

---

# Computer Architecture

Ngo Lam Trung & Pham Ngoc Hung  
Faculty of Computer Engineering  
School of Information and Communication Technology (SoICT)  
Hanoi University of Science and Technology  
E-mail: [trungnl, hungpn]@soict.hust.edu.vn

---

## Chapter 3: Arithmetic for Computers

[with materials from *Computer Organization and Design, 4<sup>th</sup> Edition*,  
Patterson & Hennessy, © 2008, MK  
and M.J. Irwin's presentation, PSU 2008]

# Content

---

- ❑ (Super) Basics of logic design
- ❑ Integer representation and arithmetic
- ❑ Floating point number representation and arithmetic

# What are stored inside computer?

---

- ❑ Data, of course!
- ❑ Data is represented as binary numbers.
- ❑ How binary numbers are treated by CPUs?
  - ❑ By logic circuits
  - ❑ Integers
    - Unsigned
    - Signed
  - ❑ Floating point numbers
    - Single precision
    - Double precision
    - Other formats

## Basics of logic design (Appendix B)

---

- ❑ Boolean logic: logic variable and operators
- ❑ Logic variable: values of 1 (TRUE) or 0 (FALSE)
- ❑ Basic operators: AND, OR, NOT
  - ❑ A AND B :  $A \cdot B$  hay  $AB$
  - ❑ A OR B :  $A + B$
  - ❑ NOT A :  $\overline{A}$
  - ❑ Order: NOT > AND > OR
- ❑ Additional operators: NAND, NOR, XOR
  - ❑ A NAND B:  $\overline{A \cdot B}$
  - ❑ A NOR B :  $\overline{A + B}$
  - ❑ A XOR B:  $A \oplus B = A \bullet \overline{B} + \overline{A} \bullet B$

# Truth tables

---

| A | B | A AND B<br>$A \bullet B$ |
|---|---|--------------------------|
| 0 | 0 | 0                        |
| 0 | 1 | 0                        |
| 1 | 0 | 0                        |
| 1 | 1 | 1                        |

| A | B | A OR B<br>$A + B$ |
|---|---|-------------------|
| 0 | 0 | 0                 |
| 0 | 1 | 1                 |
| 1 | 0 | 1                 |
| 1 | 1 | 1                 |

| A | NOT A<br>$\neg A$ |
|---|-------------------|
| 0 | 1                 |
| 1 | 0                 |

Unary operator NOT

| A | B | A NAND B<br>$A \bullet B$ |
|---|---|---------------------------|
| 0 | 0 | 1                         |
| 0 | 1 | 1                         |
| 1 | 0 | 1                         |
| 1 | 1 | 0                         |

| A | B | A XOR B<br>$A \oplus B$ |
|---|---|-------------------------|
| 0 | 0 | 0                       |
| 0 | 1 | 1                       |
| 1 | 0 | 1                       |
| 1 | 1 | 0                       |

| A | B | A NOR B<br>$A + B$ |
|---|---|--------------------|
| 0 | 0 | 1                  |
| 0 | 1 | 0                  |
| 1 | 0 | 0                  |
| 1 | 1 | 0                  |

# Laws of Boolean algebra

---

$$A \cdot B = B \cdot A$$

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

$$1 \cdot A = A$$

$$A \cdot \bar{A} = 0$$

$$0 \cdot A = 0$$

$$A \cdot A = A$$

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

$$\overline{A \cdot B} = \bar{A} + \bar{B} \text{ (DeMorgan's law)}$$

$$A + B = B + A$$

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

$$0 + A = A$$

$$A + \bar{A} = 1$$


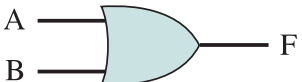
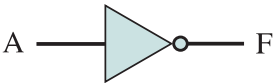

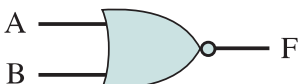
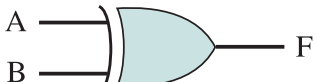
$$1 + A = 1$$

$$A + A = A$$

$$A + (B + C) = (A + B) + C$$

$$\overline{A + B} = \bar{A} \cdot \bar{B} \text{ (DeMorgan's law)}$$

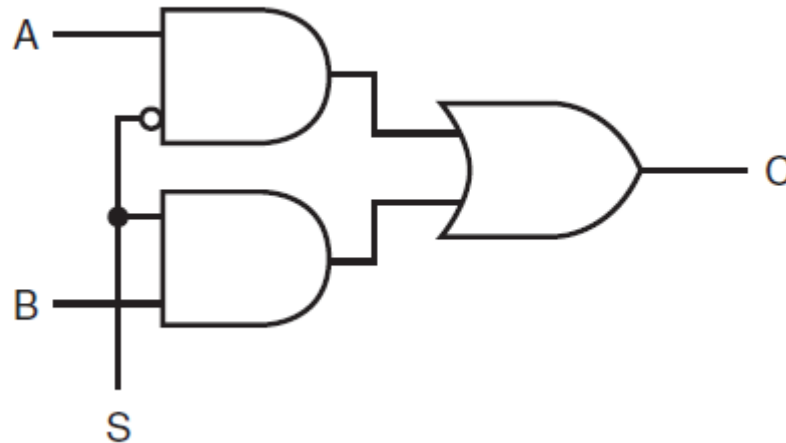
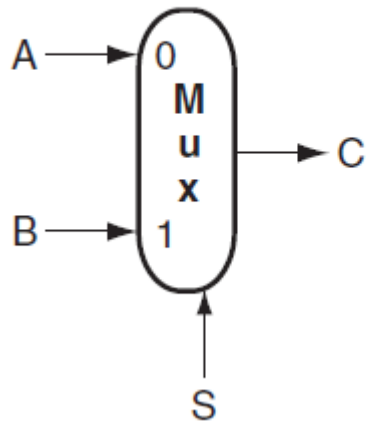
# Logic gates

| Name | Graphical Symbol  | Algebraic Function                   | Truth Table  |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
|------|---|--------------------------------------|--|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| AND  |    | $F = A \bullet B$<br>or<br>$F = AB$  | <table><tr><th>A</th><th>B</th><th>F</th></tr><tr><td>0</td><td>0</td><td>0</td></tr><tr><td>0</td><td>1</td><td>0</td></tr><tr><td>1</td><td>0</td><td>0</td></tr><tr><td>1</td><td>1</td><td>1</td></tr></table> | A | B | F | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 1 | 1 |
| A    | B   | F                                    |  |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| 0    | 0   | 0                                    |  |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| 0    | 1   | 0                                    |  |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| 1    | 0   | 0                                    |  |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| 1    | 1   | 1                                    |  |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| OR   |    | $F = A + B$                          | <table><tr><th>A</th><th>B</th><th>F</th></tr><tr><td>0</td><td>0</td><td>0</td></tr><tr><td>0</td><td>1</td><td>1</td></tr><tr><td>1</td><td>0</td><td>1</td></tr><tr><td>1</td><td>1</td><td>1</td></tr></table> | A | B | F | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 1 | 1 | 1 |
| A    | B   | F                                    |  |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| 0    | 0   | 0                                    |  |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| 0    | 1   | 1                                    |  |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| 1    | 0   | 1                                    |  |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| 1    | 1   | 1                                    |  |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| NOT  |    | $F = \overline{A}$<br>or<br>$F = A'$ | <table><tr><th>A</th><th>F</th></tr><tr><td>0</td><td>1</td></tr><tr><td>1</td><td>0</td></tr></table>   | A | F | 0 | 1 | 1 | 0 |   |   |   |   |   |   |   |   |   |
| A    | F   |                                      |  |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| 0    | 1   |                                      |  |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| 1    | 0   |                                      |  |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| NAND |    | $F = \overline{AB}$                  | <table><tr><th>A</th><th>B</th><th>F</th></tr><tr><td>0</td><td>0</td><td>1</td></tr><tr><td>0</td><td>1</td><td>1</td></tr><tr><td>1</td><td>0</td><td>1</td></tr><tr><td>1</td><td>1</td><td>0</td></tr></table> | A | B | F | 0 | 0 | 1 | 0 | 1 | 1 | 1 | 0 | 1 | 1 | 1 | 0 |
| A    | B   | F                                    |  |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| 0    | 0   | 1                                    |  |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| 0    | 1   | 1                                    |  |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| 1    | 0   | 1                                    |  |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| 1    | 1   | 0                                    |  |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| NOR  |  | $F = \overline{A + B}$               | <table><tr><th>A</th><th>B</th><th>F</th></tr><tr><td>0</td><td>0</td><td>1</td></tr><tr><td>0</td><td>1</td><td>0</td></tr><tr><td>1</td><td>0</td><td>0</td></tr><tr><td>1</td><td>1</td><td>0</td></tr></table> | A | B | F | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 1 | 0 |
| A    | B   | F                                    |  |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| 0    | 0   | 1                                    |  |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| 0    | 1   | 0                                    |  |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| 1    | 0   | 0                                    |  |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| 1    | 1   | 0                                    |  |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| XOR  |  | $F = A \oplus B$                     | <table><tr><th>A</th><th>B</th><th>F</th></tr><tr><td>0</td><td>0</td><td>0</td></tr><tr><td>0</td><td>1</td><td>1</td></tr><tr><td>1</td><td>0</td><td>1</td></tr><tr><td>1</td><td>1</td><td>0</td></tr></table> | A | B | F | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 1 | 1 | 0 |
| A    | B   | F                                    |  |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| 0    | 0   | 0                                    |  |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| 0    | 1   | 1                                    |  |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| 1    | 0   | 1                                    |  |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| 1    | 1   | 0                                    |  |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |



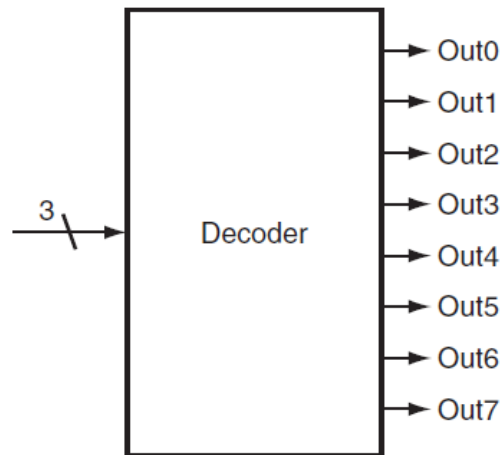
## Example: multiplexor

- ❑ Depending on  $S$ , output  $C$  is equal to one of the two inputs  $A$ ,  $B$
- ❑ Explain how this circuit works?



## Example: 3-to-8 decoder

❑ Very important in address decoder circuits



a. A 3-bit decoder

| Inputs |    |    | Outputs |      |      |      |      |      |      |      |
|--------|----|----|---------|------|------|------|------|------|------|------|
| 12     | 11 | 10 | Out7    | Out6 | Out5 | Out4 | Out3 | Out2 | Out1 | Out0 |
| 0      | 0  | 0  | 0       | 0    | 0    | 0    | 0    | 0    | 0    | 1    |
| 0      | 0  | 1  | 0       | 0    | 0    | 0    | 0    | 0    | 1    | 0    |
| 0      | 1  | 0  | 0       | 0    | 0    | 0    | 0    | 1    | 0    | 0    |
| 0      | 1  | 1  | 0       | 0    | 0    | 0    | 1    | 0    | 0    | 0    |
| 1      | 0  | 0  | 0       | 0    | 0    | 1    | 0    | 0    | 0    | 0    |
| 1      | 0  | 1  | 0       | 0    | 1    | 0    | 0    | 0    | 0    | 0    |
| 1      | 1  | 0  | 0       | 1    | 0    | 0    | 0    | 0    | 0    | 0    |
| 1      | 1  | 1  | 1       | 0    | 0    | 0    | 0    | 0    | 0    | 0    |

b. The truth table for a 3-bit decoder

# Unsigned Binary Integers

---

- ❑ Using n-bit binary number to represent non-negative integer

$$\begin{aligned} X &= X_{n-1}X_{n-2}\dots X_1X_0 \\ &= X_{n-1}2^{n-1} + X_{n-2}2^{n-2} + \dots + X_12^1 + X_02^0 \end{aligned}$$

- ❑ Range: 0 to  $+2^n - 1$

- ❑ Example

$$\begin{aligned} &0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 1011_2 \\ &= 0 + \dots + 1 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 1 \times 2^0 \\ &= 0 + \dots + 8 + 0 + 2 + 1 = 11_{10} \end{aligned}$$

- ❑ Data range using 32 bits

$$0 \text{ to } 2^{32}-1 = 4,294,967,295$$

## Eg: 32 bit Unsigned Binary Integers

| Hex        | Binary   | Decimal    |
|------------|----------|------------|
| 0x00000000 | 0...0000 | 0          |
| 0x00000001 | 0...0001 | 1          |
| 0x00000002 | 0...0010 | 2          |
| 0x00000003 | 0...0011 | 3          |
| 0x00000004 | 0...0100 | 4          |
| 0x00000005 | 0...0101 | 5          |
| 0x00000006 | 0...0110 | 6          |
| 0x00000007 | 0...0111 | 7          |
| 0x00000008 | 0...1000 | 8          |
| 0x00000009 | 0...1001 | 9          |
|            | ...      |            |
| 0xFFFFFFF0 | 1...1111 | $2^{32}-1$ |
| 0xFFFFFFF1 | 1...1110 | $2^{32}-2$ |
| 0xFFFFFFF2 | 1...1101 | $2^{32}-3$ |
| 0xFFFFFFF3 | 1...1100 | $2^{32}-4$ |

## Exercise

---

### ❑ Convert to 32-bit integers

25 = 0000 0000 0000 0000 0000 0000 0001 1001

125 = 0000 0000 0000 0000 0000 0000 0111 1101

255 = 0000 0000 0000 0000 0000 0000 1111 1111

### ❑ Convert 32-bit integers to decimal

0000 0000 0000 0000 0000 0000 1100 1111 = 207

0000 0000 0000 0000 0000 0001 0011 0011 = 307

# Signed binary integers

- Using n-bit binary number to represent integer, including negative values

$$\begin{aligned} X &= X_{n-1}X_{n-2}\dots X_1X_0 \\ &= -X_{n-1}2^{n-1} + X_{n-2}2^{n-2} + \dots + X_12^1 + X_02^0 \end{aligned}$$

- Range:  $-2^{n-1}$  to  $+2^{n-1} - 1$**

- Example**

$$\begin{aligned} &1111\ 1111\ 1111\ 1111\ 1111\ 1111\ 1111\ 1100_2 \\ &= -1 \times 2^{31} + 1 \times 2^{30} + \dots + 1 \times 2^2 + 0 \times 2^1 + 0 \times 2^0 \\ &= -2,147,483,648 + 2,147,483,644 = -4_{10} \end{aligned}$$

- Using 32 bits**

$$-2,147,483,648 \text{ to } +2,147,483,647$$

## Signed integer negation

---

- ❑ Given  $x = x_{n-1}x_{n-2} \dots x_1x_0$ , how to calculate  $-x$ ?
- ❑ Let  $\bar{x}$  = 1's complement of  $x$

$$\bar{x} = 1111 \dots 11_2 - x$$

$$(1 \rightarrow 0, 0 \rightarrow 1)$$

Then

$$\bar{x} + x = 1111 \dots 112 = -1$$

$$\rightarrow \bar{x} + 1 = -x$$

- ❑ **Example: find binary representation of -2**

$$+2 = 0000 \ 0000 \dots 0010_2$$

$$\begin{aligned} -2 &= 1111 \ 1111 \dots 1101_2 + 1 \\ &= 1111 \ 1111 \dots 1110_2 \end{aligned}$$

# Signed binary negation

$$-2^3 =$$

$$-(2^3 - 1) =$$

complement all the bits

1011

0101

and add a 1

and add a 1

0110

1010

complement all the bits

| 2'sc binary | decimal |
|-------------|---------|
| 1000        | -8      |
| 1001        | -7      |
| 1010        | -6      |
| 1011        | -5      |
| 1100        | -4      |
| 1101        | -3      |
| 1110        | -2      |
| 1111        | -1      |
| 0000        | 0       |
| 0001        | 1       |
| 0010        | 2       |
| 0011        | 3       |
| 0100        | 4       |
| 0101        | 5       |
| 0110        | 6       |
| 0111        | 7       |



## Exercise

---

❑ Find 16 bit signed integer representation of

$$16 = 0000\ 0000\ 0001\ 0000$$

$$-16 = 1111\ 1111\ 1111\ 0000$$

$$100 = 0000\ 0000\ 0110\ 0100$$

$$-100 = 1111\ 1111\ 1001\ 1100$$

## Sign extension

---

- ❑ Given n-bit integer  $x = x_{n-1}x_{n-2} \dots x_1x_0$
- ❑ Find corresponding m-bit representation ( $m > n$ ) with the same numeric value

$$x = x_{m-1}x_{m-2} \dots x_1x_0$$

- ❑ → Replicate the sign bit to the left

- ❑ Examples: 8-bit to 16-bit

+2: 0000 0010 => 0000 0000 0000 0010

-2: 1111 1110 => 1111 1111 1111 1110

# Addition and subtraction

---

## ❑ Addition

- ❑ Similar to what you do to add two numbers manually
- ❑ Digits are added bit by bit from right to left
- ❑ Carries passed to the next digit to the left

## ❑ Subtraction

- ❑ Negate the second operand then add to the first operand

$$\begin{array}{r} + \quad 0000 \ 0000 \ 0000 \ 0000 \ 0000 \ 0000 \ 0000 \ 0111_{\text{two}} = 7_{\text{ten}} \\ \quad 0000 \ 0000 \ 0000 \ 0000 \ 0000 \ 0000 \ 0000 \ 0110_{\text{two}} = 6_{\text{ten}} \\ \hline \quad 0000 \ 0000 \ 0000 \ 0000 \ 0000 \ 0000 \ 0000 \ 1101_{\text{two}} = 13_{\text{ten}} \end{array}$$

## Examples

---

- ❑ All numbers are 8-bit signed integer

$$12 + 8 =$$

$$122 + 8 =$$

$$122 + 80 =$$

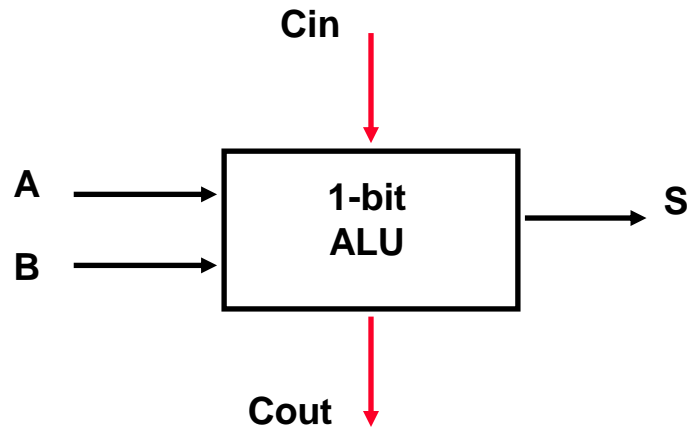
# Dealing with Overflow

- ❑ Overflow occurs when the result of an operation cannot be represented in 32-bits, i.e., when the sign bit contains a **value** bit of the result and not the proper **sign** bit
- ❑ When adding operands with different signs or when subtracting operands with the same sign, overflow can *never* occur

| Operation | Operand A | Operand B | Result indicating overflow |
|-----------|-----------|-----------|----------------------------|
| A + B     | $\geq 0$  | $\geq 0$  | $< 0$                      |
| A + B     | $< 0$     | $< 0$     | $\geq 0$                   |
| A - B     | $\geq 0$  | $< 0$     | $< 0$                      |
| A - B     | $< 0$     | $\geq 0$  | $\geq 0$                   |

# Adder implementation

## ❑ 1-bit full adder



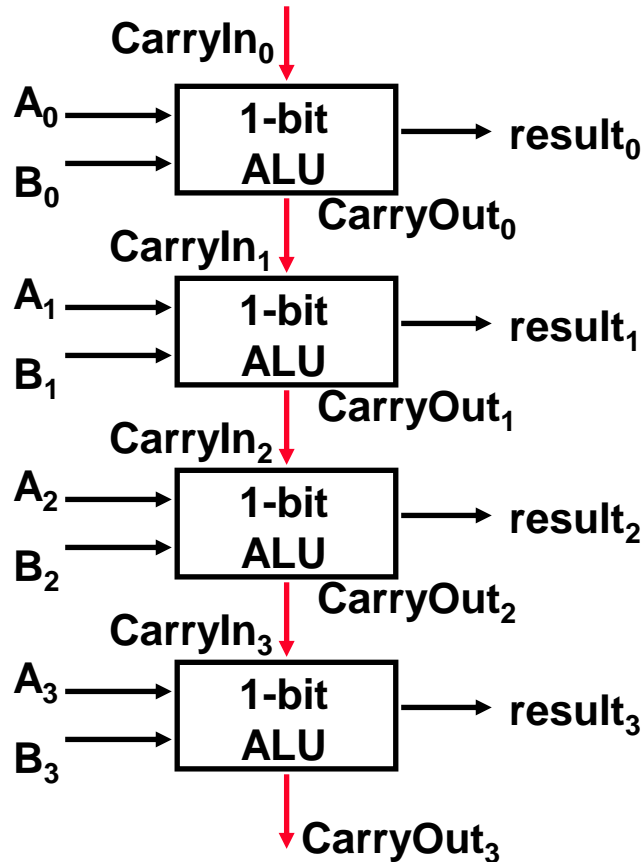
| Inputs |   |                 | Outputs |                  |
|--------|---|-----------------|---------|------------------|
| A      | B | C <sub>in</sub> | S       | C <sub>out</sub> |
| 0      | 0 | 0               | 0       | 0                |
| 0      | 0 | 1               | 1       | 0                |
| 0      | 1 | 0               | 1       | 0                |
| 0      | 1 | 1               | 0       | 1                |
| 1      | 0 | 0               | 1       | 0                |
| 1      | 0 | 1               | 0       | 1                |
| 1      | 1 | 0               | 0       | 1                |
| 1      | 1 | 1               | 1       | 1                |

❑  $S = C_{in} \oplus (A \oplus B)$

❑  $C_{out} = AB + BC_{in} + AC_{in}$

# Adder implementation

## ❑ N-bit ripple-carry adder

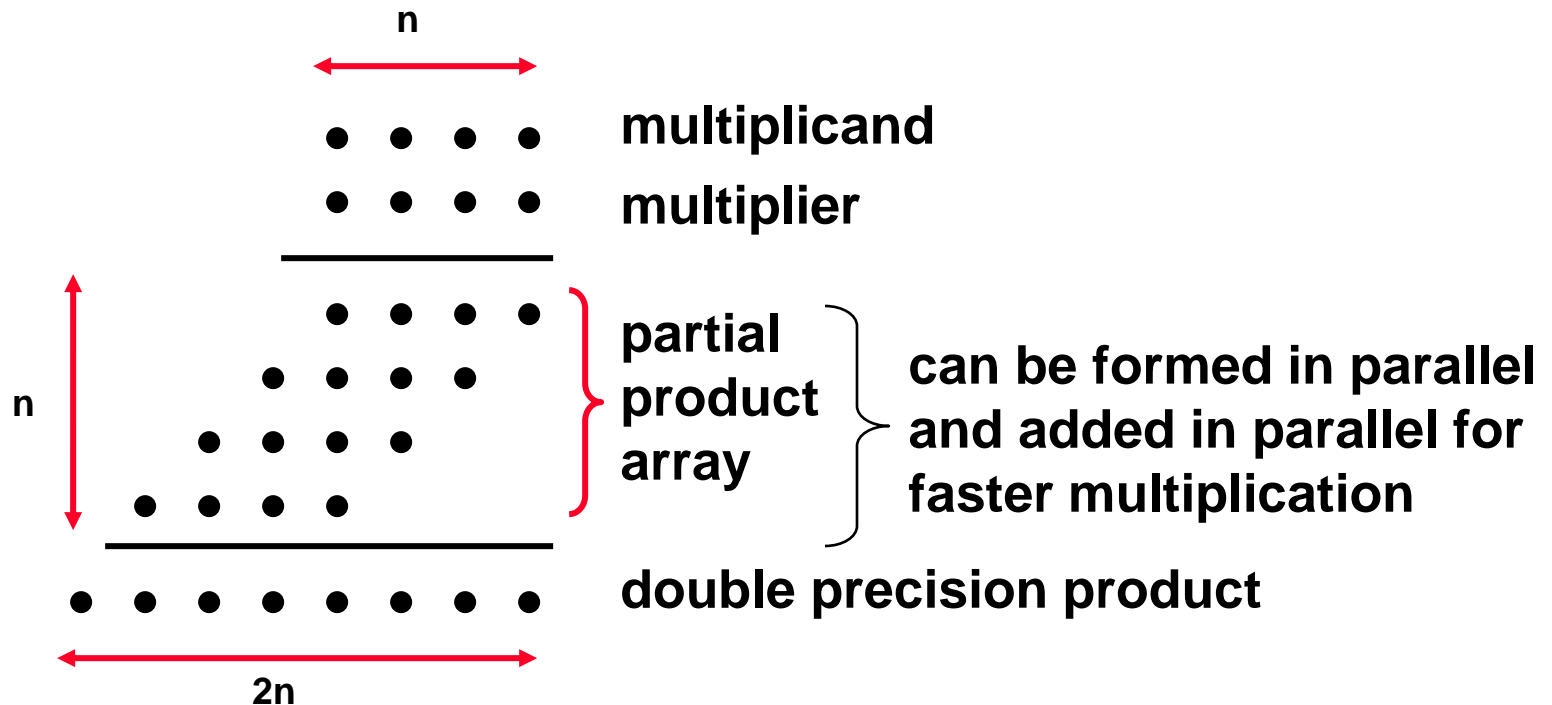


Performance depends on data length

➔ Performance is low

# Multiply

- ❑ Binary multiplication is just a *bunch* of right shifts and adds



*n-bit multiplicand and multiplier → 2n-bit product*

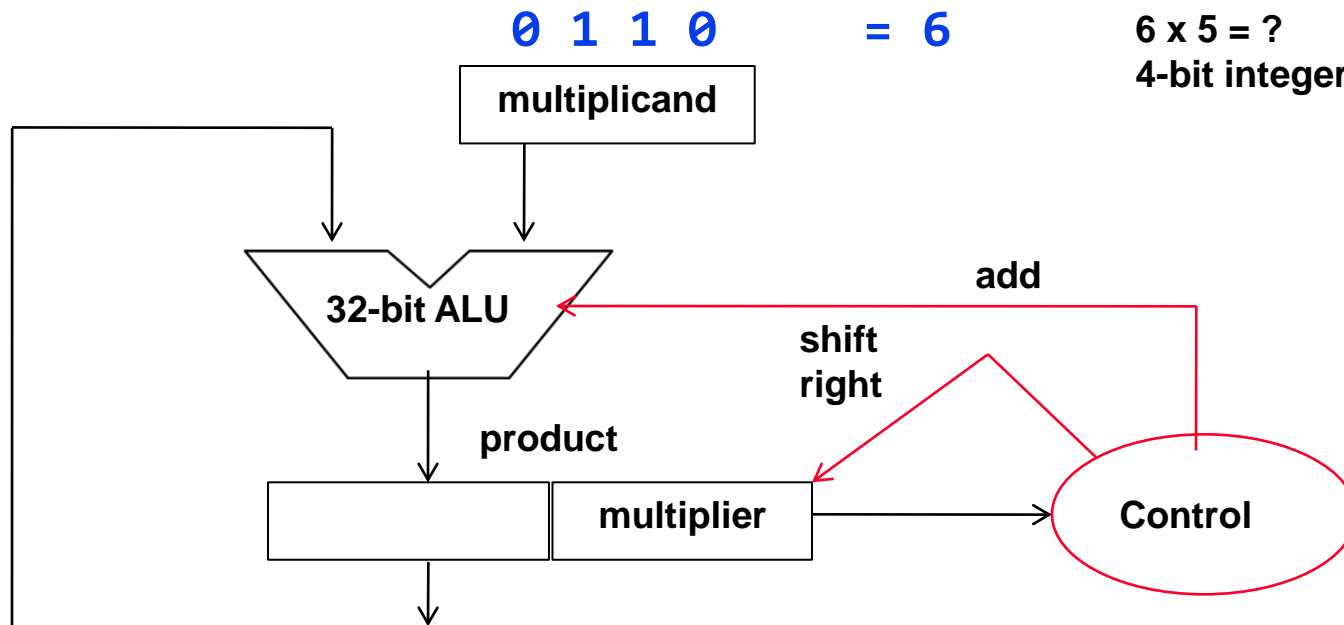


## Example

---

|              |   |  |                        |
|--------------|---|--|------------------------|
| Multiplicand |   |  | 1000 <sub>ten</sub>    |
| Multiplier   | x |  | 1001 <sub>ten</sub>    |
|              |   |  | <hr/>                  |
|              |   |  | 1000                   |
|              |   |  | 0000                   |
|              |   |  | 0000                   |
|              |   |  | 1000                   |
|              |   |  | <hr/>                  |
| Product      |   |  | 1001000 <sub>ten</sub> |

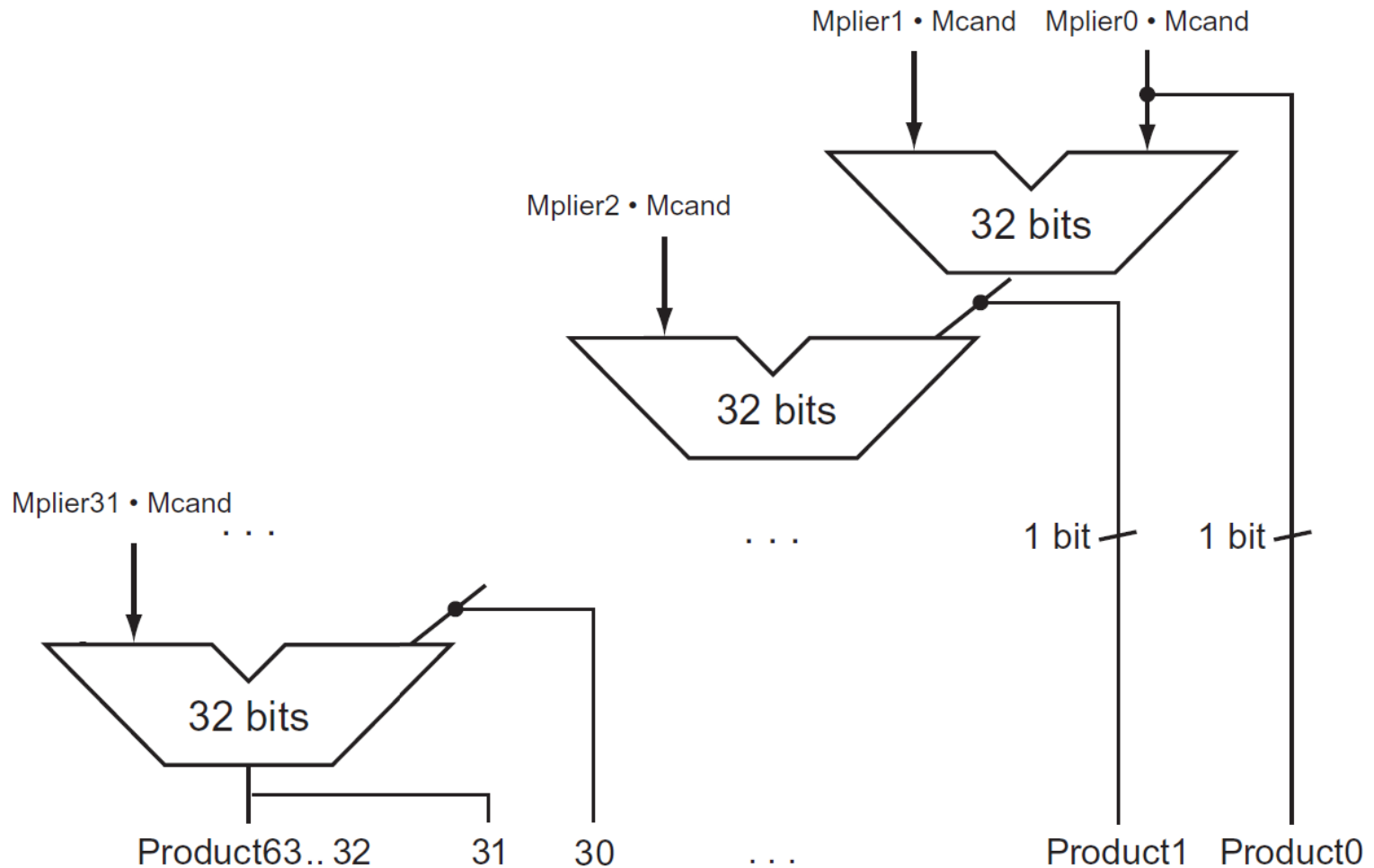
# Add and Right Shift Multiplier Hardware



|     |         |           |                          |
|-----|---------|-----------|--------------------------|
|     | 0 0 0 0 | 0 1 0 1   | = 5                      |
| add | 0 1 1 0 | 0 1 0 1   | LSB=1 → add multiplicand |
|     | 0 0 1 1 | → 0 0 1 0 | shift right              |
| add | 0 0 1 1 | 0 0 1 0   | LSB=0 → no change        |
|     | 0 0 0 1 | → 1 0 0 1 | shift right              |
| add | 0 1 1 1 | 1 0 0 1   | LSB=1 → add multiplicand |
|     | 0 0 1 1 | → 1 1 0 0 | shift right              |
| add | 0 0 1 1 | 1 1 0 0   | LSB=0 → no change        |
|     | 0 0 0 1 | → 1 1 1 0 | shift right              |
|     |         | = 30      |                          |

# Fast multiplier – Design for Moore

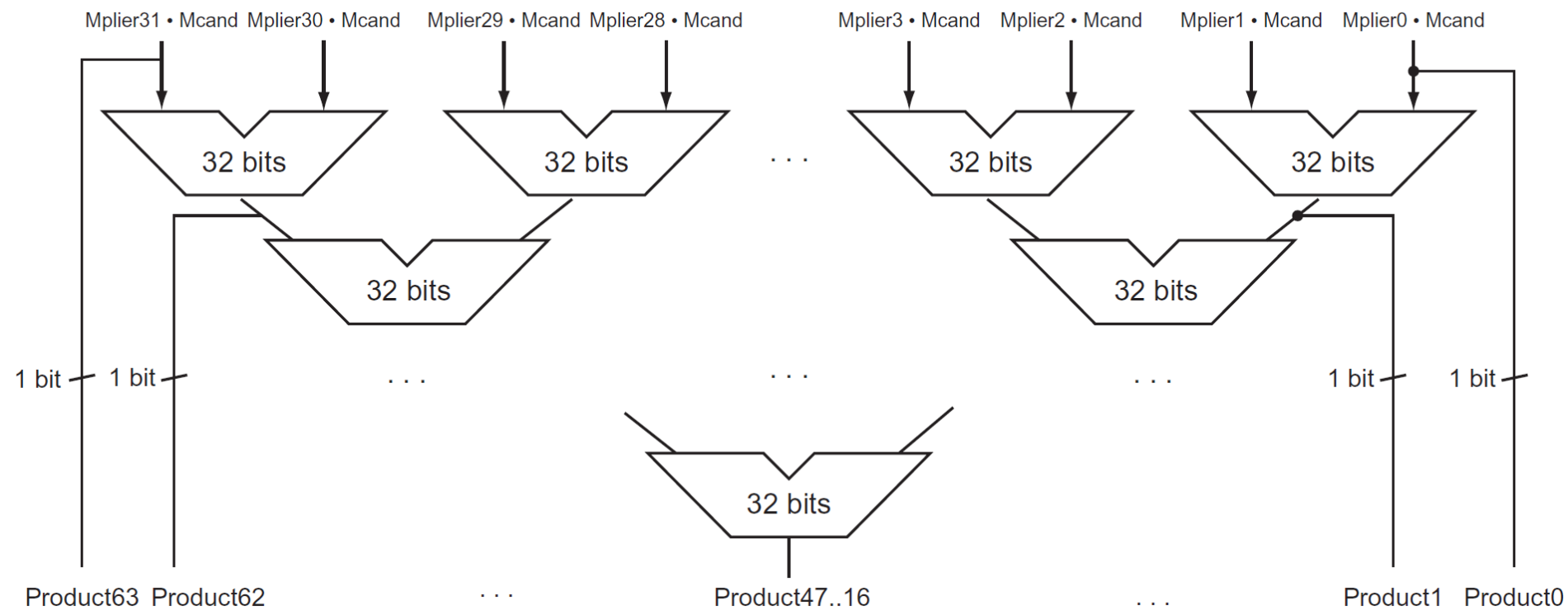
❑ Why is this fast?



# Fast multiplier – Design for Moore

❑ How fast is this?

❑ Anything wrong?



# MIPS Multiply Instruction

- ❑ Multiply (`mult` and `multu`) produces a double precision product (2 x 32 bit)

```
mult    $s0, $s1        # hi||lo = $s0 * $s1
```

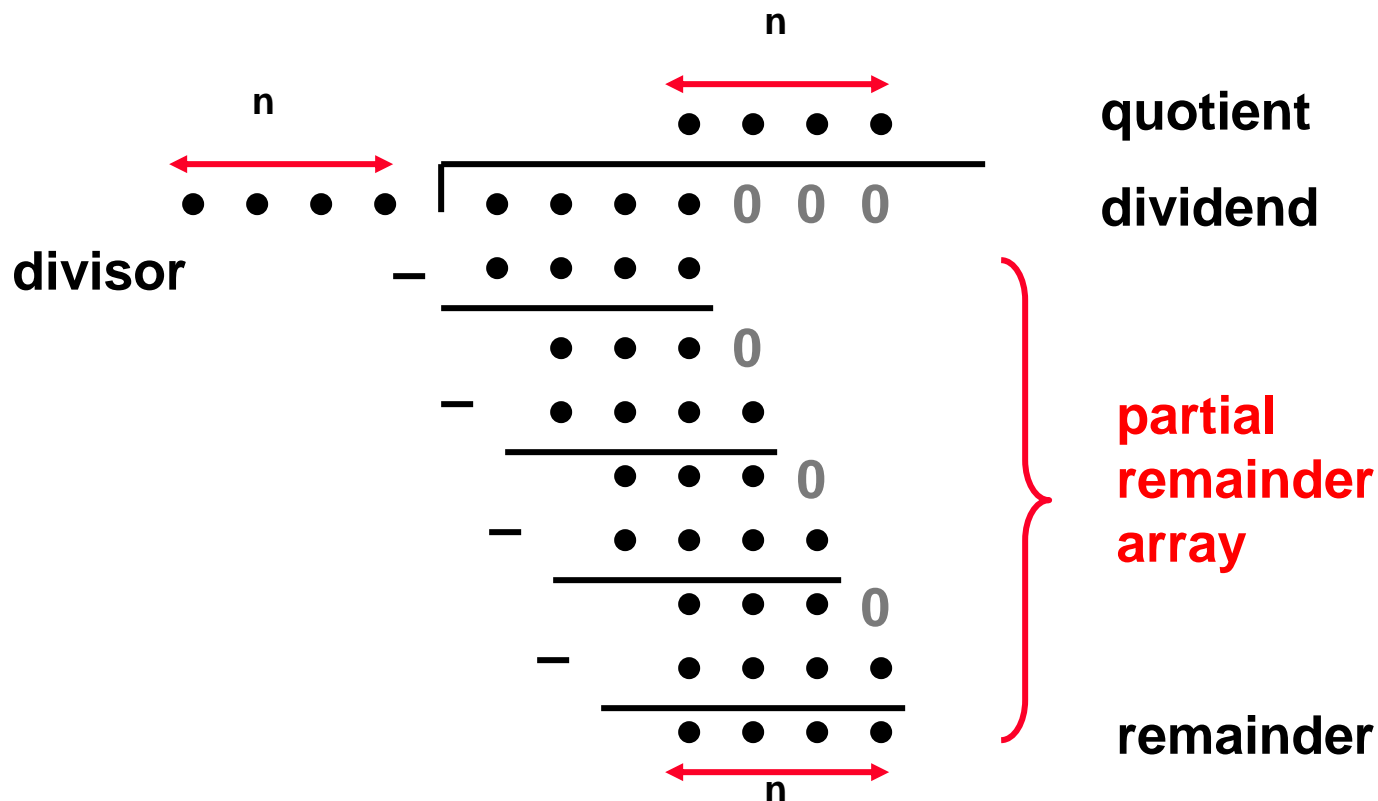
|   |    |    |   |   |      |  |
|---|----|----|---|---|------|--|
| 0 | 16 | 17 | 0 | 0 | 0x18 |  |
|---|----|----|---|---|------|--|

- ❑ Two additional registers: **hi** and **lo**
- ❑ Low-order word of the product is stored in processor register `lo` and the high-order word is stored in register `hi`
- ❑ Instructions `mfhi rd` and `mflo rd` are provided to move the product to (user accessible) registers in the register file

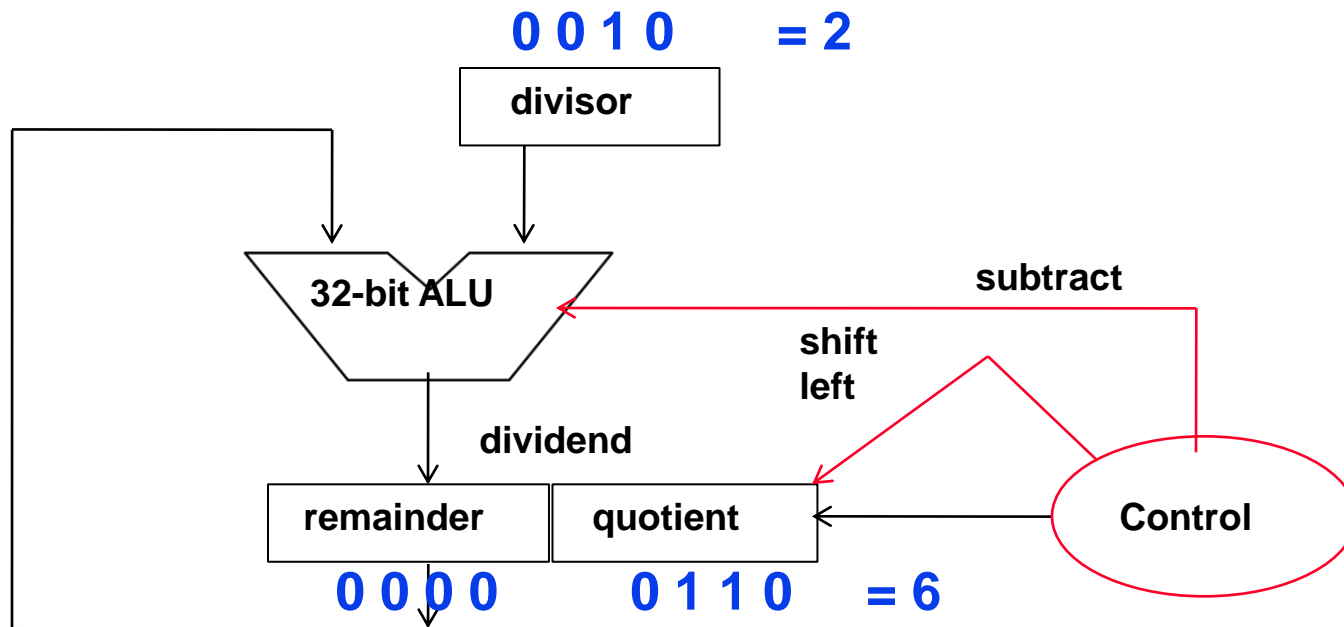
# Division

- Division is just a *bunch* of quotient digit guesses and left shifts and subtracts

$$\text{dividend} = \text{quotient} \times \text{divisor} + \text{remainder}$$



# Left Shift and Subtract Division Hardware



|     |         |   |         |                           |
|-----|---------|---|---------|---------------------------|
|     | 0 0 0 0 | ← | 1 1 0 0 |                           |
| sub | 1 1 1 0 |   | 1 1 0 0 | rem neg, so 'ient bit = 0 |
|     | 0 0 0 0 |   | 1 1 0 0 | restore remainder         |
|     | 0 0 0 1 | ← | 1 0 0 0 |                           |
| sub | 1 1 1 1 |   | 1 0 0 0 | rem neg, so 'ient bit = 0 |
|     | 0 0 0 1 |   | 1 0 0 0 | restore remainder         |
|     | 0 0 1 1 | ← | 0 0 0 0 |                           |
| sub | 0 0 0 1 |   | 0 0 0 1 | rem pos, so 'ient bit = 1 |
|     | 0 0 1 0 | ← | 0 0 1 0 |                           |
| sub | 0 0 0 0 |   | 0 0 1 1 | rem pos, so 'ient bit = 1 |
|     |         |   |         | = 3 with 0 remainder      |

## Divide Instruction

- ❑ Divide (`div` and `divu`) generates the remainder in `hi` and the quotient in `lo`

```
div    $s0, $s1      # lo = $s0 / $s1
```

```
# hi = $s0 mod $s1
```

|   |    |    |   |   |      |  |
|---|----|----|---|---|------|--|
| 0 | 16 | 17 | 0 | 0 | 0x1A |  |
|---|----|----|---|---|------|--|

- Instructions `mfhi rd` and `mflo rd` are provided to move the quotient and remainder to (user accessible) registers in the register file
- **As with multiply, divide ignores overflow so software must determine if the quotient is too large. Software must also check the divisor to avoid division by 0.**

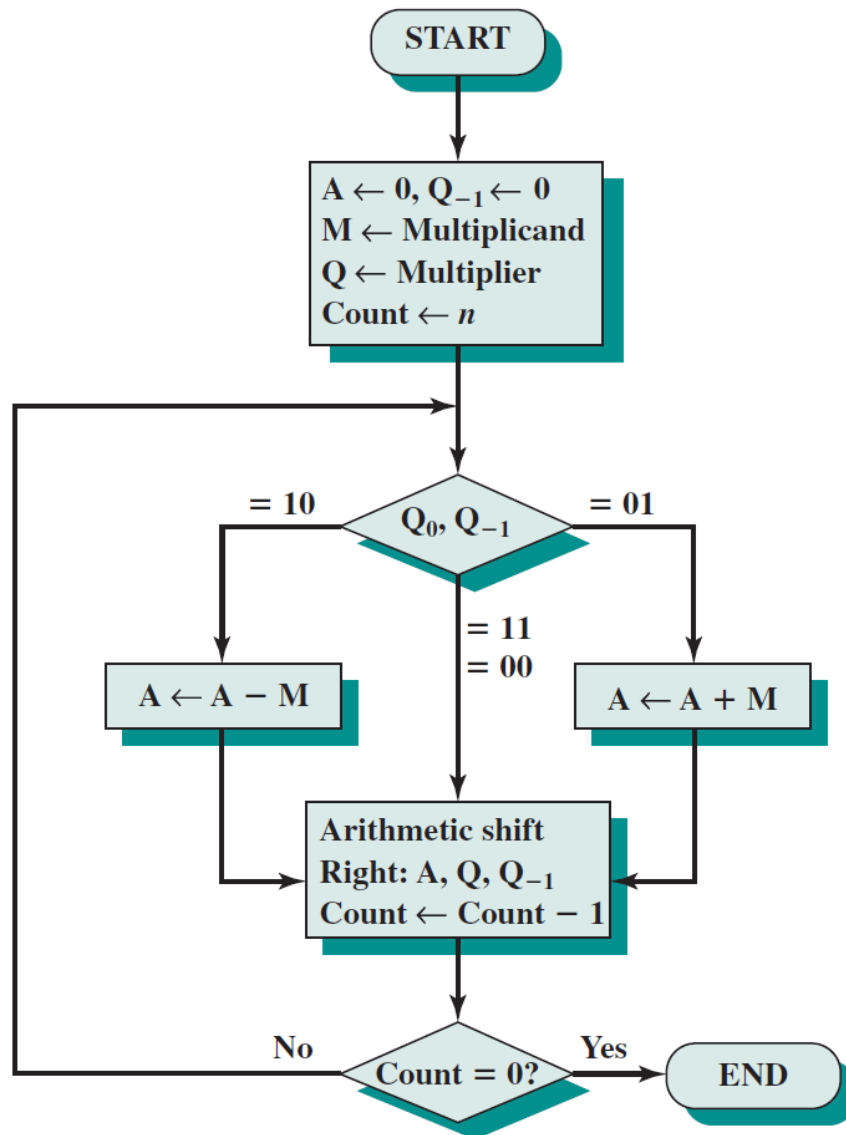


# Signed integer multiplication and division

---

- ❑ Reuse unsigned multiplication then fix product sign later
- ❑ Multiplication
  - ❑ Multiplicand and multiplier are of the same sign: keep product
  - ❑ Multiplicand and multiplier are of different sign: negate product
- ❑ Division:
  - ❑ Dividend and divisor of the same sign:
    - Keep quotient
    - Keep/negate remainder so it is of the same sign with dividend
  - ❑ Dividend and divisor of different sign:
    - Negate quotient
    - Keep/negate remainder so it is of the same sign with dividend

# Signed integer with Booth algorithm



## Representing Big (and Small) Numbers

## Encoding non-integer value?

- [illegible]

- PI number

$$\text{PI} = 3.14159\dots$$

❑ Problem: how to represent the above numbers?

➔ We need reals or floating-point numbers!

➔ Floating point numbers in decimal:

→ 1000

→  $1 \times 10^3$

→  $0.1 \times 10^4$

# Floating point number

---

- ❑ In decimal system

$$2013.1228 = 201.31228 * 10$$

$$= 20.131228 * 10^2$$

$$= 2.0131228 * 10^3$$

$$= 20131228 * 10^{-4}$$

- ❑ What is the “standard” form?

$$2.0131228 * 10^3 = \underline{2.0131228} \text{E} \underline{+03}$$

mantissa

exponent

- ❑ In binary  $X = \pm 1.xxxxx * 2^{yyyy}$

- ❑ ***Sign, mantissa, and exponent need to be represented***

# Floating point number

---

- ❑ Floating point representation in binary

$$(-1)^{\text{sign}} \times 1.F \times 2^{E-\text{bias}}$$

- ❑ Still have to fit everything in 32 bits (single precision)
- ❑ Bias = 127 with single precision floating point number



1 sign bit

8 bits

23 bits

- ❑ Defined by the IEEE 754-1985 standard
  - ❑ Single precision: 32 bit
  - ❑ Double precision: 64 bit
  - ❑ Correspond to float and double in C

## Examples

---

❑ Ex1: convert X into decimal value

$X = 1\textcolor{red}{100}\textcolor{red}{0001}\textcolor{red}{0}101\ 0110\ 0000\ 0000\ 0000\ 0000$

**sign = 1  $\rightarrow$  X is negative**

**E = 1000 0010 = 130**

**F = 10101100...00**

**$\rightarrow X = (-1)^1 \times 1.101011000..00 \times 2^{130-127}$**

**$= -1.101011 \times 2^3 = -1101.011$**

**$= -13.375$**

## Example

---

❑ Ex2: find decimal value of X

X = 0011 1111 1000 0000 0000 0000 0000 0000

**sign = 0**

**e = 0111 1111 = 127**

**m = 000...0000 (23 bit 0)**

**$X = (-1)^0 \times 1.00...000 \times 2^{127-127} = 1.0$**

## Example

---

- Ex3: find binary representation of  $X = 9.6875$  in IEEE 754 single precision

### Converting $X$ to plain binary

$$9_{10} = 1001_2$$

$$0.6875 \times 2 = 1.375 \quad \rightarrow \text{get bit } 1$$

$$0.375 \times 2 = 0.75 \quad \rightarrow \text{get bit } 0$$

$$0.75 \times 2 = 1.5 \quad \rightarrow \text{get bit } 1$$

$$0.5 \times 2 = 1.0 \quad \rightarrow \text{get bit } 1$$



$$\rightarrow 9.6875_{10} = 1001.1011_2$$



## Example

---

- Ex3: find binary representation of  $X = 9.6875$  in IEEE 754 single precision

$$X = 9.6875_{(10)} = 1001.1011_{(2)} = 1.0011011 \times 2^3$$

*Then*

$$S = 0$$

$$e = 127 + 3 = 130_{(10)} = 1000\ 0010_{(2)}$$

$$m = 001101100\dots00 \text{ (23 bit)}$$

*Finally*

$$X = 0100\ 0001\ 0001\ 1011\ 0000\ 0000\ 0000\ 0000$$

## Examples

---

□  $1.0_2 \times 2^{-1} =$

□  $100.75_{10} =$

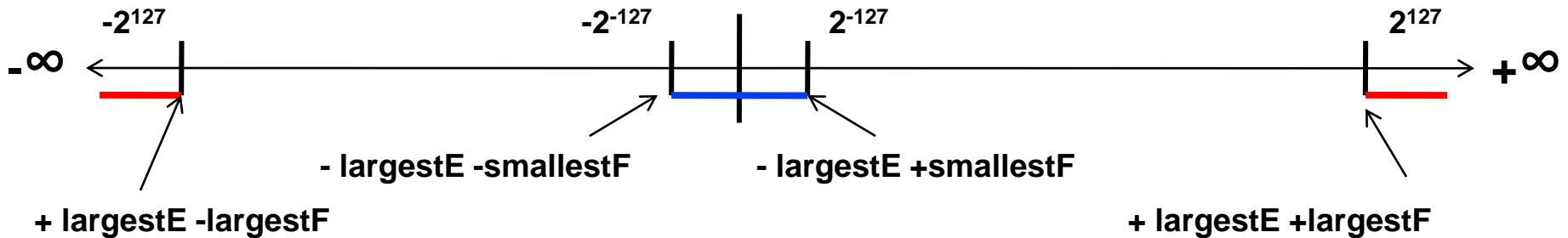
## Some special values

---

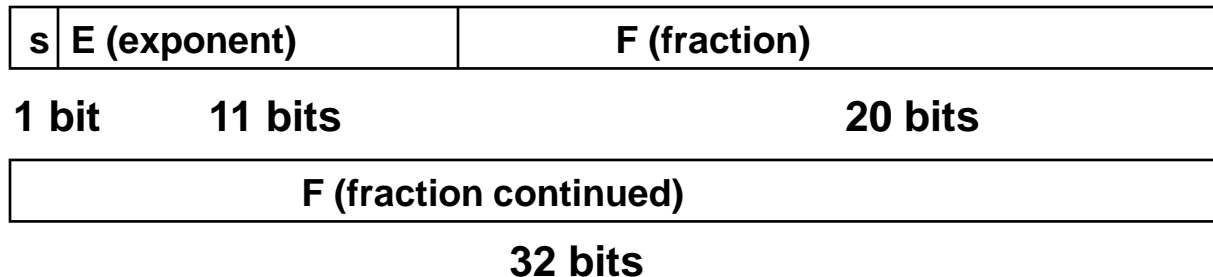
- ❑ Smallest+: 0 00000001 1.00000000000000000000000000000000  
=  $1 \times 2^{1-127}$
- ❑ Zero: 0 00000000 00000000000000000000000000000000  
= true 0
- ❑ Largest+: 0 11111110 1.11111111111111111111111111111111  
=  $(2-2^{-23}) \times 2^{254-127}$

# Too large or too small values

- ❑ **Overflow** (floating point) happens when a positive exponent becomes too large to fit in the exponent field
- ❑ **Underflow** (floating point) happens when a negative exponent becomes too large to fit in the exponent field



- ❑ Reduce the chance of underflow or overflow is to offer another format that has a larger exponent field
- ❑ **Double precision – takes two MIPS words**



## Reduce underflow with the same bit length?

---

- ❑ De-normalized number

# IEEE 754 FP Standard Encoding

- ❑ Special encodings are used to represent unusual events
  - ❑  $\pm$  infinity for division by zero
  - ❑ NaN (not a number) for invalid operations such as 0/0
  - ❑ True zero is the bit string all zero

| Single Precision       |          | Double Precision            |          | Object Represented          |
|------------------------|----------|-----------------------------|----------|-----------------------------|
| E (8)                  | F (23)   | E (11)                      | F (52)   |                             |
| 0000 0000              | 0        | 0000 ... 0000               | 0        | true zero (0)               |
| 0000 0000              | nonzero  | 0000 ... 0000               | nonzero  | $\pm$ denormalized number   |
| 0111 1111 to +127,-126 | anything | 0111 ...1111 to +1023,-1022 | anything | $\pm$ floating point number |
| 1111 1111              | + 0      | 1111 ... 1111               | - 0      | $\pm$ infinity              |
| 1111 1111              | nonzero  | 1111 ... 1111               | nonzero  | not a number (NaN)          |

# Floating Point Addition

---

## □ Addition (and subtraction)

$$(\pm F1 \times 2^{E1}) + (\pm F2 \times 2^{E2}) = \pm F3 \times 2^{E3}$$

- **Step 0: Restore the hidden bit in F1 and in F2**
- **Step 1: Align fractions by right shifting F2 by  $E1 - E2$  positions (assuming  $E1 \geq E2$ ) keeping track of (three of) the bits shifted out in G R and S**
- **Step 2: Add the resulting F2 to F1 to form F3**
- **Step 3: Normalize F3 (so it is in the form 1.XXXXXX ...)**
  - If F1 and F2 have the same sign  $\rightarrow F3 \in [1,4) \rightarrow$  1 bit right shift F3 and increment  $E3$  (check for overflow)
  - If F1 and F2 have different signs  $\rightarrow$  F3 may require *many* left shifts each time decrementing  $E3$  (check for underflow)
- **Step 4: Round F3 and possibly normalize F3 again**
- **Step 5: Rehide the most significant bit of F3 before storing the result**

# Floating Point Addition Example

---

□ Add

$$(0.5 = 1.0000 \times 2^{-1}) + (-0.4375 = -1.1100 \times 2^{-2})$$

□ **Step 0:**

□ **Step 1:**

□ **Step 2:**

□ **Step 3:**

□ **Step 4:**

□ **Step 5:**



# Floating Point Addition Example

---

□ Add:  $0.5 + (-0.4375) = ?$

$$(0.5 = 1.0000 \times 2^{-1}) + (-0.4375 = -1.1100 \times 2^{-2})$$

- **Step 0:** Hidden bits restored in the representation above
- **Step 1:** Shift significand with the smaller exponent (1.1100) right until its exponent matches the larger exponent (so once)
- **Step 2:** Add significands  
 $1.0000 + (-0.111) = 1.0000 - 0.111 = 0.001$
- **Step 3:** Normalize the sum, checking for exponent over/underflow  
 $0.001 \times 2^{-1} = 0.010 \times 2^{-2} = \dots = 1.000 \times 2^{-4}$
- **Step 4:** The sum is already rounded, so we're done
- **Step 5:** Rehide the hidden bit before storing

# Floating Point Multiplication

---

## ❑ Multiplication

$$(\pm F1 \times 2^{E1}) \times (\pm F2 \times 2^{E2}) = \pm F3 \times 2^{E3}$$

- ❑ **Step 0: Restore the hidden bit in F1 and in F2**
- ❑ **Step 1: Add the two (biased) exponents and subtract the bias from the sum, so  $E1 + E2 - 127 = E3$**   
**also determine the sign of the product (which depends on the sign of the operands (most significant bits))**
- ❑ **Step 2: Multiply F1 by F2 to form a double precision F3**
- ❑ **Step 3: Normalize F3 (so it is in the form 1.XXXXXX ...)**
  - Since F1 and F2 come in normalized  $\rightarrow F3 \in [1,4) \rightarrow$  1 bit right shift F3 and increment E3
  - Check for overflow/underflow
- ❑ **Step 4: Round F3 and possibly normalize F3 again**
- ❑ **Step 5: Rehide the most significant bit of F3 before storing the result**

# Floating Point Multiplication Example

---

## ❑ Multiply

$$(0.5 = 1.0000 \times 2^{-1}) \times (-0.4375 = -1.1100 \times 2^{-2})$$

❑ **Step 0:**

❑ **Step 1:**

❑ **Step 2:**

❑ **Step 3:**

❑ **Step 4:**

❑ **Step 5:**

# Floating Point Multiplication Example

---

## ❑ Multiply

$$(0.5 = 1.0000 \times 2^{-1}) \times (-0.4375 = -1.1100 \times 2^{-2})$$

❑ **Step 0:** Hidden bits restored in the representation above

❑ **Step 1:** Add the exponents (not in bias would be  $-1 + (-2) = -3$  and in bias would be  $(-1+127) + (-2+127) - 127 = (-1 -2) + (127+127-127) = -3 + 127 = 124$ )

❑ **Step 2:** Multiply the significands

$$1.0000 \times 1.110 = 1.110000$$

❑ **Step 3:** Normalized the product, checking for exp over/underflow

$$1.110000 \times 2^{-3} \text{ is already normalized}$$

❑ **Step 4:** The product is already rounded, so we're done

❑ **Step 5:** Rehide the hidden bit before storing

# Support for Accurate Arithmetic

---

- ❑ IEEE 754 FP rounding modes

- ❑ Always round up (toward  $+\infty$ )
- ❑ Always round down (toward  $-\infty$ )
- ❑ Truncate
- ❑ **Round to nearest even** (when the Guard || Round || Sticky are 100) – always creates a 0 in the least significant (kept) bit of F

- ❑ Rounding (except for truncation) requires the hardware to include extra F bits during calculations

- ❑ Guard and Round bit – 2 additional bits to increase accuracy
- ❑ Sticky bit – used to support **Round to nearest even**; is set to a 1 whenever a 1 bit shifts (right) through it (e.g., when aligning F during addition/subtraction)

$F = 1 . \text{xxxxxxxxxxxxxxxxxxxxxxxxxxxx} \text{ G R S}$

## Example

---

❑ Calculate:

$$0.2 \times 5 = ?$$

$$0.333 \times 3 = ?$$

$$(1.0/3) \times 3 = ?$$