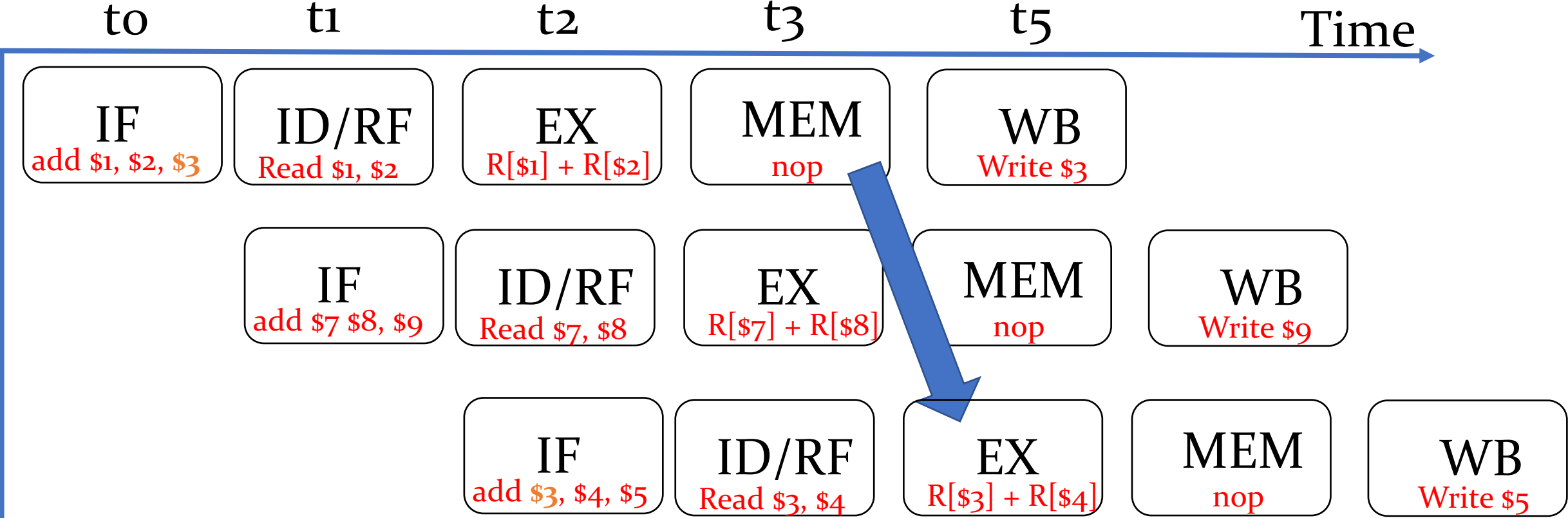
# 5-stage RAW Hazards: Add-Add Dependency



**EX-EX Forwarding (same as 4-stage pipeline)**

Instructions

# RAW Hazards: Add-Add Dependency



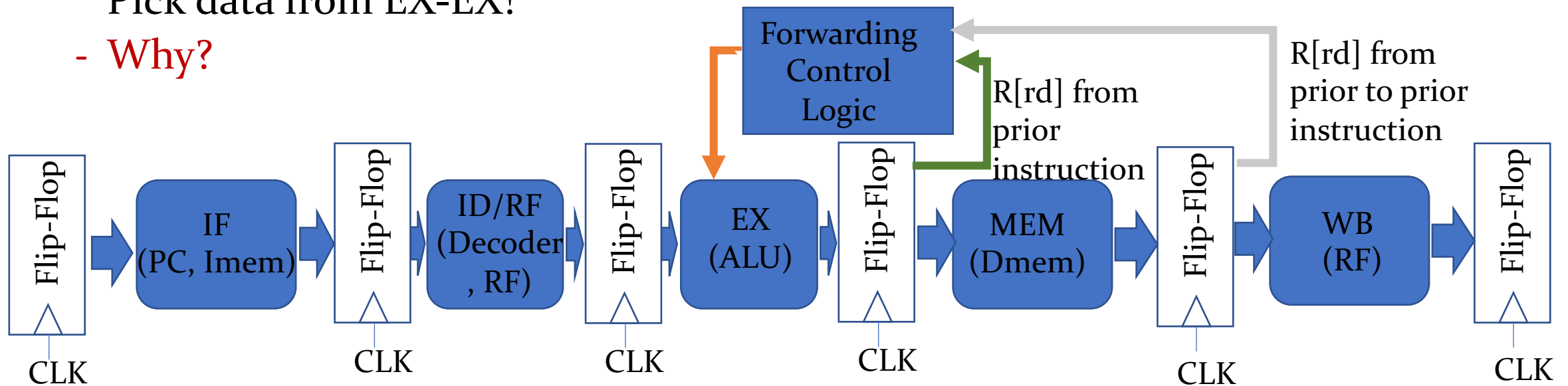
Instructions

MEM-EX Forwarding

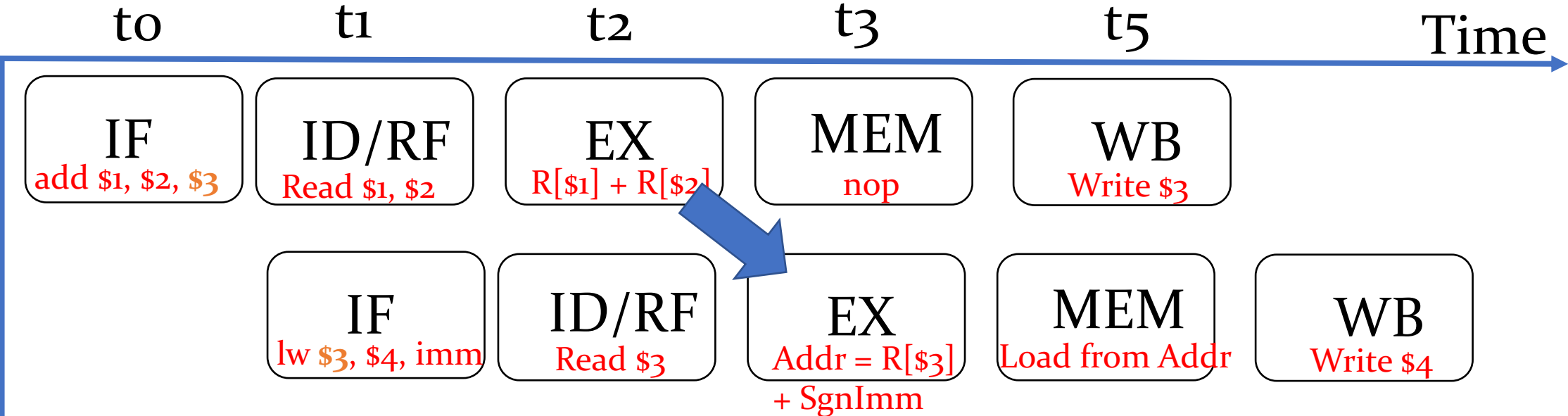
# 5 Stage Pipeline (with support for load/store)

## Two forwarding paths

- EX – green
- MEM - gray
- Same destination register in both forwarding paths?  
Pick data from EX-EX!
- Why?



# RAW Hazards: Add-Load Dependency

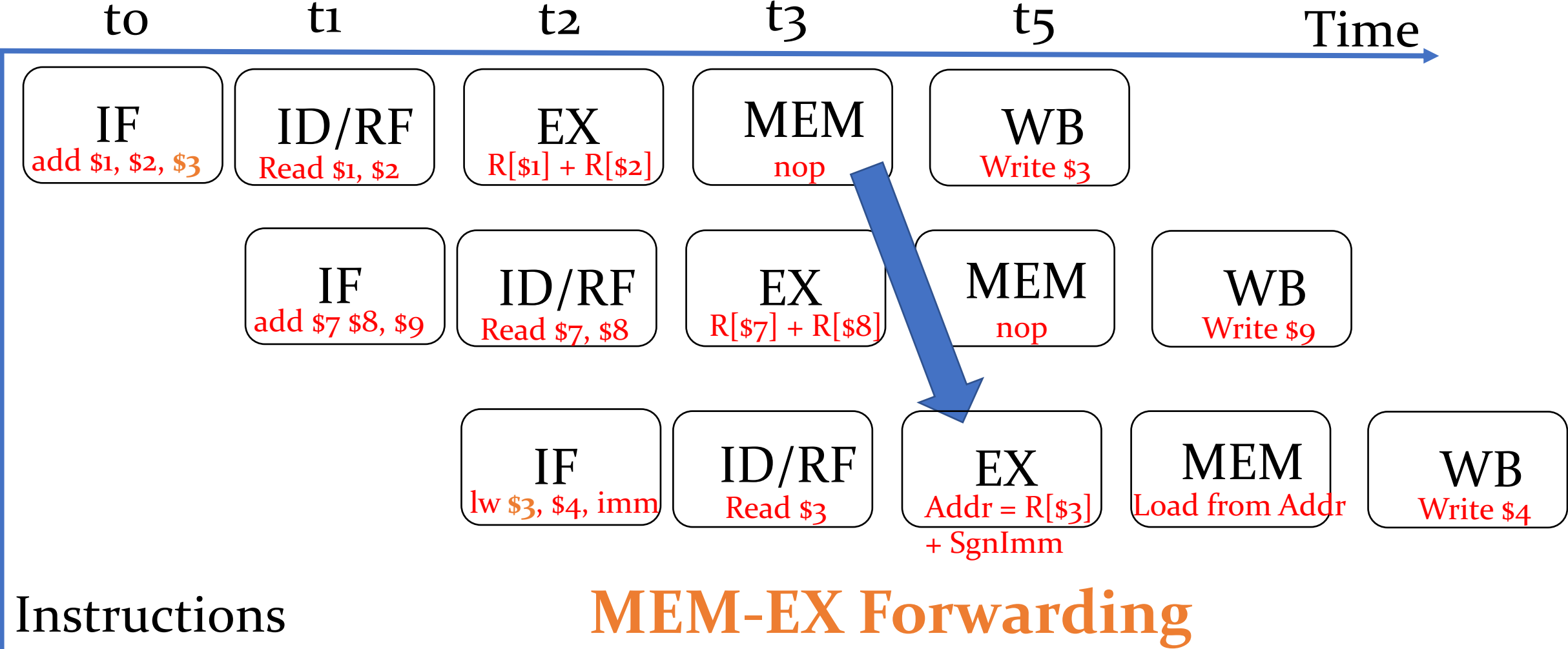


## EX-EX Forwarding

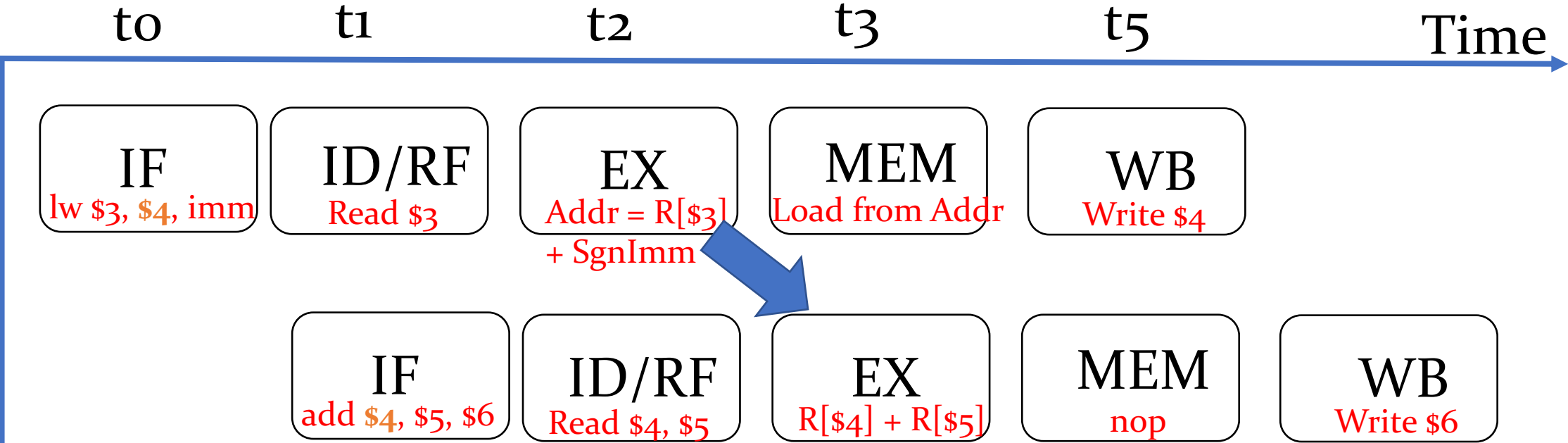
Why?

Because we calculate address with the ALU, it's the same circuit so we can use existing path

# RAW Hazards: Add-Load Dependency



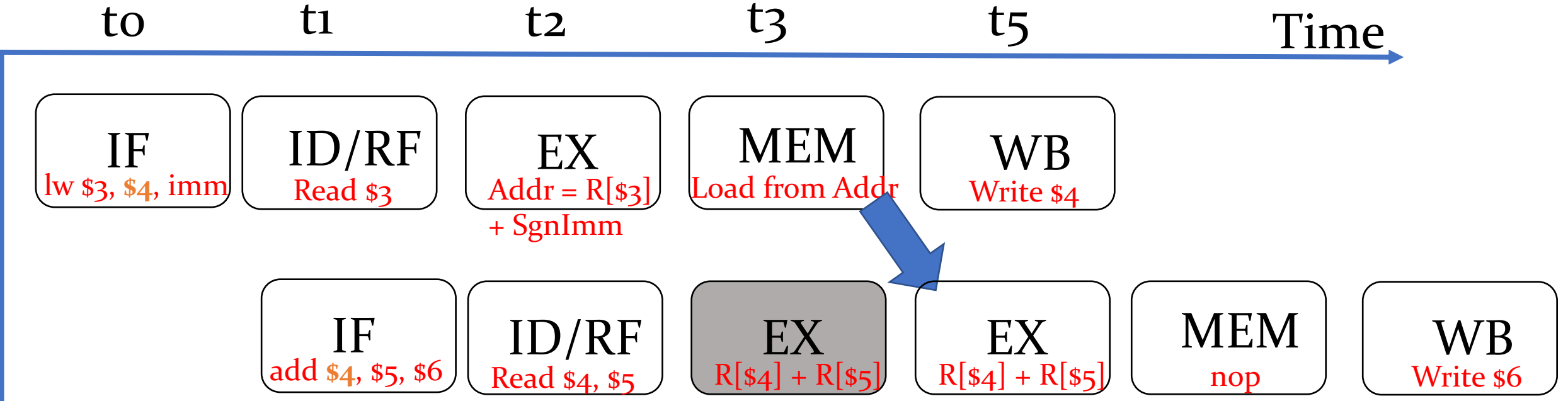
# RAW Hazards: Load-Add Dependency



Will this work?

Solution: Stall the pipeline

# RAW Hazards: Load-Add Dependency



**Solution: Stall the pipeline**

# RAW Hazards: Load-Add Dependency

SW Trick!  
(They help us for once..)  
Compilers save the day

