



University of San Carlos | Department of  
**COMPUTER ENGINEERING**

CpE 3202

Computer Organization & Architecture

# Computer Arithmetic

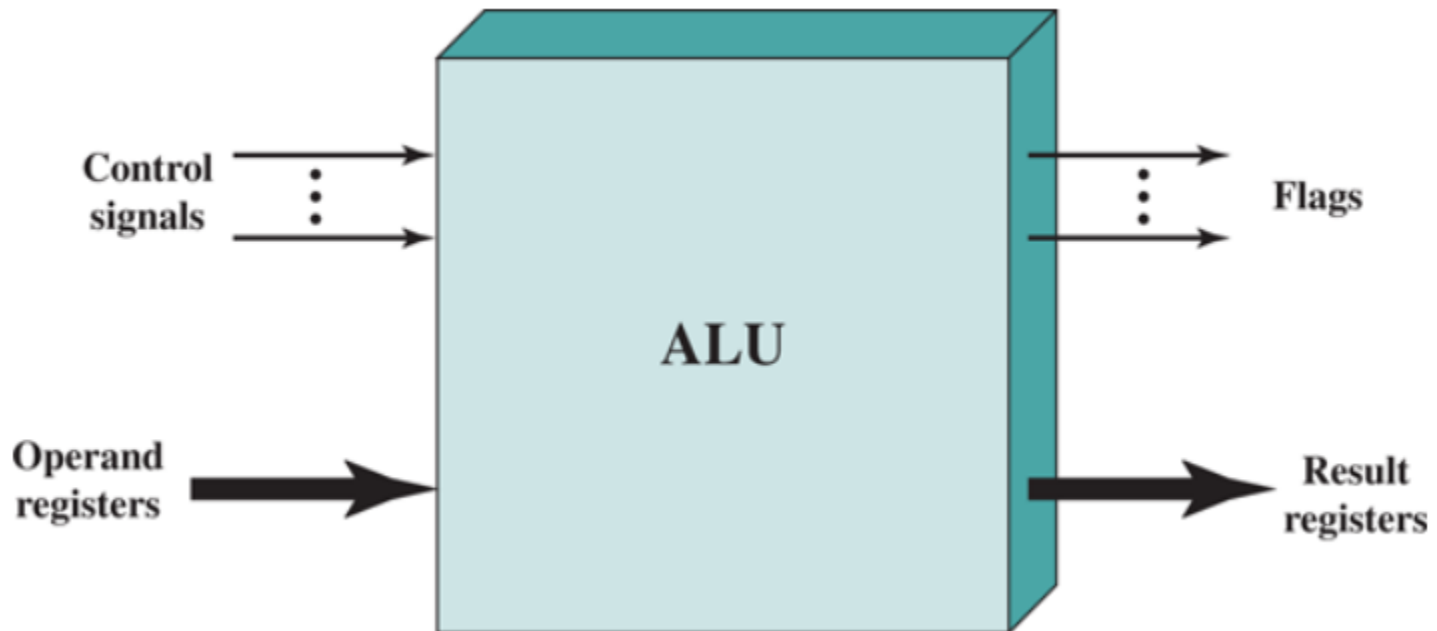
# Computer Arithmetic

- Most complex aspect of the Arithmetic and Logic Unit (ALU)
- It is performed on two very different numbers: *integer* and *floating point*

# Arithmetic and Logic Unit (ALU)

- Does the calculations
- Everything else in the computer is there to service this unit
- Handles integers
- May handle floating point (real) numbers
- May include a separate floating-point unit (FPU), as a math co-processor

# ALU Inputs and Outputs



# Integer Representation

- Only have 0 & 1 to represent everything
- Positive numbers stored in binary
  - e.g. 41 = 0010 1001

<sup>32</sup>   <sup>8</sup>   <sup>1</sup>
- No minus sign
- No radix point (e.g. *decimal point for fractional component*)
- Sign-Magnitude

54.5

1010.1110  
↑

# Sign-Magnitude

- Used to represent positive and negative integers

- Left most bit is sign bit

- 0 means positive
- 1 means negative

- Examples:  $+18 = \overset{16}{\underline{0001}} \overset{2}{\underline{0010}}$   
 $-18 = \overset{16}{\underline{1001}} \overset{2}{\underline{0010}}$

$$\underline{0000} = +0$$

$$\underline{1000} = -0$$

- Problems

- Need to consider both sign and magnitude in arithmetic
- Two representations of zero (+0 and -0)

# Twos Complement Notation

- +3 = 0000 0011<sup>2 1</sup>
- +2 = 0000 0010<sup>2</sup>
- +1 = 0000 0001<sup>1</sup>
- +0 = 0000 0000
- -1 = 1111 1111
- -2 = 1111 1110
- -3 = 1111 1101

$$\begin{array}{r}
 0000 \ 0001 \\
 + \ 1111 \ 1110 \\
 \hline
 1111 \ 1111
 \end{array}$$

$$\begin{array}{r}
 0000 \ 0010 \\
 + \ 1111 \ 1101 \\
 \hline
 1111 \ 1110
 \end{array}$$

# Benefits

- Only one representation of zero
- Arithmetic works easily (*see later*)
- Negating is fairly easy

- $3 = 0000\ 0011$  ✓

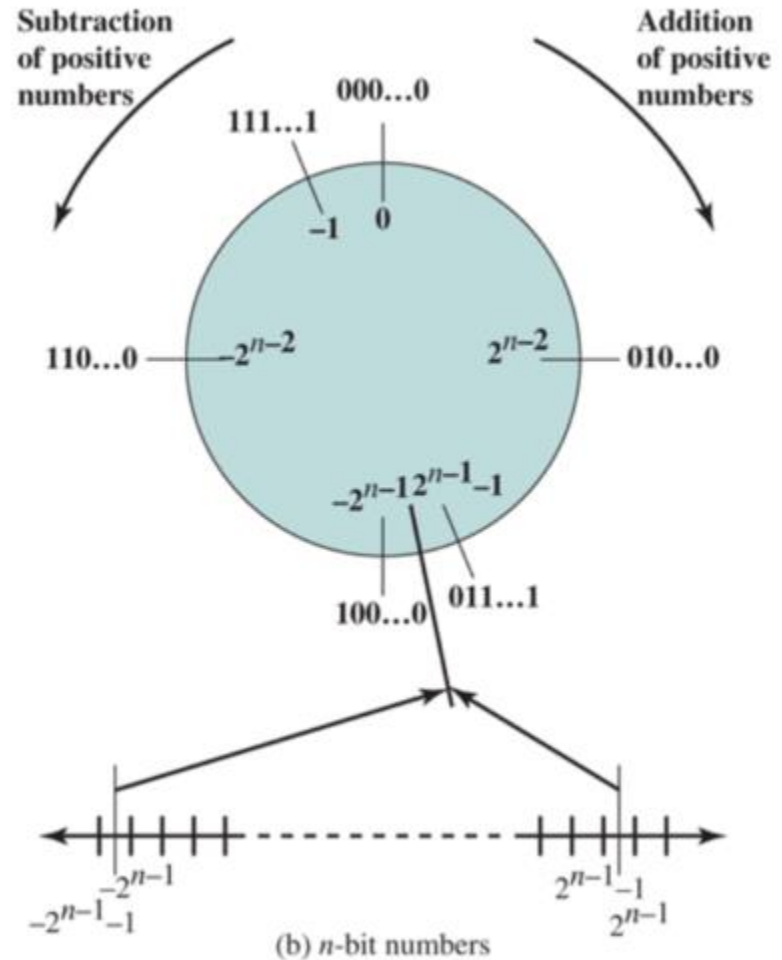
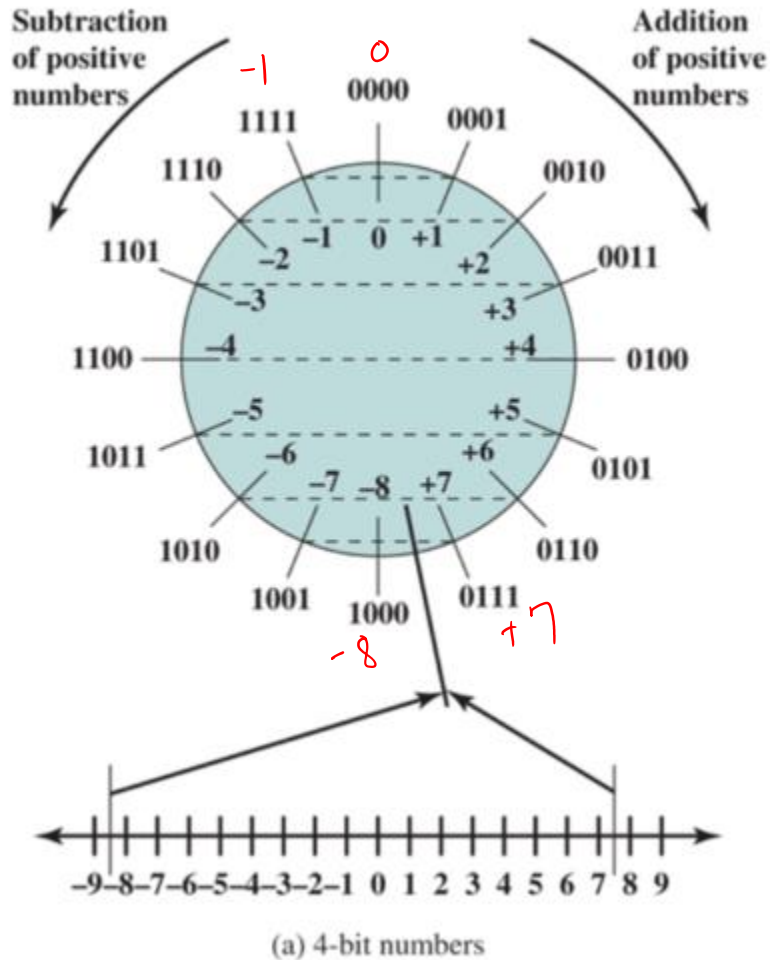
- Boolean complement gives  $1111\ 1100$

- Add 1 to LSB

$$\begin{array}{r}
 1111\ 1100 \\
 1111\ 1101 = -3 \\
 + \quad 000\ 0001 \\
 \hline
 000\ 0001 = 3
 \end{array}$$

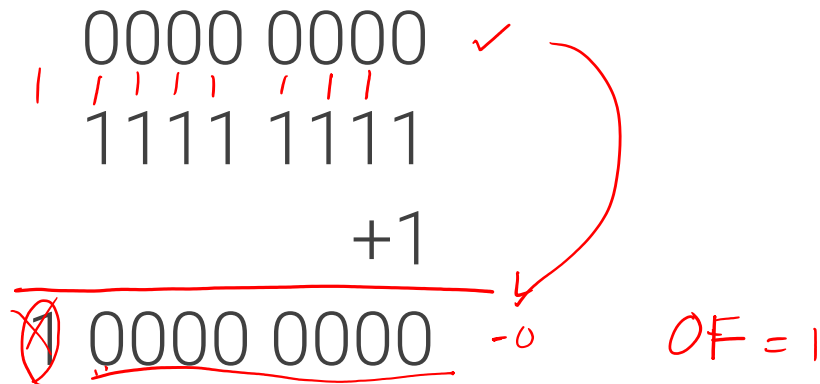


# Geometric Depiction of Two's Complement Integers



# Negation Special Case 1

- 0 =
  - Bitwise not
  - Add 1 to LSB
  - Result
- Overflow is ignored, so:
  - - 0 = 0 ✓



# Negation Special Case 2

- -128 =  $1000\ 0000$  ✓
  - Bitwise not  $0111\ 1111$
  - Add 1 to LSB  $\quad\quad\quad +1$
  - Result  $1000\ 0000$  ↩
- So:  $\uparrow$ 
  - $-(-128) = -128$  **X**
- Monitor MSB (sign bit)
  - It should change during negation

# Range of Numbers

- <sup>n=8</sup> 8 bit 2s complement <sup>n-1</sup>  $2^{n-1} - 1$

max •  $+127 = \underline{0}111\ 1111 = 2^7 - 1$

min •  $-128 = \underline{1}000\ 000\underline{0} = -2^7$  <sup>n-1</sup>  $-2^{n-1}$

- <sup>n</sup> 16 bit 2s complement

•  $+32,767 = \underline{0}111\ 1111\ 1111\ 1111 = 2^{15} - 1$

•  $-32,768 = \underline{1}000\ 0000\ 0000\ 000\underline{0} = -2^{15}$

4 bits  $n=4$   $-8$  to  $+7$

0000 0  
0001 1

⋮  
0111 7 ✓  $= 2^3 - 1 = 7$  ✓  
✓ 1000 -8  $= -2^3 = -8$   
1001  
⋮  
1111 -1

1000 = -8  
+ 0111  
-----  
1000 = 8  
8 4 2 1

# Conversion Between Lengths

- Positive numbers: pad with leading zeros

$$\bullet +18 = \quad \quad \quad \overset{16}{\text{0001}} \overset{2}{\text{0010}} = 8$$

$$\bullet +18 = \quad \text{0000 0000} \text{ 0001 0010} = 32$$

- Negative numbers: pad with leading ones

$$\bullet -18 = \quad \quad \quad \text{1} \text{001 0010} -$$

$$\bullet -18 = \quad \text{1111 1111} \text{ 1001 0010}$$

- Pad using the MSB (sign bit)

# Addition

$$\begin{array}{r}
 7 \quad 0111 \\
 + 1000 \\
 \hline
 -7 \quad 1001
 \end{array}
 \quad
 \begin{array}{r}
 1110 \\
 + 0001 \\
 \hline
 0010 = 2
 \end{array}
 \quad
 \begin{array}{r}
 4 \quad 0100 = 4 \\
 + 1011 \\
 \hline
 1100 = -4
 \end{array}
 \quad
 \begin{array}{r}
 0001 = 1 \\
 1110 \\
 + 1 \\
 \hline
 1111 = -1
 \end{array}$$

$1011 = -?$   
 $0100$   
 $+ 1$   
 $\hline$   
 $0101 = 5$

- Proceeds as if the two numbers are unsigned integers
  - Normal binary addition
- **Overflow** - result exceeds the bit-width size (extra bit is ignored)

|                                                                                                                                 |                                                                                                                               |
|---------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------|
| $  \begin{array}{r}  1001 = -7 \\  + 0101 = 5 \\  \hline  1110 = -2 \\  \text{(a) } (-7) + (+5) = -2  \end{array}  $            | $  \begin{array}{r}  1100 = -4 \\  + 0100 = 4 \\  \hline  10000 = 0 \\  \text{(b) } (-4) + (+4)  \end{array}  $               |
| $  \begin{array}{r}  0011 = 3 \\  + 0100 = 4 \\  \hline  0111 = 7 \\  \text{(c) } (+3) + (+4)  \end{array}  $                   | $  \begin{array}{r}  1100 = -4 \\  + 1111 = -1 \\  \hline  1011 = -5 \\  \text{(d) } (-4) + (-1) = -5  \end{array}  $         |
| $  \begin{array}{r}  0101 = 5 \\  + 0100 = 4 \\  \hline  1001 = \text{Overflow} \\  \text{(e) } (+5) + (+4) = 9  \end{array}  $ | $  \begin{array}{r}  1001 = -7 \\  + 1010 = -6 \\  \hline  1011 = \text{Overflow} \\  \text{(f) } (-7) + (-6)  \end{array}  $ |

5 bits  
 $00101 = 5$   
 $00100 = 4$   
 $01001 = 9$   
 not overflow

• If two numbers are added, and **they are both positive or both negative**, the overflow occurs **if and only if** the result has the opposite sign.

# Subtraction

- Take twos complement of subtrahend (S) and add to minuend (M)

•  $M - S = M (+) (-S)$

- When overflow occurs, the ALU send a signal so that no attempt is made to use the result.

$$M = 0101 = 5 \checkmark$$

$$S = 0010 = 2$$

$$-S = 1101$$

$$\begin{array}{r} 1101 \\ + 1 \\ \hline 1110 = -2 \checkmark \end{array}$$

$$M = -5$$

$$\begin{array}{r} 0101 \\ + 1010 \\ \hline 1011 = -5 \end{array}$$

$$S = 2$$

$$0010$$

$$-S = 1101$$

$$\begin{array}{r} 1101 \\ + 1 \\ \hline 1110 = -2 \end{array}$$

|                                                                                                                                                                                                                  |                                                                                                                                                                                             |
|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| $\begin{array}{r} 0010 = 2 \\ + 1001 = -7 \\ \hline 1011 = -5 \end{array}$ <p>(a) <math>M = 2 = 0010</math><br/><math>S = 7 = 0111</math><br/><math>-S = 1001</math></p>                                         | $\begin{array}{r} 0101 = 5 \\ + 1110 = -2 \\ \hline 0011 = 3 \end{array}$ <p>(b) <math>M = 5 = 0101</math><br/><math>S = 2 = 0010</math><br/><math>-S = 1110 = 3 \checkmark</math></p>      |
| $\begin{array}{r} 1011 = -5 \checkmark \\ + 1110 = -2 \checkmark \\ \hline 0001 = -7 \checkmark \end{array}$ <p>(c) <math>M = -5 = 1011</math><br/><math>S = 2 = 0010</math><br/><math>-S = 1110 = -7</math></p> | $\begin{array}{r} 0101 = 5 \\ + 0010 = 2 \\ \hline 0111 = 7 \end{array}$ <p>(d) <math>M = 5 = 0101</math><br/><math>S = -2 = 1110</math><br/><math>-S = 0010</math></p>                     |
| $\begin{array}{r} 0111 = 7 + \\ + 0111 = 7 + \\ \hline 1110 = \text{Overflow} \end{array}$ <p>(e) <math>M = 7 = 0111</math><br/><math>S = -7 = 1001</math><br/><math>-S = 0111</math></p>                        | $\begin{array}{r} 1010 = -6 - \\ + 1100 = -4 - \\ \hline 0110 = \text{Overflow} \end{array}$ <p>(f) <math>M = -6 = 1010</math><br/><math>S = 4 = 0100</math><br/><math>-S = 1100</math></p> |

$$5 - (-8) = 13$$

$$S = -8 = 1000$$

$$-S = 0111$$

$$\begin{array}{r} 1000 \\ + 0111 \\ \hline 1000 \end{array}$$

$$5 - 9 = -4$$

$$M = 5 = 0101$$

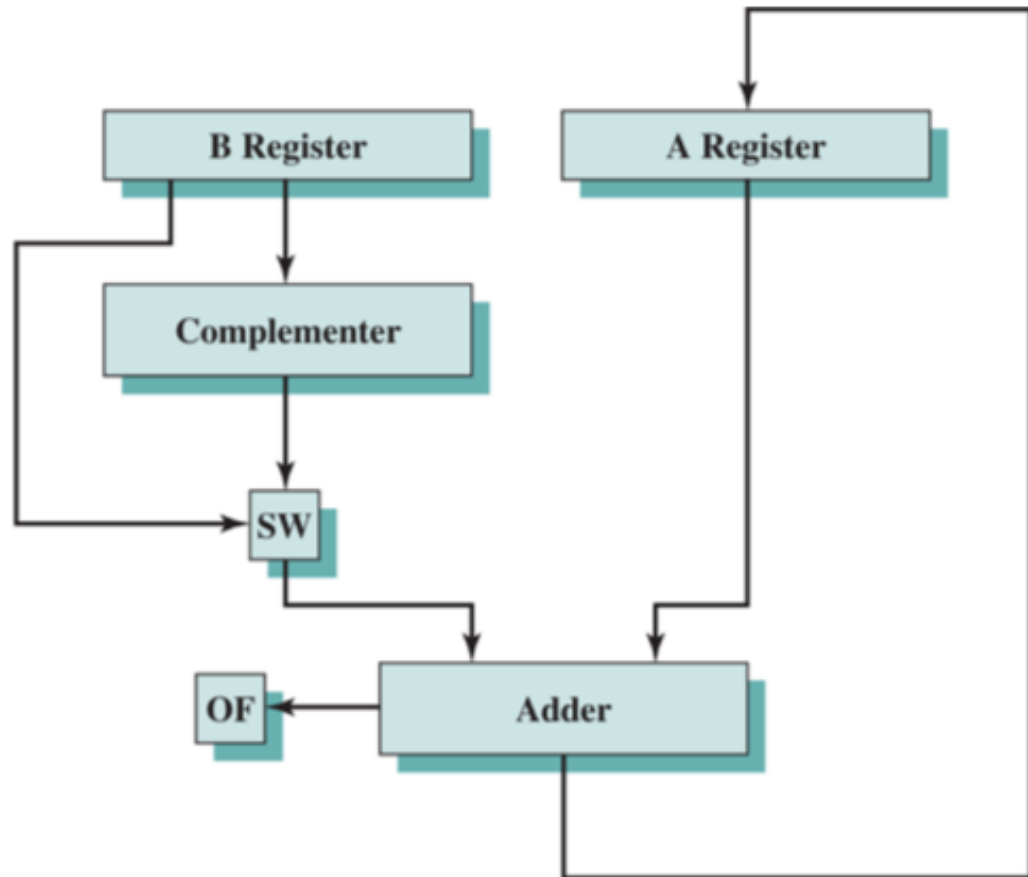
$$S = 9 = 1001$$

$$-S = 0111$$

$$\begin{array}{r} 0101 \\ + 0111 \\ \hline 0100 \end{array}$$

# Hardware for Addition and Subtraction

- For addition and subtraction, we only need addition and complement circuitry



OF = Overflow bit

SW = Switch (select addition or subtraction)



# Multiplication

- Complex
- Work out a **partial product** for each digit
- Take care with placing values (*column*)
- Add partial products

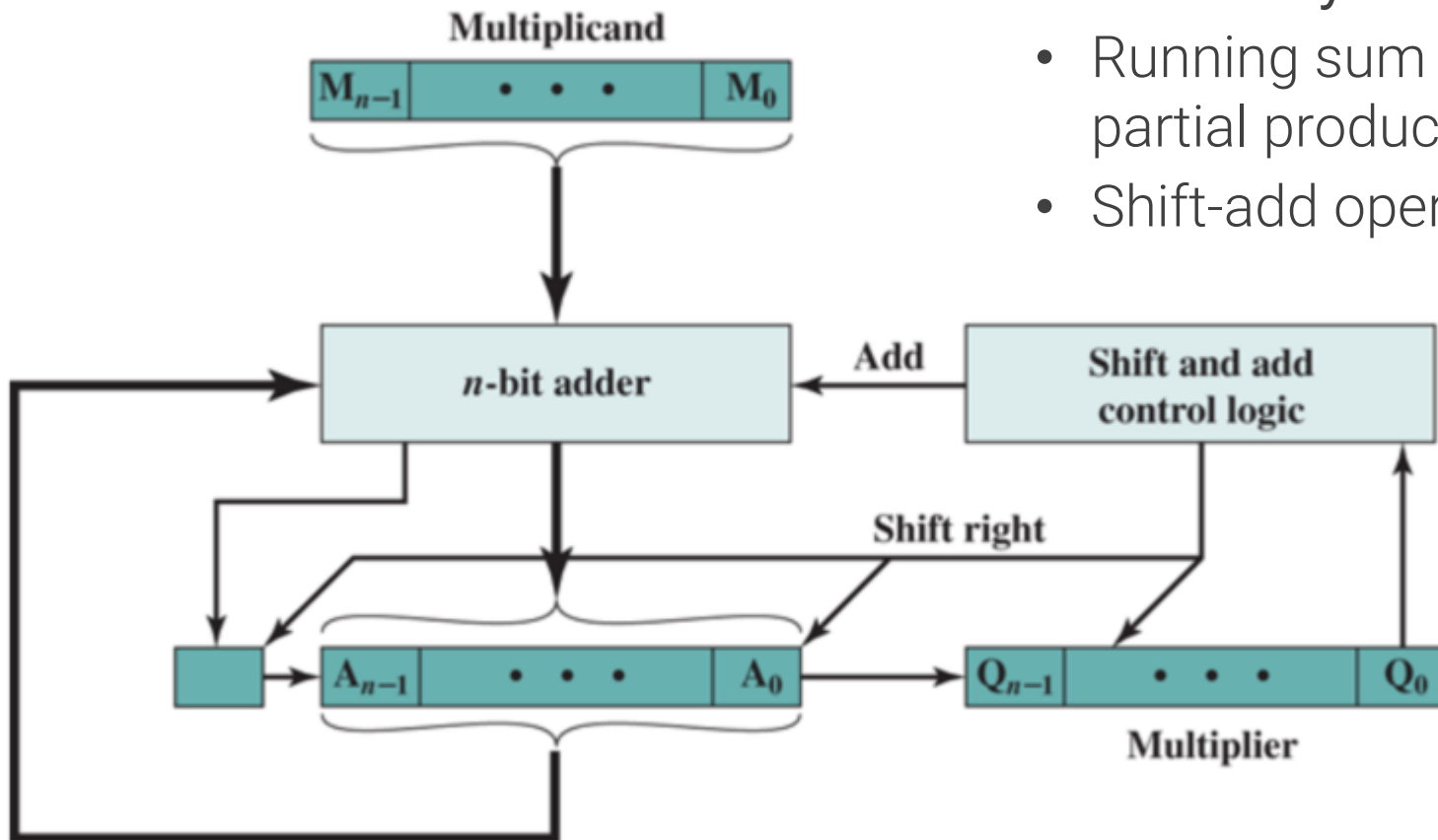
# Multiplication Example

|                                                                                                                |                                                                                                                                                                                                                                                                                                           |
|----------------------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| $\begin{array}{r} 1011 \\ \times 1101 \\ \hline 1011 \\ 0000 \\ + 1011 \\ 1011 \\ \hline 10001111 \end{array}$ | <p>Multiplicand <math>(11_{10})</math></p> <p>Multiplier <math>(13_{10})</math></p> <hr/> <p><u>Partial products</u></p> <p><b>Note:</b> IF multiplier bit is 1, copy multiplicand<br/>(<i>multiplicand is the partial product</i>);<br/>otherwise, zero</p> <hr/> <p>Product <math>(143_{10})</math></p> |
|----------------------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|

- **Note:** Need a double-length width result (*2n-bit result for n-bit binary integers*)

# Unsigned Binary Multiplication

- For efficiency:
  - Running sum on the partial products
  - Shift-add operation



(a) Block diagram

# Execution of Example

$$\begin{array}{r}
 1011 \ (11_{10}) \quad \checkmark \ M \\
 \times 1101 \ (13_{10}) \quad \checkmark \ Q \\
 \hline
 1011 \\
 0000 \\
 1011 \\
 1011 \\
 \hline
 \underline{10001111} \ (143_{10})
 \end{array}$$

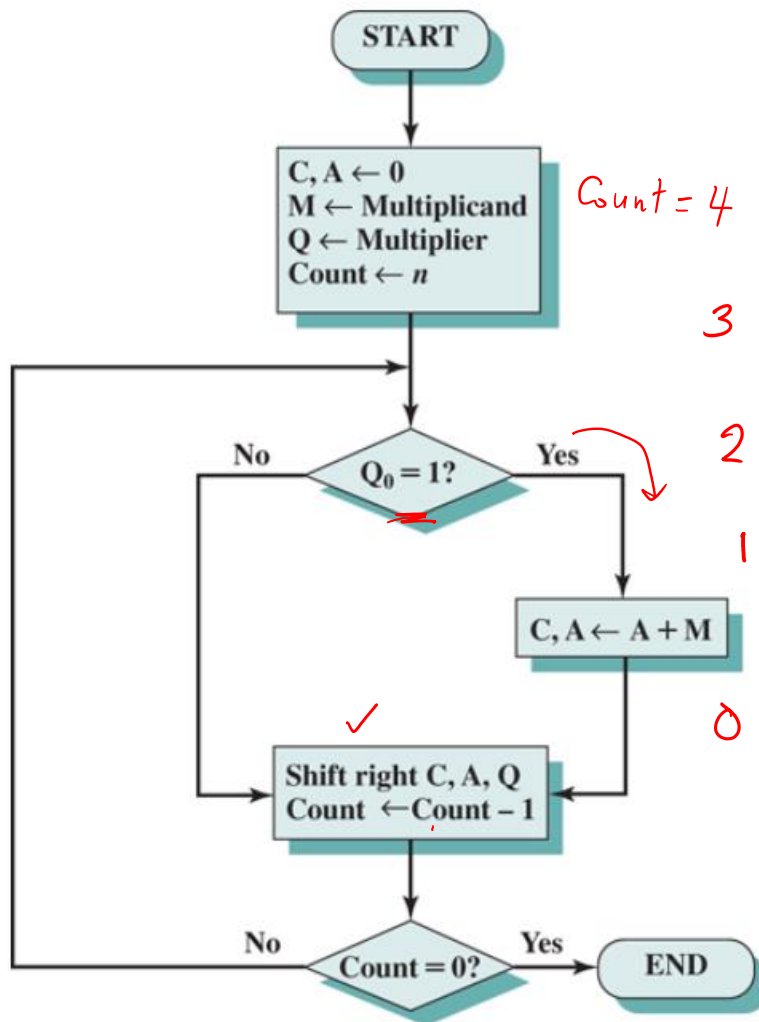
## Note:

- M = Multiplicand
- Q = Multiplier
- A = Accumulator
- C = Carry bit

| C | A    | Q    | M    | Initial values |                 |
|---|------|------|------|----------------|-----------------|
| 0 | 0000 | 1101 | 1011 |                |                 |
| 0 | 1011 | 1101 | 1011 | Add<br>Shift   | First<br>cycle  |
| 0 | 0101 | 1110 | 1011 |                |                 |
| 0 | 0010 | 1111 | 1011 | Shift          | Second<br>cycle |
| 0 | 1101 | 1111 | 1011 |                |                 |
| 0 | 0110 | 1111 | 1011 | Add<br>Shift   | Third<br>cycle  |
| 0 | 0110 | 1111 | 1011 |                |                 |
| 1 | 0001 | 1111 | 1011 | Add<br>Shift   | Fourth<br>cycle |
| 0 | 1000 | 1111 | 1011 |                |                 |

Product in {A, Q}

# Flowchart for Unsigned Binary Multiplication



| C | A    | Q    | M    |         |                |
|---|------|------|------|---------|----------------|
| 0 | 0000 | 1101 | 1011 |         | Initial values |
| 0 | 1011 | 1101 | 1011 | ✓ Add   | First cycle    |
| 0 | 0101 | 1110 | 1011 | ✓ Shift |                |
| 0 | 0010 | 1111 | 1011 | ✓ Shift | Second cycle   |
| 0 | 1101 | 1111 | 1011 | ✓ Add   |                |
| 0 | 0110 | 1111 | 1011 | ✓ Shift | Third cycle    |
| 1 | 0001 | 1111 | 1011 | ✓ Add   |                |
| 0 | 1000 | 1111 | 1011 | ✓ Shift | Fourth cycle   |
|   |      |      |      |         |                |

$$\begin{array}{r}
 0010 \\
 + 1011 \\
 \hline
 1101
 \end{array}
 \qquad
 \begin{array}{r}
 0110 \\
 + 1011 \\
 \hline
 10001
 \end{array}$$

Product  
in A, Q

# Multiplying Negative Numbers

- This does not work!
- **Solution 1**
  - Convert to positive, if required
  - Multiply as above
  - If signs were different, negate answer
- **Solution 2**
  - Booth's algorithm

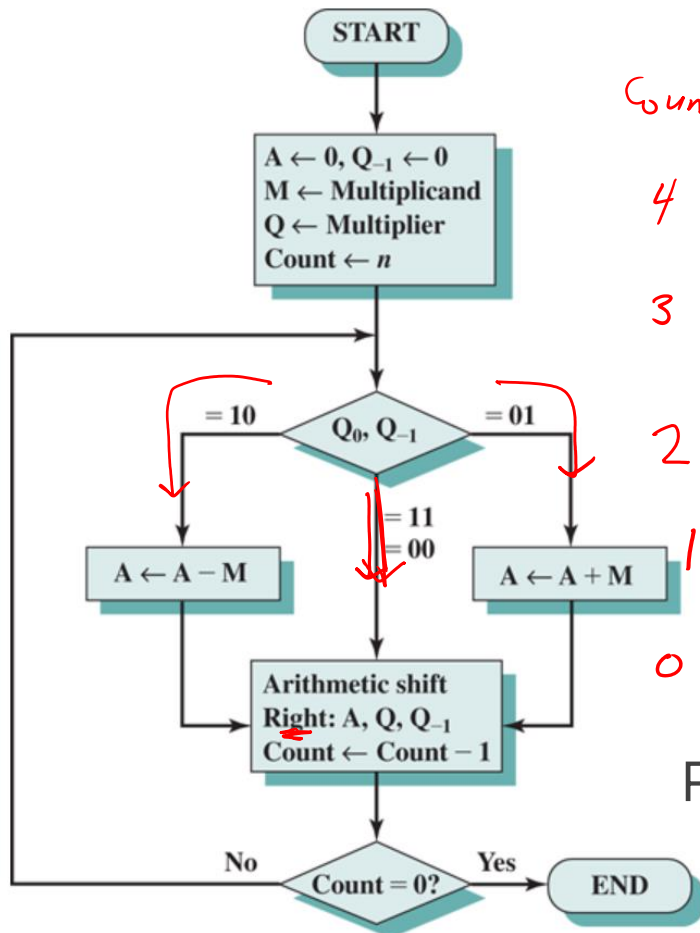
# Booth's Algorithm

$$\begin{array}{r} A \quad 0000 \\ M \quad 0111 \\ -M \quad 1000 \\ + \quad 1 \\ \hline 1001 \end{array}$$

$$\begin{array}{r} A + (-M) \\ 0000 \\ + 1001 \\ \hline 1001 \end{array}$$

$$\begin{array}{r} 1110 \\ + 0111 \\ \hline 10101 \end{array}$$

Perform:  $3_{10} \times 7_{10} = 21_{10}$



Count

4

3

2

1

0

| A    | Q <sub>0</sub> | Q <sub>-1</sub> | M    |                  |                |
|------|----------------|-----------------|------|------------------|----------------|
| 0000 | 0011           | 0               | 0111 | ✓ Initial values |                |
| 1001 | 0011           | 0               | 0111 | ✓ A ← A - M      | } First cycle  |
| 1100 | 1001           | 1               | 0111 | ✓ Shift          |                |
| 1110 | 0100           | 1               | 0111 | ✓ Shift          | } Second cycle |
| 0101 | 0100           | 1               | 0111 | ✓ A ← A + M      |                |
| 0010 | 1010           | 0               | 0111 | ✓ Shift          | } Third cycle  |
| 0001 | 0101           | 0               | 0111 | ✓ Shift          |                |
| 0001 | 0101           | 0               | 0111 | ✓ Shift          | } Fourth cycle |
|      |                |                 |      |                  |                |

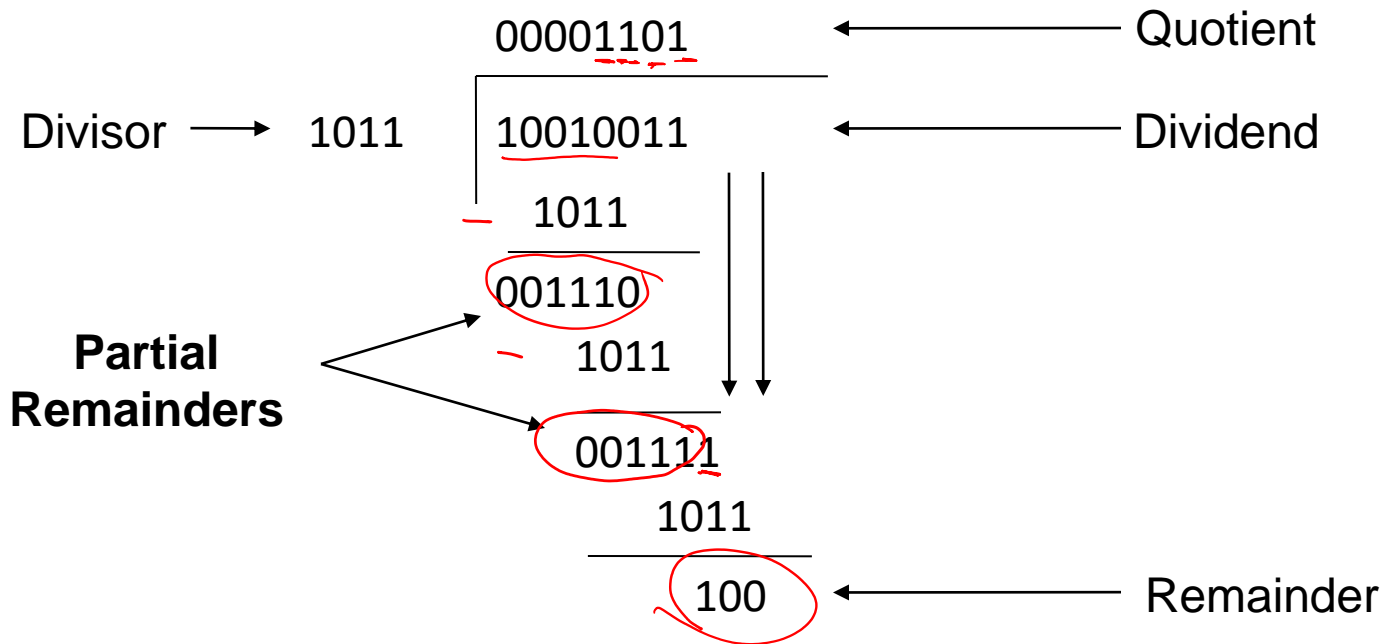
Product in {A, Q}

# Division

- More complex than multiplication
- Negative numbers are really bad!
- Based on long division

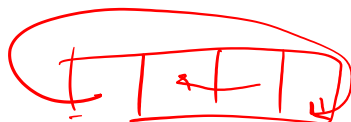
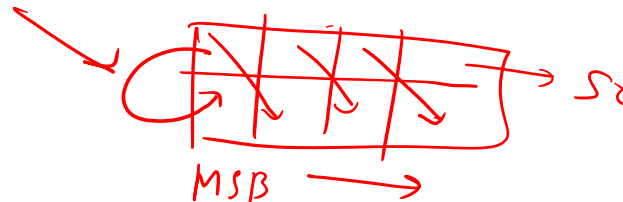
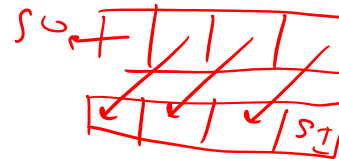
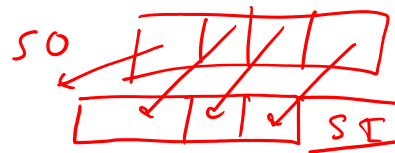


# Division of Unsigned Binary Integers



# Logical Operations

- NOT
- AND
- OR
- XOR
- Logical Shift Left
- Logical Shift Right
- Arithmetic Shift Left →
- ✓ Arithmetic Shift Right
- Rotate Left/Right



# Real Numbers

Unsigned integer

Integer

Signed integer

Sign

Integer

Unsigned fixed point

Integer

Fraction

Signed fixed point

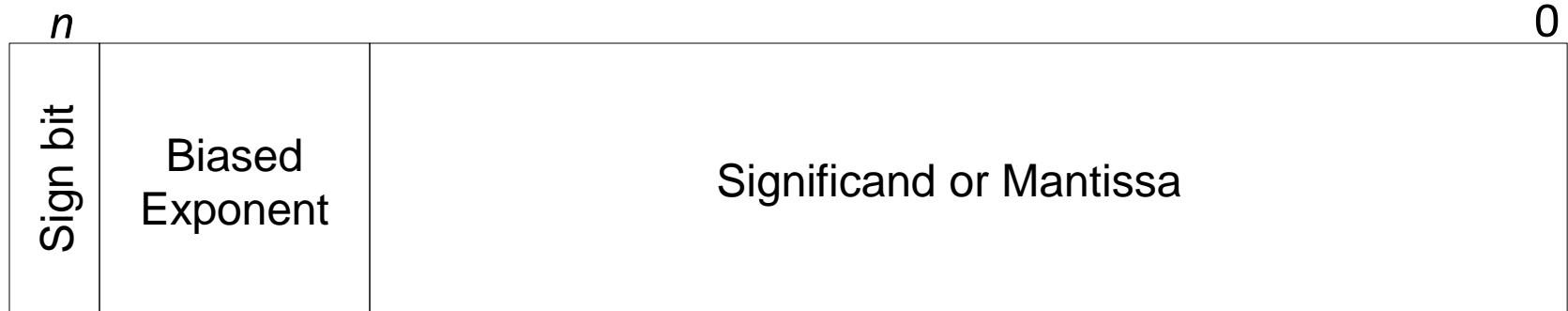
Sign

Integer

Fraction

- Numbers with fractions
- Could be done in pure binary
  - $1001.1010 = 2^4 + 2^0 + 2^{-1} + 2^{-3} = 9.625$  15.
- Where is the binary point?
  - Fixed?
    - Very limited ✓
  - Moving?
    - How do you show where it is?
- Another limitation:** Very large and very small numbers may not be represented
  - Solution in decimal is to use scientific notation
    - $976,000,000,000,000. = 9.76 \times 10^{14}$
    - $0.00000000000000976 = 9.76 \times 10^{-14}$

# Floating-Point (IEEE 754 Standard)



- The same approach used with scientific notation:  

$$\pm \text{Significand} \times 2^{\text{Exponent}}$$
- Binary point is fixed between the sign bit and body of the mantissa
- Exponent indicates the point's place value (*point position*)
- IEEE 754 Standards for floating point storage:
  - ✓ Single Precision (32 bit): 1 sign bit, 8 bit exponent, and 23 bit mantissa
  - Double Precision (64 bit): 1 sign bit, 11 bit exponent, and 52 bit mantissa

# Floating-Point (IEEE 754 Standard)

- **Bias** is a fixed value that is subtracted from the field to get the true exponent value
  - **Biased exponent** means exponent is in **excess** (*biased notation*)
  - $\text{Bias} = 2^{k-1} - 1$ 
    - $k$  = number of bits for the binary exponent
- For single precision (32 bits):  $k = 8$ 
  - $\text{Bias} = 2^7 - 1 = 127$  (8 bit represents 0 to 255)
  - True exponent values range from -127 to +128
- Leading 1 of the significand is not stored for normalized FP numbers
  - Considered a "*hidden bit*"

# Floating Point Examples (Single Precision)



(a) Format

Handwritten notes: *32 bits* (under the entire format), *8* (under the biased exponent), and *23* (under the significand).

|                                |     |          |          |                          |     |                             |
|--------------------------------|-----|----------|----------|--------------------------|-----|-----------------------------|
| $1.1010001 \times 2^{10100}$   | $=$ | <u>0</u> | 10010011 | 101000100000000000000000 | $=$ | $1.6328125 \times 2^{20}$   |
| $-1.1010001 \times 2^{10100}$  | $=$ | <u>1</u> | 10010011 | 101000100000000000000000 | $=$ | $-1.6328125 \times 2^{20}$  |
| $1.1010001 \times 2^{-10100}$  | $=$ | 0        | 01101011 | 101000100000000000000000 | $=$ | $1.6328125 \times 2^{-20}$  |
| $-1.1010001 \times 2^{-10100}$ | $=$ | 1        | 01101011 | 101000100000000000000000 | $=$ | $-1.6328125 \times 2^{-20}$ |

(b) Examples

# Conversion Example 1

- Convert  $100_{10}$  to single precision floating-point
- Convert to binary

$$100_{10} = 0110\overset{6}{0}\overset{4}{1}\overset{3}{0}\overset{2}{0}_2$$

|     |    |    |    |   |   |   |   |
|-----|----|----|----|---|---|---|---|
| 128 | 64 | 32 | 16 | 8 | 4 | 2 | 1 |
| 0   | 1  | 1  | 0  | 0 | 1 | 0 | 0 |

- In a binary representation form of 1.xxx:

$$\cancel{0}110\cancel{0}100 = 1.1001\cancel{00} \times 2^6 = \boxed{1.1001} \times 2^6$$

- Exponent = 6 ✓

- Biased Exponent = 6 + 127 = 133 =  $\overset{128}{1}000\overset{4}{0}101\overset{1}{1}$  (8 bits) ✓

- Sign = 0 (positive) ✓

- Significand/mantissa = 100100 (23 bits: pad with trailing 0s)

- Normalized - leading 1 of the significand is not stored

- Final answer:

- $100_{10} = \underline{0100} \underline{0010} \underline{1100} \underline{1000} \dots$  (32 bits in total)

- $100_{10} = \underline{0x42C80000}$  (in hex)

# Conversion Example 2

- Convert  $-175_{10}$  to single precision floating-point
  - $-175_{10} = 10101111_2$  (with bit weights 128, 64, 32, 16, 8, 4, 2, 1 written above)
  - $10101111 = 1.0101111 \times 2^7$  (the leading 1 is circled in red)
  - Exponent = 7 ✓
  - Biased Exponent =  $7 + 127 = 134 = 10000110$  (8 bits) (134 is circled in red, and 128, 64, 32, 16, 8, 4, 2 are written above)
  - Sign = 1 (negative) ✓
  - Significand/mantissa = 0101111 (23 bits: pad with trailing 0s)
    - Normalized - leading 1 of the significand is not stored
- Final answer:
  - $-175_{10} = 11000011001011110000 \dots$  (32 bits in total)
  - $-175_{10} = 0xC32F0000$



# Example 3: Converting Back

- Convert 0xC32F0000 to decimal
- Extract components from 1100 0011 0010 1111 ...
  - ✓ Sign = 1 (negative) ✓
  - Biased Exponent = 1000 0110 (8 bits) = 134<sub>10</sub>
  - Unbias: 134 - 127 = 7
  - Exponent = 7 ✓
  - Significand/mantissa = 0101111 (23 bits: remove trailing 0s)
  - Denormalized significand = 1.0101111
- Move binary point by 7 (exponent) places = 1010 1111.
- Convert to decimal
  - <sup>128 32 8 4 2 1</sup> 1010 1111<sub>2</sub> = 175<sub>10</sub>
  - Sign is negative, so -175<sub>10</sub>

# Conversion Example 4

- Convert **0x41C8000** to decimal
  - 0100 0001 1100 1000 0000 ...
  - Sign = 0 (positive) ✓
  - Biased Exponent = 1000 0011 (8 bits) =  $131_{10}$  ✓
  - Unbias:  $131 - 127 = 4$  ✓
  - Exponent = 4 ✓
  - Significand/mantissa = 1001 ✓ (23 bits: remove trailing 0s)
  - Denormalized significand = 1.1001
- Move binary point by 4 (exponent) places = 11001 <sup>1/c 8 1</sup>
- $11001_2 = 25_{10}$  ✓

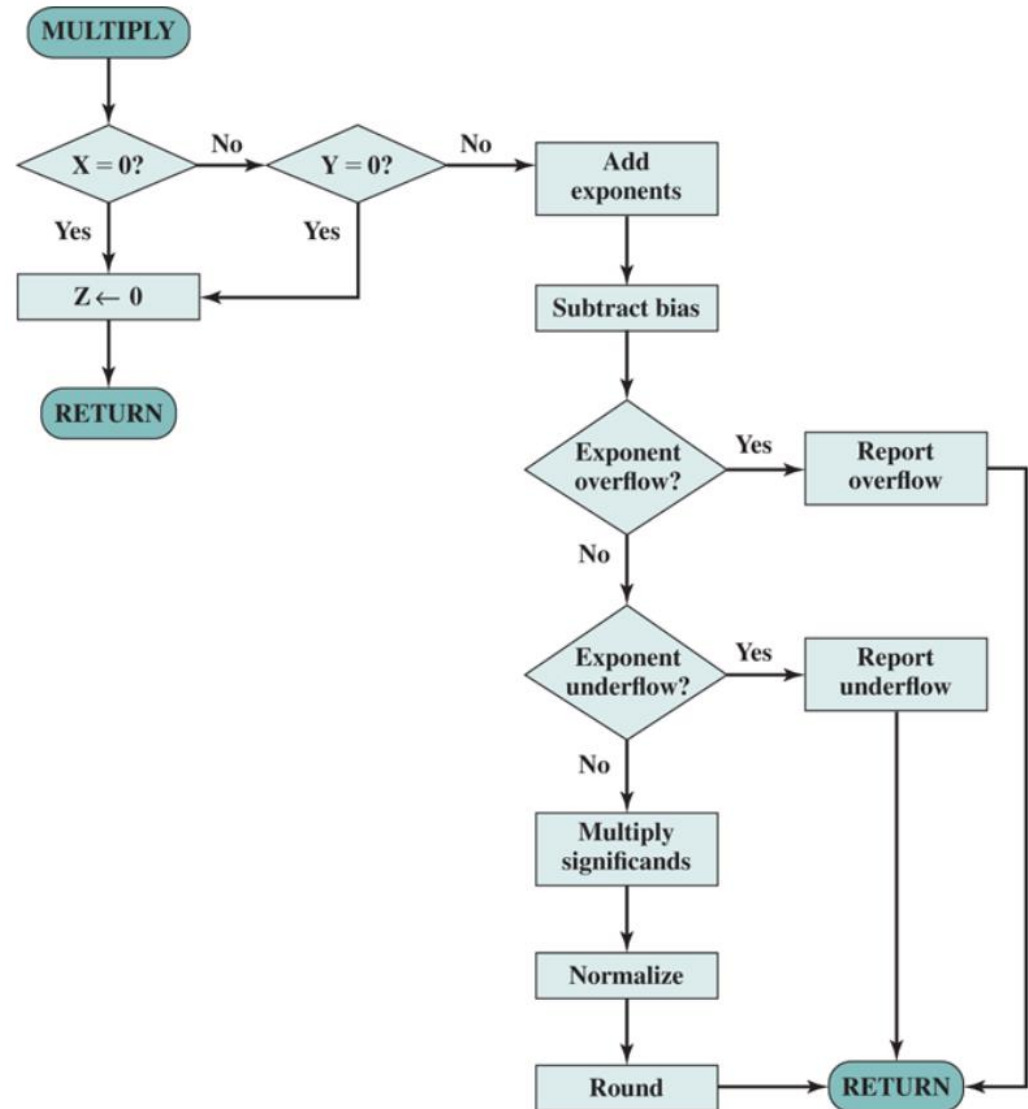
# FP Arithmetic: ADD and SUB

- Check for zeros
- Align significands (*adjusting exponents*)
- Add or subtract significands
- Normalize result

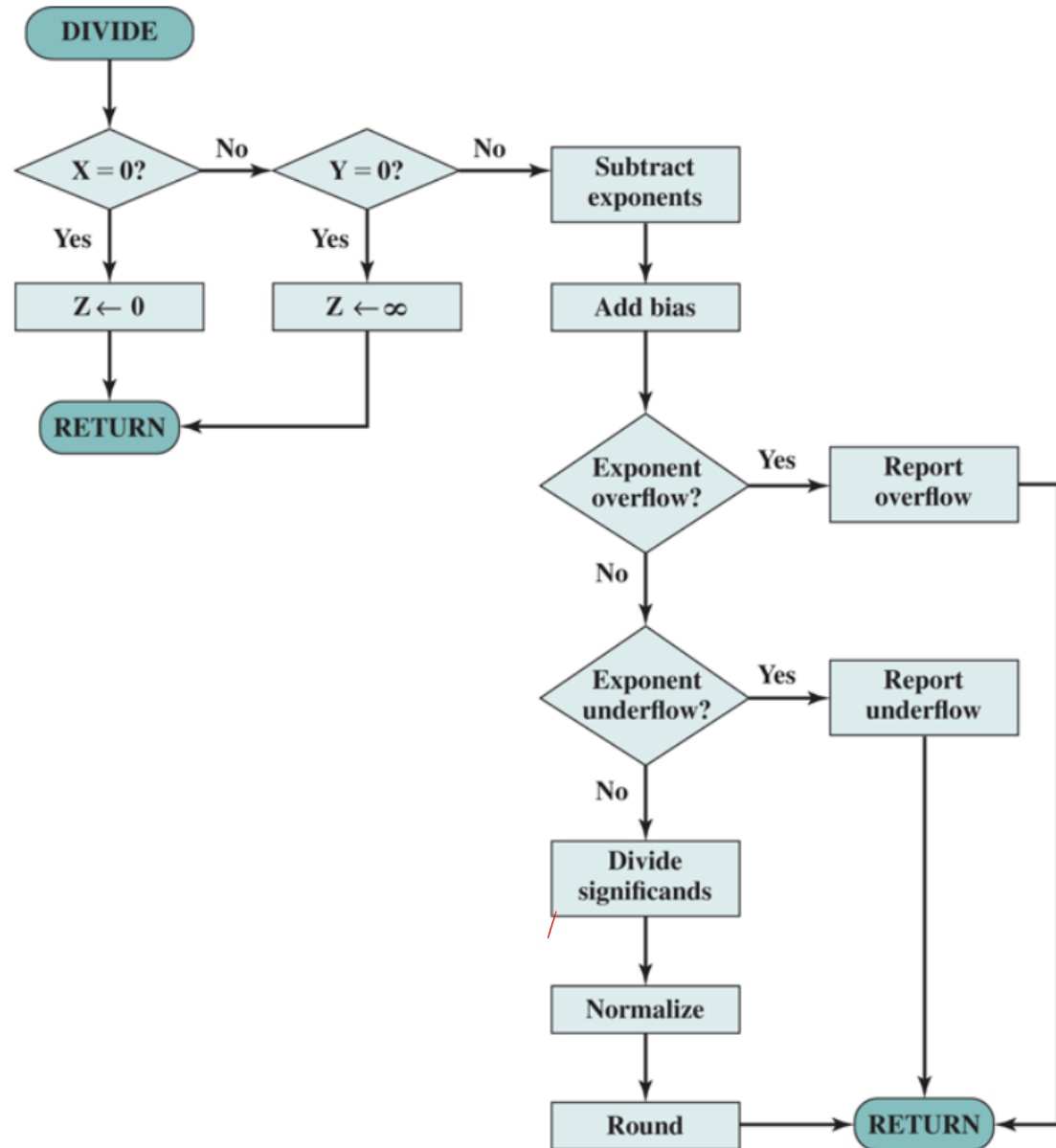
# FP Arithmetic: MUL and DIV

- Check for zeros
- Add/subtract exponents
- Multiply/divide significands (*watch sign*)
- Normalize
- Round
- All intermediate results should be in double length storage

# Floating-Point Multiplication



# Floating-Point Division





CpE 3202  
Computer Organization & Architecture

# End of Lecture

Note: The lecture slide was based on the accompanying lecture material from the Computer Organization and Architecture textbook by William Stallings. All copyright belong to the latter. For instructional use only. **Do not distribute.**

## References:

- Stallings, W. Computer Organization and Architecture, 6th edition, Pearson Education, Inc. (2003).
- Stallings, W. Computer Organization and Architecture, 11<sup>th</sup> edition, Pearson Education, Inc. (2019).
- <https://www.tutorialspoint.com/fixed-point-and-floating-point-number-representations>
- DeGroat, Joanne. "IEEE Floating Point", Lecture Slides.