

CS 110
Computer Architecture
Review Midterm II

<http://shtech.org/courses/ca/>

School of Information Science and Technology SIST

ShanghaiTech University

Slides based on UC Berkley's CS61C

Midterm II

- Date: Tuesday, May 9
- Time: 10:15 - 12:15 (similar to last time)
- Venue: Teaching Center 201 + 203
- One empty seat between students
- Closed book:
 - You can bring one A4 page with notes (both sides; English preferred; Chinese is OK): Write your Chinese and **Pinyin** name on the top!
 - This time **HANDWRITTEN** only!
 - You will be provided with the MIPS "green sheet"
 - No other material allowed!

Midterm II

- Switch cell phones **off!** (not silent mode – off!)
 - Put them in your bags.
- Bags to the front. Nothing except paper, pen, 1 drink, 1 snack on the table!
- No other electronic devices are allowed!
 - No ear plugs, music, ...
- Anybody touching any electronic device will FAIL the course!
- Anybody found cheating (copy your neighbors answers, additional material, ...) will FAIL the course!



Midterm II

- Ask questions today!
- And in next weeks Q&A session
 - Suggest topics for review in piazza!
- This review session does not/ can not cover all possible topics!
- Topics: SDS till Data-Level-Parallelism

Synchronous Digital Systems

Hardware of a processor, such as the MIPS, is an example of a Synchronous Digital System

Synchronous:

- All operations coordinated by a central clock
 - “Heartbeat” of the system!

Digital:

- Represent all values by discrete values
- Two binary digits: 1 and 0
- Electrical signals are treated as 1's and 0's
 - 1 and 0 are complements of each other
- High /low voltage for true / false, 1 / 0

CMOS Circuit Rules

- Don't pass weak values => Use Complementary Pairs
 - N-type transistors pass weak 1's ($V_{dd} - V_{th}$)
 - N-type transistors pass strong 0's (ground)
 - Use N-type transistors only to pass 0's (N for negative)
 - Converse for P-type transistors: Pass weak 0s, strong 1s
 - Pass weak 0's (V_{th}), strong 1's (V_{dd})
 - Use P-type transistors only to pass 1's (P for positive)
 - Use pairs of N-type and P-type to get strong values
- Never leave a wire undriven
 - Make sure there's always a path to V_{dd} or GND
- Never create a path from V_{dd} to GND (ground)
 - This would short-circuit the power supply!

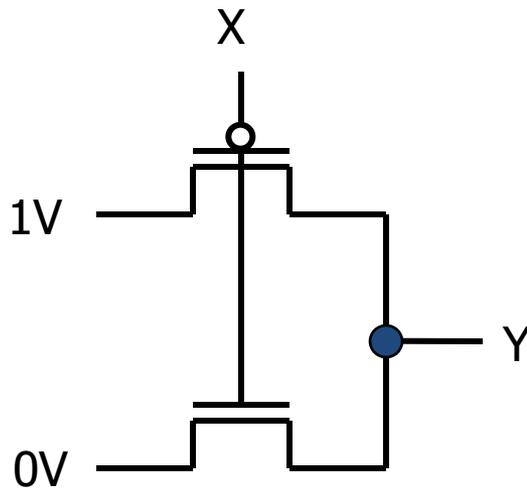
CMOS Networks

p-channel transistor

on when voltage at Gate is low

off when:

voltage(Gate) > voltage (Threshold)



n-channel transistor

off when voltage at Gate is low

on when:

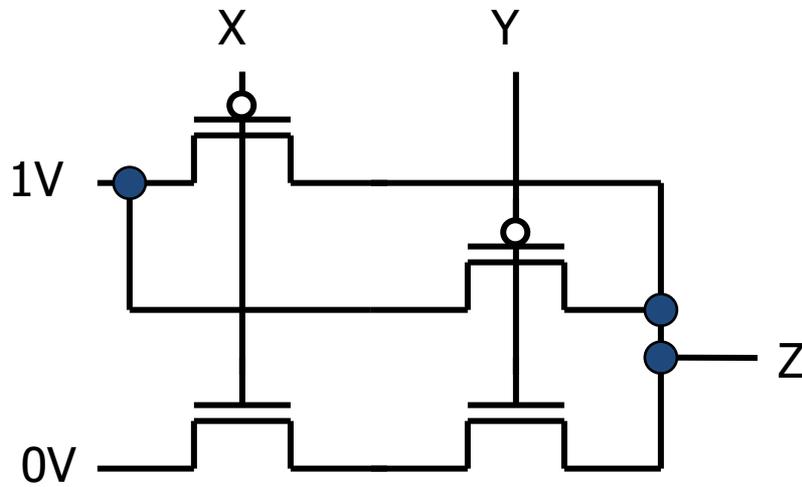
voltage(Gate) > voltage (Threshold)

what is the relationship between x and y?

x	y
0 Volt (GND)	1 Volt (Vdd)
1 Volt (Vdd)	0 Volt (GND)

Called an *inverter* or *not gate*

Two-Input Networks



what is the relationship between x, y and z?

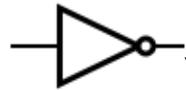
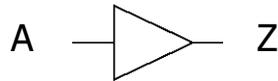
x	y	z
0 Volt	0 Volt	1 Volt
0 Volt	1 Volt	1 Volt
1 Volt	0 Volt	1 Volt
1 Volt	1 Volt	0 Volt

Called a *NAND gate*
(*NOT AND*)

Combinational Logic Symbols

- Common combinational logic systems have standard symbols called logic gates

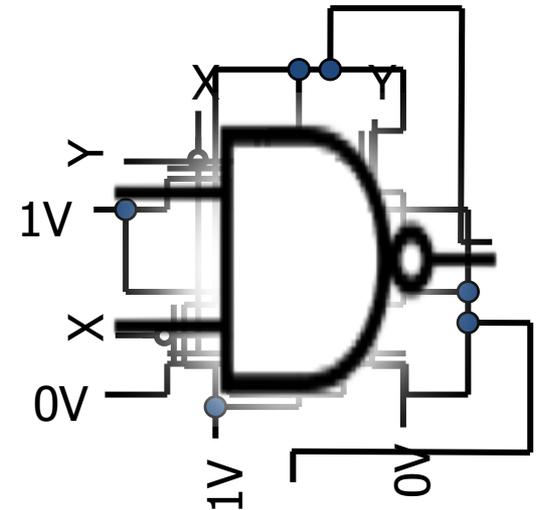
– Buffer, NOT



– AND, NAND



– OR, NOR



Inverting versions (NOT, NAND, NOR) easiest to implement with CMOS transistors (the switches we have available and use most)

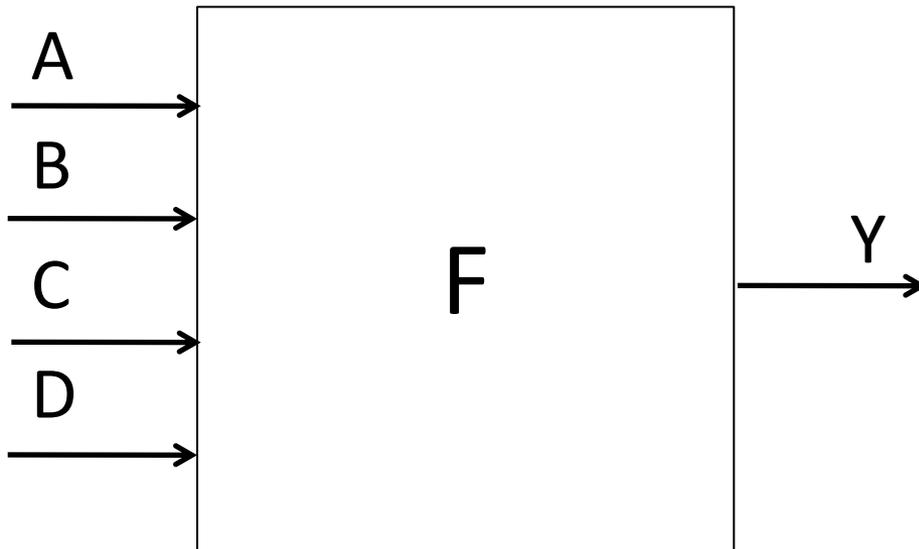
Boolean Algebra

- Use plus “+” for OR
 - “logical sum” $1+0 = 0+1 = 1$ (True); $1+1=2$ (True); $0+0 = 0$ (False)
- Use product for AND ($a \bullet b$ or implied via ab)
 - “logical product” $0*0 = 0*1 = 1*0 = 0$ (False); $1*1 = 1$ (True)
- “Hat” to mean complement (NOT)
- Thus

$$\begin{aligned} & ab + a + \bar{c} \\ = & a \bullet b + a + \bar{c} \\ = & (a \text{ AND } b) \text{ OR } a \text{ OR } (\text{NOT } c) \end{aligned}$$



Truth Tables for Combinational Logic

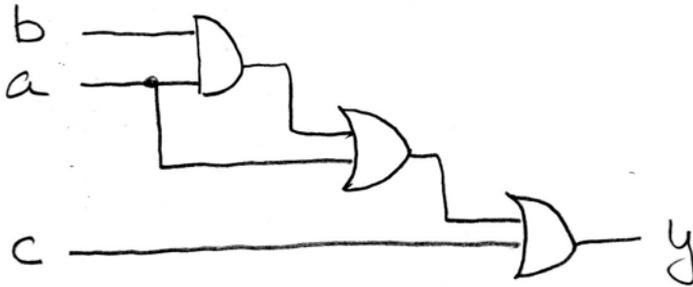


Exhaustive list of the output value
generated for each combination of inputs

How many logic functions can be defined
with N inputs?

a	b	c	d	y
0	0	0	0	F(0,0,0,0)
0	0	0	1	F(0,0,0,1)
0	0	1	0	F(0,0,1,0)
0	0	1	1	F(0,0,1,1)
0	1	0	0	F(0,1,0,0)
0	1	0	1	F(0,1,0,1)
0	1	1	0	F(0,1,1,0)
0	1	1	1	F(0,1,1,1)
1	0	0	0	F(1,0,0,0)
1	0	0	1	F(1,0,0,1)
1	0	1	0	F(1,0,1,0)
1	0	1	1	F(1,0,1,1)
1	1	0	0	F(1,1,0,0)
1	1	0	1	F(1,1,0,1)
1	1	1	0	F(1,1,1,0)
1	1	1	1	F(1,1,1,1)

Boolean Algebra: Circuit & Algebraic Simplification



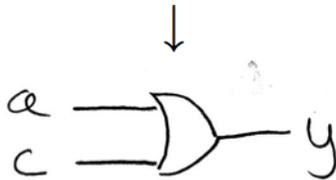
original circuit

$$y = ((ab) + a) + c$$

equation derived from original circuit

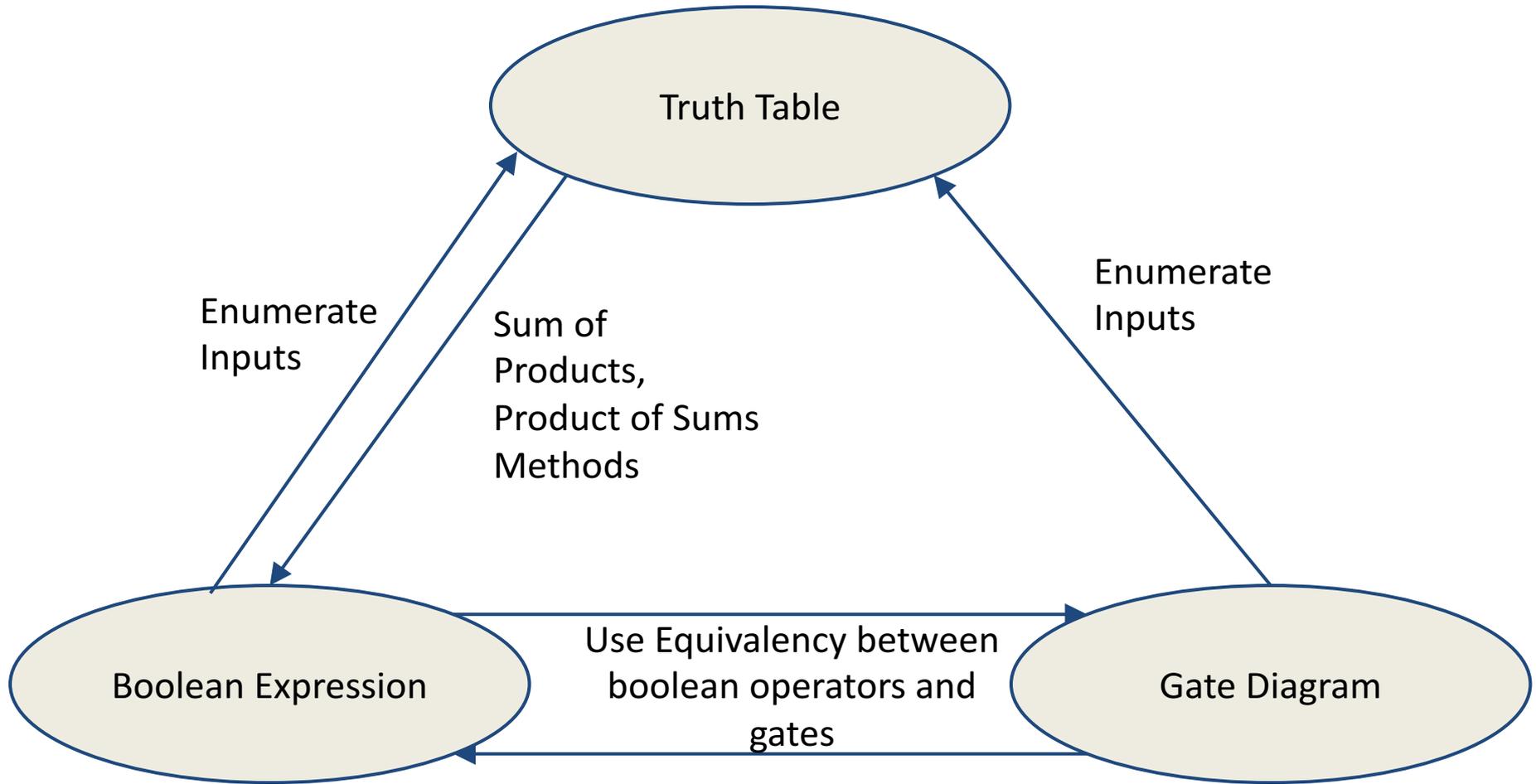
$$\begin{aligned} &\downarrow \\ &= ab + a + c \\ &\downarrow \\ &= a(b + 1) + c \\ &= a(1) + c \\ &= a + c \end{aligned}$$

algebraic simplification



simplified circuit

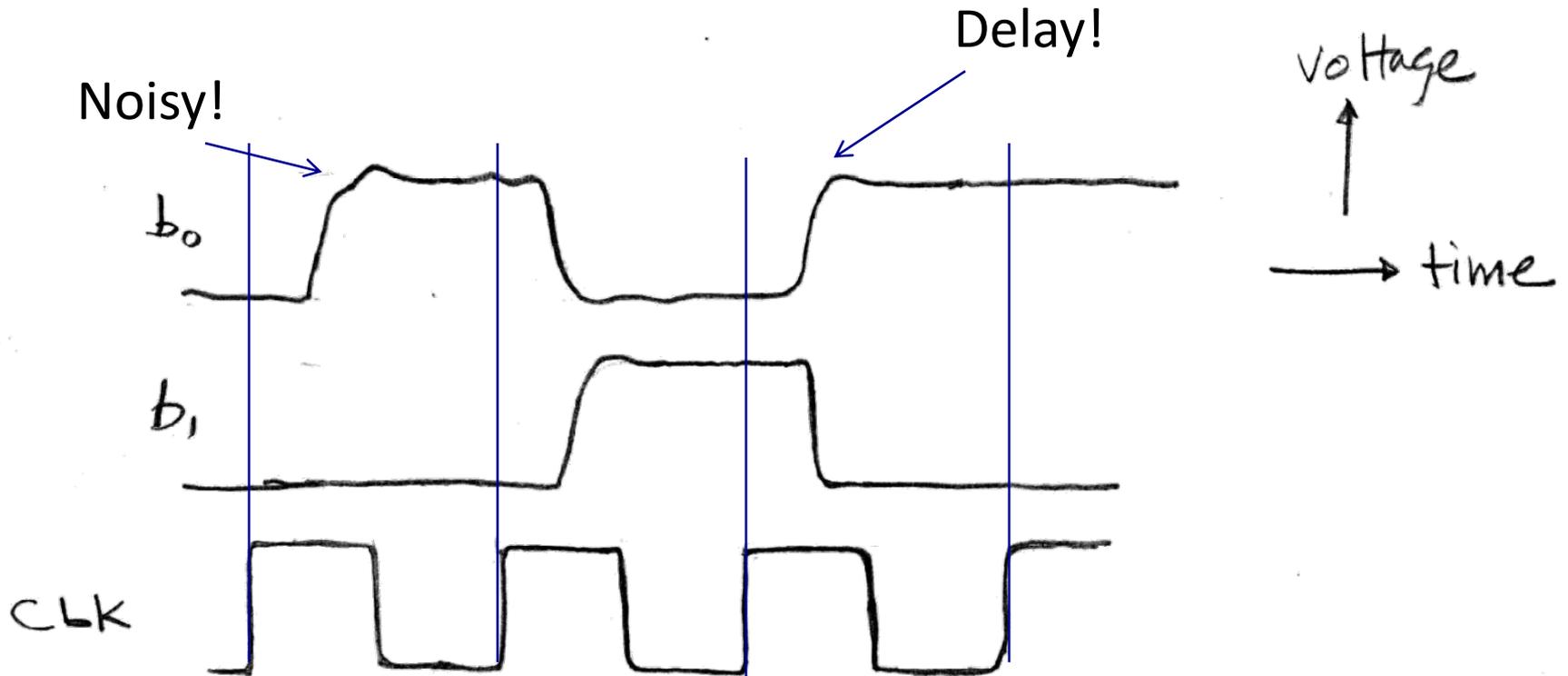
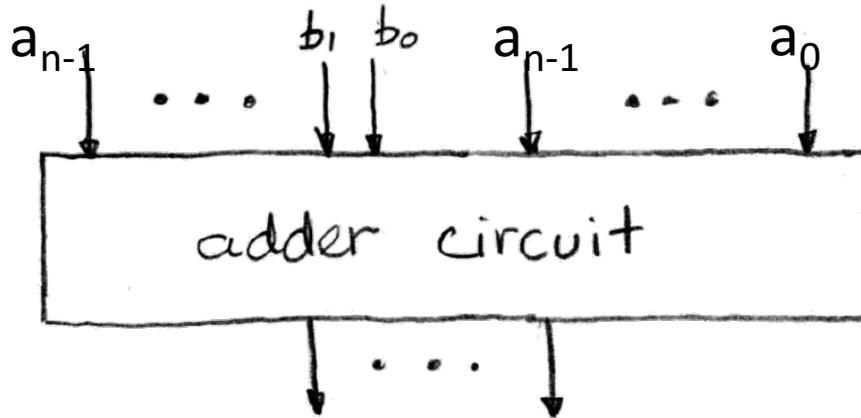
Representations of Combinational Logic (groups of logic gates)



Question

- Simplify $Z = A + BC + \overline{A}(\overline{BC})$
- A: $Z = 0$
- B: $Z = \overline{A(1 + BC)}$
- C: $Z = (A + BC)$
- D: $Z = BC$
- E: $Z = 1$

Signals and Waveforms



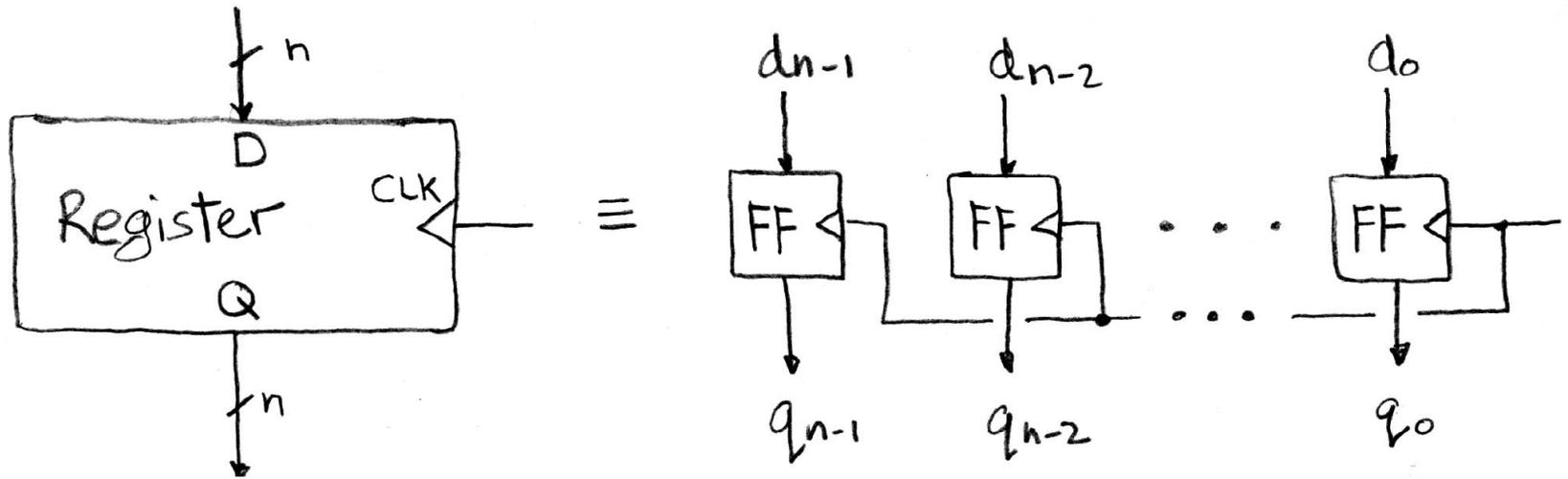
Type of Circuits

- *Synchronous Digital Systems* consist of two basic types of circuits:
 - Combinational Logic (CL) circuits
 - Output is a function of the inputs only, not the history of its execution
 - E.g., circuits to add A, B (ALUs)
 - Sequential Logic (SL)
 - Circuits that “remember” or store information
 - aka “State Elements”
 - E.g., memories and registers (Registers)

Uses for State Elements

- Place to store values for later re-use:
 - Register files (like \$1-\$31 in MIPS)
 - Memory (caches and main memory)
- *Help control flow of information between combinational logic blocks*
 - State elements hold up the movement of information at input to combinational logic blocks to allow for orderly passage

Register Internals



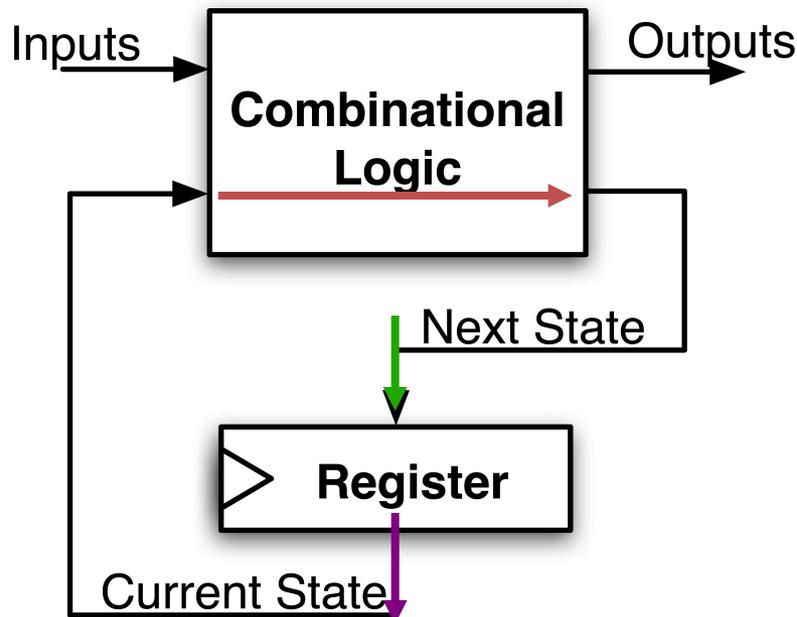
- n instances of a “Flip-Flop”
- Flip-flop name because the output flips and flops between 0 and 1
- D is “data input”, Q is “data output”
- Also called “D-type Flip-Flop”

Recap of Timing Terms

- **Clock (CLK)** - steady square wave that synchronizes system
- **Setup Time** - when the input must be stable before the rising edge of the CLK
- **Hold Time** - when the input must be stable after the rising edge of the CLK
- **“CLK-to-Q” Delay** - how long it takes the output to change, measured from the rising edge of the CLK
- **Flip-flop** - one bit of state that samples every rising edge of the CLK (positive edge-triggered)
- **Register** - several bits of state that samples on rising edge of CLK or on LOAD (positive edge-triggered)

Maximum Clock Frequency

- What is the maximum frequency of this circuit?



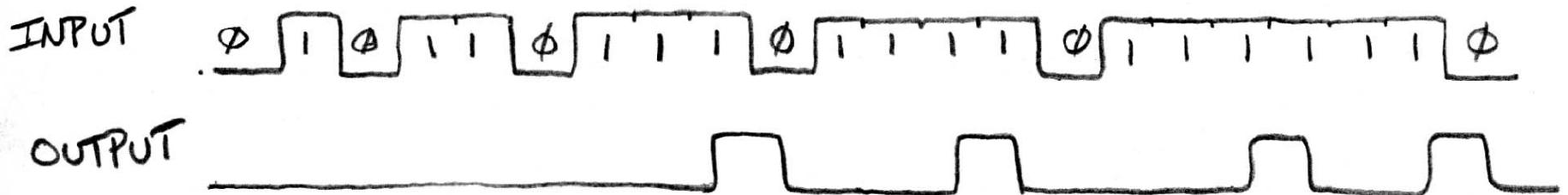
Hint:

Frequency = $1/\text{Period}$

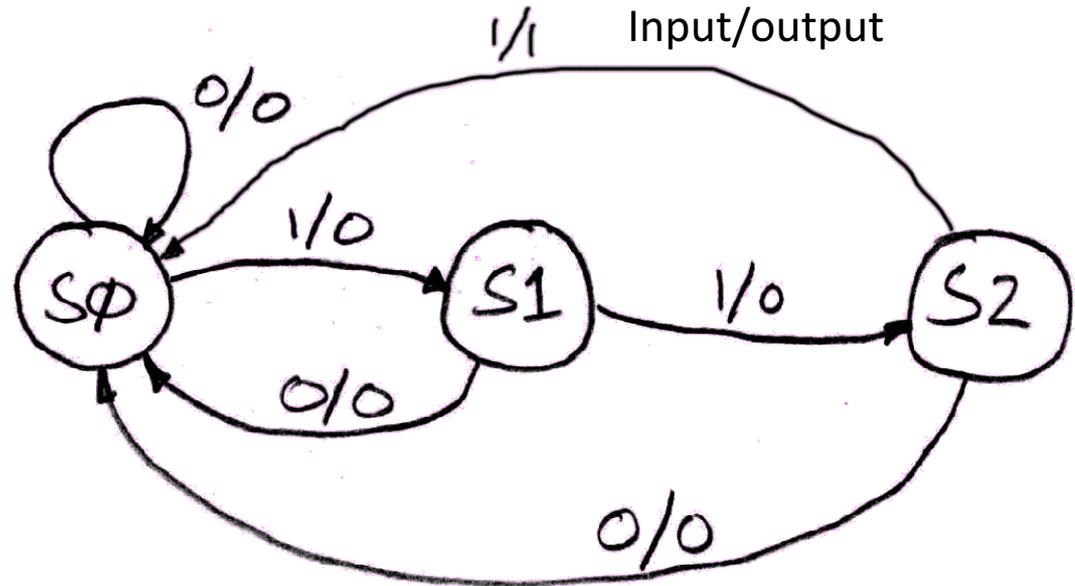
Max Delay = Setup Time + CLK-to-Q Delay + CL Delay

FSM Example: 3 ones...

FSM to detect the occurrence of 3 consecutive 1's in the input.



Draw the FSM...



Assume state transitions are controlled by the clock: on each clock cycle the machine checks the inputs and moves to a new state and produces a new output...

Question

Convert the truth table to a boolean expression
(no need to simplify):

A: $F = xy + x(\sim y)$

B: $F = xy + (\sim x)y + (\sim x)(\sim y)$

C: $F = (\sim x)y + x(\sim y)$

D: $F = xy + (\sim x)y$

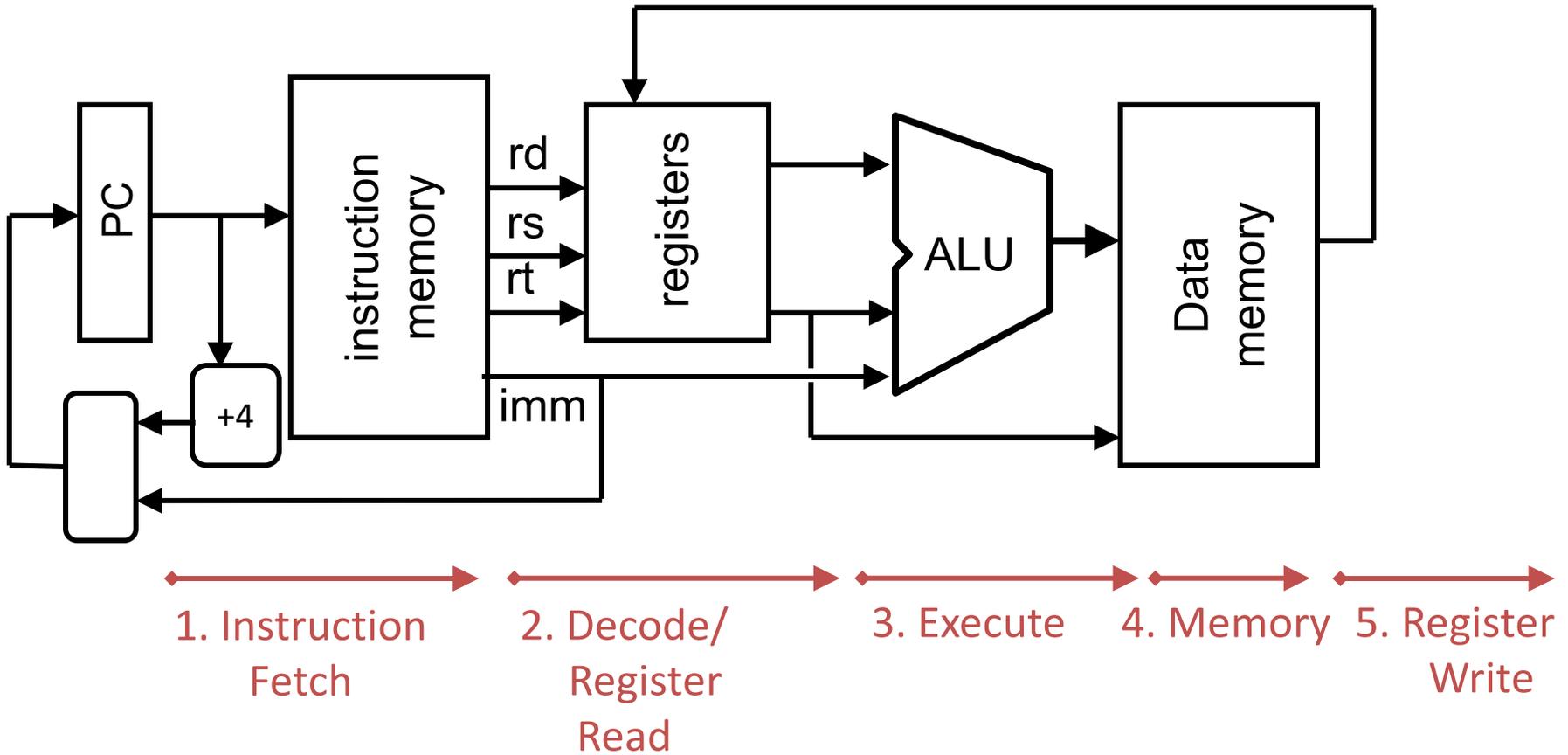
E: $F = (x+y)(\sim x+\sim y)$

x	y	F(x,y)
0	0	0
0	1	1
1	0	0
1	1	1

Datapath: Five Stages of Instruction Execution

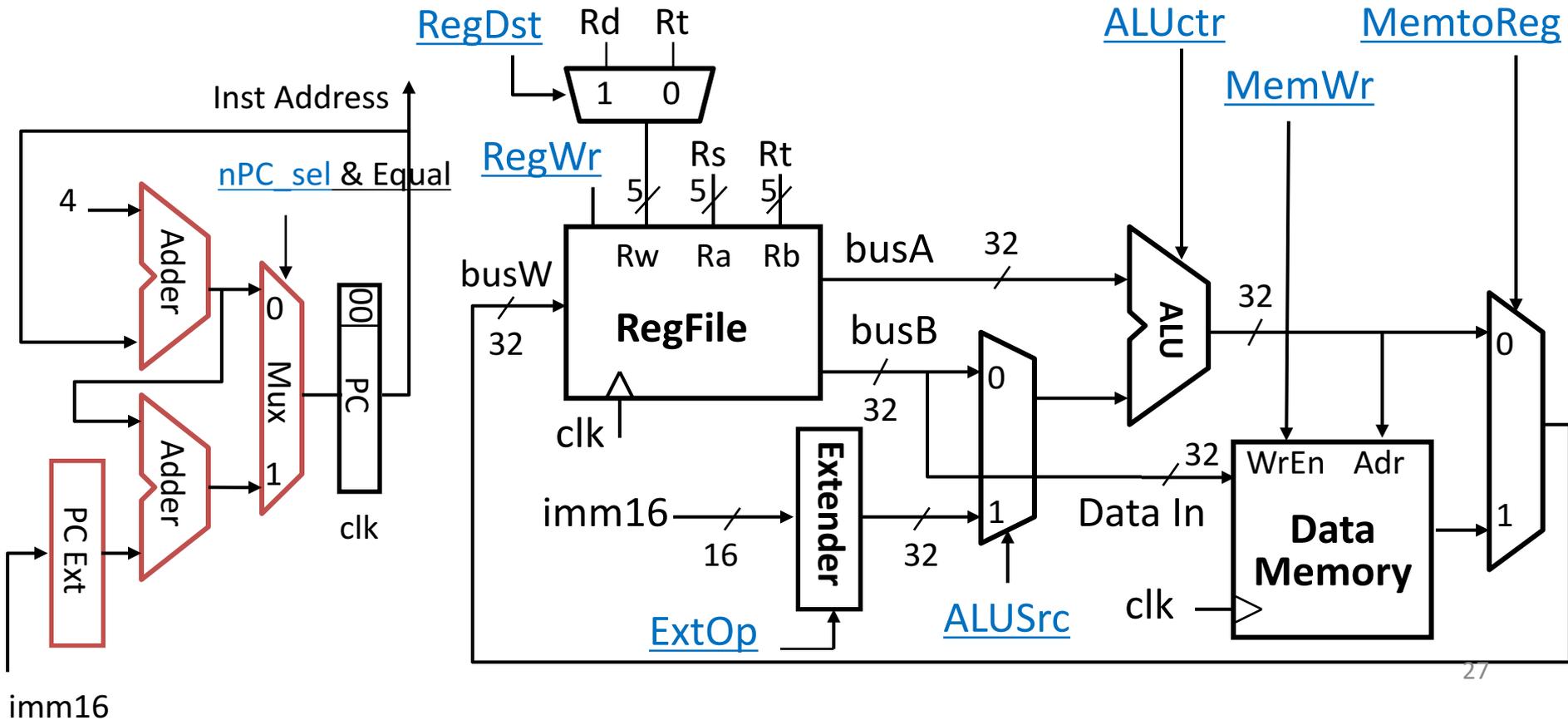
- Stage 1: Instruction Fetch
- Stage 2: Instruction Decode
- Stage 3: ALU (Arithmetic-Logic Unit)
- Stage 4: Memory Access
- Stage 5: Register Write

Stages of Execution on Datapath

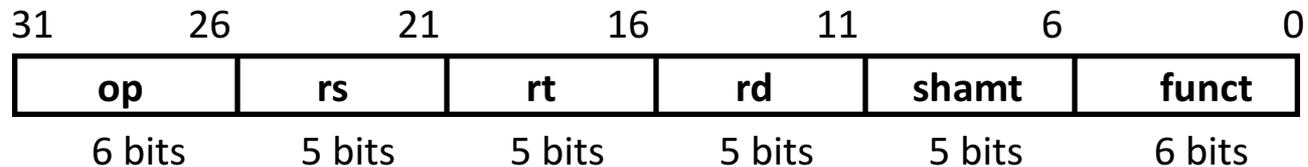


Datapath Control Signals

- ExtOp: “zero”, “sign”
- ALUsrc: 0 => regB; 1 => immed
- ALUctr: “ADD”, “SUB”, “OR”
- nPC_sel: 1 => branch
- MemWr: 1 => write memory
- MemtoReg: 0 => ALU; 1 => Mem
- RegDst: 0 => “rt”; 1 => “rd”
- RegWr: 1 => write register



RTL: The Add Instruction



`add rd, rs, rt`

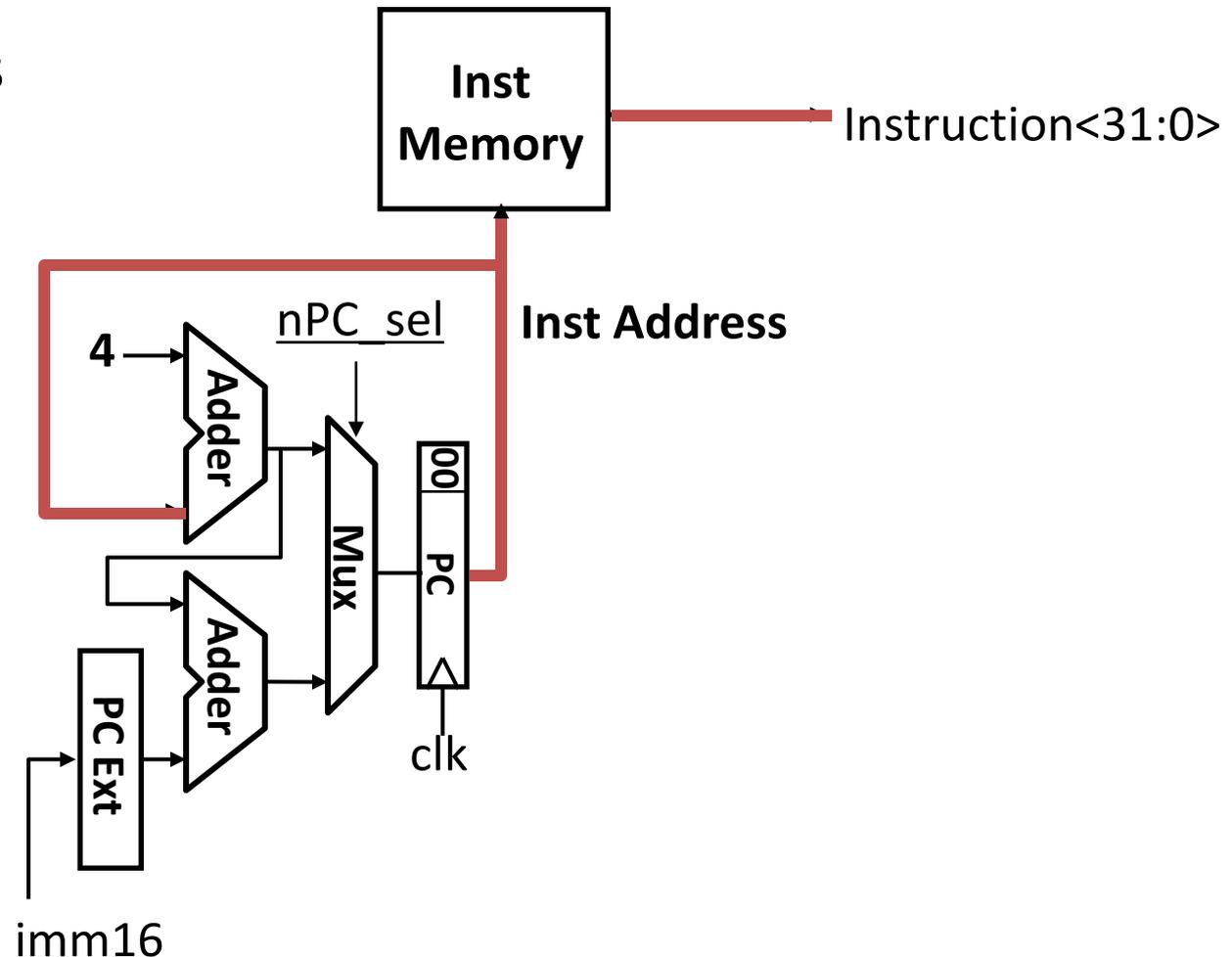
- $\text{MEM}[\text{PC}]$ Fetch the instruction from memory
- $\text{R}[\text{rd}] = \text{R}[\text{rs}] + \text{R}[\text{rt}]$ The actual operation
- $\text{PC} = \text{PC} + 4$ Calculate the next instruction's address

Instruction Fetch Unit at the Beginning of Add

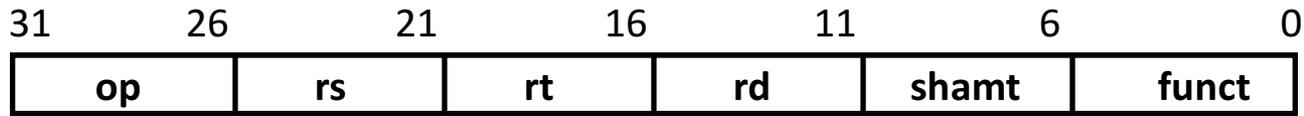
- Fetch the instruction from Instruction memory:

Instruction = MEM[PC]

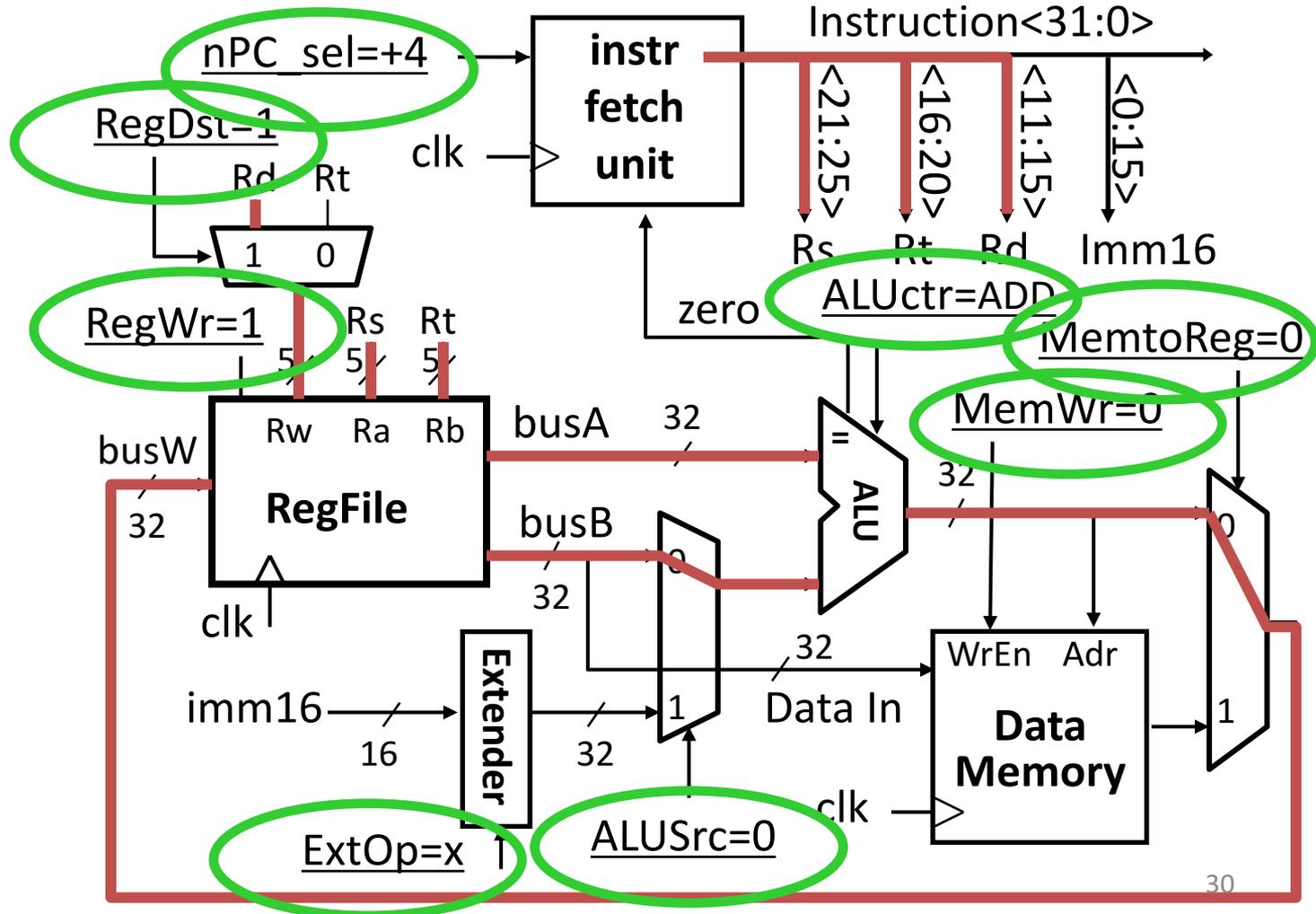
- same for all instructions



Single Cycle Datapath during Add



$$R[rd] = R[rs] + R[rt]$$



Summary of the Control Signals (1/2)

inst Register Transfer

add $R[rd] \leftarrow R[rs] + R[rt]; PC \leftarrow PC + 4$
ALUSrc=RegB, ALUctr="ADD", RegDst=rd, RegWr, nPC_sel="+4"

sub $R[rd] \leftarrow R[rs] - R[rt]; PC \leftarrow PC + 4$
ALUSrc=RegB, ALUctr="SUB", RegDst=rd, RegWr, nPC_sel="+4"

ori $R[rt] \leftarrow R[rs] + \text{zero_ext}(\text{Imm16}); PC \leftarrow PC + 4$
ALUSrc=Im, Extop="Z", ALUctr="OR", RegDst=rt, RegWr, nPC_sel="+4"

lw $R[rt] \leftarrow \text{MEM}[R[rs] + \text{sign_ext}(\text{Imm16})]; PC \leftarrow PC + 4$
ALUSrc=Im, Extop="sn", ALUctr="ADD", MemtoReg, RegDst=rt, RegWr,
nPC_sel = "+4"

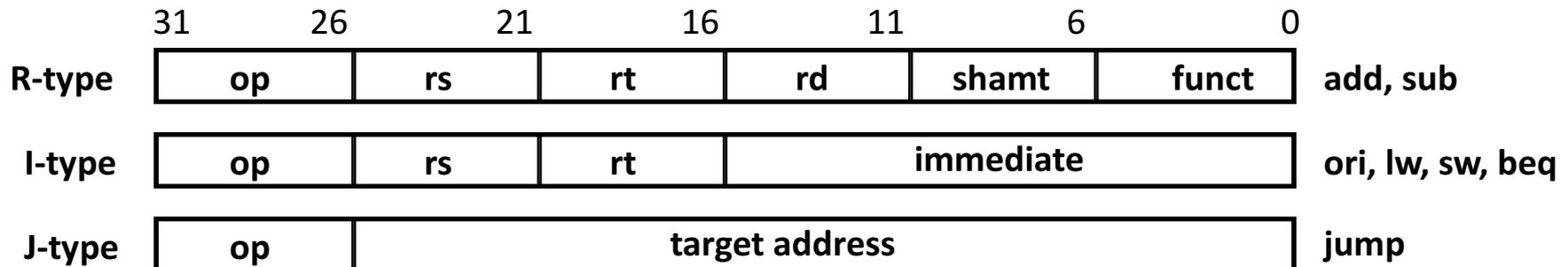
sw $\text{MEM}[R[rs] + \text{sign_ext}(\text{Imm16})] \leftarrow R[rs]; PC \leftarrow PC + 4$
ALUSrc=Im, Extop="sn", ALUctr = "ADD", MemWr, nPC_sel = "+4"

beq if (R[rs] == R[rt]) then PC $\leftarrow PC + \text{sign_ext}(\text{Imm16})$ || 00
else PC $\leftarrow PC + 4$
nPC_sel = "br", ALUctr = "SUB"

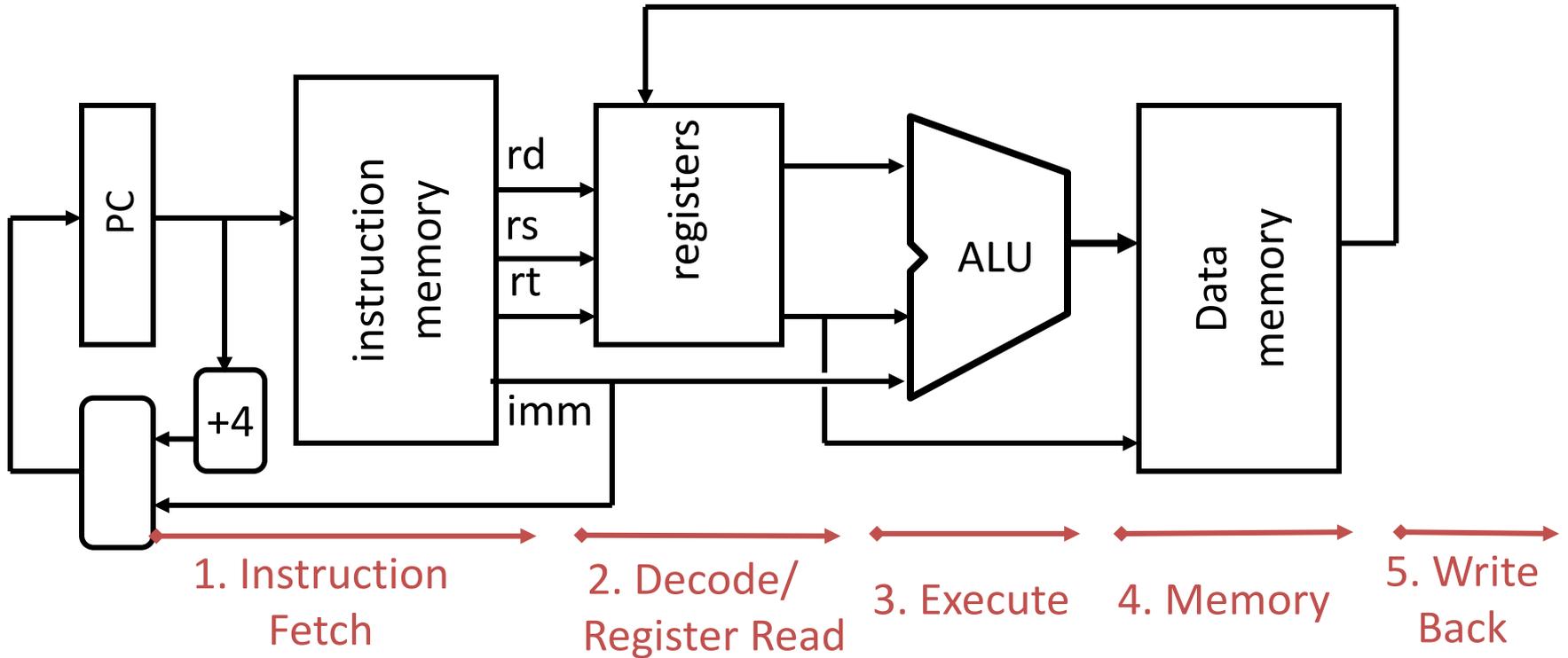
Summary of the Control Signals (2/2)

See Appendix A → **func**
 See Appendix A → **op**

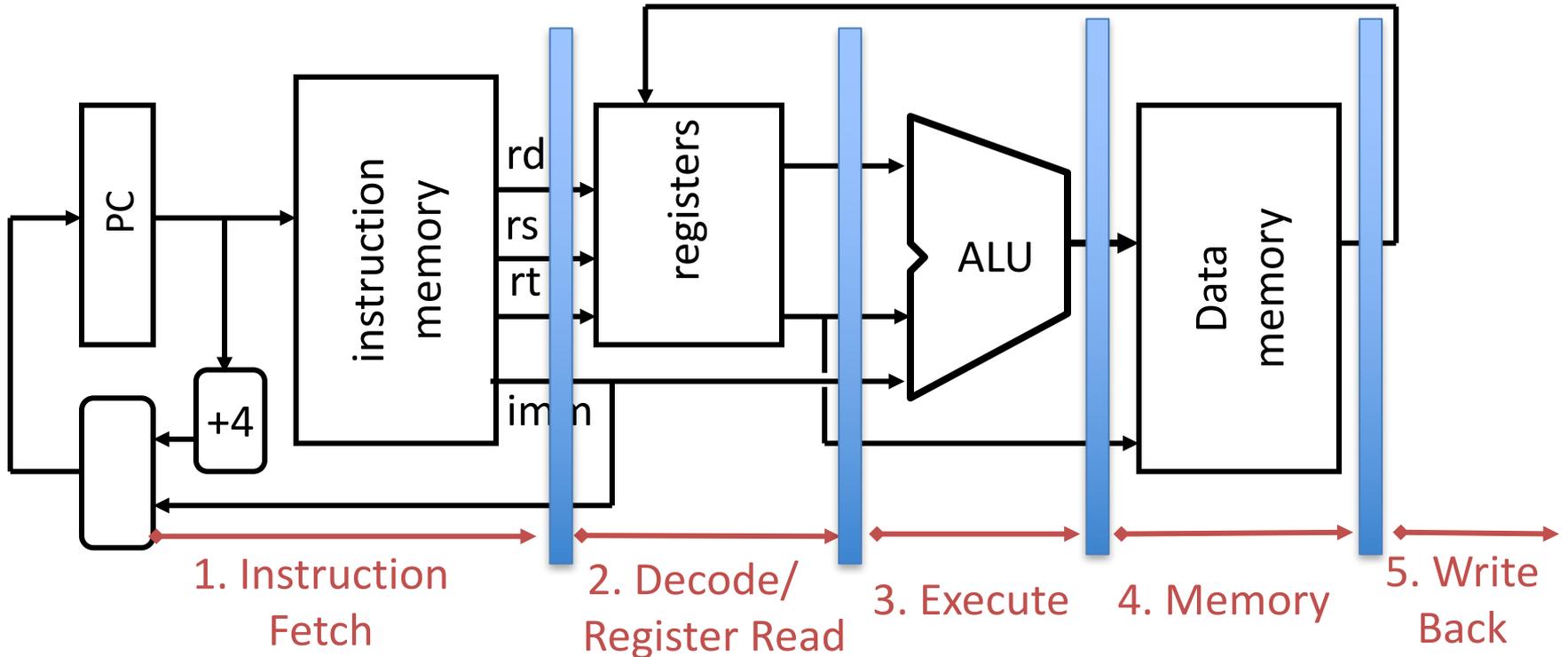
	10 0000	10 0010	We Don't Care :-)				
	00 0000	00 0000	00 1101	10 0011	10 1011	00 0100	00 0010
	add	sub	ori	lw	sw	beq	jump
RegDst	1	1	0	0	x	x	x
ALUSrc	0	0	1	1	1	0	x
MemtoReg	0	0	0	1	x	x	x
RegWrite	1	1	1	1	0	0	0
MemWrite	0	0	0	0	1	0	0
nPCsel	0	0	0	0	0	1	x
Jump	0	0	0	0	0	0	1
ExtOp	x	x	0	1	1	x	x
ALUctr<2:0>	Add	Subtract	Or	Add	Add	Subtract	x



Pipelining: Single Cycle Datapath



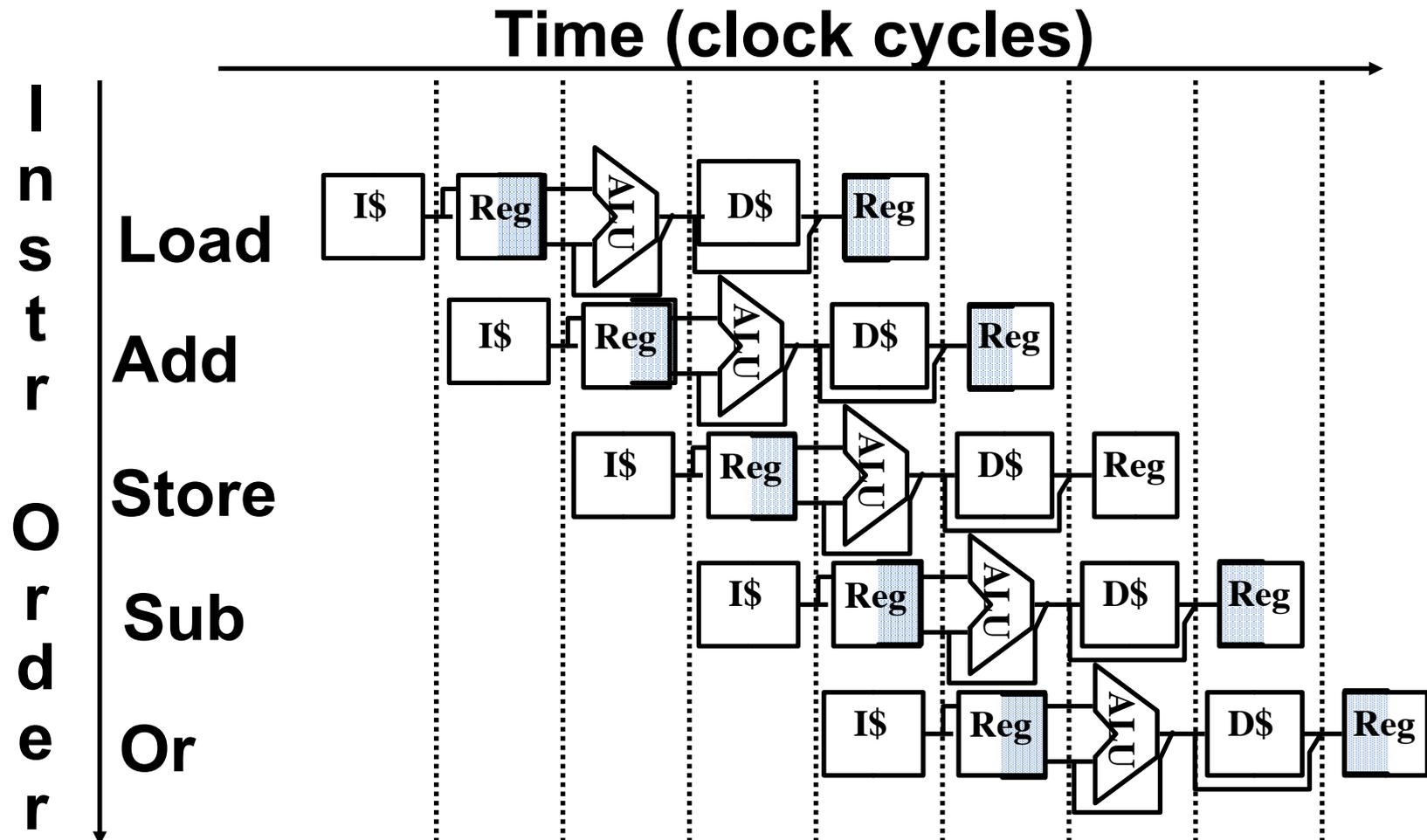
Pipeline registers



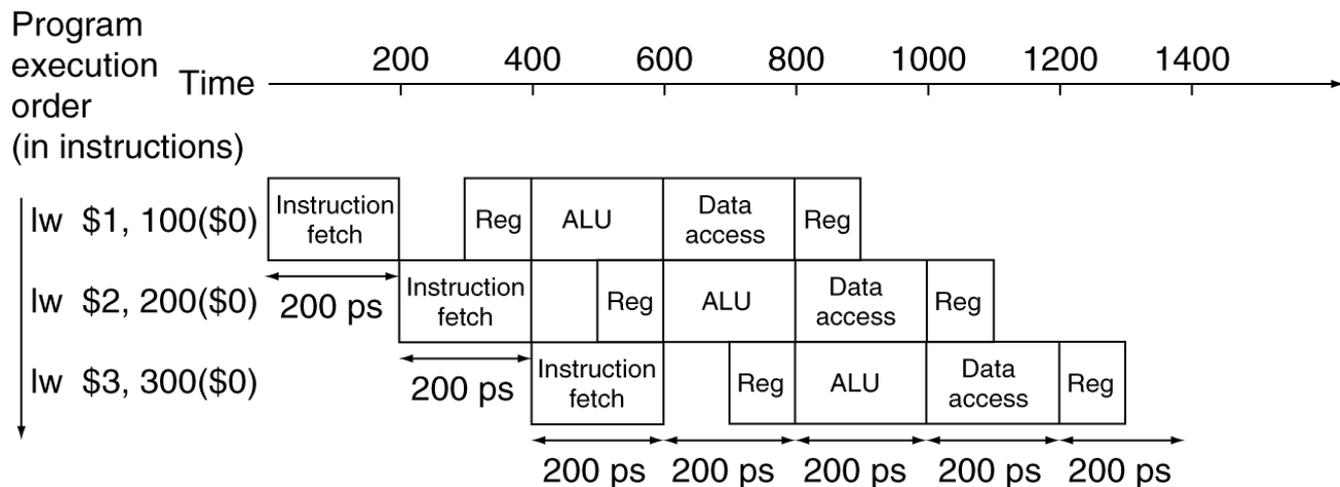
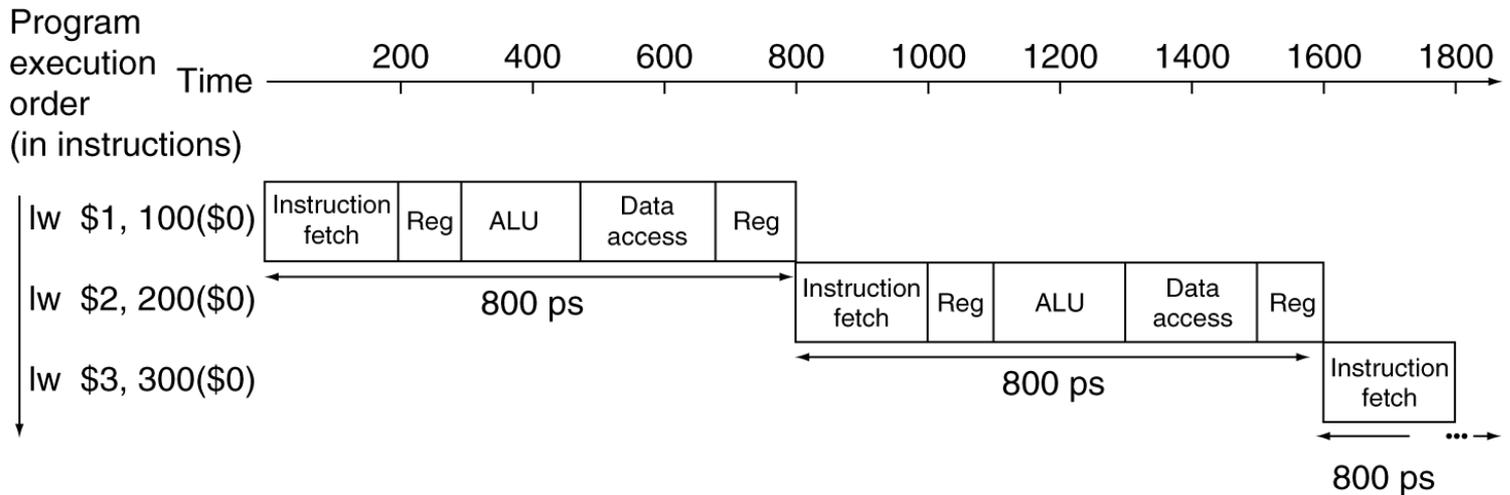
- Need registers between stages
 - To hold information produced in previous cycle

Graphical Pipeline Representation

- RegFile: left half is write, right half is read



Pipelining Performance (3/3)



Pipelining Hazards

A *hazard* is a situation that prevents starting the next instruction in the next clock cycle

1) *Structural hazard*

- A required resource is busy (e.g. needed in multiple stages)

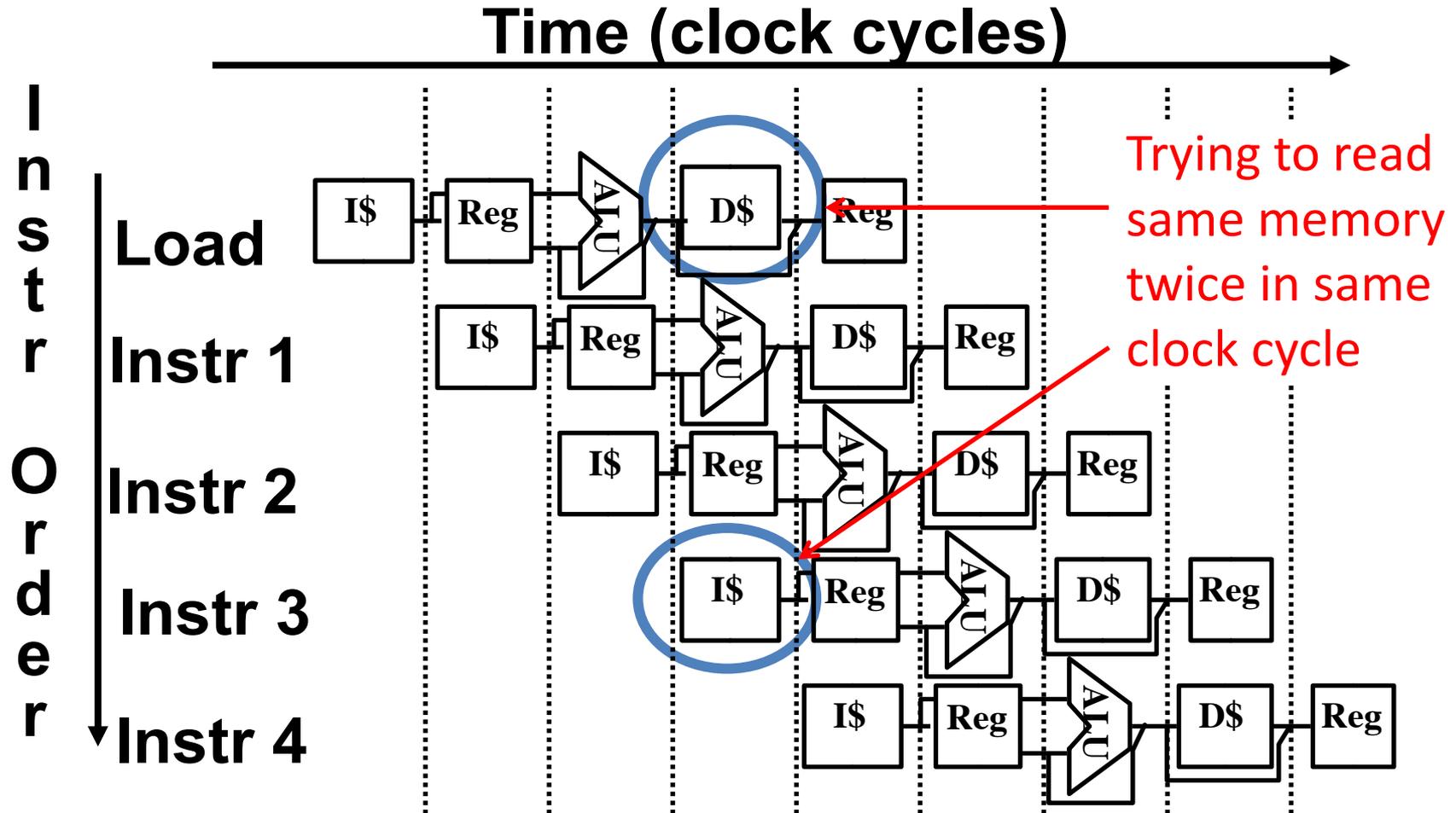
2) *Data hazard*

- Data dependency between instructions
- Need to wait for previous instruction to complete its data read/write

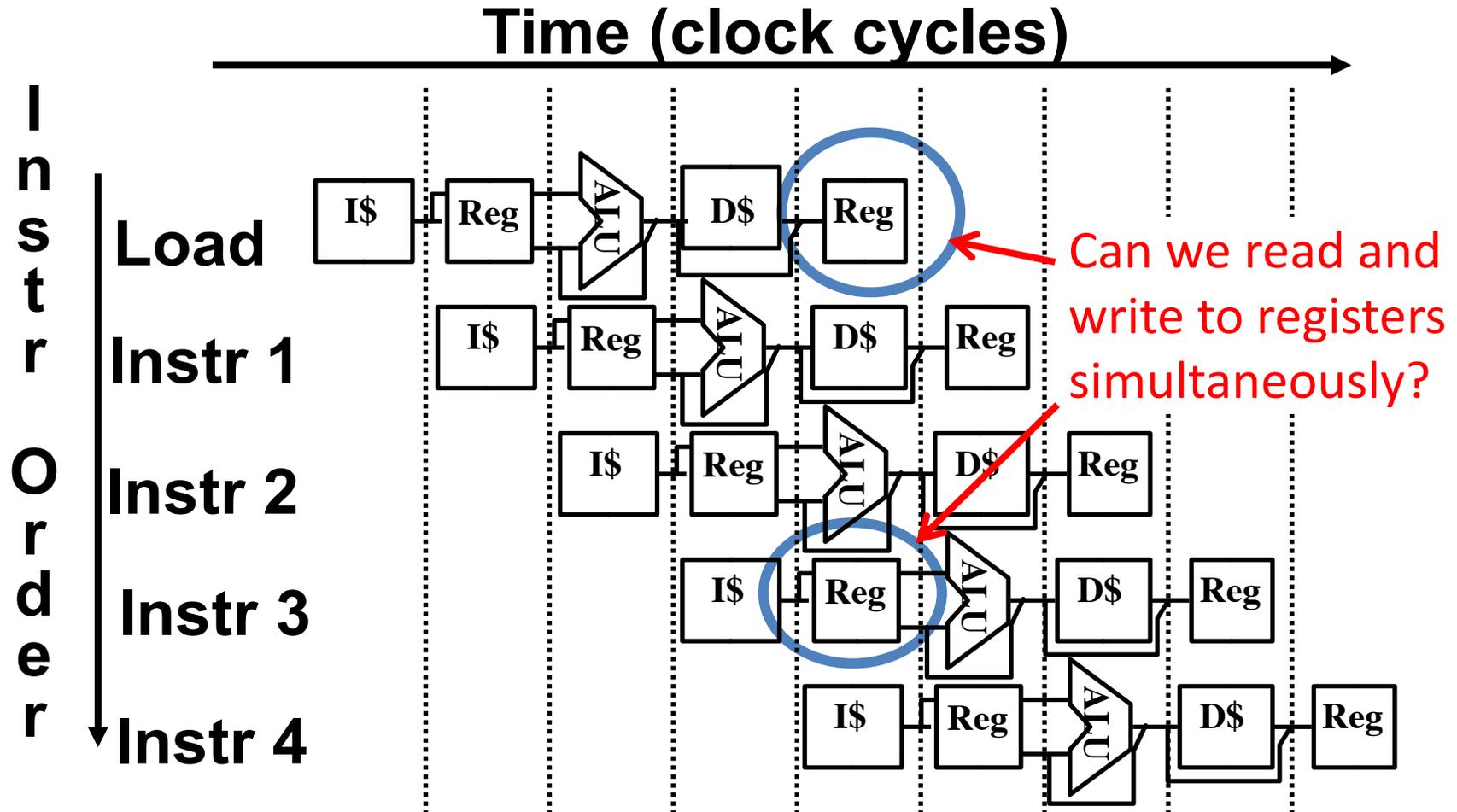
3) *Control hazard*

- Flow of execution depends on previous instruction

Structural Hazard #1: Single Memory



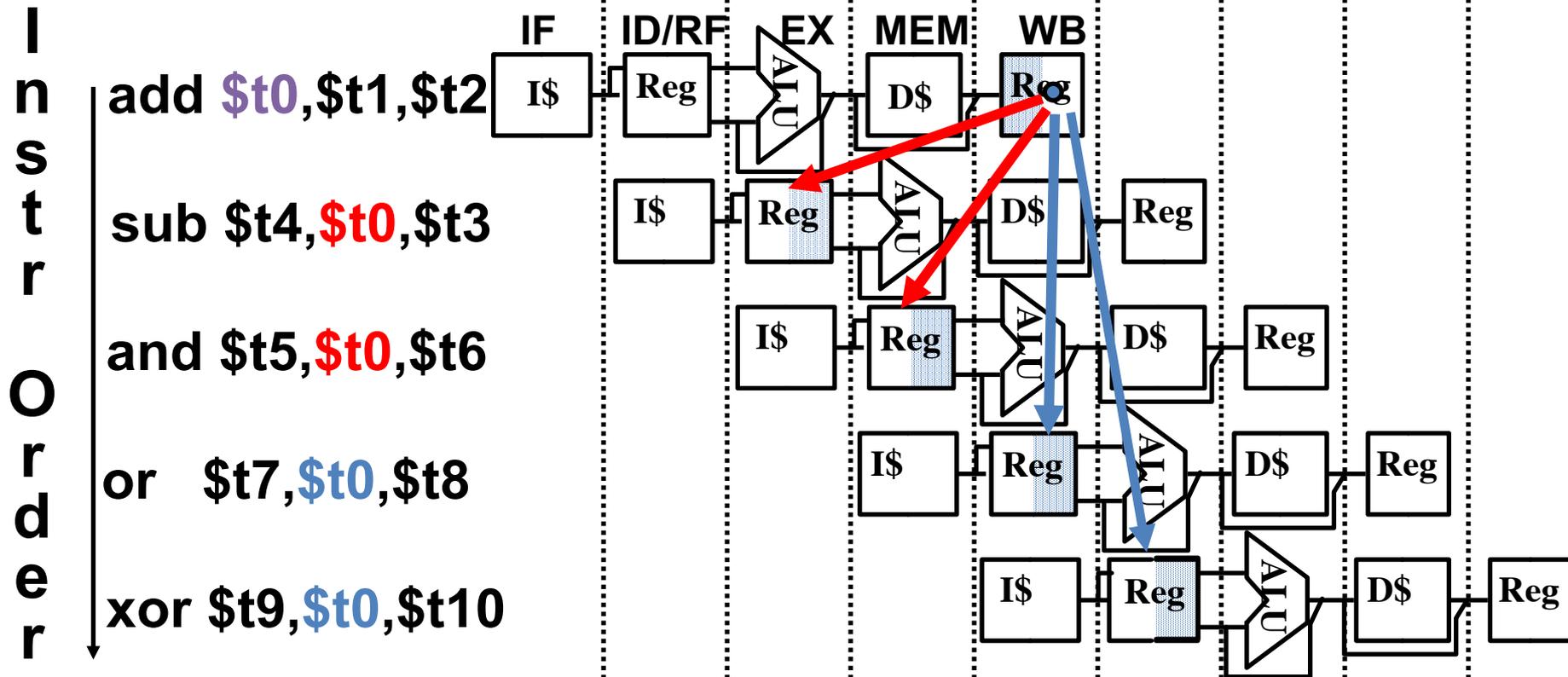
Structural Hazard #2: Registers (1/2)



2. Data Hazards (2/2)

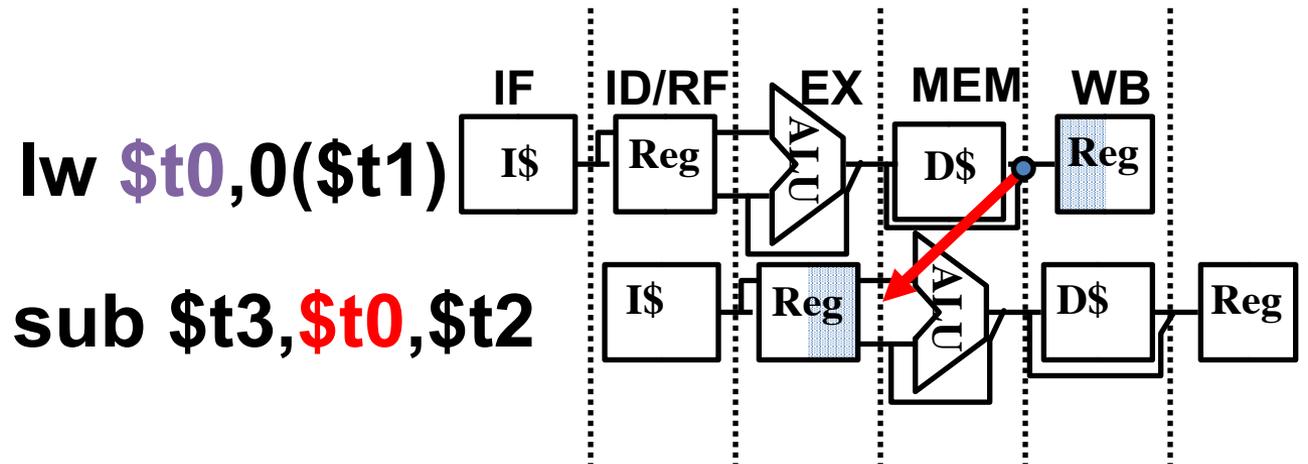
- Data-flow *backwards* in time are hazards

Time (clock cycles)



Data Hazard: Loads (1/4)

- **Recall:** Dataflow backwards in time are hazards



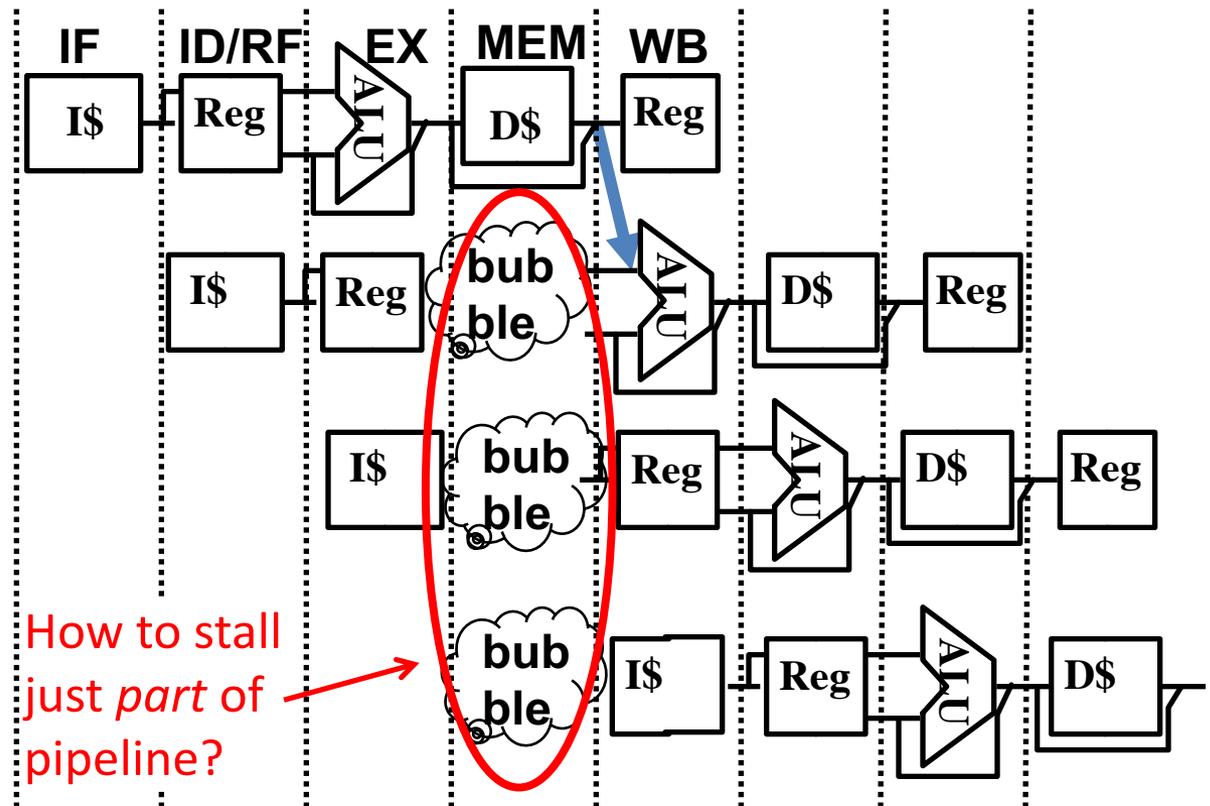
- Can't solve all cases with forwarding
 - Must *stall* instruction dependent on load, then forward (more hardware)

Data Hazard: Loads (2/4)

- *Hardware stalls pipeline*
 - Called “hardware interlock”

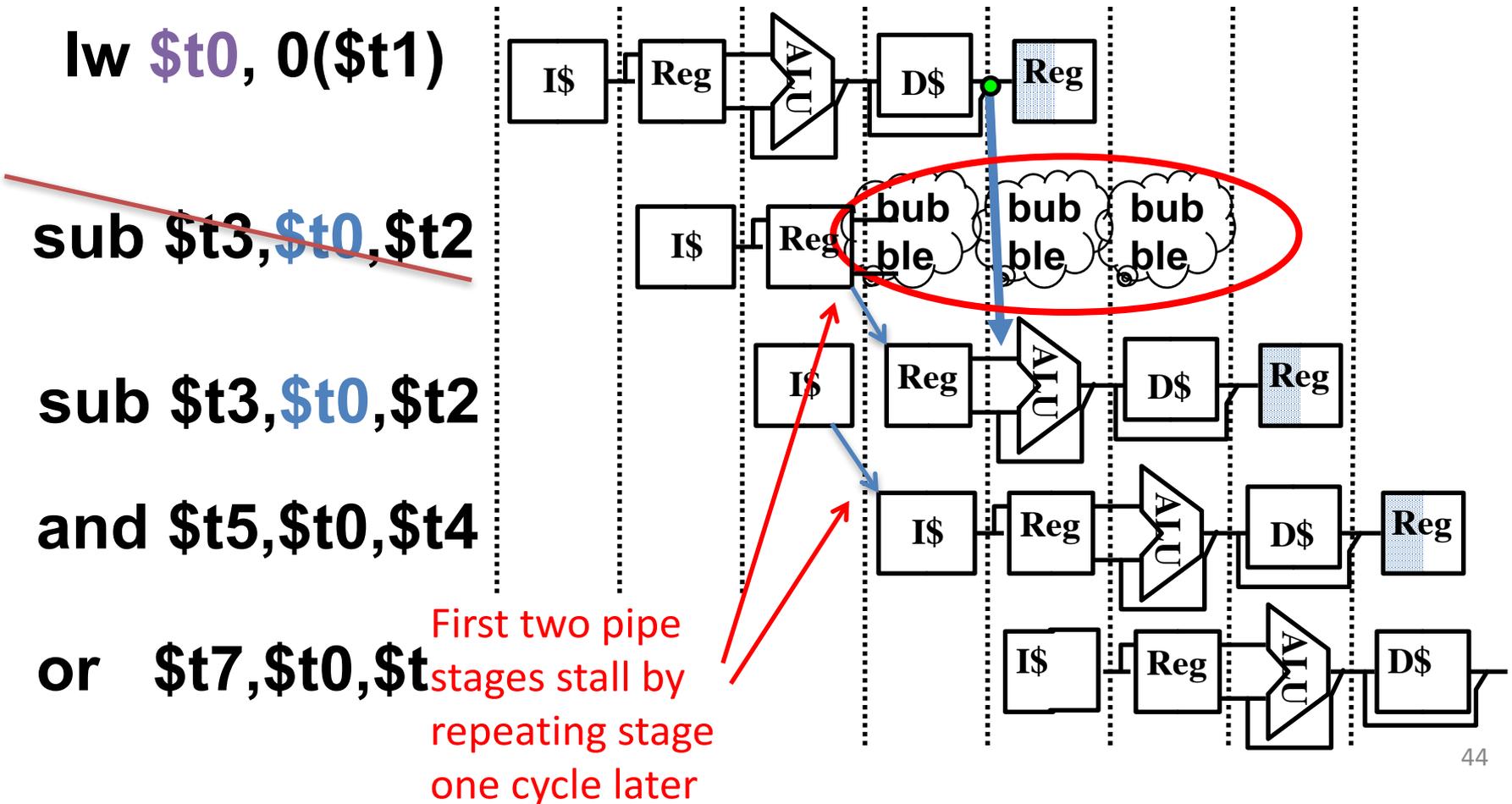
Schematically, this is what we want, but in reality stalls done “horizontally”

lw \$t0, 0(\$t1)
sub \$t3,\$t0,\$t2
and \$t5,\$t0,\$t4
or \$t7,\$t0,\$t6



Data Hazard: Loads (3/4)

- Stalled instruction converted to “bubble”, acts like nop



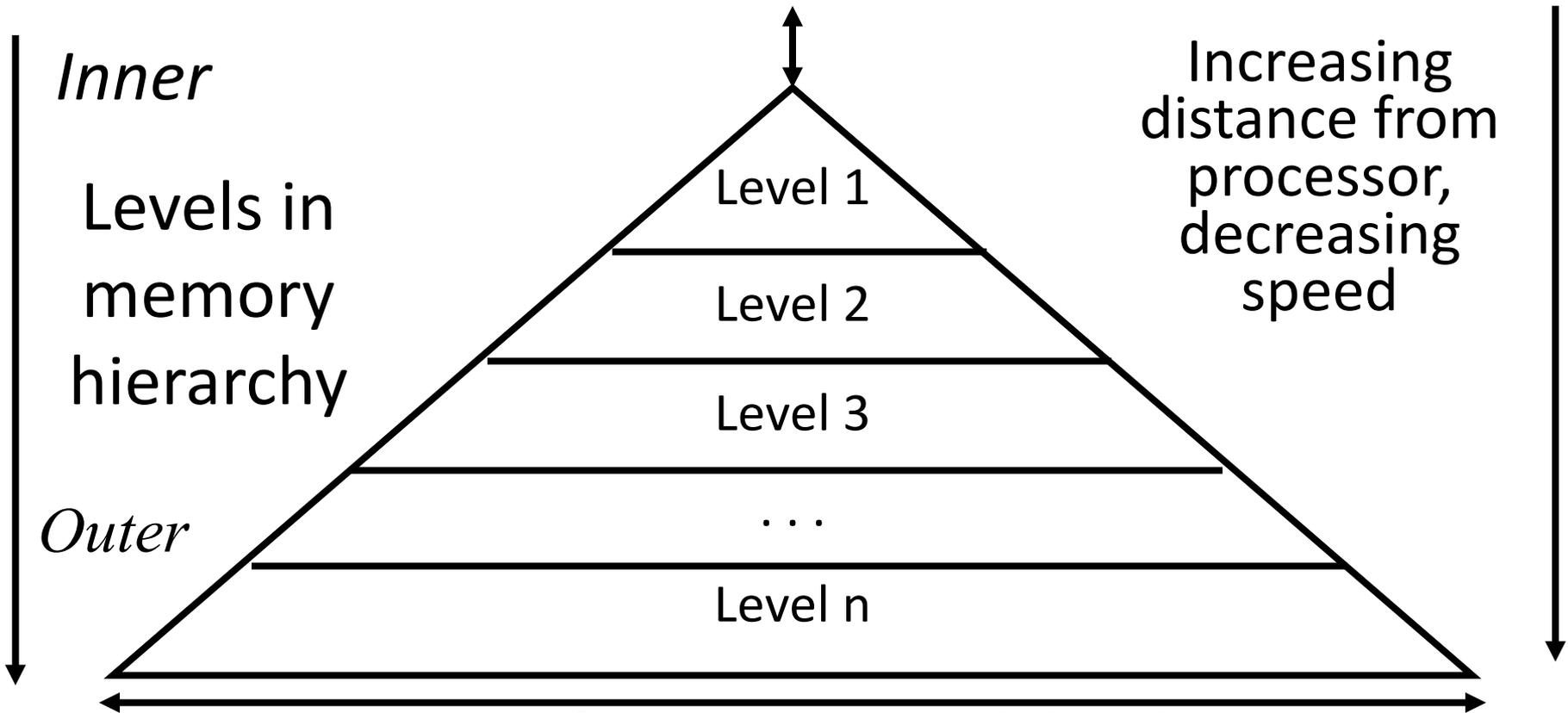
3. Control Hazards

- Branch determines flow of control
 - Fetching next instruction depends on branch outcome
 - Pipeline can't always fetch correct instruction
 - Still working on ID stage of branch
- BEQ, BNE in MIPS pipeline
- Simple solution Option 1: *Stall* on every branch until branch condition resolved
 - Would add 2 bubbles/clock cycles for every Branch! (~ 20% of instructions executed)

Caches

Big Idea: Memory Hierarchy

Processor



Inner

Levels in
memory
hierarchy

Outer

Increasing
distance from
processor,
decreasing
speed

Level 1

Level 2

Level 3

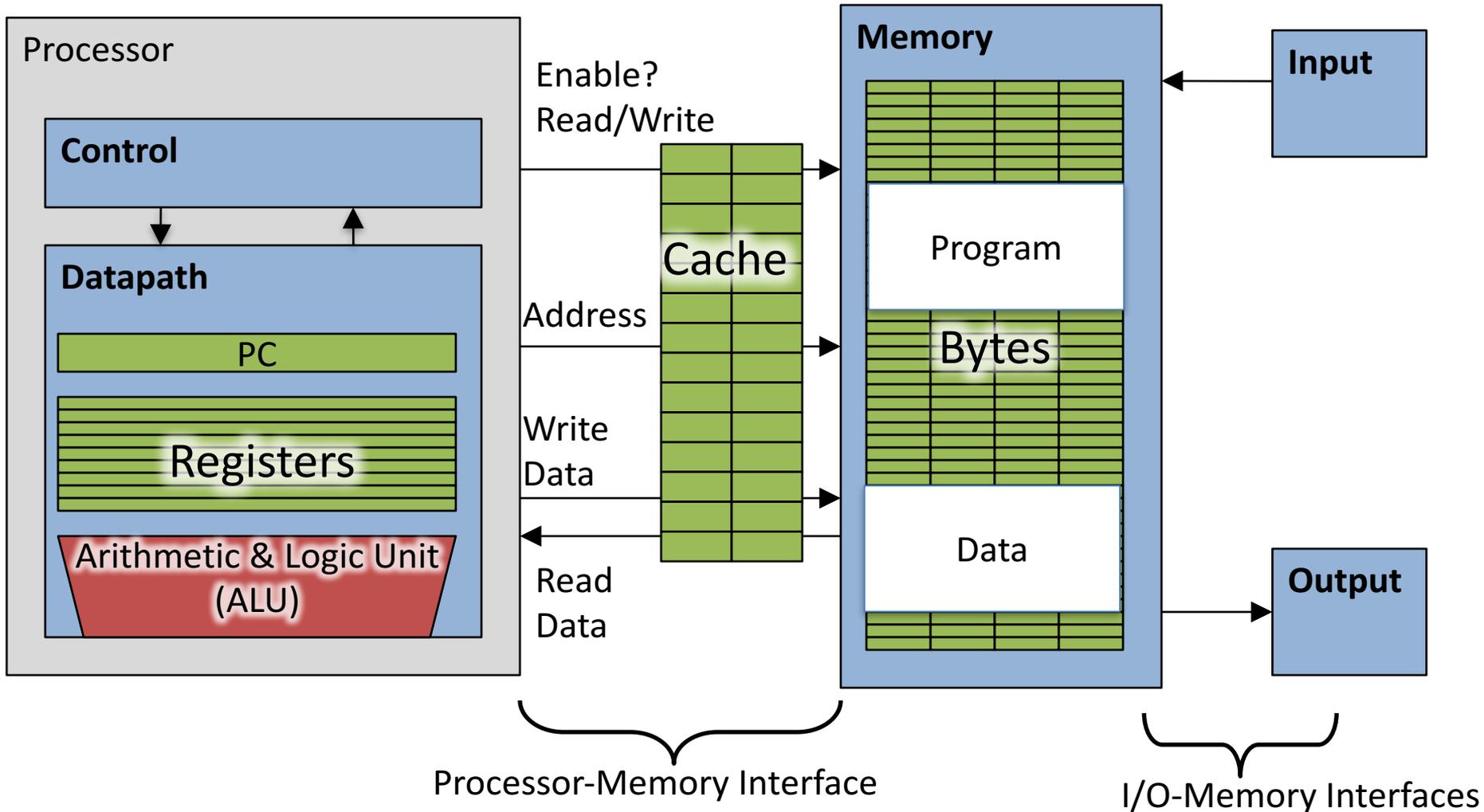
...

Level n

Size of memory at each level

*As we move to outer levels the latency goes up
and price per bit goes down. Why?*

Adding Cache to Computer



Total Cash Capacity =

Associativity * # of sets * block_size

*Bytes = blocks/set * sets * Bytes/block*

$$C = N * S * B$$



$$\begin{aligned} \text{address_size} &= \text{tag_size} + \text{index_size} + \text{offset_size} \\ &= \text{tag_size} + \log_2(S) + \log_2(B) \end{aligned}$$

Clicker Question: C remains constant, S and/or B can change such that

$$C = 2N * (SB)' \Rightarrow (SB)' = SB/2$$

$$\begin{aligned} \text{Tag_size} &= \text{address_size} - (\log_2(S) + \log_2(B)) = \text{address_size} - \log_2(SB) \\ &= \text{address_size} - (\log_2(SB) - 1) \end{aligned}$$

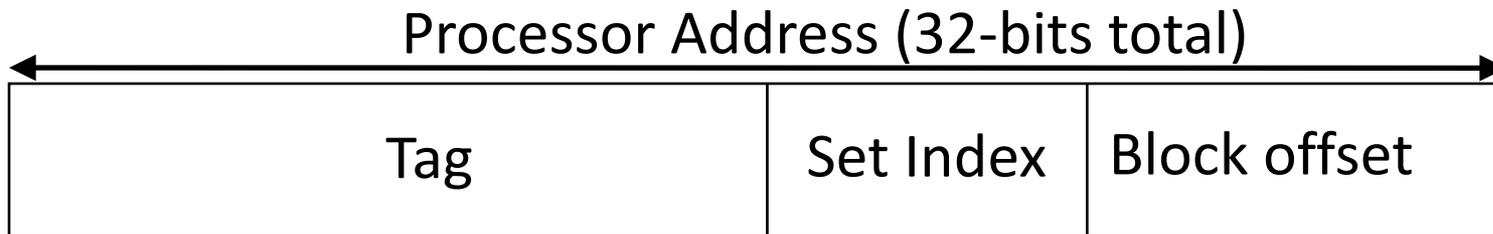
- Principle of Locality for Libraries /Computer Memory
- Hierarchy of Memories (speed/size/cost per bit) to Exploit Locality
- Cache – copy of data lower level in memory hierarchy
- Direct Mapped to find block in cache using Tag field and Valid bit for Hit
- Cache design choice:
 - Write-Through vs. Write-Back

Cache Organizations

- “Fully Associative”: Block can go anywhere
 - First design in lecture
 - Note: No Index field, but 1 comparator/block
- “Direct Mapped”: Block goes one place
 - Note: Only 1 comparator
 - Number of sets = number blocks
- “N-way Set Associative”: N places for a block
 - Number of sets = number of blocks / N
 - N comparators
 - **Fully Associative: $N = \text{number of blocks}$**
 - **Direct Mapped: $N = 1$**

Processor Address Fields used by Cache Controller

- **Block Offset:** Byte address within block
- **Set Index:** Selects which set
- **Tag:** Remaining portion of processor address



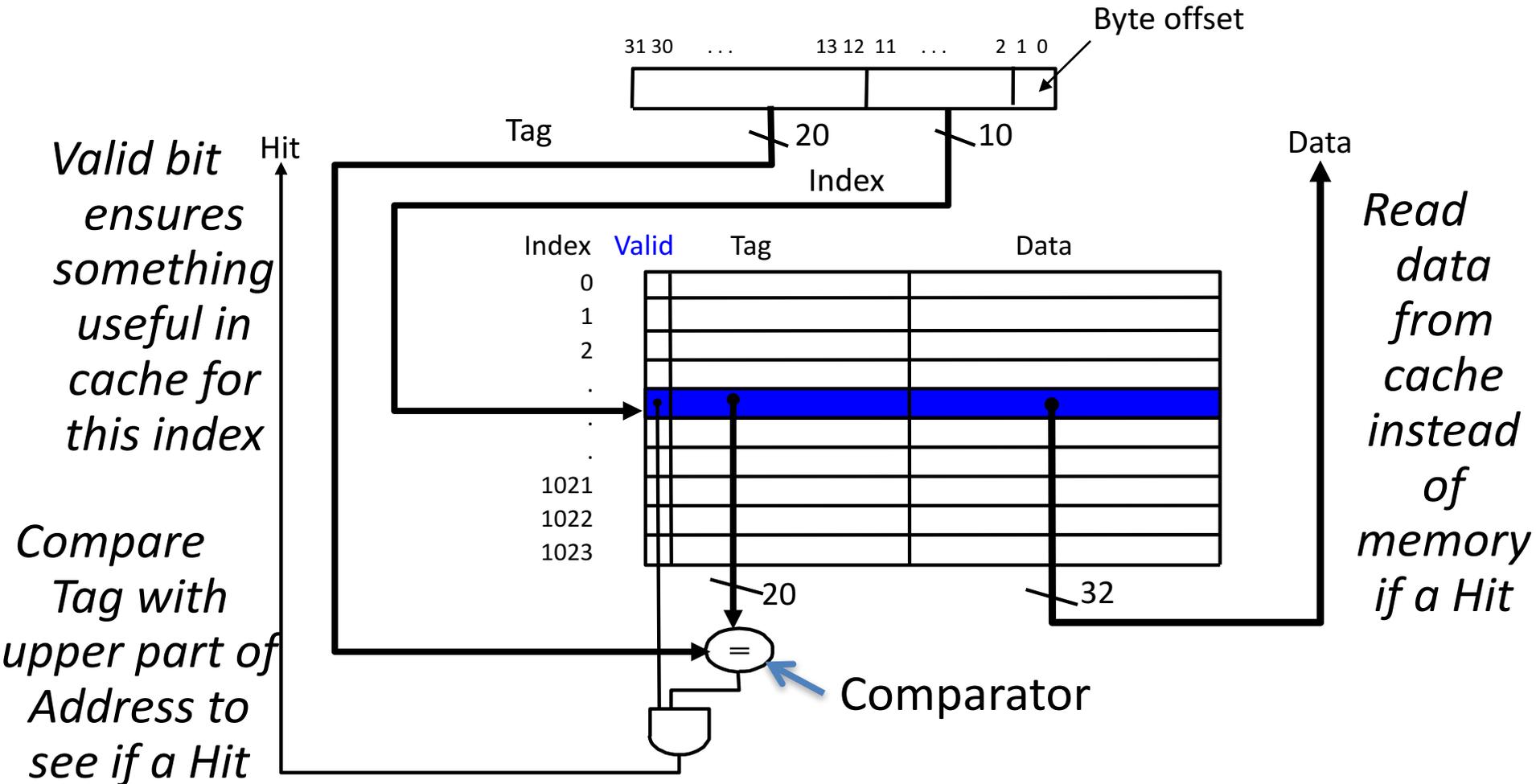
- Size of Index = \log_2 (number of sets)
- Size of Tag = Address size – Size of Index – \log_2 (number of bytes/block)

Write Policy Choices

- Cache hit:
 - **write through**: writes both cache & memory on every access
 - Generally higher memory traffic but simpler pipeline & cache design
 - **write back**: writes cache only, memory `written only when dirty entry evicted
 - A dirty bit per line reduces write-back traffic
 - Must handle 0, 1, or 2 accesses to memory for each load/store
- Cache miss:
 - **no write allocate**: only write to main memory
 - **write allocate** (aka fetch on write): fetch into cache
- Common combinations:
 - write through and no write allocate
 - write back with write allocate

Direct-Mapped Cache Review

- One word blocks, cache size = 1K words (or 4KB)



Sources of Cache Misses (3 C's)

- *Compulsory* (cold start, first reference):
 - 1st access to a block, “cold” fact of life, not a lot you can do about it.
 - If running billions of instructions, compulsory misses are insignificant
- *Capacity*:
 - Cache cannot contain all blocks accessed by the program
 - Misses that would not occur with infinite cache
- *Conflict* (collision):
 - Multiple memory locations mapped to same cache set
 - Misses that would not occur with ideal fully associative cache

Impact of Cache Parameters on Performance

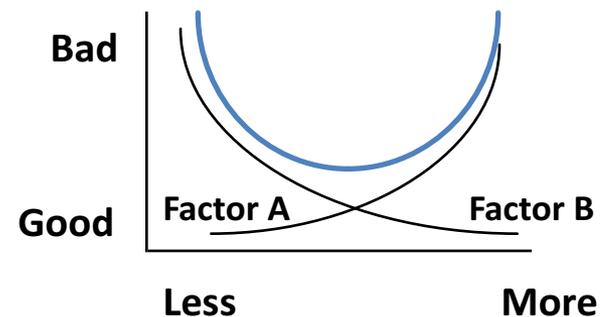
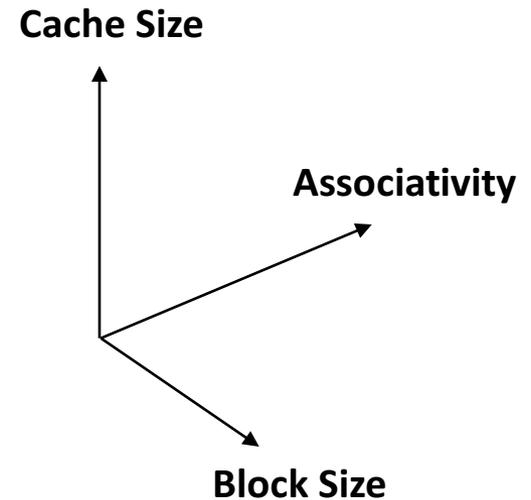
- $AMAT = \text{Hit Time} + \text{Miss Rate} * \text{Miss Penalty}$
 - Note, we assume always first search cache, so must charge hit time for both hits and misses!
- For misses, characterize by 3Cs

Local vs. Global Miss Rates

- *Local miss rate* – the fraction of references to one level of a cache that miss
- Local Miss rate L2\$ = $\frac{\$L2 \text{ Misses}}{\$L1 \text{ Misses}}$
- *Global miss rate* – the fraction of references that miss in all levels of a multilevel cache
 - L2\$ local miss rate \gg than the global miss rate
- Global Miss rate = $\frac{\$L2 \text{ Misses}}{\text{Total Accesses}}$
= $\left(\frac{\$L2 \text{ Misses}}{\$L1 \text{ Misses}}\right) \times \left(\frac{\$L1 \text{ Misses}}{\text{Total Accesses}}\right)$
= Local Miss rate L2\$ \times Local Miss rate L1\$
- AMAT = Time for a hit + Miss rate \times Miss penalty
- AMAT = Time for a L1\$ hit + (local) Miss rate L1\$ \times (Time for a L2\$ hit + (local) Miss rate L2\$ \times L2\$ Miss penalty)

In Conclusion, Cache Design Space

- Several interacting dimensions
 - Cache size
 - Block size
 - Associativity
 - Replacement policy
 - Write-through vs. write-back
 - Write-allocation
- Optimal choice is a compromise
 - Depends on access characteristics
 - Workload
 - Use (I-cache, D-cache)
 - Depends on technology / cost
- Simplicity often wins



Iron Law of Performance

- A program executes instructions
- CPU Time/Program
= Clock Cycles/Program x Clock Cycle Time
= Instructions/Program
x Average Clock Cycles/Instruction
x Clock Cycle Time
- 1st term called *Instruction Count*
- 2nd term abbreviated *CPI* for average *Clock Cycles Per Instruction*
- 3rd term is 1 / Clock rate

Workload and Benchmark

- *Workload*: Set of programs run on a computer
 - Actual collection of applications run or made from real programs to approximate such a mix
 - Specifies programs, inputs, and relative frequencies
- *Benchmark*: Program selected for use in comparing computer performance
 - Benchmarks form a workload
 - Usually standardized so that many use them

IEEE 754 Floating Point Standard (2/3)

- IEEE 754 uses “biased exponent” representation
 - Designers wanted FP numbers to be used even if no FP hardware; e.g., sort records with FP numbers using integer compares
 - Wanted bigger (integer) exponent field to represent bigger numbers
 - 2’s complement poses a problem (because negative numbers look bigger)
 - Use just magnitude and offset by half the range

IEEE 754 Floating Point Standard (3/3)

- Called **Biased Notation**, where bias is number subtracted to get final number
 - IEEE 754 uses bias of 127 for single prec.
 - Subtract 127 from Exponent field to get actual value for exponent

- **Summary (single precision):**



- $(-1)^S \times (1 + \text{Significand}) \times 2^{(\text{Exponent}-127)}$

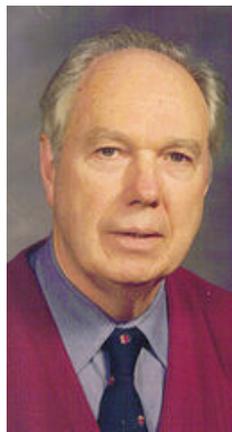
- Double precision identical, except with exponent bias of 1023 (half, quad similar)

Flynn* Taxonomy, 1966

		Data Streams	
		Single	Multiple
Instruction Streams	Single	SISD: Intel Pentium 4	SIMD: SSE instructions of x86
	Multiple	MISD: No examples today	MIMD: Intel Xeon e5345 (Clovertown)

- Since about 2013, SIMD and MIMD most common parallelism in architectures – usually both in same system!
- Most common parallel processing programming style: Single Program Multiple Data (“SPMD”)
 - Single program that runs on all processors of a MIMD
 - Cross-processor execution coordination using synchronization primitives
- SIMD (aka hw-level *data parallelism*): specialized function units, for handling lock-step calculations involving arrays
 - Scientific computing, signal processing, multimedia (audio/video processing)

*Prof. Michael Flynn, Stanford



Big Idea: Amdahl's Law

$$\text{Speedup} = \frac{1}{(1 - F) + \frac{F}{S}}$$

Non-speed-up part \rightarrow $(1 - F)$ $\frac{F}{S}$ \leftarrow Speed-up part

Example: the execution time of half of the program can be accelerated by a factor of 2.

What is the program speed-up overall?

$$\frac{1}{\frac{0.5 + 0.5}{2}} = \frac{1}{0.5 + 0.25} = 1.33$$

Strong and Weak Scaling

- To get good speedup on a parallel processor while keeping the problem size fixed is harder than getting good speedup by increasing the size of the problem.
 - *Strong scaling*: when speedup can be achieved on a parallel processor without increasing the size of the problem
 - *Weak scaling*: when speedup is achieved on a parallel processor by increasing the size of the problem proportionally to the increase in the number of processors
- **Load balancing** is another important factor: every processor doing same amount of work
 - Just one unit with twice the load of others cuts speedup almost in half

Data Level Parallelism

- Loop Unrolling
- Intel SSE SIMD Instructions
 - Exploit data-level parallelism in loops
 - One instruction fetch that operates on multiple operands simultaneously
 - 128-bit XMM registers
- SSE Instructions in C
 - Embed the SSE machine instructions directly into C programs through use of intrinsics
 - Achieve efficiency beyond that of optimizing compiler