

211 Review: Storing Data in Memory

How much data can be stored in 1 mem. address?

→ 1 byte ; aka 8 bits

What if you want to store an integer?

→ The top row is "byte address 0"
 → An int is size = 4 bytes, so it will span multiple addresses in memory!

→ 4 bytes → 4 addresses in memory.

→ Ex: $329,421 = 0b000b\ 0000\ 0000\ 0101\ 0000\ 0110\ 1100\ 1101$

→ The leftmost byte of a number, e.g. $0000\ 0000$

Byte Address	Value
0	0b00000101 (5)
1	0b00000100 (4)
2	0b11111110 (-2)

What is the Most Significant Byte (MSB)?

→ The leftmost byte of a number, e.g. $0000\ 0000$

Little Endian

What is the Least Significant Byte (LSB)?

→ The rightmost byte, e.g. $1100\ 1101$

Byte Address	Value
0	1100 1101
1	0000 0110
2	0000 0101
3	0000 0000

What is Little Endian Ordering?

→ Where the LSB of an int goes in the first byte address

→ AKA, fill the table from "right to left";

- the last 2 nibbles of num → 1st address
- nibbles 5 & 6 → 2nd address
- nibbles 2 & 2 → 4th address

Big Endian

What is Big Endian Ordering?

→ Where the MSB of an int goes in the largest byte address

→ Little Endian

Byte Address	Value
0	0000 0000
1	0000 0101
2	0000 0110
3	1100 1101

Which do we use in this class?

→ $(\text{Index}[\text{element}] \times \text{size of (element type)}) + \text{base_address}$

How do you compute the address of an element in an array?

→ Ex: An int array with 10 elements

- Each element requires 4 bytes/addresses
- element 0 → address $0 \times 4 = \text{addr } 0$ (or whatever the base addr. is)
- element 8 → address $8 \times 4 = \text{addr } 32$

Diagrams to visualize arrays?

→ Ex: $\text{arr} = [5, 90, 100, 40]$

1 row = 1 byte

Addr.	Data
0	0000 0101
1	0000 0000
2	0000 0000
3	0000 0000
4	0101 1010
5	0000 0000
6	0000 0000
7	0000 0000
8	0110 0100
9	0000 0000
10	0000 0000
11	0000 0000
12	0010 1000
13	0000 0000
14	0000 0000
15	0000 0000

1 row = 4 bytes

Addr.	Data
0x00000000	5 0000 0000 0000 0000 0000 0000 0101
0x00000004	90 0000 0000 0000 0000 0000 0101 1010
0x00000008	100 0000 0000 0000 0000 0000 0110 0100
0x0000000C	40 0000 0000 0000 0000 0000 0010 1000
0x00000010	
0x00000014	
0x00000018	
0x0000001C	
0x00000020	

How do we know how many bits are needed to address our memory?

→ It depends on the size of our total memory!

→ $\text{bits_per_addr} = \log_2(\text{size of memory in bytes})$

→ Ex:

• Memory = 16 bytes ... each address is $\log_2(16) = 4$ bits (eg, addr 1 = 0010)

• Memory = 4 GB ...

• 1 GB = 2^{30} bytes

• each address is $\log_2(4 \cdot 2^{30}) = 32$ bits

• eg, addr 1 is 0000 0000 0000 0000 0000 0000 0000 0010.

Digital Logic

What is a logic gate?

- A device that performs a boolean function to produce a single binary output.
- Acts as a building block for digital circuits.
- Several types of logic gates, each of which represent a single logical operator.

What is an "Inverter" or "Not" Gate?

- Represents NOT logical operator
- Can be represented by a symbol, truth table, or equation:

Symbol

$$A \rightarrow \neg \rightarrow Y$$

Equation

$$Y = \bar{A}$$

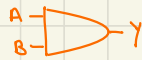
Truth Table

A	Y
0	1
1	0

Gate

What is an AND Gate?

Symbol



Equation

$$Y = AB$$

or

$$Y = A \cdot B$$

Truth Table

A	B	Y
0	0	0
0	1	0
1	0	0
1	1	1

What is an OR Gate?



$$Y = A + B$$

A	B	Y
0	0	0
0	1	1
1	0	1
1	1	1

What is an XOR Gate?

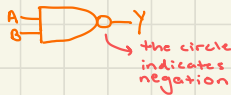


$$Y = A \oplus B$$

A	B	Y
0	0	0
0	1	1
1	0	1
1	1	0

- Y is true when either A OR B is true, but not both. "XOR" = exclusive OR

What is a NAND Gate?



$$Y = \overline{A \cdot B}$$

or

$$Y = \overline{AB}$$

A	B	Y
0	0	1
0	1	1
1	0	1
1	1	0

- Y is true when (A and B) is false ... just negate the result of AB.

What is a NOR Gate?



$$Y = \overline{A + B}$$

A	B	Y
0	0	1
0	1	0
1	0	0
1	1	0

- Negate the result of (A OR B)

What is a XNOR Gate?



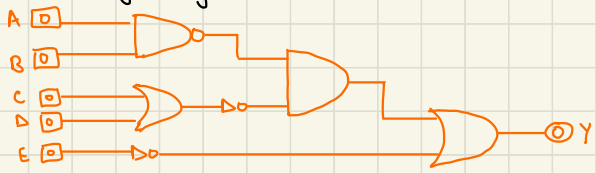
$$Y = \overline{A \oplus B}$$

A	B	Y
0	0	1
0	1	0
1	0	0
1	1	1

- Negate the result of (A XOR B)

Example: How do you form an equation from a logic diagram?

Logic Diagram

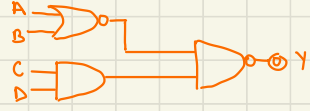


Equation

$$(\overline{A \cdot B} \cdot \overline{C + D}) + \overline{E} = (\overline{A \cdot B})(\overline{C + D}) + \overline{E}$$

How do you do the reverse?

→ Equation: $Y = \overline{(A+B) \cdot C} \cdot D$



How do you form an equation from a truth table?

→ Go through each line of the table where $Y=1$ & create an equation that describes it using AND statements.

→ Then, your final equation is the addition (aka OR) of all those equations

• "sum of products" form

Example:

A	B	C	Y
0	0	0	1
0	0	1	0
0	1	0	1
0	1	1	0
1	0	0	0
1	0	1	1
1	1	0	0
1	1	1	0

→ $\overline{A} \overline{B} \overline{C} = Y$

→ $\overline{A} B \overline{C} = Y$

→ $A \overline{B} C = Y$

$$Y = \overline{A} \overline{B} \overline{C} + \overline{A} B \overline{C} + A \overline{B} C$$

Transistors

- Basic Circuits -

What is voltage?

→ The force that makes currents flow!

→ Measured in Volts (V)

What is current?

→ The rate of flow of electrons.

→ Measured in Ampères (A)

Analogy to understand voltage & current?

→ A Dam with water & the top and a reservoir at the bottom.

→ **Voltage** \approx The desire of the water to flow downhill.

• NOT the actual movement of the water. **Voltage itself is static.**

→ **Current** \approx The actual flow of the water downstream.

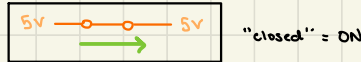
• The size of the opening affects the amount of water (current) that can flow.

• The wider the "opening", the higher the current

What is a closed circuit?

→ A circuit that is fully connected & allows electricity to flow uninterrupted.

→ When a switch (and circuit) is closed, we know the voltage on both ends of the circuit.

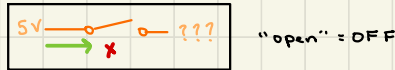


What is an open circuit?

→ A circuit that contains a broken connection.

→ Electricity stops flowing @ point where connection was lost.


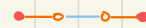
→ When a switch is open, we don't necessarily know the voltage at the end of the circuit - We'd need to see the full circuit.



What are switches?

→ Controlled by physical contact. **SYMBOL:**

What is the symbol?

→  or 

- Transistors -

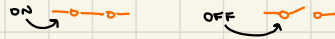
What are transistors?

→ Like switches, but **controlled by a voltage, rather than physical contact.**

What "states" do transistors operate in?

→ Just 2 states - ON or OFF ... its binary.

→ **This is why all machines communicate in binary (1s and 0s)!**



Wait... I thought you said transistors are controlled by voltage?

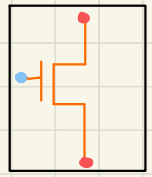
→ They are, but we can use binary terms rather than exact voltage levels for discussing them.

→ **high voltage** = "logic high" = logic 1 = **1**

→ **low voltage** = "logic low" = logic 0 = **0**

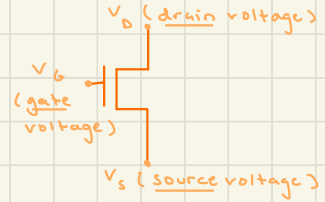
What is the symbol for transistors?

→ The terminal labeled with the **blue** dot determines whether current can flow between the terminals labeled with the **red** dot



What are the terminals in an **NMOS** transistor?

1. **Gate**: controls whether transistor is on or off
2. **Source**: endpoint
3. **Drain**: endpoint

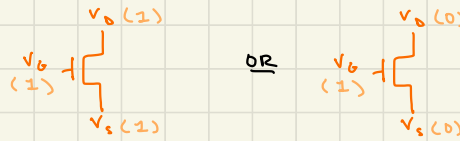


What does it mean if the **gate voltage is high**?

→ a.k.a. "logic 1"
→ Current can flow between the source and the drain... they will be connected & will have the same voltage.

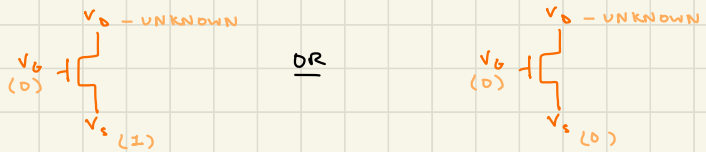
Example?

→ e.g., if $V_S = \text{logic } 0$, V_D also = $\text{logic } 0$. And vice versa.



What does it mean if the **gate voltage is low**?

→ a.k.a. "logic 0"
→ Source & drain are not connected; current cannot flow.
→ The voltage of the drain will be unknown.



What is a **PMOS** transistor?

→ Same 3 terminals as with NMOS transistor.

→ The **gate voltages are reversed**:

- Gate voltage LOW (logic=0) = current can flow
- Gate voltage HIGH (logic=1) = current cannot flow

How do we define the **behavior** of a transistor?

- | | |
|---|---|
| <p><u>nmos</u></p> <ul style="list-style-type: none">• behaves as an open switch when V_G is low.• behaves as a closed switch when V_G is high. | <p><u>pmos</u></p> <ul style="list-style-type: none">• behaves as an open switch when V_G is high.• behaves as a closed switch when V_G is low. |
|---|---|

What is Moore's Law?

→ A trend describing the reduction in transistor size in machines/computers

→ Moore's Law: The # of transistors in a given area on a chip doubles every 2 years

• ~2,000 transistors in 1970 → > 10 billion today!

→ The trend is coming to an end now, as the industry develops alternative technologies to continue improving chip performance.

Building Logic Gates with Transistors & CMOS

RECALL: EX of a logic gate?

→ Inverter gate: $A \rightarrow Y$ ($Y = \bar{A}$)

→ See pg 3 of notes.

What is power?

→ A component (of a circuit?) that produces a logic high (1) value

→ Symbol: \top

→ Always connected to output via a pmos transistor.

• RECALL: pmos transistor is closed ON when input = 0

→ A component that produces a logic low (0) value.

→ Symbol: \downarrow

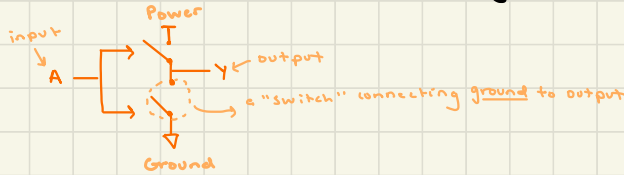
→ Always connected to output via an nmos transistor.

• RECALL: nmos transistor is closed ON when input = 1

What is ground?

→ RECALL: Switches are a way to physically control the opening & closing of circuits, while transistors operate on logic values.

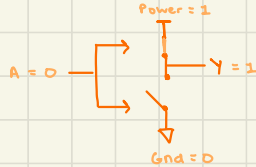
Visual example of an Inverter with switches?



How does this inverter work when input = 0?

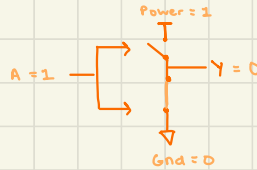
→ When input = 0, we manually close the top/power switch.

→ When input = 1, we manually close the bottom/ground switch.



→ Since power produces logic = 1 & power is connected to output, output = 1

How does it work when input = 1?

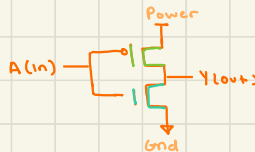


→ Since ground produces logic = 0 & ground is connected to output, output = 0.

How would we build this inverter with transistors?

→ Replace the power switch with a pmos transistor. Replace ground switch w/ an nmos transistor.

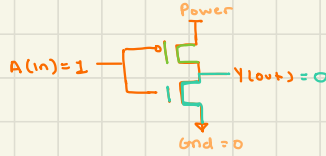
→ Then, invert the output of the power transistor (with the \circ symbol)



How does this inverter work when input = 1?

→ If input = 1, pmos tr. is OPEN, so power is not connected to output.

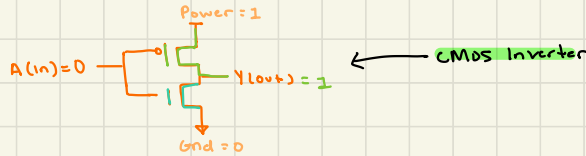
→ If input = 1, nmos tr. is CLOSED, so ground is connected to output.



How does this inverter work when input = 0?

→ If input = 0, pmos tr. is CLOSED, so power is connected to output.

→ If input = 0, nmos tr. is OPEN, so ground is not connected to output.



What is a CMOS?

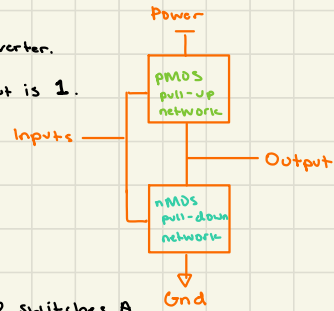
→ CMOS = "Complementary MOS" ... the pMOS & nMOS transistors complement each other to form the logic gate.

→ The inverter described above is a CMOS Inverter.

What is the pull-up network?

→ What "drives" the output whenever the output is 1.

→ ALWAYS associated with PMOS transistor



What is the pull-down network?

→ What "drives" the output whenever output is 0.

→ ALWAYS associated with nMOS transistor.

→ Only one network will be on at a time.

What are "switches in series"?

→ Think of them as an AND function. When 2 switches A and B are connected in series, current will only flow if A AND B are on/closed

→ Current = switch A AND switch B

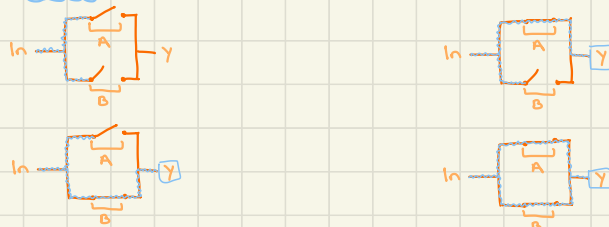
How do they look?



What are switches in parallel?

→ Think of them as an OR function. Current will flow if A OR B is closed.

→ Current = switch A OR switch B



How do we build a CMOS for the NAND logic gate?

→ RECALL: NAND $Y = \overline{A \times B} / Y = \overline{AB}$

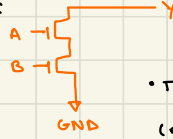
A	B	Y
0	0	1
0	1	1
1	0	1
1	1	0

→ When AB is true, we want output = 0, so

we can build "switches in series" (aka

A AND B) and connect them to GROUND, since GROUND always sends

output of 0:



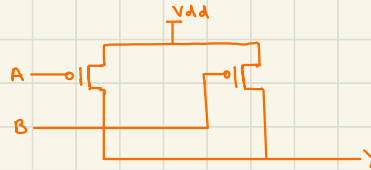
• This part of the CMOS will result in something (namely, logic = 0) being sent to Y whenever (AB) is true. $A=1, B=1 \Rightarrow \text{Ground}=1 \Rightarrow Y=0$

What do we build for when A and B aren't both logic=1?

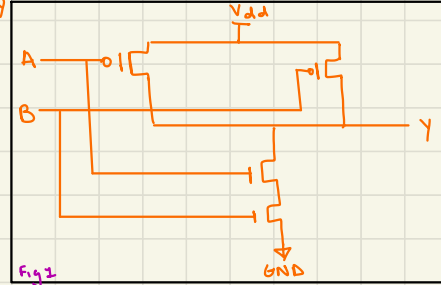
→ When $A=0$ OR $B=0$, we want output = 1. So we can build "switches in parallel" (aka A OR B), but invert the values being sent to the A & B transistors.

• if $A=0$ and $B=0$, send 1 and 1 to the 2 switches. Since at least one is 1, output of 1 is received by power. $A=0, B=0 \Rightarrow \text{Power}=1 \Rightarrow Y=1$

• if $A=0$ and $B=1$, send 1 and 0 to the 2 switches. Since at least one is 1, output of 1 is received by power. $A=0, B=1 \Rightarrow \text{Power}=1 \Rightarrow Y=1$



How do we combine these 2 pieces to form the CMOS NAND?



How do we build a CMOS for the NOR logic gate?

→ RECALL: NOR $Y = \overline{A+B}$

A	B	Y
0	0	1
0	1	0
1	0	0
1	1	0

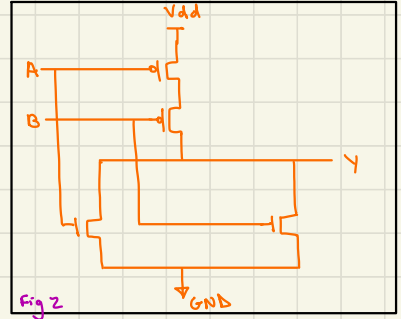
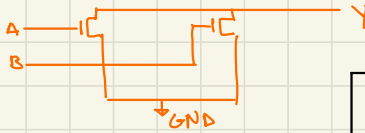
• Lets build "switches in series" & invert the inputs so that result = 1 if and only if $A=0$ and $B=0$

• Since we actually want the output to be 1 (output = result) when $A=0$ and $B=0$, we should connect these switches to power.



A	B	Y
0	0	1
0	1	0
1	0	0
1	1	0

- We want output = 0 whenever $A+B$ (A OR B is 1) = TRUE
- To do this, let's build switches in parallel so that whenever $A=1$ or $B=1$, result = 1
- Since we actually want output = 0 when the above "result" = 1, we should connect these switches to ground.



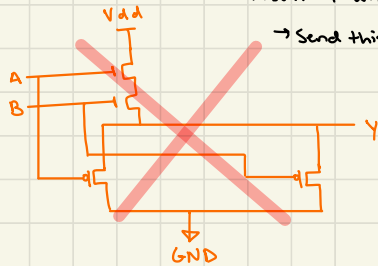
How do we combine these 2 pieces to form the CMOS NOR?

How do we build a CMOS for the AND logic gate?

→ RECALL: AND = $Y = AB$

A	B	Y
0	0	0
0	1	0
1	0	0
1	1	1

- When $A \& B$ is true, send output to POWER; use switches in series.
- Use switches in parallel; send inverted inputs so that the result = 1 whenever these are true



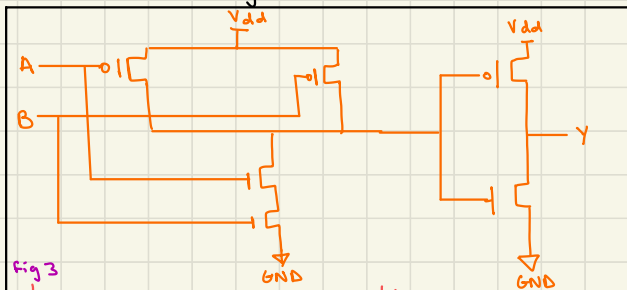
But actually, this does NOT work!

→ Send this output to GROUND

Why aren't the diagram & method above correct for CMOS "AND"?

→ because we can't use nmos transistors in the pull-up network?

→ SOLUTION: Use a NAND gate & add an inverter:



1) When the NAND gate sends an output of 0 through GND, Vdd receives $\bar{0} = 1$, & thus $Y = 1$

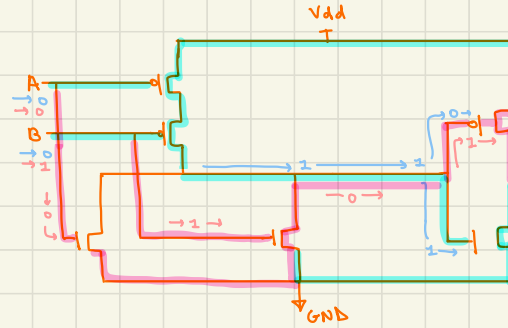
2) When NAND sends output = 1 through Vdd, GND receives output 1, & thus sends output = 0 to Y.

NAND Gate

Inverter to invert $\bar{\bar{NAND}} = \text{AND}$

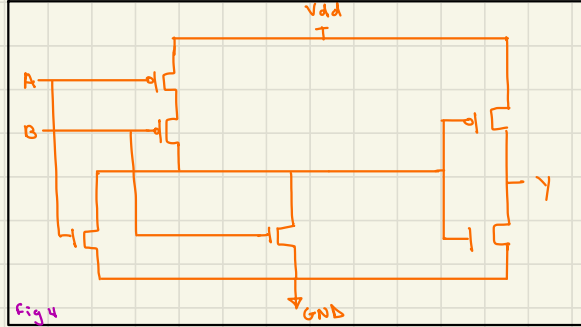
How do we build a CMOS for the OR gate?

→ Similar concept: A NOR Gate + an inverter:



1) When $A=0$ & $B=0$, V_{dd} sends output = 1 to inverter. Inverter sends output = 1 to GND and op = 0 to V_{dd} . Thus, GND sends output 0 to Y

2) When $A=0$ and $B=1$, GND sends output = 0 to inverter. Inverter sends output = $\bar{0} = 1$ to V_{dd} , and output = 0 to GND. Thus, V_{dd} sends output 1 to Y.



How many transistors are required to build each type of gate?

Gate	# of transistors
NOT	2 (-DO-)
AND	6 (Fig. 3)
OR	6 (Fig. 4)
NAND	4 (Fig. 1)
NOR	4 (Fig. 2)
XOR	12
XNOR	12

Why do AND and OR require 6 transistors?

→ Note that AND and OR gates require

$$4 \text{ (NAND/NOR gate)} + 2 \text{ (Inverter)} = \underline{6} \text{ transistors!}$$

Digital Logic Pt. 2

Why do we want to minimize transistor count when building logic gates?

→ Smaller amt. of transistors Reduces:

- Delay from input to output
- The area of the circuit taken up ; this is good b/c we can pack more logic (i.e. transistors) in a given space.
- Power consumption
- The cost of the circuit

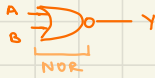
What is a logically equivalent circuit to an inverter ?

→ The "bubble" on an AND or OR symbol :



6 (OR) + 2 (NOT) = 8 transistors

is equivalent to



4 (NOR) transistors

But obviously, one uses much less transistors than the other!

What is a logically equivalent circuit for NAND ?



is equivalent to



This is the better, conventional way to draw NAND

What about NOR ?



is equivalent to

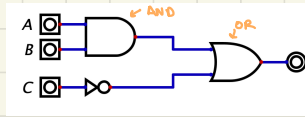


This is the better, conventional way to draw NOR

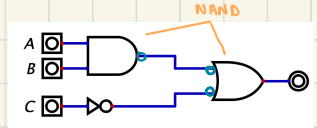
How do we minimize transistor count on logic gate diagrams?

1. Replace AND and OR gates with NAND gates, or NOR gates.
2. Add bubbles to make the circuit logically equivalent.
3. Cancel out the free bubbles
4. Draw the NAND and NOR gates in their conventional form.

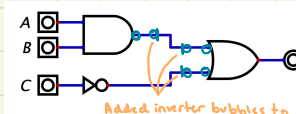
Example?



1)

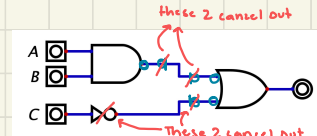


2)

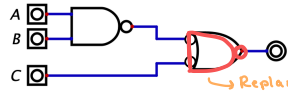


Added inverter bubbles to make the circuit equivalent after replacing the AND and OR with NANDs

3)



4)



→ Replace with

- Gates with more than 2 inputs -

What does the symbol for a 3-input AND gate look like?



OR (logical equivalent)



Symbol for 3-input OR gate?



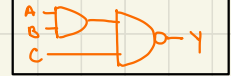
OR (logical equivalent)



Symbol for 3-input NAND gate?



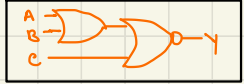
OR (logical equivalent)



Symbol for 3-input NOR gate?



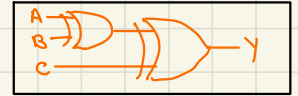
OR (logical equivalent)



Symbol for 3-input XOR gate?



OR (logical equivalent)



Symbol for 3-input XNOR gate?



OR (logical equivalent)



The Laws of Boolean Algebra

What do we use boolean algebra for?

→ To create logically equivalent circuits that use fewer transistors ←
(MOTIVATION: Recall notes on Digital Logic pt. 2)

What are boolean algebra laws?

→ Laws about boolean statements where:
• Letters (A, B, C, etc.) represent literals • "+" = OR
• "." (algebra multiplication) = AND

→ Every Boolean alg law comes in a pair of statements: the AND version and the OR version.

What is the Identity Law?

→ $A + 0 = A$ (Any literal A or'd with 0 will always output A).
→ $A \cdot 1 = A$ (Any literal A and'd with 1 will always output A).

What is the Null Law?

→ $A + 1 = 1$ (A OR 1 always = 1)
→ $A \cdot 0 = 0$ (A AND 0 always = 0)

What is the Idempotent Law?

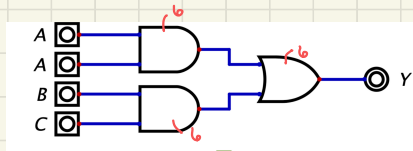
→ $A + A = A$ (A OR A = A)
→ $A \cdot A = A$ (A AND A = A)

What is the Complement Law?

→ $A + \bar{A} = 1$ (A OR NOT A = 1)
→ $A \cdot \bar{A} = 0$ (A AND NOT A = 0)

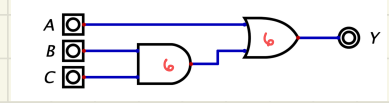
What is an example of optimizing a circuit?

→ The equation $Y = AA + BC$ produces this circuit:
• it has $6 + 6 + 6 = 18$ transistors

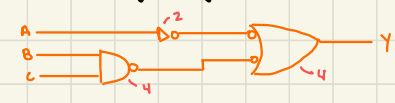


→ Using the idempotent law $AA = A$, we can reduce the equation to $Y = A + BC$, producing the circuit

• it has $6 + 6 = 12$ transistors



→ By converting the gates into NAND gates (RECALL: Digital Logic 2), we can further reduce the circuit to 10 transistors:



What is the Commutative Law?

→ $A + B = B + A$
→ $A \cdot B = B \cdot A$

What is the Associative Law?

→ $(A + B) + C = A + (B + C)$
→ $(A \cdot B) \cdot C = A \cdot (B \cdot C)$

What is the Distributive Law?

→ $A \cdot (B + C) = (A \cdot B) + (A \cdot C)$
↳ (A AND (B OR C)) = (A AND B) OR (A AND C)
→ $A + (B \cdot C) = (A + B) \cdot (A + C)$
↳ (A OR (B AND C)) = (A OR B) AND (A OR C)

Proof for the distributive law (OR version) ?

→ Lets use the other laws to prove the truth of $A + (B \cdot C) = (A+B) \cdot (A+C)$

$$\begin{aligned} 1. \text{ Distribute out the terms: } & A + (BC) = AA + AC + AB + BC \\ & \downarrow \\ & = A + AC + AB + BC \\ & = A(1 + C + B) + BC \\ & = A(2) + BC \\ & = A + BC \end{aligned}$$

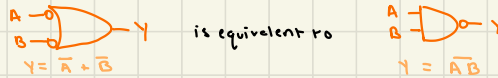
What is De Morgan's Theorem?

→ Used to simplify large NOT bars (e.g. $\overline{A+B}$)

$$\rightarrow \overline{A+B} = \bar{A} \cdot \bar{B} \quad (\text{NOT}(A \text{ OR } B) = \text{NOT}(A) \text{ AND } \text{NOT}(B))$$

$$\rightarrow \overline{AB} = \bar{A} + \bar{B} \quad (\text{NOT}(A \text{ AND } B) = \text{NOT}(A) \text{ OR } \text{NOT}(B))$$

→ This theorem is actually how we created the "conventional" NAND gate:



→ De Morgan's: flipping the inputs ($\overline{AB} \rightarrow \bar{A} \cdot \bar{B}$), operations ($AB \rightarrow A+B$), and the outputs ($A+B \rightarrow \bar{A} + \bar{B}$)... Thus $\overline{AB} \rightarrow \bar{A} + \bar{B}$

When would we need to use these laws for simplification?

→ When converting a truth table into a boolean equation!

1. Create an equation in "sum of products" form (RECALL)
2. Use the laws to simplify!

Karnaugh Maps (K Maps)

What are K Maps?

→ A tool used to simplify boolean expressions

• An alternate tool to using boolean algebra laws.

→ K-Maps provide a graphical method to find simplified boolean expressions.

→ Benefit: If used properly, K Maps always derive the most simplified equation!

Unlike boolean alg, where you aren't always 100% sure if you've maximized simplification.

How do K Maps work?

→ A diagram/grid that makes it easy to identify combinable terms.

→ Each box in the K-map corresponds to a single row/entry in the truth table.

Example of a K Map for

a 3-Input Function?

→ Given this

truth table:

A	B	C	Y
0	0	0	0
0	0	1	1
0	1	0	0
0	1	1	1
1	0	0	0
1	0	1	0
1	1	0	0
1	1	1	0

A		BC			
		00	01	11	10
0	0	R ₀	R ₁	R ₃	R ₂
1	0	R ₄	R ₅	R ₇	R ₆

Because in Row 2, A=0, B=1, C=0

→ The "0" & "1" on y-axis indicate values of A

→ The "00", "01", "11", "10" are values of B & C, respectively.

B=0, C=0 B=0, C=1 B=1, C=0 B=1, C=1

How do we fill in the K-Map?

→ With the Y (output) of each row!

A		BC			
		00	01	11	10
0	0	0	1	1	0
1	0	0	0	0	0

What is gray code?

→ The ordering that we used for the B and C values in the ex above.

→ In Gray code, successive values differ by only 1 bit - unlike binary:

Gray Code: 00, 01, 11, 10

Binary: 00, 01, 10, 11

00
01
11
10

00
01
10
11

eg L, R, above, or below

Why is gray code useful for K-Maps?

→ Each square in the K-Map differs from its adjacent squares by one literal:

A		BC			
		00	01	11	10
0	0	0	1	1	0
1	0	0	0	0	0

• Entry 0 & Entry 1 both share " $\bar{A}\bar{B}$ ", and only differ by the value of C literal

• Entries 1 & 3 both differ only in their B literal

So how do you use a K-Map to find simplified Boolean equations?

Example of using a K-Map?

1. Make the K-Map "table", & fill in all squares (according to truth table)

	BC			
A	00	01	11	10
0	0	1	1	0
1	0	0	0	0

2. Group the terms together -- all terms where the entry is 1 must get grouped.

• We'll learn later about how to create the groups.

	BC			
A	00	01	11	10
0	0	1	1	0
1	0	0	0	0

• "Grouping" Entries 1 & 3

3. Find the literals that each group has in common, and "OR" them in the final eq.

• Entry 1: $\bar{A}\bar{B}C$

• Entry 3: $\bar{A}BC$

• Both entries share \bar{A} and C ! ANS: $Y = \bar{A}C$

→ EX:

	BC			
A	00	01	11	10
0	0	0	0	0
1	1	1	1	0

1. For each group, analyze each entry in the group. Extract the terms that all entries have in common (there should be 2?). AND these terms together

• Left group: $A=1, B=0 \rightarrow A\bar{B}$

• Right group: $A=1, C=1 \rightarrow AC$

2. OR the terms from step 1 together to create the final equation.

$$Y = A\bar{B} + AC$$

What are the Grouping Rules for K-Maps?

1. All squares in each group must contain only 1s, and every cell containing a 1 must be in at least one group.

	B		
A	0	1	1
0	0	0	0
1	1	1	1

	B		
A	0	1	1
0	0	0	0
1	1	1	1

2. Groups may be horizontal or vertical, but NOT diagonal.

3. Groups must contain 2^n cells, where $n=0, 1, 2$, etc. aka, must be a power of 2. e.g. 1, 2, 4, 8, ... cells.

	BC			
A	00	01	11	10
0	1	1	0	0
1	1	1	1	0

	BC			
A	00	01	11	10
0	0	0	0	0
1	1	1	1	0

4. Each group must be as large as possible.

	BC	00	01	11	10
A	0	1	1	1	1
	1	0	0	1	1

✓

	BC	00	01	11	10
A	0	1	1	1	1
	1	0	0	1	1

X

5. Groups may wrap around the table.

	BC	00	01	11	10
A	0	1	0	0	1
	1	1	0	0	1

Here, the group is $\bar{A}\bar{B}\bar{C}$, $A\bar{B}\bar{C}$, $\bar{A}B\bar{C}$, and $AB\bar{C}$
 $(Y = \bar{C})$

6. When multiple groups, every group should contain at least one unique 1 (aka unique entry).

	BC	00	01	11	10
A	0	1	1	1	1
	1	0	0	1	1

✓

	BC	00	01	11	10
A	0	1	1	1	1
	1	0	0	1	1

X

the bottom group is redundant

How do we create a K-Map from a truth table with 4 variables?

→ Truth table:

Row #	A	B	C	D	Y
0	0	0	0	0	0
1	0	0	0	1	0
2	0	0	1	0	1
3	0	0	1	1	1
4	0	1	0	0	1
5	0	1	0	1	1
6	0	1	1	0	0
7	0	1	1	1	0
8	1	0	0	0	0
9	1	0	0	1	0
10	1	0	1	0	1
11	1	0	1	1	0
12	1	1	0	0	1
13	1	1	0	1	1
14	1	1	1	0	0
15	1	1	1	1	0

→ K-Map (Orange numbers indicate corresponding Row):

	CD	00	01	11	10
AB	00	0	1	3	2
	01	4	5	7	6
	11	12	13	15	14
	10	8	9	11	10

using Gray code order so that adjacent cells differ by only 1 literal

this cell is 0110 ... therefore it corresponds to Row 6.

Example of reading a 4-variable K-Map?

	CD	00	01	11	10
AB	00	1	0	0	0
	01	1	1	1	0
	11	1	1	1	1
	10	1	0	0	0

How do you put a boolean equation in Sum of Products (SOP) Form?

How do we simplify a boolean equation using a K-Map?

→ EXAMPLE: $Y = ABC + AB\bar{C} + \bar{A}CB + ACB + \overline{AB\bar{C}\bar{D}} + \overline{\bar{A}B + \bar{C}}$

→ Step 1: Put the equation in SOP form.

$$ABC + AB\bar{C} + \bar{A}CB + ACB + \overline{AB\bar{C}\bar{D}} + \overline{\bar{A}B + \bar{C}}$$

$$\rightarrow ABC + AB\bar{C} + \bar{A}B + \bar{C}B + ACB + AB + C + D + AC + \bar{B}C$$

→ Step 2: Fill out the K-Map.

$$Y = ABC + AB\bar{C} + \bar{A}B + AC + ACB + AB + C + D + AC + \bar{B}C$$

	CD			
AB	00	01	11	10
00	0	1	1	1
01	1	1	1	1
11	1	1	1	1
10	0	1	1	1

→ each term in the SOP-form equation can tell you where to place a 1 in the K-Map.
 e.g. $AB\bar{C} \rightarrow A=1, B=1, C=0$ then $Y=1$... entry 15 has a 1

→ Step 3: Group the entries & find the simplified equation!

	CD			
AB	00	01	11	10
00	0	1	1	1
01	1	1	1	1
11	1	1	1	1
10	0	1	1	1

- Group 1: $D = 1$
- Group 2: $C = 1$
- Group 3: $B = 1$

$ANS: Y = B + C + D$

Quiz 0 Review

Logic Gates

Name	Equation	Symbol	Transistor count
AND	$Y = AB$		6
OR	$Y = A + B$		6
NOT	$Y = \bar{A}$		2
NAND	$Y = \overline{AB}$	OR	4
NOR	$Y = \overline{A+B}$	OR	4
XOR	$Y = A \oplus B$		12
XNOR	$Y = \overline{A \oplus B}$		12

211 Review

Unit	bits	example
byte	8	00001101
nibble	4	1101
word	16 or 32?	

- With n bits,
 - you can represent 2^n distinct values.
 - 2's complement most negative num: $-(2^{n-1})$
- LARGEST NUM: $2^n - 1$
- 2's complement most positive num: $(2^{n-1}) - 1$

→ LSB = RIGHTMOST 8 bits. Stored in the 1st mem address (little Endian)

→ MSB = LEFTMOST 8 bits. Stored in the last mem address

$$42,185 = \underbrace{060000 \ 0000 \ 0000 \ 0000}_{\text{MSB}} \ 1010 \ 0100 \ \underbrace{1100 \ 1001}_{\text{LSB}}$$

→ Bits needed to address memory = \log_2 (size in bytes of total memory)

Address	Value
0	1100 1001
1	1010 0100
2	0000 0000
3	0000 0000

Quiz 0 Review

211 Review

Decimal	0	1	2	3	4	5	6	7	8	9	10	11
Binary	0000	0001	0010	0011	0100	0101	0110	0111	1000	1001	1010	1011
Hex	0x0	0x1	0x2	0x3	0x4	0x5	0x6	0x7	0x8	0x9	0xA	0xB

Decimal	12	13	14	15	16	17	18	19
Binary	1100	1101	1110	1111	10000	10001	10010	10011
Hex	0xC	0xD	0xE	0xF	0x10	0x11	0x12	0x13

→ Converting decimal to hex:

1. Fill out this chart $\overline{\quad} \overline{\quad} \overline{\quad} \overline{\quad}$, where every blank can have a max value of F (aka 15)

Ex 50 → 0 0 3 2 because $3(16) + 2(1) = 50$

→ Finding address of element in array: $(\text{index} [i] \times \text{Size of element type}) + \text{base addr}$

Ex • Element 10 in int array; base addr = $0x0000300C$

1. $10 \times \text{sizeof(int)} = 4 = 40$

2. Add decimal num from step 1 to decimal version of last relevant digits of base address:

• $0C \rightarrow 12$

• $40 + 12 = 52$

3. Convert ans to decimal: $52 \rightarrow 16(3) + 1(4) \rightarrow 0x34$

4. Add result to base addr (replacing last digits w/ result):

$$\begin{array}{r} 0x00003000 \\ + 0x00000034 \\ \hline 0x00003034 \quad \checkmark \end{array}$$

Transistors

→ nMOS: $V_g \text{ low} = \text{open}$
 $V_g \text{ high} = \text{closed} = \text{ON}$

→ pMOS: $V_g \text{ low} = \text{closed} = \text{ON}$
 $V_g \text{ high} = \text{open}$

nMOS → "normal" MOS → the normal thinking would be that "high" power results in it flowing through. Therefore, in nMOS, $V_g \text{ high} \Rightarrow \text{circuit closed} \Rightarrow \text{power on}$.

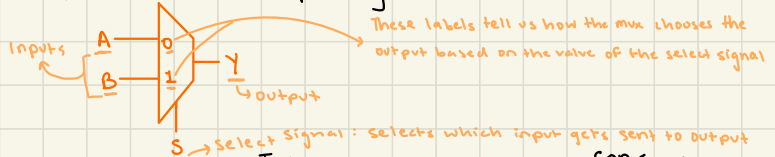
Multiplexers

What is a 2:1 multiplexer?
(2:1 mux)

→ A circuit that chooses an output from among 2 inputs, based on the value of a select signal.

→ Multiplexers are useful for implementing if-statements.

What is the 2:1 mux symbol?



Example of a regular schematic for an if-statement?

Statement

if (S == 0) : Y = A
else if (S == 1) : Y = B

Truth Table

S	A	B	Y
0	0	0	0
0	0	1	0
0	1	0	1
0	1	1	1
1	0	0	0
1	0	1	1
1	1	0	0
1	1	1	1

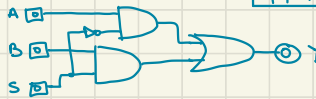
SOP Equation

$$Y = \bar{S}A\bar{B} + \bar{S}AB + S\bar{A}B + SAB$$

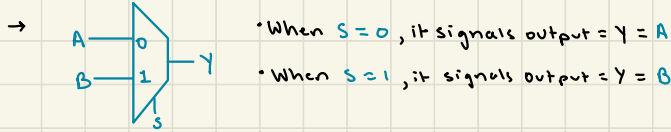
$$= A\bar{S}(\bar{B} + B) + BS(\bar{A} + A)$$

$$= A\bar{S} + BS$$

Diagram



How could we replace this diagram with a 2:1 mux?



→ This mux diagram replaces/avoids having to draw that bigger more complicated gate diagram.

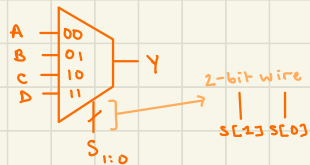
What is a 4:1 multiplexer?
(4:1 mux)

→ A circuit that chooses an output from among 4 inputs, based on the value of a select signal.

→ Since we need 4 signaling options, S must be a 2-bit wire that can send a 2-bit signal ($n = 2, 2^2 = 4$ "bit combinations")

How does the 4:1 mux diagram look?

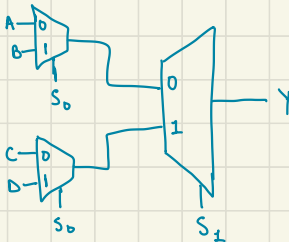
→ EX: For the Truth Table:



S ₁	S ₀	Y
0	0	A
0	1	B
1	0	C
1	1	D

- When S = 0000, Y = A
- When S = 0001, Y = B
- When S = 0010, Y = C
- When S = 0011, Y = D

How can you implement the same truth table with a 4:1 mux made out of 2 2:1 muxes?



Definitions in Digital Logic

What is a literal?

→ A single variable. May be complemented.

→ eg A, B, \bar{A}

What is a product term?

→ An AND of individual literals

→ eg $A\bar{B}C$. But NOT $A\bar{BC}$ (because " \bar{BC} " isn't a literal)

What is a minterm?

→ A product term in which all variables appear once.

→ eg $\bar{A}\bar{B}\bar{C}, \bar{A}BC, \bar{A}\bar{B}C, \text{ or } ABC$

→ but NOT $A, \bar{A}C, BC, \text{ etc.}$

How do you derive an equation using minterms?

Truth Table

A	B	C	F
0	0	1	1
0	0	1	0
0	1	0	1
0	1	1	0
1	0	0	0
1	0	1	1
1	1	0	0
1	1	1	1

Minterm	Minterm name
$\bar{A}\bar{B}C$	m_0
$\bar{A}BC$	m_1
$A\bar{B}C$	m_2
ABC	m_3
$A\bar{B}\bar{C}$	m_4
$A\bar{B}C$	m_5
$AB\bar{C}$	m_6
ABC	m_7

1. Write the minterm for each row.

2. Take the sum of the minterms for rows where output = 1

Equation

$$F = \bar{A}\bar{B}C + \bar{A}BC + A\bar{B}C + ABC$$

OR

$$F(A,B,C) = \sum(m_0, m_2, m_5, m_7)$$

What is sum of products?

→ When 2 or more product terms are summed (OR) together.

→ e.g. $Y = AB + A\bar{C}$

→ but NOT $Y = (A+B)(C+A)$

What is Canonical Sum of Products form?

→ An SOP form in which each product contains all literals, e.g.

$$F = \bar{A}\bar{B}\bar{C} + \bar{A}BC + A\bar{B}C + ABC$$

What is Simplified Sum of Products form?

→ When you use boolean algebra to simplify the canonical SOP form, e.g.

$$F = \bar{A}\bar{B}\bar{C} + C(\bar{A}B + A(\bar{B} + B)) \Rightarrow \bar{A}\bar{B}\bar{C} + C(\bar{A}B + A)$$

$$\Rightarrow \bar{A}(\bar{B}\bar{C}) + AC + AB \Rightarrow \boxed{\bar{A}\bar{C} + AC}$$

What is a maxterm?

→ A term in which all variables appear once, as literals OR'd together

→ eg $\bar{A} + B + \bar{C}$

→ To write the maxterm for a truth table row, sum the complement of each literal's value together.

→ Eg:

A	B	C	Y
0	1	0	1

 → $A + \bar{B} + C = Y$

How do you derive an equation using maxterms?

1. Write the maxterm for each row

2. Take the product of the maxterms for rows where output = 0

Truth Table

A	B	C	F	Maxterm	Maxterm name
0	0	0	1	$\bar{A} + \bar{B} + \bar{C}$	M_0
0	0	1	0	$A + B + \bar{C}$	M_1
0	1	0	1	$A + \bar{B} + C$	M_2
0	1	1	0	$A + \bar{B} + \bar{C}$	M_3
1	0	0	0	$\bar{A} + B + C$	M_4
1	0	1	1	$\bar{A} + B + \bar{C}$	M_5
1	1	0	0	$\bar{A} + \bar{B} + C$	M_6
1	1	1	1	$\bar{A} + \bar{B} + \bar{C}$	M_7

Equation

$$F = (A + B + \bar{C})(A + \bar{B} + \bar{C})(\bar{A} + B + C)(\bar{A} + \bar{B} + C)$$

OR

$$F = \prod (M_1, M_3, M_4, M_6)$$

Why does this derived equation work?

→ Creating the canonical SOP equation for \bar{F}

$$\bar{F} = \bar{A}\bar{B}C + \bar{A}BC + A\bar{B}\bar{C} + AB\bar{C}$$

And negating both sides + applying DeMorgans

$$F = \overline{\bar{A}\bar{B}C + \bar{A}BC + A\bar{B}\bar{C} + AB\bar{C}}$$

$$F = (A + B + \bar{C})(A + \bar{B} + \bar{C})(\bar{A} + B + C)(\bar{A} + \bar{B} + C)$$

Actually yields the same term!

K-map Definitions

What is an implicant?

→ Any product term (RECALL: $AB, \bar{B}C$, etc.) whose output for a given Boolean equation is 1.

Example?

→ aka, in a K-Map, the terms that define any group of 1s.

A \ BC		00	01	11	10
		0	0	0	1
1		0	1	1	0

Implicants

- $\bar{A}BC$ BC
- $A\bar{B}C$ AC
- ABC

What is a prime implicant?

→ An implicant that is not a subset of any other implicant.

→ aka, in a K-Map, an implicant that corresponds to a group which can NOT be covered by any other group.

Example?

AB \ CD		00	01	11	10
		00	0	1	1
01		0	1	1	0
11		1	1	1	0
10		1	1	0	0

- red = implicants
- prime implicants: $\bar{C}D, \bar{B}\bar{C}$

What is an essential prime implicant?

→ A prime implicant where:

- At least 1 element is not covered by 1 or more other prime implicants.

→ aka, in a K-Map, a group that is necessary to use in the final solution to cover all 1s.

Example?

AB \ CD		00	01	11	10
		00	1	1	0
01		0	1	1	0
11		1	1	1	0
10		1	1	0	0

- red → essential prime implicants, aka the final groups (in this case)
- orange → non-essential prime implicant.

What is a non-essential prime implicant?

→ A prime implicant that:

- contains NO elements which can't be covered by another prime implicant group

What is the formal procedure for using K-Maps to derive equations?

1. Convert truth table to K-Map.

	CD			
AB	00	01	11	10
00	0	0	1	1
01	1	1	1	1
11	1	1	0	0
10	0	0	1	1

2. Include all essential prime implicants.

$$B\bar{C}, \bar{B}C$$

	CD			
AB	00	01	11	10
00	0	0	1	1
01	1	1	1	1
11	1	1	0	0
10	0	0	1	1

3. Include non-essential prime implicants as needed to cover all 1s.

	CD			
AB	00	01	11	10
00	0	0	1	1
01	1	1	1	1
11	1	1	0	0
10	0	0	1	1

choosing either of the non-essential primes (green or orange) will give us the most simplified equation.

$$Y = B\bar{C} + \bar{B}C + \bar{A}C = B\bar{C} + C(\bar{B} + \bar{A})$$

OR

$$Y = B\bar{C} + \bar{B}C + \bar{A}B = \bar{B}C + B(\bar{A} + \bar{C})$$

Lab 0

What is a tunnel?

Example?

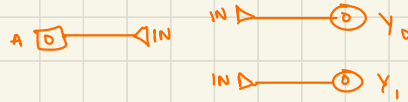
→ A way to draw an "invisible wire" to bind 2 points together

→ Tunnels are grouped by case sensitive labels

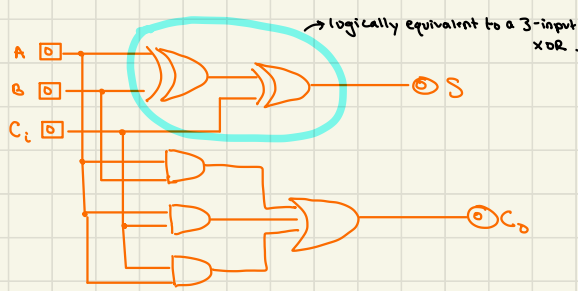
→ Normal circuit:



→ With Tunnels:

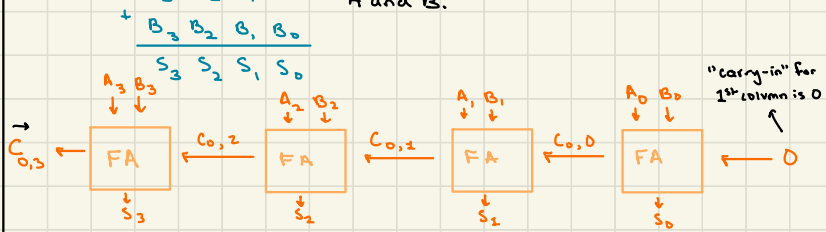


What would the full adder diagram look like?



How can we create a circuit to add 2 4-bit values?
(4-bit Adder)

→ For $C_{0,3}$ $C_{0,2}$ $C_{0,1}$ $C_{0,0}$, we can use 4 Full-adders to add together A and B.



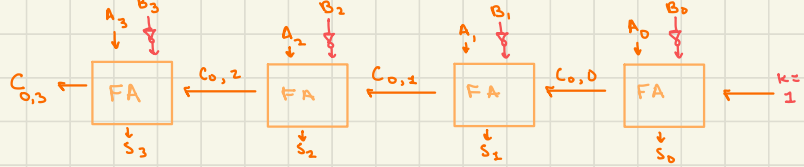
How do we create a circuit for binary subtraction?

→ **RECALL:** with 2's complement, $A - B \approx A + (-B)$. To binary subtract 2 numbers A and B, we have to negate B and then add it to A.

→ Instead of building a separate circuit, we can modify the ripple carry adder (diagram above) s.t. it can perform subtraction AND addition.
• How? By modifying it s.t. it can negate B.

How can we modify the ripple carry adder to create the Subtractor circuit?

→ **RECALL:** to negate a binary num, negate each bit & then add 1
→ Solution: negate each input of B, and send in $k=1$ as the carry-in for column 0:



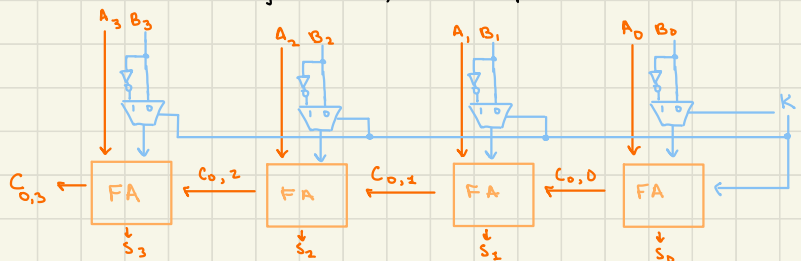
How do we create an adder-subtractor circuit?

→ Goal: Same circuit to perform both addition and subtraction

→ Solution: Let $C_{in,0} = K$, and use a 2:1 mux with K as the signal, to determine whether to send B or \bar{B} to the Full adder!

• if $K=0$: binary addition, 2:1 mux outputs B

• if $K=1$: binary subtraction, 2:1 mux outputs \bar{B}



What is an alternate way to implement Adder-Subtractor?

→ Using XOR instead of a 2:1 Mux.

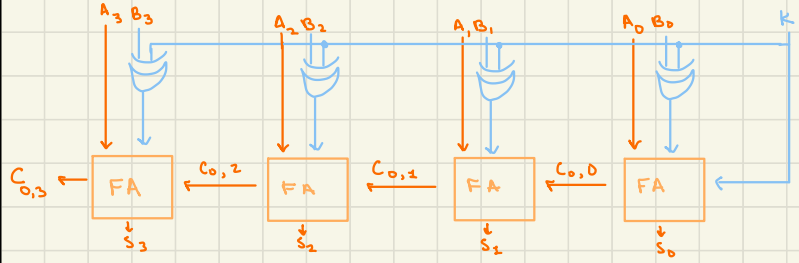
• RECALL: $N \oplus 1 = \bar{N}$ $N \oplus 0 = N$ (boolean rules)

→ If we XOR B_i with K ,

• when $K = 0$ (addition), $B_i \oplus K = B_i$

• when $K = 1$ (subtraction), $B_i \oplus K = \bar{B}_i$

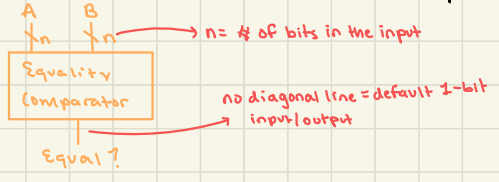
How would this diagram look?



Comparator and Shifter

What is an equality comparator circuit?

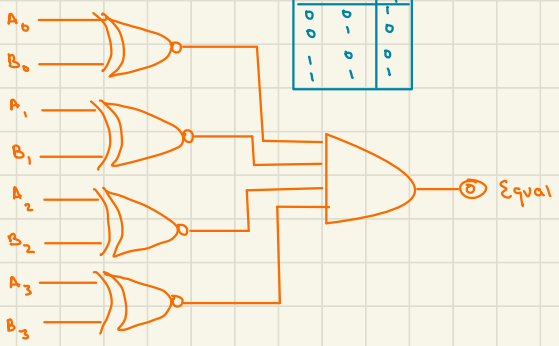
→ A circuit that determines whether the 2 inputs are equal



How would you design a circuit to compare 4-bit inputs and output 1 if they are equal?

→ Notice that when comparing 2 1-bit inputs, $\overline{A \oplus B}$ (XNOR) = 1 when $A=B$ and 0 when $A \neq B$:

A	B	Y
0	0	1
0	1	0
1	0	0
1	1	1



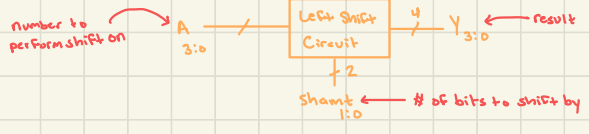
→ Alternatively, we can save transistors by converting the AND gate to a NOR gate and subsequently turning the XNOR gates into XOR gates.

How do we design a circuit to perform a left-shift on a 4-bit input?

→ Since $n=4$, shifting an input A by more than $n-1=3$ bits will always produce 00000. So we only need to support the $A \ll 1$, $A \ll 2$, and $A \ll 3$ operations in our circuit.

- Ex: $A = 0b1101$
- $A \ll 1 = 1010$
- $A \ll 2 = 0100$
- $A \ll 3 = 1000$
- $A \ll 4 = 0000$

→ Solution: Have a 2-bit "shamt" input that acts as a select signal to the output Y . Its 2 bits, so can represent up to 4 values.

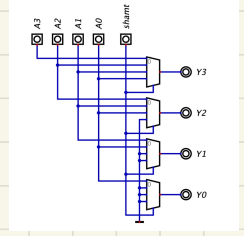


What would be the truth table?

→ Let A_3, A_2, A_1, A_0 depict the value of A , and same for Y .

Output Value		Y_3	Y_2	Y_1	Y_0
Shift Amount	0	A_3	A_2	A_1	A_0
1	A_2	A_1	A_0	0	0
2	A_1	A_0	0	0	0
3	A_0	0	0	0	0

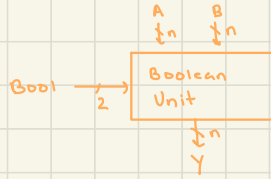
The diagram?



Arithmetic Logic Unit (ALU)

What is a boolean unit?

- A circuit/logic unit that can perform 4 boolean operations (AND, OR, XOR, NOR)
- Encompasses the logic for all 4 operations, and uses a 2-bit control signal (RECALL: multiplexers!) to tell the circuit which operation to perform.



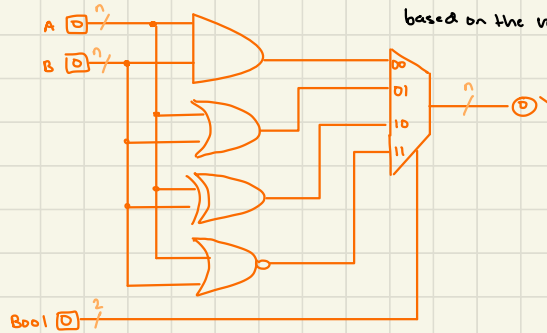
How does the Boolean unit operate?

- The value of the bool select signal indicates which operation to perform:

Bool	Operation	Expression
00	AND	$A \cdot B$
01	OR	$A + B$
10	XOR	$A \oplus B$
11	NOR	$\overline{A + B}$

- The unit will perform all of the operations, and then output the one that was requested based on the value of "bool".

What is the diagram for a boolean unit?



What is the Bidirectional Shifter?

- A circuit/logic unit that can perform 3 shift operations: left shift, logical right shift, or arithmetic right shift.
- Encompasses the logic for all 3 operations, and uses a 2-bit control signal to tell the circuit which operation to perform.

RECALL: What is a "logical" vs "arithmetic" right shift?

Operation	Expression	Meaning	Example
logical right shift	$A \gg B$	Shift B to the right by A bits, and pad the left side with 0s	$1100 \gg 2 =$
arithmetic right shift	$A \ggg B$	Shift B to the right by A bits, and pad the left side with the MSB of B! • if MSB(B) = 1, B is neg., so we pad with 1s. • if MSB(B) = 0, then $B \ll A \approx B \lll A$	$1100 \ggg 2 =$

How does the bidirectional shifter operate?

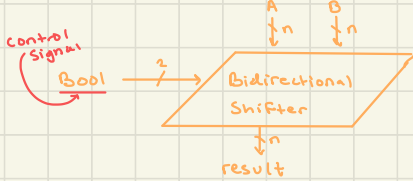
→ The 2-bit select signal "bool" indicates which operation to perform:

Bool	Operation	Expression
0 0	Left shift	$B \ll A$
0 1	NONE	N/A
1 0	logical R-shift	$B \gg A$
1 1	arithmetic R-shift	$B \ggg A$

The left bit (Bool[1]) tells you the direction of the shift (L or R)

we only need to perform 3 diff operations

The right bit tells you the type of shift (Logic or arithmetic)

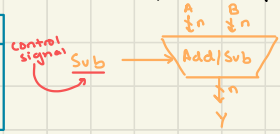


What is the add/sub?

→ RECALL: The adder/subtractor logic unit!

→ Uses a 1-bit control signal to tell the circuit which operation to perform.

Sub	Operation	Expression
0	Add	$A + B$
1	Sub	$A - B$



What is an Arithmetic Logic Unit?

→ A circuit/logic unit that can perform add/sub, bidirectional shift, AND boolean operations.

→ The value of the Shift and Math select signals are used to indicate which operation to perform.

- Shift (2-bit): if Shift = 1, perform a shift operation based on the value of Bool. If Shift = 0, perform a boolean operation based on the value of Bool.
- Math (1-bit): if Math = 1, perform an add/sub op. based on the value of Sub. if Math = 0, perform a boolean or shift op. based on the values of Shift and Bool.

What are "Shift" and "Math"?

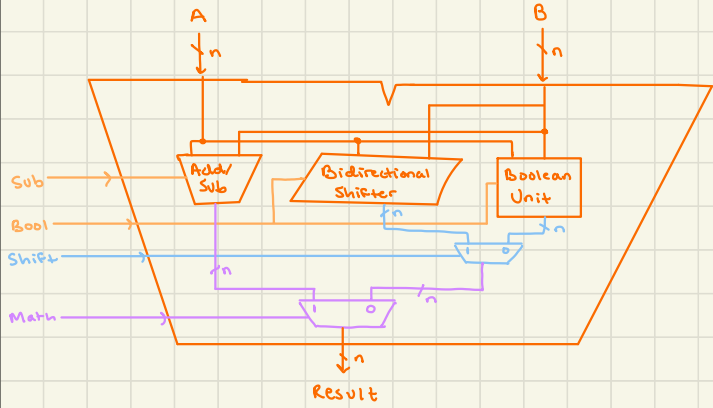
Table of all operations in the ALU?

Math	Shift	Sub	Bool	Operation
1	X	0	XX	$A + B$
1	X	1	XX	$A - B$
0	1	X	00	$B \ll A$
0	1	X	10	$B \gg A$
0	1	X	11	$B \ggg A$
0	0	X	00	$A \text{ AND } B$
0	0	X	01	$A \text{ OR } B$
0	0	X	10	$A \text{ XOR } B$
0	0	X	11	$A \text{ NOR } B$

When Math = 1, we don't care about the values of Shift and Bool; we just need to know if doing ADD or SUB

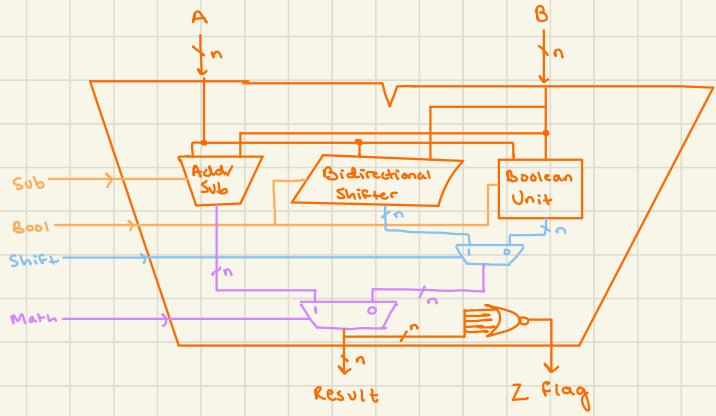
When Math = 0, we don't care about the value of Sub.

What is the ALU diagram?



What is the Z Flag?

→ A component of the ALU that outputs logic 1 if & only if the Result = 0



Waveform Diagrams

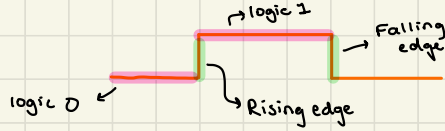
What are waveform diagrams?

→ A way to represent the values of our signals over time.

→ Each one-bit signal will have its own waveform.

How does a waveform look?

→ Like a square wave, where a low value \approx logic 0 and a high value \approx logic 1:



What are rising and falling edges?

→ Rising edge: When a signal "rises" from logic low (0) to logic high (1).

→ Falling edge: When a signal "falls" from logic high (1) to logic low (0).

Example of a waveform diagram?

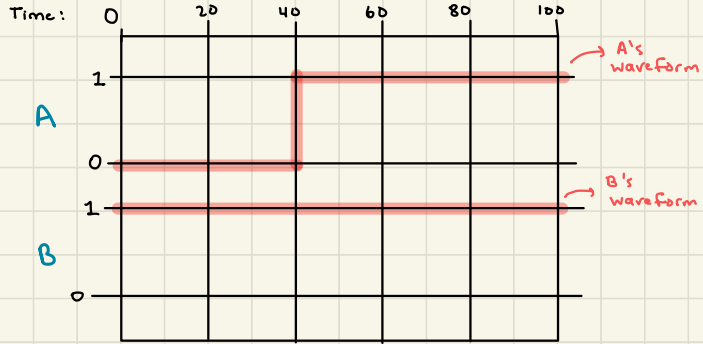
→ Take the following circuit:



• At time = 0 ps, we will

set input $A=0$ and $B=1$

• At time $t=40$ ps, we will change input A to $A=1$.



What is propagation delay?

→ The time taken for a signal to travel from input to output

• It doesn't happen immediately, as we've been assuming thus far.

→ REASON: Time for electricity to travel through a wire; capacitors charging & discharging, etc.

→ Denoted: t_{pd} = time of propagation delay.

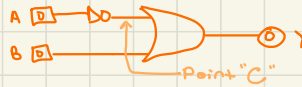
How does propagation delay affect our waveform diagrams?

→ When we change an input (e.g. A or B), the change in the output's waveform diagram won't be reflected immediately - b/c first there is a delay of \approx picoseconds.

→ The exact value of t_{pd} for different logic gates will be provided by the manufacturer; don't need to calculate.

Example of computing propagation delay?

→ Let $t_{pd}(\text{NOT gate}) = 10$ and $t_{pd}(\text{OR gate}) = 20$.



→ IF the value of input A changes at $t = n$, the value won't be inverted until $n + 10 \text{ ps}$.

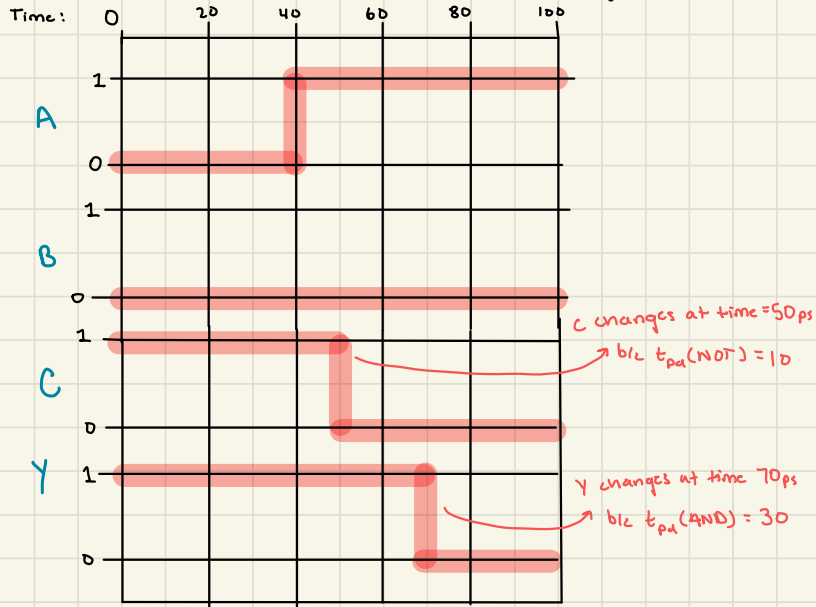
→ IF the value of input A changes at $t = n$, the value of Y won't change until $n + 10 \text{ ps}(\text{NOT gate}) + 20 \text{ ps}(\text{OR gate}) = n + 30 \text{ ps}$

• Therefore, delay from A → Y = 30 ps

→ Delay from B to Y = 20 ps

→ GIVEN: A and B start at 0. at $t = 40 \text{ ps}$, we will set A = 1.

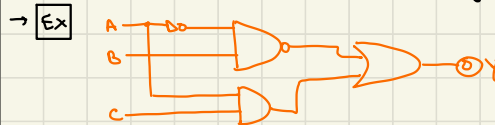
→ To visualize the p.d., we also track a "point C" placed right after the inverter.



How do we draw the waveform diagram for this circuit?

What is the longest combinatorial path?

→ Given a circuit and the t_{pd} values for each logic gate, the LCP is the path from an input to Y that has the greatest propagation delay.



Gate	$t_{pd}(\text{ps})$
NOT	15
NAND	20
AND	30
OR	40

→ LCP = A → Y :

$$15 \text{ ps} + 20 \text{ ps} + 40 \text{ ps} = 75 \text{ ps}$$

→ Shortest path = B → Y : $20 \text{ ps} + 40 \text{ ps} = 60 \text{ ps}$

Flip-Flops

What is a D Flip Flop?

→ A digital electronic circuit that is used to delay ("D"=delay) the change of state of its output signal using clock timing.

→ Diagram:



- D = the value we want to store
- Q = the stored value (from D) at any given point.
- \bar{Q} = Inverted stored value
- C = the clock; determines when D is stored.

What is a clock?

→ A one-bit signal that oscillates (or "toggles") back and forth between 0 and 1 at a constant/consistent pace (e.g. "clock period = 100 ps")

→ The purpose of the clock is to ensure that our circuits stay in sync.

What is a period?

→ RSCALL Calculus: The period of a clock is the time between one rising edge and the next:



→ The period is how we define/quantify the speed of a clock; e.g. "A clock with clock period = 700 ps".

→ The period determines how fast our circuit runs; the shorter the period, the faster our circuit.

What are the 2 types of D Flip-Flops?

1. Positive-Edge triggered:

- Q stores the value of D when the clock goes from 0 → 1
- Q updates on every rising edge of the clock

2. Negative-Edge triggered:

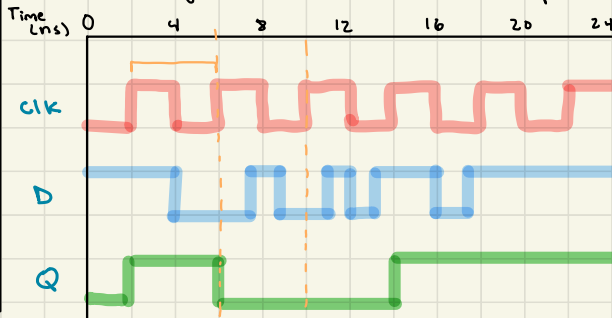
- Q stores the value of D when the clock goes from 1 → 0
- Q updates on every falling edge of the clock

How does a positive-edge D Flip-Flop work?

→ Basically, we have 1 or more inputs and some logic circuit that outputs a value D based on the input(s) — e.g., a NAND circuit, a full adder, an ALU, etc.

→ The Flip-Flop adds a variable Q which starts at some specified value (e.g. 0) and then, every time the clock has a rising edge — aka every p seconds, where p = the period length —, the value of Q is updated to equal the value of D at that point in time.

Example of a waveform diagram for a pos. edge-triggered Flip-Flop?



- The clock has a period of 4 ns
- Every time clock has a rising edge, check the value of D and change Q (if needed).
- Notice that D=0 at both indicated clock rising edges

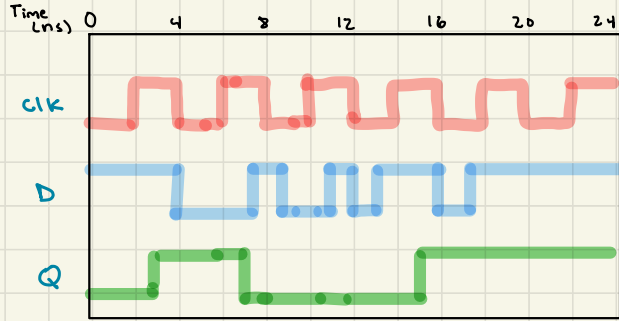
What is Clock-to-Q delay?

Example?

→ The amount of time that it takes for the "update" of D's value to appear on the output (Q) after the clock trigger.

→ Aka, the propagation delay for Q — it doesn't immediately update when CLK goes from 0 → 1 ... there is a t_{pd} in between.

→ The same waveform diagram from the previous ex, but with $clk-to-q\text{-delay} = 1\text{ ns}$:

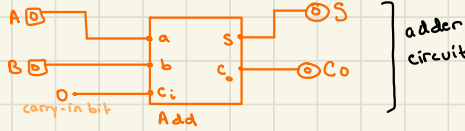


Registers Motivation: Back-to-back additions

How would we perform back-to-back additions with one adder circuit?

→ Let t_{pd} for our adder = 600 ps — meaning that after we input A and B, it will take 600 ps for a value and a carry-out bit to be output.

→ Ex: Performing $7+7$, $18+7$, and $30+8$, back-to-back.



What are the steps to doing this?

1. At time $t = 0\text{ps}$, set $A = 7$ and $B = 7$. Initially, S and $C_o = 0$.
2. We can't start addition #2 until the circuit outputs $S = 14$ — 600 ps bit of t_{pd} .
At $t = 600\text{ps}$, we can now change $A = 18$ and $B = 7$.
3. At $t = 1200\text{ps}$, S has updated to $= 25$. Now we are ready to perform the next op. So we set $A = 30$ and $B = 8$.
4. Once we are done w/ this last one, we can set $A = 0$ and $B = 0$.

Time (ps)	A	B	S	C _o
0	7	7	0	0
600	18	7	14	0
1200	30	8	25	0
1800	0	0	38	0

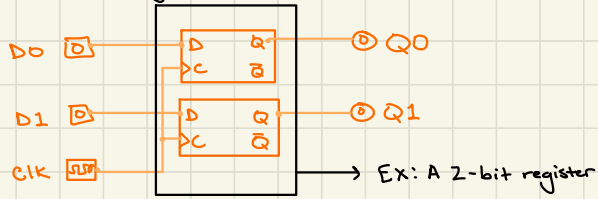
→ In total, performing these 3 operations took 1800 ps.

- Registers -

What is a register?

→ A combination of multiple Flip Flops together

→ Used for storing multi-bit values



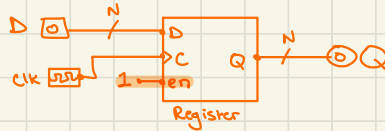
Why do we use registers?

→ If we wanted to do back-to-back operations like in the prev. addition example, we can't just manually set the input values every 600 ps.

• Instead, we use a **CLOCK** and we use registers to update the inputs at the proper times.

What is an N-bit

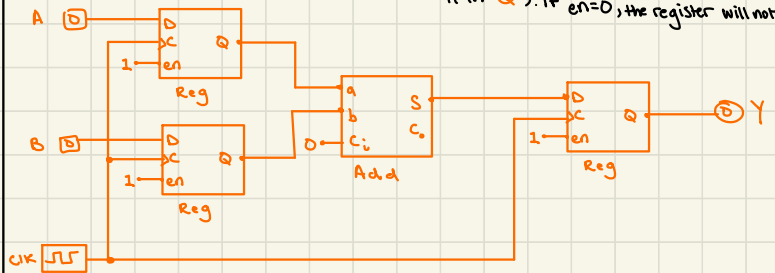
positive-edge triggered register?



• "en" = Enable

• if $en = 1$, the register will read D on the rising edge of the clock (and store it in Q). If $en = 0$, the register will not read D .

Ex: How would we use registers to perform the additions?



How would the timing work?

→ Each register's stored value will update on the rising edge of the clock.

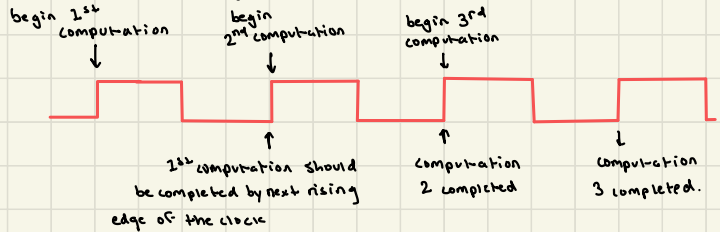
This means that when we begin (e.g. 0ps), the 1st reg will have nothing, and A & B will be storing "1" and "1".

• Upon the first rising edge (+ the CLK-to-Q delay), the register will output the vals of A and B that they store - meaning we can begin the 1st computation.

How will we utilize the clock cycle?

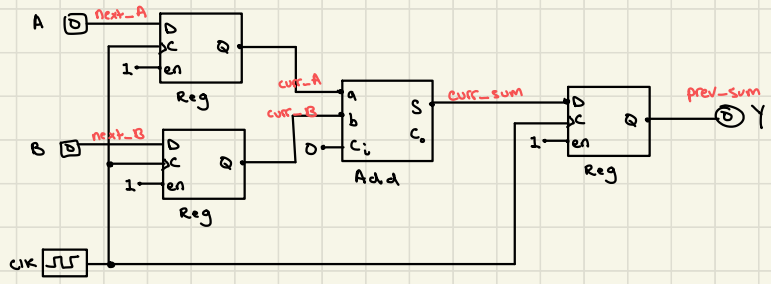
→ Due to the schedule of when the registers update, we can perform one addition per clock cycle.

→ We should set the length of the clock cycle s.t. each computation can be done in 1 cycle:



Timing

Example of performing back-to-back additions with registers?



→ In this Ex, we will perform $30+8$, $18+7$, and $7+7$. We will complete one addition per clock cycle.

What is the starting state?

→ All values = 0

What happens in Cycle 0?

→ On the Rising Edge, the registers will update to hold values for A and B:

- next-A = 30
- next-B = 8

→ After the clk-to-q delay, the values stored in registers will get stored in Q:

- next-A = 30
- next-B = 8
- curr-A = 30
- curr-B = 8

→ After the adder propagation delay (Add pd), the adder will output the sum of Q(A) and Q(B):

- next-A = 30
- next-B = 8
- curr-A = 30
- curr-B = 8
- curr-sum = 38
- prev-sum = 0

→ The sum will not be sent to output Y until the next rising edge (plus clk-to-q delay) because curr-sum first has to go through that final register.

What happens in Cycle 1?

→ Cycle 1 begins on the next Rising Edge of the clock, when the new operands are sent in:

- next-A = 18
- next-B = 7
- curr-A = 30
- curr-B = 8
- curr-sum = 38
- prev-sum = 0

→ After clk-to-q delay:

- next-A = 18
- next-B = 7
- curr-A = 18
- curr-B = 7
- curr-sum = 38
- prev-sum = 38

(Now, the value of the first add operation is "available")

→ After Adder pd:

- next-A = 18
- next-B = 7
- curr-A = 18
- curr-B = 7
- curr-sum = 25
- prev-sum = 38

(the new sum is output by the adder)

What happens in **Cycle 2**?

→ On next Rising edge:

- next_A = 7
- curr_A = 18
- curr_sum = 25
- next_B = 7
- curr_B = 7
- prev_sum = 38

→ After clk-to-q delay:

- next_A = 7
- curr_A = 7
- curr_sum = 25
- next_B = 7
- curr_B = 7
- prev_sum = 25

→ After Adder pd:

- next_A = 7
- curr_A = 7
- curr_sum = 14
- next_B = 7
- curr_B = 7
- prev_sum = 26

What happens in **Cycle 3**?

→ After the next rising edge & clk-to-q delay, the sum $7+7=14$ will be sent to Y. i.e., **prev_sum = 14**.

How long is the clock cycle/period?

→ We can adjust the length of the clock cycle - aka time between one rising edge and the next - based on our needs! Its up to us.

→ There is a minimum clock period at which the circuit will work properly, so our clock period can be set to any value \geq min. clock period.

What is the **minimum clock period**?

→ The most optimized (i.e. shortest) time within which the circuit can complete one operation.

- We don't include the clk-to-q delay of the final output register when calculating the min. clock period.

How do we calculate the min clock period in this example circuit?

→ Min clock period = clk-to-q delay of the input registers + adder propagation delay

↳ b/c the longest path has \pm adder

How do we calculate the min clock period in general?

→ **Min clock period = clk-to-q delay + longest combinational delay**

What is the "Longest combinational delay"?

→ The length (in time) of the longest path between 2 registers - AKA longest combinational path

→ Basically, the amt of times it takes for each operation unit in a path between 2 registers

→ Does NOT include the clk-to-q delay time of the registers.

What is the **max clock frequency**?

→ $1 / \text{min clock period}$

Pipelining

What is a single-cycle implementation?

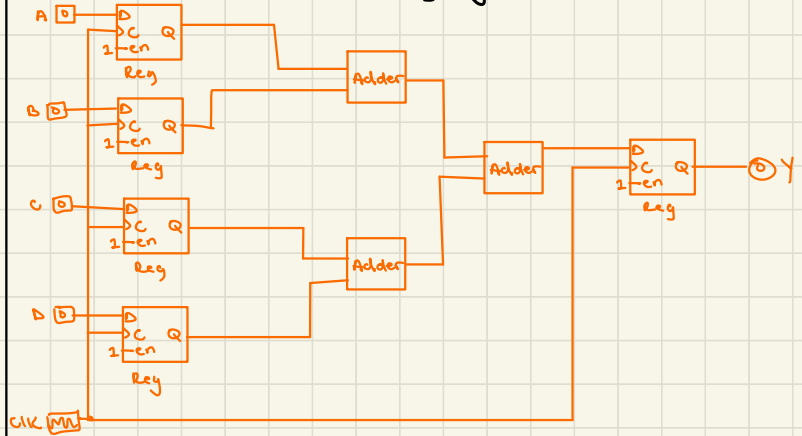
→ An implementation of a circuit such that it completes one operation in a **single clock cycle**.

What is pipelining?

→ Breaking our circuit down into multiple stages that can operate simultaneously in order to make more efficient use of our hardware.

• We create "stages" by adding registers.

Example of a single-cycle circuit?



→ This circuit performs $Y = A + B + C + D$

→ Let $clk\text{-to-}q$ delay = $20ps$ & adder $pd = 480ps$

→ The longest path btwn 2 registers has 2 adders, so $480 + 480 = 960ps$

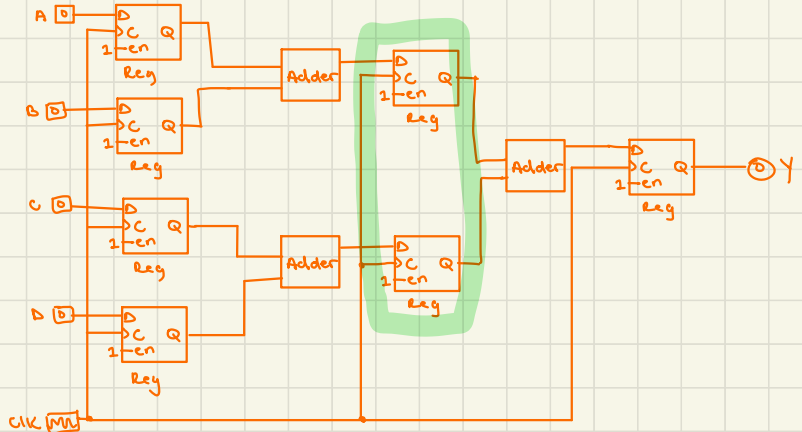
→ $20ps$ ($clk\text{-to-}q$ delay) + $960ps$ (LCD) = $980ps$

→ If we run this circuit at the min clock period, it will take $980ps \cdot 10 = 9800ps$ to complete 10 operations.

Ex: What is the longest combinational delay?

Ex: What is the min clock period?

How do we implement pipelining on this circuit?



How are these 2 circuits different?

→ The single-cycle impl. performs $A + B$ and $C + D$ simultaneously, then immediately does the sum of the 2.

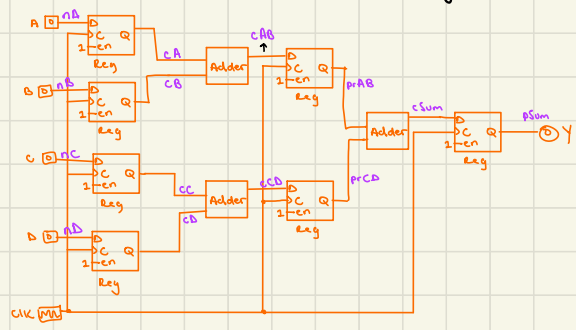
→ The pipelined impl. performs $A + B$ and $C + D$ simultaneously, but then sends the outputs into 2 registers.

EX: How does the pipelined circuit work?

→ GIVEN: $clk-to-q$ delay = 20ps and Adder pd = 480ps

→ We will set our clock period to 500ps & do these operations, starting at $t=0$ ps.

1. $10 + 9 + 8 + 7$
2. $5 + 4 + 3 + 2$
3. $10 + 10 + 20 + 20$



What will happen at...

$t = 20$ ps?

→ At the beginning, nA, nB, nC, nD will hold the operands for the 1st addition.

→ The registers will update to reflect the input values:

- $nA, cA = 10$
- $nB, cB = 9$
- $nC, cC = 8$
- $nD, cD = 7$

$t = 500$ ps?

→ After adder pd, the outputs are updated.

→ Additionally, since we broke our circuit into 2 stages of addition, we can now start operation 2!

- $nA = 5$
- $nB = 4$
- $nC = 3$
- $nD = 2$
- $cA = 10$
- $cB = 9$
- $cC = 8$
- $cD = 7$
- $cAB = 19$
- $cCD = 15$

$t = 520$ ps?

→ The inputs for op2 move past the registers and the inputs to the 2nd adder also move past the registers.

→ Why? B/c $t = 500$ ps was the rising edge of the clock, and then it took 20ps of $clk-to-q$ delay for the registers to update.

- $nA = 5$
- $nB = 4$
- $nC = 3$
- $nD = 2$
- $cA = 5$
- $cB = 4$
- $cC = 3$
- $cD = 2$
- $cAB = 19$
- $cCD = 15$
- $prAB = 19$
- $prCD = 15$

$t = 1000$ ps?

→ 480ps later, all 3 of the adders have performed & output values.

→ Now, we can also start operation 3!

- $nA = 10$
- $nB = 10$
- $nC = 20$
- $nD = 20$
- $cA = 5$
- $cB = 4$
- $cC = 3$
- $cD = 2$
- $cAB = 9$
- $cCD = 5$
- $prAB = 19$
- $prCD = 15$
- $cSum = 34$

$t = 1020 \text{ ps} ?$

→ After the R.E. at $t = 1000 \text{ ps}$, the registers update to hold the values of A, B, C, D for op. 3, $A+B$ and $C+D$ for op. 2, and $A+B+C+D$ for op. 1.

→ Operation 1 is now complete.

- $nA = 10$
- $nB = 10$
- $nC = 20$
- $nD = 20$
- $cA = 10$
- $cB = 10$
- $cC = 20$
- $cD = 20$
- $cAB = 9$
- $cCD = 5$
- $prAB = 9$
- $prCD = 5$
- $cSum = 34$
- $pSum = 34$

$t = 1500 \text{ ps} ?$

→ All 3 adders perform addition operations & output values.

- $nA = 0$
- $nB = 0$
- $nC = 0$
- $nD = 0$
- $cA = 10$
- $cB = 10$
- $cC = 20$
- $cD = 20$
- $cAB = 20$
- $cCD = 40$
- $prAB = 9$
- $prCD = 5$
- $cSum = 14$
- $pSum = 34$

$t = 1520 \text{ ps} ?$

→ The values of $A+B$ and $C+D$ for op. 3 and $A+B+C+D$ for op. 2 get sent through the registers.

→ Operation 2 is now complete.

→ It took 2 clock cycles (op. 1 was done at 1000 ps, and our clock period was 500)

→ The longest circuit between 2 registers has 1 adder, so LCP = 480 ps

→ $20 \text{ ps (clk-to-q delay)} + 480 \text{ ps (LCD)} = 500 \text{ ps}$

→ Op 1 finished at 1000 ps, but operations 2 and beyond will finish at 1500, 2000, 2500... etc. ps, so for 10 operations:

$$500 \text{ ps} * (10 + 1) = 5,500 \text{ ps}$$

	Single Cycle	Pipelined
longest comb. path	960 ps	480 ps
min clock period	980 ps	500 ps
# of cycles to complete 1 op.	1	2
Time to complete 10 ops at MCP	9800 ps	5,500 ps

The pipelined circuit is much faster!

Ex: How many clock cycles did it take to complete 1 operation?

Ex: What is the longest combinational path/delay?

Ex: What is the min clock period?

Ex: At the MCP, how long will it take to do 10 operations?

Comparison of the single cycle & pipelined impls of $Y = A+B+C+D$?

What are the execution times for both implementations?

→ To perform n operations:

- Single Cycle: $(n) \cdot (\text{clock period})$
- Pipelined: $(n+1) \cdot (\text{clock period})$

Why is the pipelined impl. usually faster?

→ Because it has a significantly smaller MCP.

→ $n(x)$ v.s. $(n+1)(x \div 2)$... the 2nd expression will OFTEN yield a smaller value.

- Metrics -

What is latency?

→ Let SC-Y denote the single-cycle example circuit performing $Y = A + B + C + D$ and let PL-Y denote the pipelined example.

→ The amount of time it takes to complete a single operation, from beginning to end.

• Ex: Latency of SC-Y = 1 clock cycle = 980 ps

• Latency of PL-Y = 2 clock cycles = $2 \times 500 = 1000$ ps

What is throughput?

→ The number of operations that can be completed in a given amount of time.

• PL-Y has a higher throughput than SC-Y.

What is speedup?

→ The measure of how much faster the pipelined impl. is in comparison to the single-cycle impl.:

$$\text{SPEEDUP} = \frac{\text{Time to complete } n \text{ ops on SC impl.}}{\text{Time to complete } n \text{ ops on Pipelined impl.}}$$

• Ex: The speedup of operation $Y = A + B + C + D$ for 10 operations =

$$9800 \text{ ps} / 5500 \text{ ps} = 1.78$$

311 Quiz 1 Review

→ TOPICS: K-maps - including "terms", Multiplexers, Digital Logic terms - Canonical SOP, maxterms, minterms, Adder, comparator, shifter, ALU, Wave form diagrams, timings on T-F diagrams.

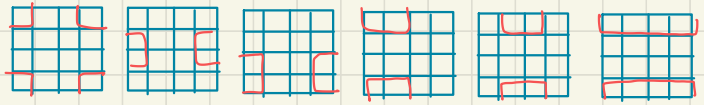
K-Maps

→ Find the literals that each group has in common, and OR them for the final term.

Grouping Rules:

- Each group must contain only 1s, and all 1s have to get grouped
- Groups must be 2^n cells (e.g. 1, 2, 4, 8, ...)
- Each group must be as large as possible
- If > 2 groups, each group must have at least 1 unique cell.
- Ways to wrap around the table:

AB \ CD	00	01	11	10
00	1	0	1	1
01	0	0	0	x
11	0	0	1	1
10	1	0	x	1



STEPS:

1. Create K-map
2. Include all essential prime implicants
3. Include all non-essential prime implicants as needed to cover all 1s.

Implicants

→ Implicant = any product/minterm in the SOP form of an equation.

• Ex: $Y = MN + MND + ND$ = implicant ... aka every square with a "1" in the K-map.

→ Prime Implicant = All possible groups that can be formed in a K-map. The p.i. itself is the term that defines a given group

→ Essential Prime Implicant = groups that cover at least one minterm that cannot be covered by another prime implicant.

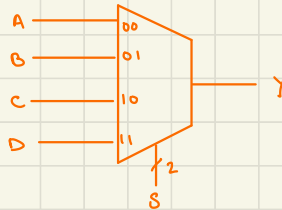
Multiplexers

→ Chooses output based on a select signal:

→ If # of inputs = n , select signal S

must be at least $\log_2(n)$ bits.

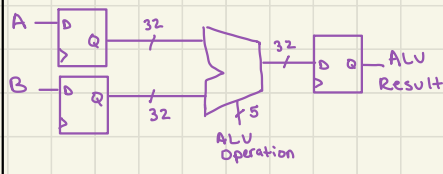
• aka, $2^{(\# \text{ of bits})} \geq n$



• if $S = 00$, $Y = A$; if $S = 01$, $Y = B$;
if $S = 10$, $Y = C$; if $S = 11$, $Y = D$.

What are the control signals associated with each operation?

ALU Operation	Function
0b 00000	AND
0b 00100	OR
0b 01000	XOR
0b 01100	NDR
0b 00001	add
0b 10001	sub
0b 00011	set on less than



* assume all registers are connected to the same clock.

What is the "set on less than" (SLT) operation?

→ Performs the operation "A < B".
 • If $A < B$, ALU Result = 1. If $A \geq B$, ALU Result = 0.

How does the ALU work?

→ It performs one operation per cycle.
 → For each cycle, we must set 3 inputs: A, B, and ALUOp

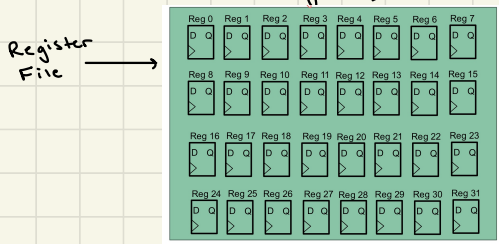
MIPS Processor

RECALL: What is a register?

→ an operand which has to do with memory/hardware
 → MIPS defines 32 general purpose registers, \$0 through \$31, each of which have their own meanings. Each one can hold a 32-bit value. \$0 always holds val=0.

What is the Register File?

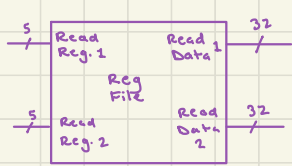
→ A small piece of memory for storing intermediate results of computations.
 → stored in hardware (aka the registers)
 → Contains the 32 32-bit registers, \$0 to \$31:



What is the MIPS processor?

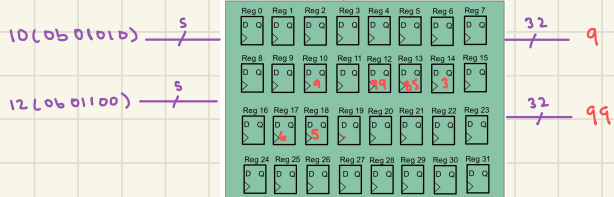
→ A piece of logic that allows us to perform dependent operations!
 • e.g. $a = 5 + 4$, $b = 7 + 8$, and $c = a + b$
 → How? By combining the RegFile with the ALU Unit so that we can read from and write to registers.

How do we read from the register File?



→ Set the "read reg." inputs to the 5-bit val. of the number of the register you want to read.
 → The value stored at that register will appear on the corresponding output.

Example of reading from the RegFile?



How do you write to the Register File?

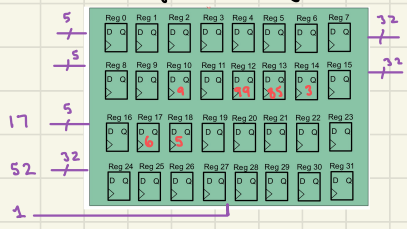


- Set the "write register" to the num. of the register where you want to store data.
- Set the "write data" input to the value you want to store.
- **RegWrite** = a control signal that we set to 1 when writing, and 0 when not writing.

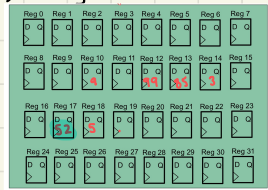
→ When **RegWrite** = 1, we can still simultaneously read data on output if we want to.

Example of writing to the Reg File?

→ EX: Writing 52 to register \$17:



After the next rising edge of the clock, the RegFile will update:

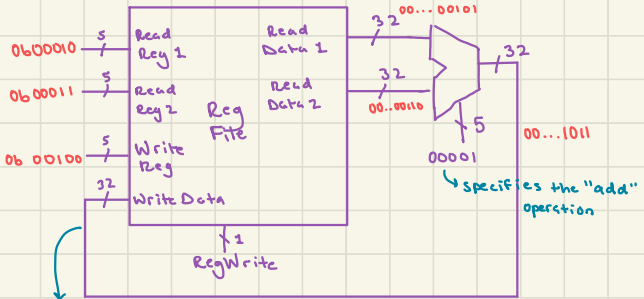


Putting it all Together

How are the RegFile & ALU unit combined in the MIPS processor?

- The "read data" outputs are given as inputs A and B to the ALU
- The output of the ALU operation is wired as the input to the "write data"!
- Lets perform $\$4 = \$2 + \$3$, where the vals stored at \$2 and \$3 are 5 and 6

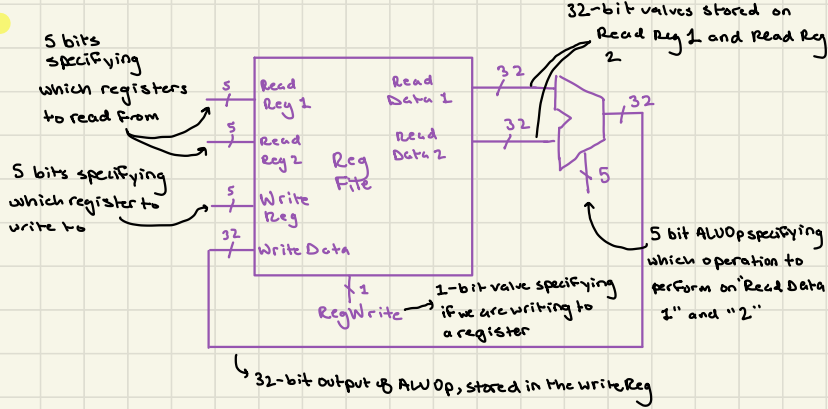
Example?



the ALU output is sent to the "write data" input

R-Format Instructions

RECAP: What are the inputs & outputs of a MIPS processor?



How is assembly code turned into binary code?

→ An assembly instruction can be translated into a 32-bit binary string that can then be input into the MIPS processor!

→ The 32-bit data is broken up into several "sections"

→ Assembly instructions to perform an operation where the registers are the operands.

• e.g., not adding a constant val, but rather 2 vals already stored at 2 registers.

What are r-format instructions?

[Operation] [rd], [rs], [rt]

What is the syntax of r-format instructions?

1. "operation": the operation being performed (e.g. AND, add, etc.)

2. **\$rd**: the destination register; where the output of the operation will be stored.

3. **\$rs, \$rt**: the source & target registers whose stored values will be the

inputs to the operation (e.g. "A" and "B")

→ EX:

What are the fields of the 32-bit binary instruction?

Field:	opcode	rs	rt	rd	shamt	funct
Bits:	31:26	25:21	20:16	15:11	6:10	5:0
	(6 bits)	(5 bits)	(5 bits)	(5 bits)	(5 bits)	(5 bits)

→ The 5-bit binary nums of the 3 registers in the operation (rs, rt, rd) are stored in the "rs", "rt", and "rd" fields.

What are the "opcode" and "funct" fields?

→ They are specified by the ISA, and they tell the processor which operation to perform.

→ The values at opcode & funct are used to determine the value of the ALUOp control signal.

What are the opcode & Funct fields for common operations?

Instruction	Opcode	Funct
add	000000 ₂	10 0000 ₂
sub	000000 ₂	10 0010 ₂
and	000000 ₂	10 0100 ₂
or	000000 ₂	10 0101 ₂
nor	000000 ₂	10 0111 ₂
slt	000000 ₂	10 1010 ₂

ALUOp:
 00001
 10001
 00000
 00100
 01100
 00011

Example converting assembly to binary?

→ EXAMPLE: the instruction `add $t2, $t7, $t0`

1) Interpret the instruction and identify register fields:

• `$t2 = $t7 + $t0`
 ↓ ↓ ↓
 rd rs rt

2) Fill out fields:

opcode	rs	rt	rd	shamt	Funct
10 0000	00111	01010	01100	00000	100000

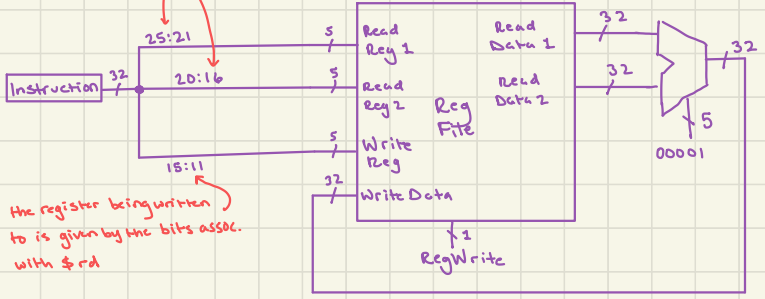
ANS: 00 1000 0000 1110 1010 0110 0000 0010 0000

→ The 32-bit inst. is taken as an input, and then split into substrings corresponding to the different fields!

How is the 32-bit r-instruction interpreted by the MIPS processor?

field	bits of inst
funct	5:0
shamt	10:6
rd	15:11
rt	20:16
rs	25:21
opcode	31:26

the bits specifying the rt & rs register numbers are sent to read reg 1 and 2



the register being written to is given by the bits assoc. with \$rd

Summary: What are all of the R-format operations?

→ Syntax: `OP $rd, $rs, $rt`, $R[rd] = R[rs] \text{ OP } R[rt]$

• $R[rd]$ denotes "RegFile[rd]", aka, the contents stored at register rd.

NAME	OPERATION
add	$R[rd] = R[rs] + R[rt]$
sub	$R[rd] = R[rs] - R[rt]$
and	$R[rd] = R[rs] \text{ AND } R[rt]$
or	$R[rd] = R[rs] \text{ OR } R[rt]$
nor	$R[rd] = R[rs] \text{ NOR } R[rt]$
slt	$R[rd] = (R[rs] < R[rt]) ? 1 : 0$

if $R[rs] < R[rt]$, $R[rd] = 1$. else, $R[rd] = 0$

I-Format Instructions

What are i-format instructions?

- Where the instruction operands are a combination of a register and a 16-bit constant (an "immediate" operand)
- Used when you want to perform an op between a register and an immediate value (i.e. a constant that is not in a register)
 - As opposed to R-type instructions, which perform operations between 2 registers.

[operation] [rt], [rs] , [immediate]

↗ result stored here

What is the syntax of an i-format instruction?

- In i-format instructions, we need the last 16 bits of the inst. to be designated for storing the constant val, so we don't have an rd field. **Instead, we store the result in the reg specified by rt.** We write to rt.

→ Ex: $5 = 4 + 12 \rightarrow \text{addi } \$5, \$4, 12$

→ Syntax: $R[rt] = R[rs] \text{ OP } (\text{immediate})$

How do we set a register to a certain constant value?

- Using i-format instructions with the source operand register \$0, since it always holds value = 0!

→ Ex:

$a = 5 \xrightarrow{\hspace{2cm}} \text{addi } \$10, \$0, 5$ (\$10 now holds var a)
 $b = 3 \xrightarrow{\hspace{2cm}} \text{addi } \$11, \$0, 3$ (\$11 holds var b)
 $c = a + b \xrightarrow{\hspace{2cm}} \text{add } \$12, \$10, \11 (\$12 holds var c = a + b)

What are the fields of the 32-bit i-instruction?

Field:	opcode	rs	rt	immediate
Bits:	31:26	25:21	20:16	15:0
	(6 bits)	(5 bits)	(5 bits)	(16 bits)

- The rd, shamt, and funct fields are excluded to make room for the immediate value.

→ INST: $\text{addi } \$5, \$0, 2$

opcode	rs	rt	immediate
001000	00000	00101	0000 0000 0000 1100

• ANS: 0b 0010 0000 0000 0101 0000 0000 0000 1100

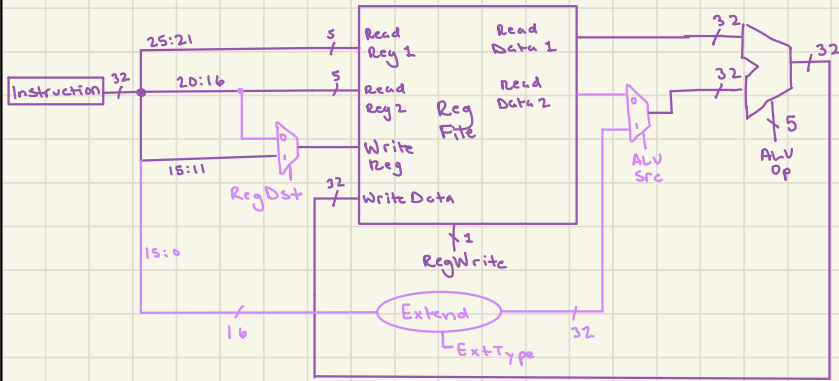
Example converting assembly → binary?

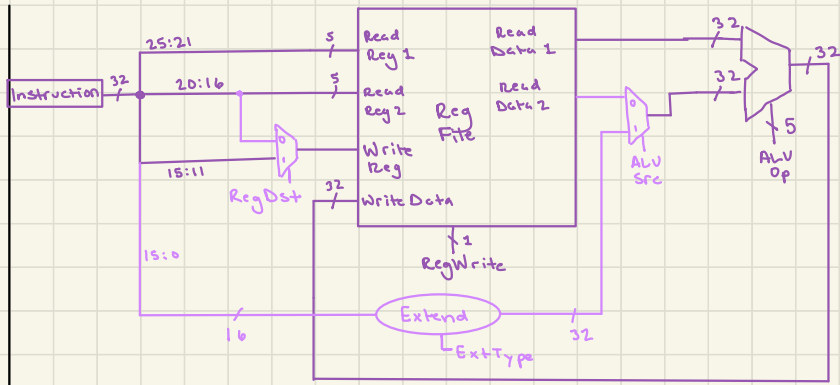
What are the opcodes for i-format inst. operations?

Instruction	Opcode	ALU Op
addi	00 1000 ₂	00 001
andi	00 1100 ₂	00 000
ori	00 1101 ₂	00 100
slli	00 1010 ₂	00 011

How do we modify the MIPS processor hardware to execute both R and I instructions?

- 1) Add a 2:1 mux with a control signal that takes as inputs:
 - The 5-bit value of rd (aka bits 15:11, aka the write register for r-instructions), AND
 - The 5-bit value of rt (aka bits 20:16, aka the write register for i-instructions)And uses a $RegDst$ control signal to determine which value will get sent to the Write Register input.
- 2) Add a 2:1 mux that takes as inputs:
 - The 32-bit value $R[rt]$ (aka "Read Data 2" aka the 2nd operand for r-instructions, AND
 - The 32-bit value $Extended_immediate$ (aka [extension + bits 15:0] aka the 2nd operand for i-instructions)And uses an $ALUSrc$ control signal to determine which value will get sent to the second ALU input.
- 3) A piece of logic that takes the $immediate$ value (aka bits 15:0) as an input, and outputs the extended 32-bit imm. value.
 - It has an $ExtType$ control signal that specifies whether the value should be zero- or sign- extended.





Summary: Control signals for all r & i inst operations?

→ Whenever we are doing a boolean or logic operation (and, or, nor, xor, add, sub, slt), $RegWrite = 1$

→ R-instructions:

- $ALUSrc == 0$ (taking R[rt] as the 2nd source operand)
- $RegDst == 1$ (storing the result in \$rd)
- $ExtType == X$ (DON'T CARE) (no value that requires extension)

→ I-instructions:

- $ALUSrc == 1$ (taking bits 15:0 (extended), aka imm., as 2nd source operand)
- $RegDst == 0$ (storing the result in \$rt)
- $ExtType == 0$ for boolean operations
 - andi
 - ori
 - nori
- $ExtType == 1$ for logic operations
 - addi
 - sli

MIPS Programming

What are Labels?

→ A way to "mark" sections of assembly code that you may want the program to be able to jump to or repeatedly execute. Like the beginning & end of a loop.

• EX: For-loops.

Example?

→ loop.c:

```
sum = 0;
for (i = 0; i < 5; i++) {
    sum += 2;
}
sum += 10
```

→ loop.o:

• we'll store sum and i in \$10, \$8
 • we store the value representing the "terminating point" of i, aka 5, in \$9

→ We enclose the instructions in loop.c's for-loop in a Label:

```

addi $10 $0 0
addi $8 $0 0
addi $9 $0 5
LoopStart:
addi $10 $10 2
addi $8 $8 1
LoopEnd:
addi $10 $10 10
```

initialize sum, i, and terminating point

Label "LoopStart"

sum = sum + 2

i++ (i = i + 1)

Label "LoopEnd"

sum = sum + 10

How do we actually utilize the labels?

→ With branch instructions that check whether the terminating condition (in this case $i \geq 5$) has been reached!

→ Instructions that evaluate some condition. If the condition is true, they tell the program to "branch" to a different spot in the code. Specifically, they define a label that should be branched to if the condition is true.

• EX: If $i = 5$, end the for-loop & go to the label for the code to be executed after the loop ends.

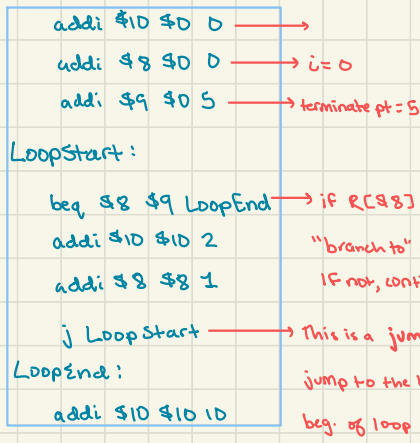
What are branch instructions?

Instruction	Format	Meaning: Branch to Label if...
Branch on equal	beq \$rs \$rt Label	$\$rs == \rt
Branch on not equal	bne \$rs \$rt Label	$\$rs != \rt
B.O. greater than	bgt \$rs \$rt Label	$\$rs > \rt
B.O. greater than or equal	bge \$rs \$rt Label	$\$rs \geq \rt
B.O. less than	blt \$rs \$rt Label	$\$rs < \rt
B.O. less than or equal	ble \$rs \$rt Label	$\$rs \leq \rt

What are all of the branch instructions?

→ Note that for Branch instructions, RegWrite = 0!

Example loop.o using branch instructions?



What are Native Instructions?

→ Instructions that are supported by the datapath (the MIPS hardware diagram thing) AND instr. that have actual hardware support.

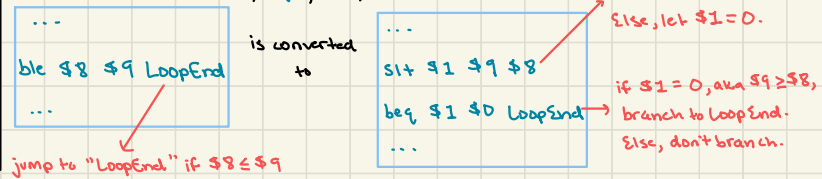
- Includes: all the instructions we learned up until now
- Includes: `beq`, `bne`

What are Pseudo Instructions?

→ Ones that aren't supported by the datapath; they only exist to make programs easier to read & write.
→ When the program is assembled, pseudo instructions are converted to one or more native instructions.

- Includes: `bgt`, `bge`, `blt`, `ble`

Example of a pseudo instruction being converted to a native one?



- Register Usage -

What is register 0 (`$0`) for?

→ ALWAYS holds the value 0. Cannot be written to.
→ Useful for when you need the value 0 (like initializing variables with `addi $X $0 0`)

What is register 1 (`$1`) for?

→ Used by the assembler to resolve pseudo instructions.
→ You CAN write to this register, but the assembler may override the value when it needs to resolve a pseudo-inst.
• Shouldn't rely on `$1` to store values you might need later.

MIPS Memory Model

RECALL: What is the structure of a computer's main memory?

→ Data stored in a "primary storage" component like RAM, accessible to the CPU via a BUS (communication system between diff pieces of hardware)

- Unlike Register data, which is stored in the register file which is actually located in the CPU.

1. Stack: Used to manage Function calls & local variables.

- STACK grows down

2. Dynamic Data aka Heap: Used for dynamic memory allocation.

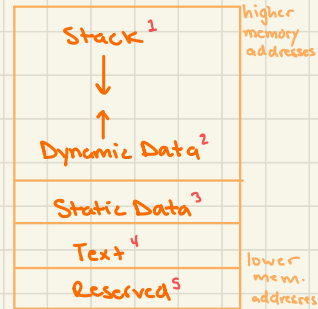
- HEAP grows up

3. Static Data: variables allocated at compile-time.

- e.g., statically allocated arrays.

4. Text: Programs (aka actual code).

5. Reserved memory for the OS



How does Memory compare to the Register File?

→ Memory:

- large
- takes longer to access
- NOT part of the CPU

→ Register File:

- small
- much, much, much faster to access
- part of the CPU

→ We need MM b/c we can't store all of our data in the RF

→ We store values/data that we are currently working with in the RF.

What is the load word instruction?

→ lw ; used to read 4 bytes (aka 1 word) of data from memory and store it in a register.

→ `lw $rt offset($rs)`

What do \$rt, offset, and \$rs mean in the lw inst.?

→ \$rs : a register that (already) holds a memory address.

→ offset : an immediate value that is added to R(\$rs) to obtain the address

- e.g., `lw $rt 4($rs)` ≈ loading a word from `4 + [addr stored at $rs]`

→ \$rt : Where we store the 4 byte (32-bit) data that we read from addr `[offset + R($rs)]`

Formula for the lw instruction?

→ $R[\$rt] = M[R[\$rs] + \text{offset}]$, where "M" means Main Memory.

Example of the lw instruction?

- Let $\$10 = 0x00004000$, and let $A = [5, 90, 100, 40, 11]$ be an int array stored at base addr. $0x00004000$ (aka $A[0]$). The size of an int is 4 bytes.
- EX: `lw $11 8($10)` will store the value 100 in register $\$11$.
 - $\$10$ holds $0x00004000$ & offset = 8, so we want to read the data at address $0x00004000 + 8 = 0x00004008$
 - Since $\text{sizeof}(\text{int}) = 4$ bytes, $0x00004008$ holds the 3rd element of A .

What is the store word instruction?

- `sw`; Used to store 4 bytes (1 word) of data from the RF into memory!
- `sw $rt offset($rs)`
- $\$rs$: A register that (already) holds a memory address.
- `offset`: An immediate value that is added to $R[\$rs]$ to obtain the address where we want to store the data.
- $\$rt$: Holds the data/value that we want to store at the MM address.
- $M[R[\$rs] + \text{offset}] = R[\$rt]$

What do $\$rt$, `offset`, and $\$rs$ mean in the sw inst.?

What is the formula?

Example using sw?

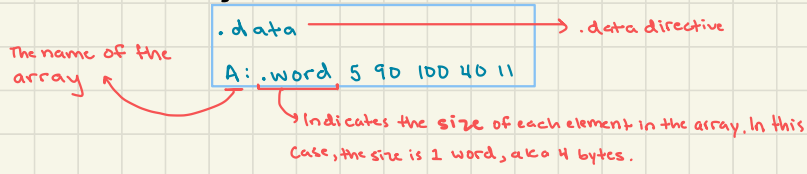
- Let $\$10 = 0x00004000$, $\$7 = 8$, and let $A = [5, 90, 100, 40, 11]$ be an int array stored at base addr. $0x00004000$ (aka $A[0]$).
- EX: `sw $7 12($10)` will modify array A into:

$$A = [5, 90, 100, 8, 11]$$

How do we store arrays in memory?

- We store statically allocated arrays in the static data segment (below the heap).
- To store ANY data in the static data segment, we use the `.data` directive in our assembly code.

Example of using .data to store data in static memory?



How do we actually get a memory address into a register?

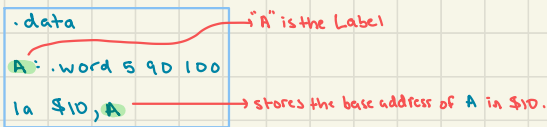
- E.g., the value stored in $\$rs$ for load-word & store-word instructions.
- ANS: with the load address instruction, `la`.

What is the syntax of the la inst.?

```
la $rt Label
```

- Gets the start address of the array/variable specified by `Label` and stores it in $\$rt$.
- For an array, `la` gets the start address of the array.

Example using la?

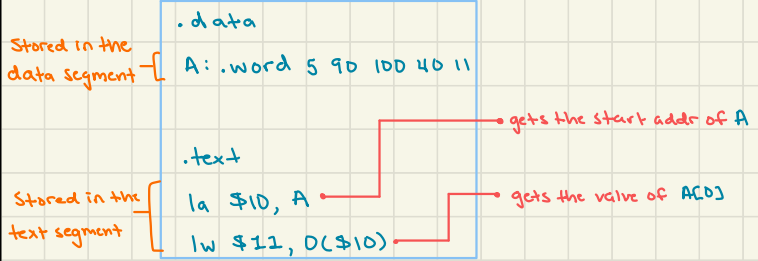


How do we store data in the text segment of memory?

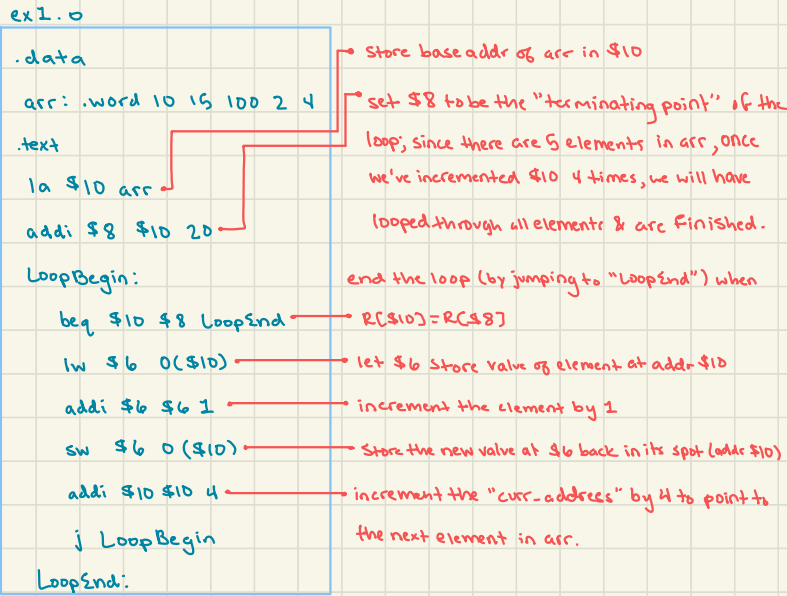
Example of using the `.text` directive?

→ Using the `.text` directive

→ We will store our actual assembly code in the text segment, using the `.text` directive.



→ Let `arr = [10, 15, 100, 2, 4]`



What is the Static instruction count?

→ The number of native instructions that a program has.

→ AKA, if you convert every pseudo-instruction into a native one, & then count up all the lines of assembly code written. Labels don't count.

→ Ex: `ex1.o` has 9 static instructions.

• every inst. is native except `la`

• `la $6, arr` resolves to `lui $1, 4097`
`ori $10, $1, 0`

What does the static inst. count tell us?

→ How much space the program takes up in memory.

What is the dynamic instruction count?

- The # of native instructions that actually get executed.
- E.g., if there is a for-loop that has 5 inst. & 3 iterations, the dynamic inst. count (DIC) = $5 \times 3 = 15$
- Ex: ex1.o has $3 + 6(5) + 1 = 34$ dynamic instructions.
- Gives us an idea of the runtime.
- The DIC & SIC are 2 metrics to look at when trying to improve program efficiency.

What does the DIC tell us?

How can we reduce the instruction count of ex1.o?

<pre> ex1.o .data arr: .word 10 15 100 2 4 .text la \$10 arr addi \$8 \$10 20 LoopBegin: beq \$10 \$8 LoopEnd lw \$6 0(\$10) addi \$6 \$6 1 sw \$6 0(\$10) addi \$10 \$10 4 j LoopBegin LoopEnd: </pre>	→	<pre> ex1.o .data arr: .word 10 15 100 2 4 .text la \$10 arr addi \$8 \$10 20 LoopBegin: lw \$8 0(\$6) addi \$8 \$8 1 sw \$8 0(\$6) addi \$6 \$6 4 bne \$6 \$7 LoopBegin </pre>
---	---	---

beq \$rs \$rt Label is native

the assembler always uses \$1 for storing temp. values when resolving instructions.

What are some pseudo instructions & their translations?

- `beq $rs imm Label` → `addi $1 $0 imm`
`beq $rs $1 Label`
- `bge $8 $9 Label` → `slt $1 $8 $9`
`beq $1 $0 Label` → `$1 will equal 1 if $8 < 9$, and 0 if $8 \geq 9$`
- `ble $9 $14 Label` → `slt $1 $14 $9`
`beq $1 $0 Label`

What is the load immediate instruction?

- Sets `$rd` to the immediate value:
`li $rd immediate` $R[rd] = \text{immediate}$

What is the move instruction?

- Copies the value of `$rs` into `$rd`:
`move $rd $rs` $R[rd] = R[rs]$

Example using these?

- `addi $8 $0 4` → `li $8 4`
- `addi $7 $6 0` → `move $7 $6`

Load Word and Store Word: Encoding & Hardware Support

RECALL: What are the "load word" and "store word" instructions?

→ LOAD WORD: $lw\ \$rt\ immediate\ (\$rs)$

- Store the data value at mem. address $[\$rs + imm]$ in $\$rt$

→ STORE WORD: $sw\ \$rt\ immediate\ (\$rs)$

- Store the data inside $\$rt$ at mem. address $[\$rs + imm.]$

How are lw and sw instructions encoded in binary?

→ Same as i -format!

Field:	opcode	rs	rt	immediate
Bits:	31:26	25:21	20:16	15:0
	(6 bits)	(5 bits)	(5 bits)	(16 bits)

- lw : opcode = 100011

- sw : opcode = 101011

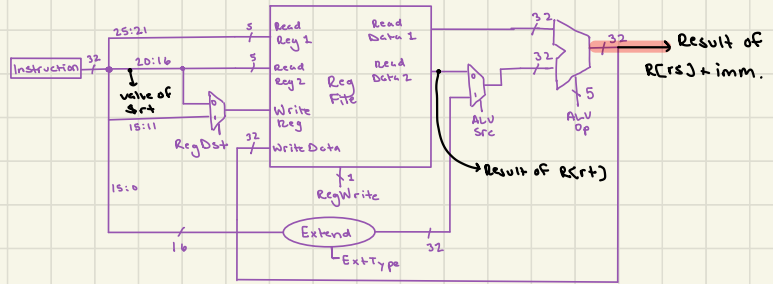
What "operation" is being done for lw and sw ?

→ Computing the address, aka $R[\$rs] + \text{SignExtend}(immediate)$

- We can use the ALU for this!

- Since it is addition, always sign extend the immediate.

How do we extend the Datapath to support lw and sw ?

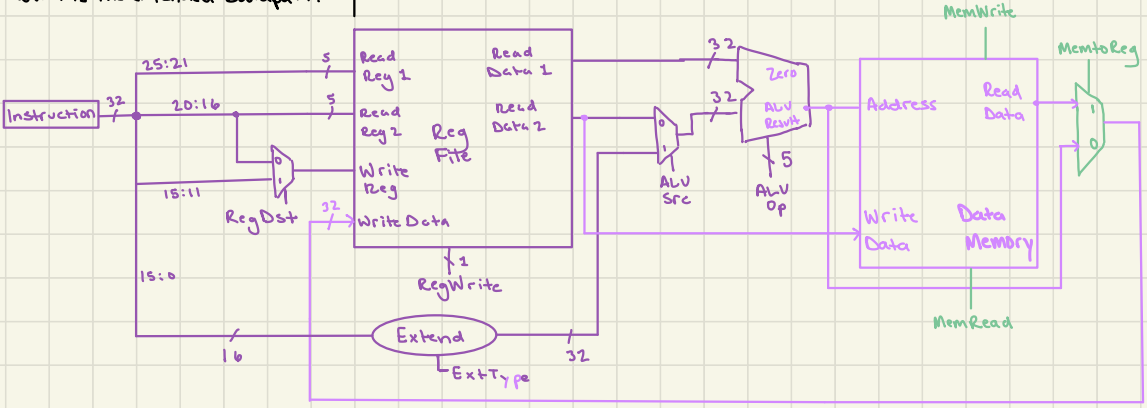


→ From the current DP, we can obtain 3 things:

- $R[rt]$, aka the value to place in memory for sw instructions
- $\$rt$, aka the register to store the mem. data for lw instructions
- $R[\$rs] + imm$, aka the memory address in question

→ To perform operations involving memory, we add a new component & new control signals!

What is the extended Datapath?



What is the MemRead control signal?

- Tells us if we have to read from memory
- lw = reading from Mem = MemRead = 1
- sw = writing to Mem = MemRead = 0
- When MemRead = 1, the "memory unit" obtains the value at "Address" in Main Mem, and places it in the "Read Data" slot that gets output

What is the MemWrite control signal?

- Tells us if we have to write to memory
- lw : MemWrite = 0
- sw : MemWrite = 1
- When MemWrite = 1, the memory unit takes the data placed in "Write Data" & puts it in Main Mem. at the address.

What is the MemtoReg control signal?

- Tells us what output to send to the Write Register: the output of the memory unit processes - aka the data from some addr. in MM - OR the output of the ALU.
- For add, sub, etc. ops we learned so far, we don't care about the memory unit, so MemtoReg = 0
- For lw, MemtoReg = 1

The Program Counter (PC)

What is the Stored Program Concept?

- RISCAL: Instructions (like `add $8 $7 $6`) are stored in memory as 32-bit binary numbers.
- Since inst are 32 bits, each take up 4 bytes in memory. In the Text section.
- The instructions for a given program are stored in subsequent memory addresses that increment by 4.

How does the code that we write get "resolved" for execution?

- 1) Pseudo-instructions are converted into native ones
- 2) Addresses for labels and offsets get resolved.

Example of the stored program concept?

Address in memory where the inst. is stored

machine code representation of the inst.

Assembled program (code that is actually executed)

og program (the code that you write)

Address	Code	Basic	Source
0x0003000	0x20060000	<code>addi \$6,\$0,0</code>	<code>la \$6 arr</code>
0x0003004	0x20c70014	<code>addi \$7,\$6,20</code>	<code>addi \$7,\$6,20</code>
0x0003008	0x8cc80000	<code>lw \$8,0(\$6)</code>	<code>lw \$8,0(\$6)</code>
0x000300c	0x10800001	<code>addi \$8,\$8,1</code>	<code>addi \$8,\$8,1</code>
0x0003010	0xacc80000	<code>sw \$8,0(\$6)</code>	<code>sw \$8,0(\$6)</code>
0x0003014	0x20c60004	<code>addi \$6,\$6,4</code>	<code>addi \$6,\$6,4</code>
0x0003018	0x14c7ffffb	<code>bne \$6,\$7,-5</code>	<code>bne \$6 \$7 LoopBegin</code>

notice that addr increment by 4

What is the program Counter (PC)?

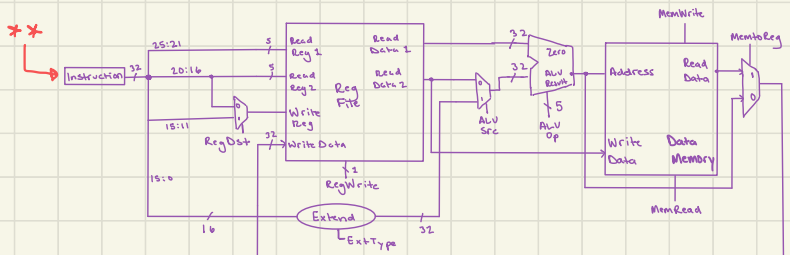
- A register that holds the address of the instruction being currently executed!
- Holds the 32-bit address in a 32-bit register

How is the value of the PC set?

- The PC is not inside the register file.
- For EX, when executing `addi $8 $8 1` from the prog. above, `PC=0x000300c`
- Our datapath executes one instruction per clock cycle. So on every clock cycle, the datapath will increment the PC value by 4 to fetch the next instruction.

WAIT so... what does this have to do with anything?

- The PC is what automates the process of sending 32-bit instructions as inputs to our data path!

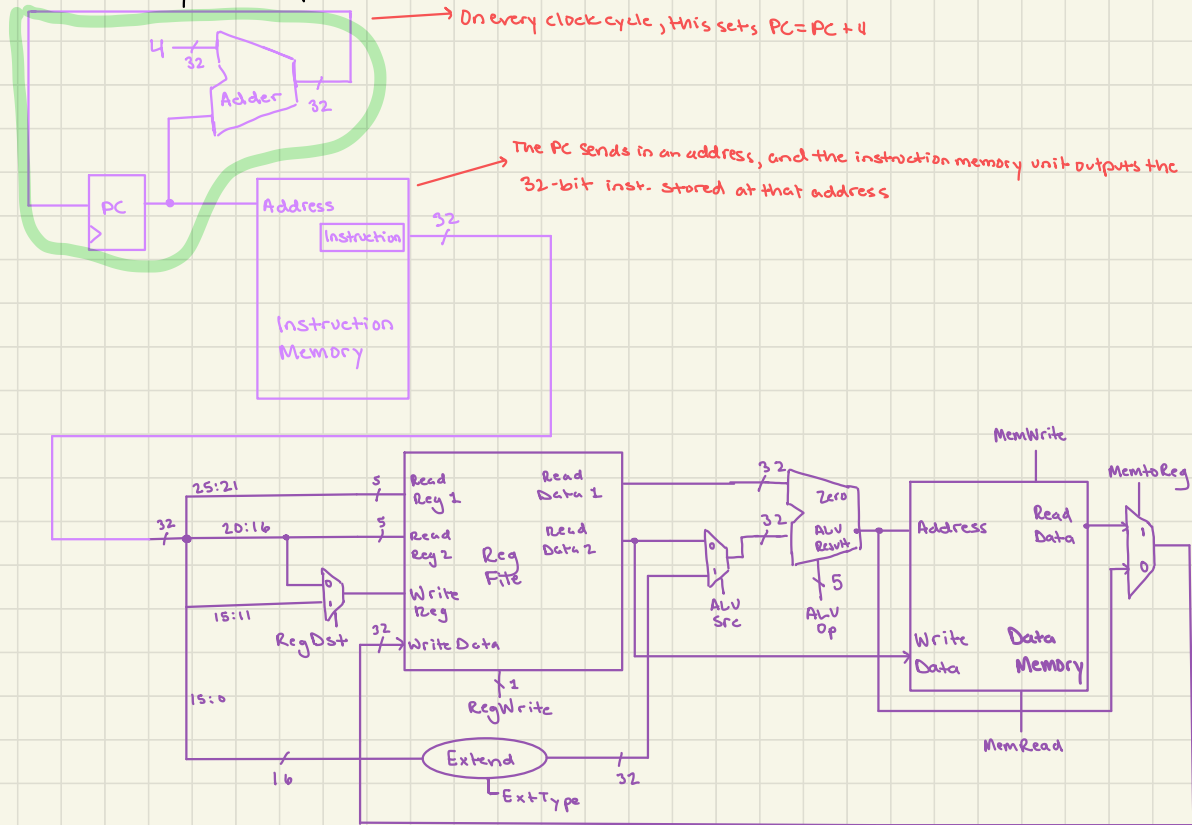


How do we use the PC to do this?

→ The PC holds the address of the instruction that should be executed in. On every clock cycle, we need to do the following:

- 1) Increment $PC = PC + 4$
- 2) Send the PC val into an "instruction memory" unit to extract the 32-bit instruction at that address
- 3) Send the 32-bit instruction into our datapath

How do we add the PC and Instruction Memory to our datapath?



Branch Instruction Encoding

Instruction	Opcode
beq	000100
bne	000101

How do we encode branch instructions in binary?

→ In i-format!

beq rs rt Label

→ For Ex, beq \$t0 \$t1 LoopBegin

Field:	opcode	rs	rt	immediate
Bits:	000100	01010	01011	????

What is the target instruction?

→ The instruction that the program will jump to if the branch is taken

What is the "byte offset to target instruction"?

→ The distance, in bytes, between the current branch instruction itself, and the target instruction.

→ The byte offset will always be a multiple of 4, which means that the last 2 bits of the byte offset will always be 0. Because 0x0, 4, 8, C, etc. all have 00 at the end.

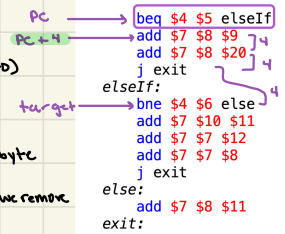
→ We actually calculate the byte offset relative to the next instruction, e.g. PC+4.

How do you generate the value of the immediate to encode a branch instruction?

1. Obtain the value of PC+4 - aka, the next inst. below the branch instruction, b/c the PC is currently on the branch instruction.

EX: beq \$t4 \$t5 elseif

2. Compute the byte offset between PC+4 and the target instruction (the inst. that we would branch to)



-AKA, (target addr) - (PC+4)

3. Since each inst. are 4 addresses away from each other, the byte

offset is always a multiple of 4, and its last 2 bits are 0. So we remove the last 2 bits.

$$\text{immediate} = (\text{target addr} - (\text{PC} + 4)) \gg 2$$

ANS = 0000 0000 0000 0100
remove these

How do we compute the branch target address given the immediate?

→ EX: The instruction 0x15320007 translates to:

opcode	rs	rt	immediate
000101	01001	10010	0000 0000 0000 0111

→ Assume the current PC = 0x000301C

1. Extract the immediate value: 0b0000 0000 0000 0111

2. Sign-extend it to 30 bits: 0b00 0000 0000 0000 0000 0000 0111

3. Append 2 0s to the end of it: 0b00 0000 0000 0000 0000 0000 0001 1100

4. Add PC+4 to it:

$$\begin{array}{r}
 0x000001C \\
 + 0x000301C \\
 + 0x0000004 \\
 \hline
 \text{The branch target address} \rightarrow 0x000303C
 \end{array}$$

this syntax means concatenate

$$\text{Branch Target Addr} = \text{PC} + 4 + \{30\text{-bit sign-ext imm.}, 0b00\}$$

- Datapath Support for Branching -

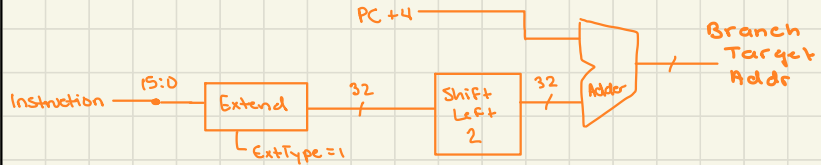
How do we obtain the branch target address?

→ RECALC: to generate the branch target address:

1. Sign Extend the imm. value to 32 bits
2. Append 2 0s
3. Add it to PC+4

How do we do this on the datapath?

→ We can achieve step 2 by left-shifting the 32-bit imm. by 2!

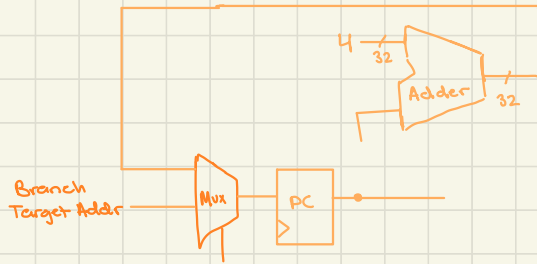


What do we do if we shouldn't branch upon a branch inst.?

→ The branch target addr. = the value we should set PC to **IF we want to branch.**

→ If not, set $PC = PC + 4$, like usual

→ **We will need a mux to make this decision:**



How will we set the control signal for the mux?

→ Using a piece of logic that interprets the **beq** and **bne** branch instructions.

→ This logic will output 1 to the Mux if we should branch, and 0 otherwise.

How do we support **beq** & **bne** on the datapath?

→ Using 2 new 1-bit control signals to indicate them:

• **beq**: 1 if inst. is beq, 0 otherwise

• **bne**: 1 if inst. is bne, 0 otherwise

→ **beq \$rs \$rt Label**: if $\$rs = \rt , we want to branch.

→ **bne \$rs \$rt Label**: if $\$rs \neq \rt , we want to branch.

→ SUMMARY: Output **1** to the Mux control signal IF:

• **beq** = 1 AND $\$rs = \rt OR

• **bne** = 1 AND $\$rs \neq \rt

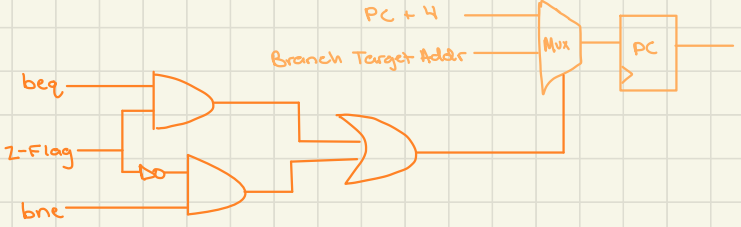
How do we implement this logic?

→ Using the Z-Flag/Zero Output of the ALU!!

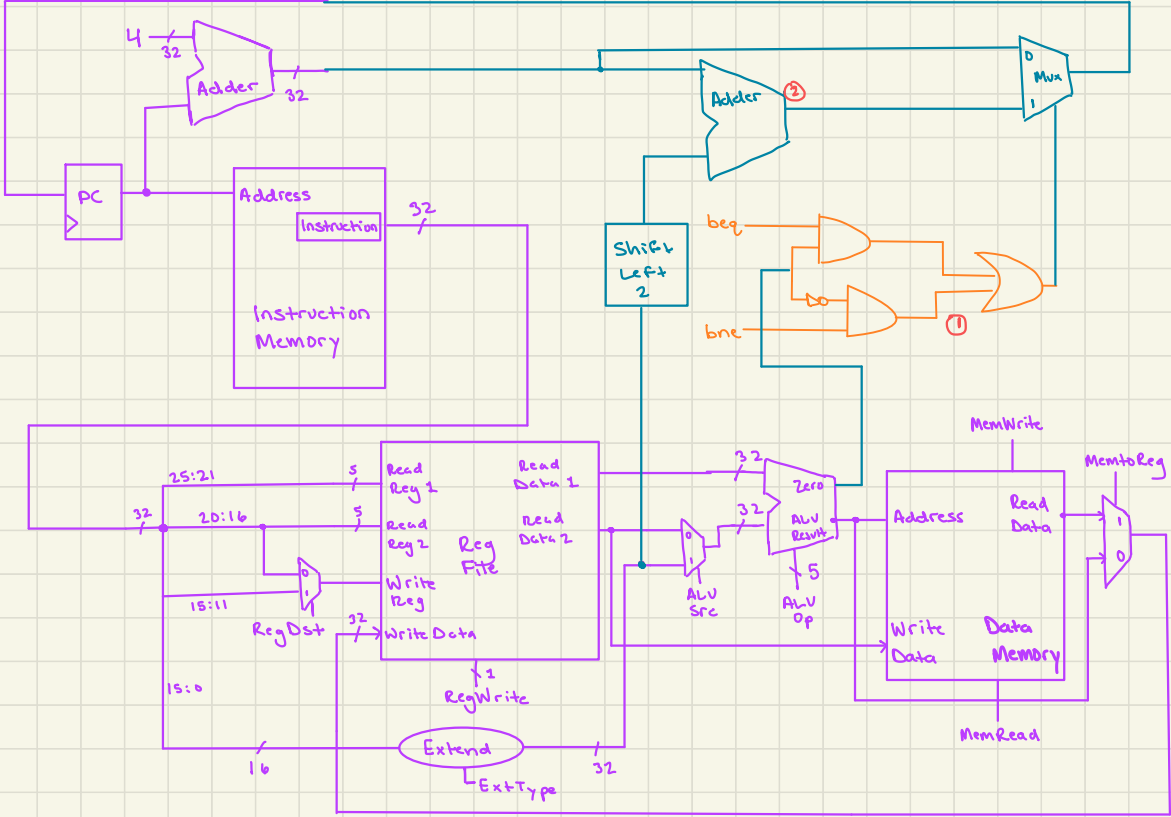
• RECALL: The ALU has a 1-bit "Zero Output" that is set to 1 whenever the ALU Output is 0.

→ If we set the ALUop to perform subtraction, the Z-Flag will tell us whether $\$rs = \rt

- If $\$rs = \rt , ALU Output = $\$rs - \$rt = 0$ and Zero Output = 1
- If $\$rs \neq \rt , ALU Output = $\$rs - \$rt \neq 0$ and Zero Output = 0



How do we implement all of this logic on the whole datapath?



① Uses beg & bne control signals to determine whether we should jump

② Computes the branch target addr. using the immediate.

Jump Instructions

RECALL: What are jump instructions?

→ An instruction that tells the program to jump to an address (specified by the label) to continue execution.

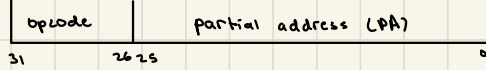
j Label

• area, set PC to the address of Label

opcode = 000010

How are jump instructions encoded in binary?

→ In a new instruction format called **j-format**:



What is the problem?

→ We can only fit a 26-bit address in the instruction, but the PC needs a 32-bit address.

How can we reduce the size of the address to 30 bits?

→ We can chop off/not store the last 2 bits, since the last 2 bits of any instruction address are always 0 — due to word alignment:

- 0b0...0000 (0)
- 0b0...0100 (4)
- 0b0...1000 (8)
- 0b0...1100 (12)
- 0b0...11100 (24)

How do we get rid of the last 4 extra bits?

→ By limiting the range of address values we can jump to

→ When setting the PA field of a jump instruction, we will chop off the top 4 bits of the target address.

→ When constructing the new PC addr given the PA field of a jump instruction, we will set the top 4 bits to be the top 4 bits of (PC+4).

Why do we do this?

→ Setting the top 4 bits to be those of (PC+4) gives us an address range that is nearby to where we are currently located.

• To jump outside this range, we'd have to use a different instruction.

What is the JumpAddress given a binary j-instruction?

$$\text{JumpAddr} = \underbrace{[(PC+4)[31:28]]}_{\text{bits 31:28 of PC+4}}, \underbrace{\text{partial_address}}_{\text{bits 25:0 of j-inst}}, 0b000$$

Given a program with a jump inst, do we compute the binary encoding of the partial address?

→ Ex, where the program begins at address 0x0003000:

```

1 .data
2 A: .word 8, 9, 14, 15
3
4 .text
5
6
7 la $5 A # address of A
8 addi $6 $0 0 # sum
9 addi $7 $0 4 # size of array
10 addi $8 $0 0 # i: loop counter
11
12 LoopBegin:
13 beq $8 $7 LoopEnd # exit loop when we reach the end of the array
14 sll $9 $8 2 # compute byte offset
15 add $9 $9 $5 # compute address of A[i]
16 lw $9 0($9) # get the value of A[i]
17 add $6 $6 $9 # sum += A[i]
18 addi $8 $8 1 # i++
19 j LoopBegin # jump back to start of array
20 LoopEnd:
    
```

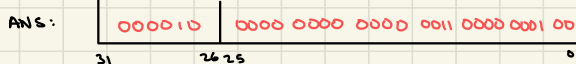
- Address:
- 0x0003000
 - 0x0003004
 - 0x0003008
 - 0x000300C
 - 0x0003010
 - 0x0003014
 - 0x0003018
 - 0x000301C
 - 0x0003020
 - 0x0003024
 - 0x0003028

1. Determine the address we want to jump to:

• The address of the first inst. after the label → 0x0003010

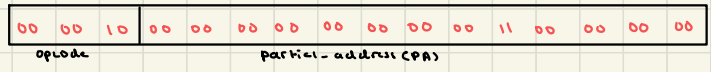
2. Convert it to binary: 0b0000 0000 0000 0000 0011 0000 0001 0000

3. Remove top 4 bits & bottom 2 bits: 0b0000 0000 0000 0011 0000 0001 00



SUMMARY: Given a 32-bit jump instruction, how do we calculate the jump address - aka, the address we have to set PC to?

→ EX: Instruction $0x08000300$ where $PC = 0x01001000$

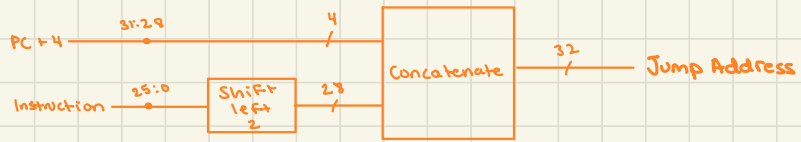


1. Extract the partial address bits $0b000000000000001100000000$
 2. Shift the value left by 2 (in order to append 2 zeroes) $0b00000000000000110000000000$
 3. Concatenate the upper 4 bits of $(PC+4)$ $PC+4 = 0x01001004$
 $(PC+4)[31:28] = 0b0000$
- ANS: jump address = $0b000000000000000000000000110000000000$

- Hardware Support for jump instructions -

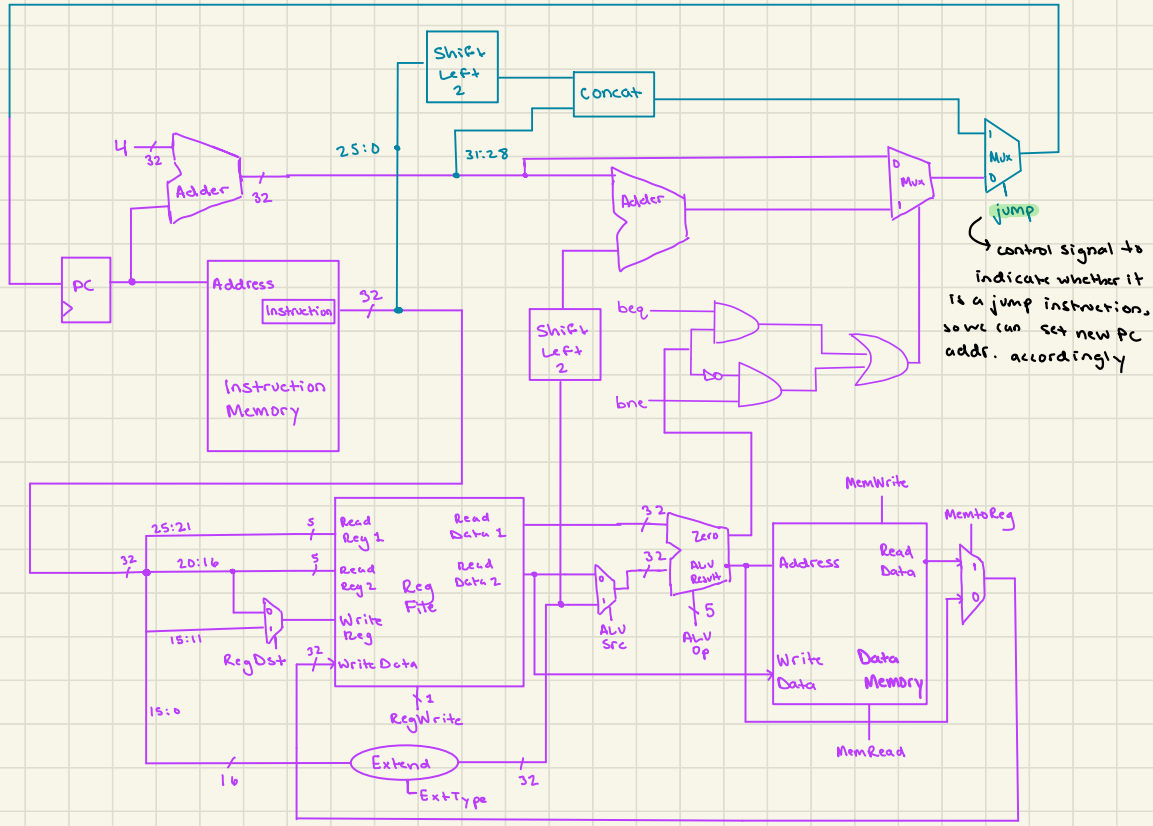
How do we perform the process described above, with our hardware?

1. Use a splitter to extract the PA bits of the 32-bit instruction.
2. Use a shift logic unit to append 2 Os by left shifting by 2.
3. Use a concatenate logic unit to concatenate the upper 4 bits of $(PC+4)$



How do we add this logic on the datapath?

→ Using a new 1-bit jump control signal:



What are all of the control signal values for a **jump** instruction?

Control Signal	Value
RegDst	X
Jump	1
beq	X
bne	X
MemRead	0*
MemtoReg	X
ALUOp	XXXXX
MemWrite	0*
ALUSrc	X
RegWrite	0*
ExtType	X

These 3 should never be "don't care"s. If we aren't writing to a reg, or reading or writing from MM, set these to 0.

What are all of the control signal values for a **beq** instruction?

Control Signal	Value
RegDst	X
Jump	0
beq	1
bne	0
MemRead	0*
MemtoReg	X
ALUOp	Subtraction
MemWrite	0*
ALUSrc	0
RegWrite	0*
ExtType	1

→ Why? Because we want to perform $\$rs - \rt .

What is the range of addresses we can jump to for a jump inst.?

→ For the given PC at the time of inst, we can jump to any address that has the same top 4 bits as current PC!

• Ex: If PC = $0x0D700000$, our range is $0x0D000000$ to $0x0DFFFFFFC$ (not $0x0DFFFFFF$ b/c of word alignment). That's a total of ~ 65, 273, 855 addresses!

What is the range of addresses we can jump to for a branch inst.?

→ If PC = $6x0D700000$, range is $0x0D6E0004$ to $0x0D720000$

Function Calls

- What is the **calling convention**? → A scheme for how functions receive arguments & return values
- **Caller**: A function that calls another function
 - **Callee**: A function called by another function
- What is the protocol for callers and callees? → Before making the function call, the caller places the parameter args in registers $\$a0 - \$a3$ (4-7)
- When finished, the callee places the return values in registers $\$v0 - \$v1$ (2-3)
- How does the callee know where to return when it finishes execution? → **Return Address**: the address of the inst. where the callee will return.
- When the caller calls the callee, they place the ret. addr. in a special register ($\$ra = \31) so the callee knows where to return.
- So how does a caller actually call a function? → Using the **jump-and-link (JAL)** instruction! To jump to a function.
- Calling a function means 2 things on the 'hardware' side:
 1. Setting PC to the address of the function being called
 2. Setting reg $\$ra$ to the address we want the program to resume executing at after the function is over.
- What does the **JAL** instruction do? → Syntax: `jal Label`
- Sets the PC to be the address of "Label" - aka the callee's label
 - AND automatically sets register $\$ra$ to be the return address
 - return addr = address of the inst. in main right after the JAL instruction
- How does the callee return execution to the caller when its done? → Using the **jump register** instruction!
- `jr $ra`
- sets the PC to be the address value stored at $\$ra$

Summary: What are some special registers and their uses?

Number	Name	Uses
$\$0$	$\$zero$	• always holding the value 0
$\$1$	$\$at$	• Resolving pseudo instructions
$\$2$	$\$v0$	• Where callee places return values
		• Where we place the "operation number" when performing <code>syscall</code>
		• Where the OS places value read from user input after <code>syscall</code>
$\$3$	$\$v1$	• Where callee places return values
$\$4$	$\$a0$	• Where caller places function parameter args
		• Where we place the value we want to print when performing <code>syscall</code>
$\$5$	$\$a1$	} Where caller places function parameter args
$\$6$	$\$a2$	
$\$7$	$\$a3$	
$\$31$	$\$ra$	• Stores return address upon a function call

Example program performing a function call?

→ C program:

```
int sum3(int a, int b, int c) {
    return a + b + c;
}

int main() {
    int a = 1;
    int b = 5;
    int c = 8;

    int y = sum3(a, b, c);
    printf("%d\n", y);
}
```

→ Assembly program (MIPS):

PC Address	.text	
	main:	
0x00003000	addi \$a0 \$0 1	} store int a, b, and c in the argument registers a0, a1, a2
0x00003004	addi \$a1 \$0 5	
0x00003008	addi \$a2 \$0 8	
0x0000300c	jal sum3	→ jump to Label sum3 AND set \$ra to return instruction addr (PC+4): \$ra = 0x00003010
0x00003010	addi \$a0 \$v0 0	→ We know that the return value y is stored in \$v0. Move it to \$a0 (\$4) so we can print
0x00003014	addi \$v0 \$0 1	→ Set \$v0 (\$2) to 1 to indicate "print an integer from \$4"
0x00003018	syscall	
0x0000301c	addi \$v0 \$0 10	→ set \$2 to 10 for "exit program".
0x00003020	syscall	NECESSARY! otherwise prog will keep executing the stuff below
0x00003024	sum3:	
0x00003028	add \$v0 \$a0 \$a1	} add the values in all arg. registers and store ans in \$v0
0x0000302c	add \$v0 \$v0 \$a2	
0x00003030	jr \$ra	→ return to main; resume execution at PC = \$ra = 0x00003010

Shift Instructions

What is the "shift left logical" instruction?

`sll $rd, $rt, $amt`

- Shift the contents of $\$rt$ to the LEFT by $\$amt$ & store result in $\$rd$
- $R[rd] = R[rt] \ll \$amt$

Why is `sll` useful?

→ For computing the byte offset when we need to index an array!

→ RECALL: $x \ll y = x \cdot 2^y$

→ Given an array, if we want to access index k , we shift k left by 2 and add the result to our base address!

Example?

`.data`

`arr: .word 1 5 10 15 20`

`.text`

`addi $a1, $0, 3` → we want `arr[3]`

`la $5, arr` → stores address of `arr[0]` in $\$5$

`sll $a1, $a1, 2` → $\$a1 = 3 \ll 2 = 3 \cdot 4 = 12$, the byte offset

`add $6, $5, $a1` → stores addr of `arr[0] + 12 = addr of arr[3]` in $\$6$

`lw $7, 0($6)` → $\$7 = arr[3] = 15$

What are the 2 right shift instructions?

→ Shift right logical:

`srl $rd, $rt, $amt`

- Pad left side with 0s

$$R[rd] = R[rt] \gg \$amt$$

→ Shift right arithmetic:

`sra $rd, $rt, $amt`

- Pad left side with MSB to preserve sign

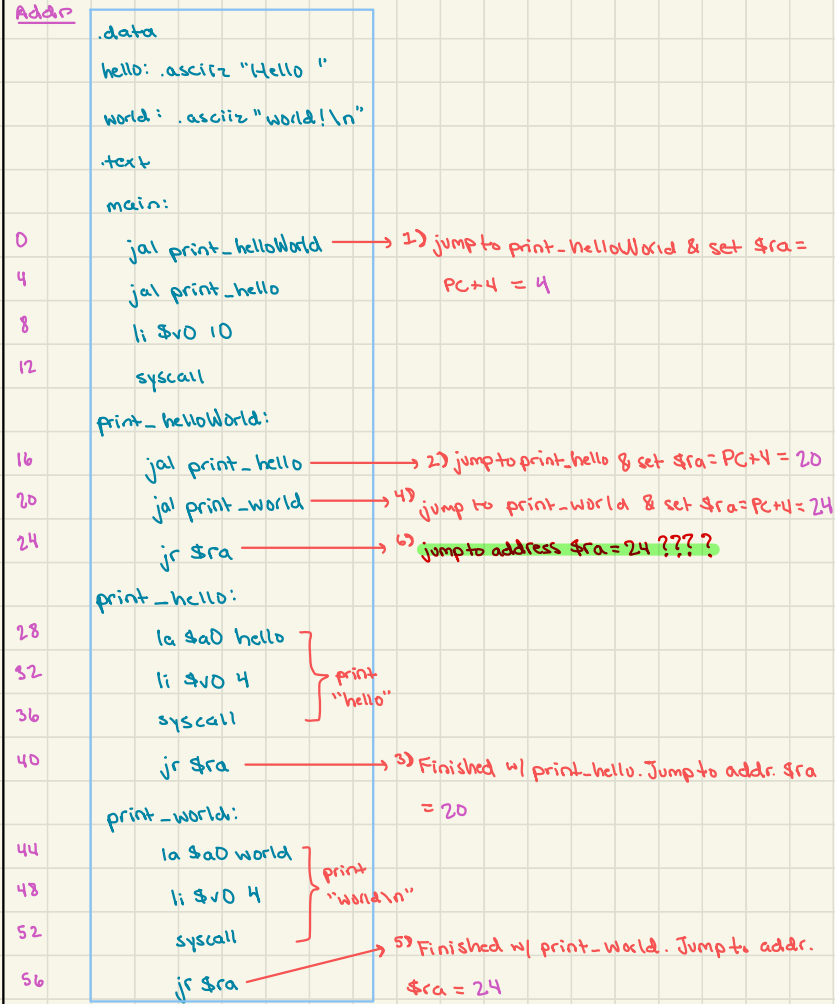
Storing Variables on the Stack

When is the register $\$ra$ not enough to help return from function calls?

→ When we have nested function calls.

→ **RECALL:** $\$ra$ is used to store the return address when a function call is made using `jal`.

→ Ex:



What has happened in this example?

→ At step 6, we are ready to return to main, as it was the one that called `print_helloWorld`

• We expected `jr $ra` to take us to $PC = 20$

→ BUT, at this point $\$ra$ was overwritten, so we have lost the return address to main and are stuck infinitely jumping to address 24 !!

How can we fix this problem?

→ By saving the ret. addr. somewhere else before we overwrite it. Saving it to another register is not a sustainable solution b/c it can still get overwritten.

→ Instead, we will store the RA on the stack

When would we save the RA on the stack?

- When the callee begins, save the value of \$ra to the stack
- When the callee is ready to return, get the val. from the stack & place it back into \$ra, before the jr \$ra instruction.

```
print_helloWorld:
    # save $ra=20 on stack
    jal print_hello
    jal print_world
    # restore ra=20 from stack
    jr $ra
```

→ Now, step 6 (from EX on prev page) will take us back to main by setting PC=20!

- The stack pointer -

RECALL: What is the stack?

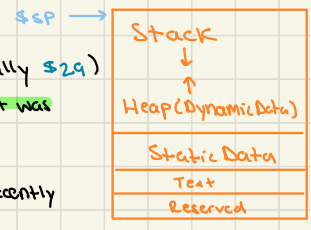
→ An area of memory used to store temporary function data, such as local variables.

What is the stack pointer?

→ The stack grows DOWN (high → low mem. addresses)

→ A special register \$sp (which is actually \$29) that is used to point to the last item that was placed on the stack

- \$sp holds the address of the most recently placed val on the stack.



How do we store a new item on the stack?

→ The stack grows down, so we must decrement the addr. value stored in \$sp by 4 to move to a new spot on the stack.

Example?

① Current stack: \$sp = 0x00001034 and points to the last val added, 12.

Address	Data
\$sp → 0x00001034	12
0x00001030	
0x0000102c	
0x00001028	
0x00001024	

→ To store a new item, the value of \$ra:

- 1) Decrement sp by 4
- 2) store value in memory

```
addi $sp $sp -4
sw $ra 0($sp)
```

② Updated stack:

Address	Data
0x00001034	12
\$sp → 0x00001030	0x00003004
0x0000102c	
0x00001028	
0x00001024	

How do we remove an item from the stack?

- "Remove" aka to fetch back the value we stored
- When we want to restore a val from the stack, we simply
 1. read it into a register (with lw)
 2. and then move the stack pointer back up
- We don't have to "clear" the value from memory; by moving sp back up, we've "unallocated" that stack space. Blk next time we want to store on the stack, we will just overwrite this value.

Example?

→ ① Current stack: \$sp = 0x00001030

→ ② Read / "remove" / restore value of ra:

```
lw $ra 0($sp)
addi $sp $sp 4
```

Address	Data
0x00001034	12
\$SP → 0x00001030	0x00003004
0x0000102c	
0x00001028	

→ ③ Updated stack:

Address	Data
\$SP → 0x00001034	12
0x00001030	0x00003004
0x0000102c	
0x00001028	
0x00001024	

- Calling Convention with Function - saved values -

RECALL: Why do we have a calling convention?

- Imagine each function in a program was written by a different person, & they didn't know how the other person wrote their func, or what registers they used.
- Then, we can't necessarily rely on a val. saved in a reg. to not be overwritten, so we must have rules (conventions for preserving all parts of the program (for ex, param args in \$a0-\$a3 & return values in \$v0-\$v1)

What is another example of when we might need to save values to the stack?

→ If a callee overwrites caller-saved values! Ex:

```
main:
    addi $s1 $0 0
    la $s2 array-1
    la $s3 array-2
    lw $a0 0($s2)
    jal fun
    Fun:
    addi $s1 $a0 1
    sub $s2 $0 $a0
    mul $v0 $s1 $s2
    jr $ra
```

fun overwrites the value in \$s1 and \$s2!



How can we fix this problem?

→ By having the callee save the values in the caller-saved registers to the stack before overwriting them.

Example of the fix?

→ Then, before returning to caller, the callee restores the values from the stack!

→ The callee "Fun":

```

Fun:
  addi $sp $sp -8
  sw $s2 4($sp)
  sw $s1 0($sp)

  addi $s1 $a0 1
  sub $s2 $b $a0
  mul $v0 $s1 $s2

  lw $s1 0($sp)
  lw $s2 4($sp)
  addi $sp $sp 8

  jr $ra
  
```

Annotations:

- allocate space to store 2 reg. values on the stack by decrementing \$sp by 8.
- Store contents of \$s2 at addr [sp addr + 4]
- Store contents of \$s1 at addr [sp addr]
- Now, callee can use registers \$s1 and \$s2 w/o worrying
- Before returning to the caller, restore the values of \$s1 and \$s2 by reading from the stack
- move sp back up to "deallocate"

What is a caller saved register in the calling convention?

→ Registers that the caller uses, BUT expects that the callee may overwrite them.

• The callee CAN overwrite these registers w/o saving them

→ If caller intends to use them after a func. call, it is the caller's responsibility to save them.

What are callee saved registers?

→ Registers that the caller expects the callee to NOT overwrite - caller expects the vals in these registers to be saved.

→ If the callee wants to use these registers, it is the callee's responsibility to save their values to the stack and then later restore them.

What are the designated caller- and callee- saved registers in MIPS?

Registers	Name	Use	Type
\$2-3	\$v0-\$v1	callee stores return values	caller-saved
\$4-7	\$a0-\$a3	caller stores param. args	caller-saved
\$8-15	\$t0-\$t7	registers for callee to use	caller-saved
\$16-23	\$s0-\$s7	registers for callee to use	callee-saved
\$24-25	\$t8-\$t9	registers for callee to use	caller-saved
\$31	\$ra	return address pointer	caller-saved

If the caller is ALSO a callee, they are responsible for saving their own RA before calling another func.