

# CHAPTER 3-2

Arithmetic for Computers

---

By Pattama Longani  
Collage of arts, media and Technology

# REAL NUMBER

- programming languages support numbers with fractions, which are called *reals* in mathematics.
- Here are some examples of reals:
  - $3.14159265 \dots_{\text{ten}}$  (pi)
  - $2.71828 \dots_{\text{ten}}$  (*e*)
  - $0.0000000001_{\text{ten}}$  or  $1.0_{\text{ten}} \times 10^{-9}$
  - $3,155,760,000_{\text{ten}}$  or  $3.15576_{\text{ten}} \times 10^9$

# FLOATING POINT

- **Scientific notation:** a single digit to the left of the **decimal point**.

- $-2.34 \times 10^{56}$



normalized

- $+0.002 \times 10^{-4}$



not normalized

- $+987.02 \times 10^9$



- In binary

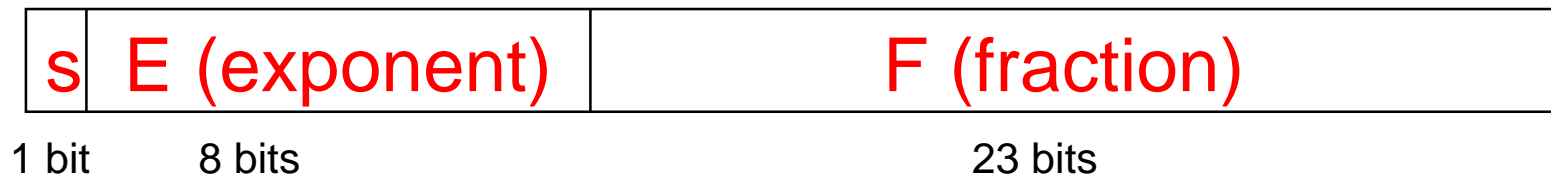
- $\pm 1.XXX_2 \times 2^{yyy}$

- $x = \textit{fraction}, y = \textit{exponent}$

- such numbers is called **floating point** because it represents numbers in which the binary point is not fixed.

# FLOATING POINT

- A designer of a floating-point representation must find a compromise between the size of the **fraction** and the size of the **exponent**.
- This representation is called **sign and magnitude**, since the sign is a separate bit from the rest of the number.



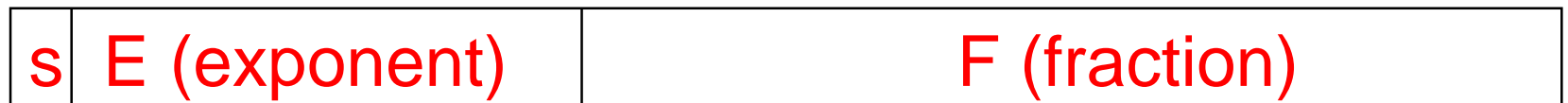
**Single Precision Floating Point**

# EXCEPTION IN FLOATING POINT

- 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**
- One way to reduce the chance of underflow or overflow is to offer another format that has a larger exponent field

# FLOATING POINT

- Double precision – takes two MIPS words



1 bit

11 bits

20 bits



32 bits

- These formats go beyond MIPS. They are part of the **IEEE 754 floating-point standard**, found in virtually every computer invented since 1980.

# FLOATING POINT STANDARD

- Defined by IEEE Std 754-1985
- Developed in response to divergence of representations
  - Portability issues for scientific code
- Now almost universally adopted
- Two representations
  - Single precision (32-bit)
  - Double precision (64-bit)

# IEEE FLOATING-POINT FORMAT

single: 8 bits

double: 11 bits

single: 23 bits

double: 52 bits



$$x = (-1)^S \times (1 + \text{Fraction}) \times 2^{(\text{Exponent} - \text{Bias})}$$

- S: sign bit (0  $\Rightarrow$  non-negative, 1  $\Rightarrow$  negative)
- Exponent: excess representation: actual exponent + Bias
  - Ensures exponent is unsigned
  - Single: Bias = 127 =  $(2^8/2) - 1$
  - Double: Bias = 1023 =  $(2^{11}/2) - 1$

$$(-1)^S \times (1 + (s_1 \times 2^{-1}) + (s_2 \times 2^{-2}) + (s_3 \times 2^{-3}) + (s_4 \times 2^{-4}) + \dots) \times 2^E$$



# EXAMPLE : Represent $-0.75$

- $0.75 = 75/100_{\text{ten}} = 3/4_{\text{ten}} = 3/2^2_{\text{ten}}$   
 $= 11_2 / 2^2_{\text{ten}} = 1.1_2 \times 2^{-1}$  ← Exponent-Bias

- $(-1)^1 \times 1.1_2 \times 2^{-1}$   $x = (-1)^S \times (1 + \text{Fraction}) \times 2^{(\text{Exponent} - \text{Bias})}$

- $S = 1$

- Fraction =  $1000...00_2$

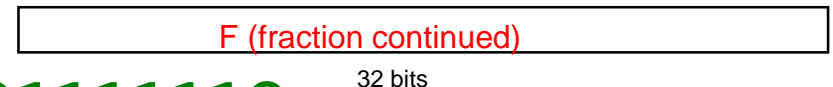
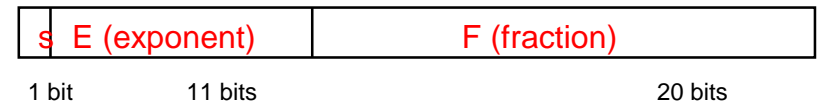
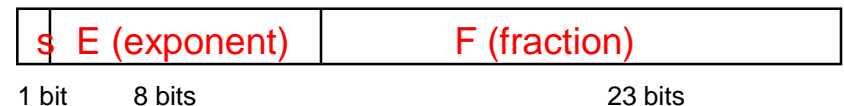
- Exponent =  $-1 + \text{Bias}$

- Single:  $-1 + 127 = 126 = 01111110_2$

- Double:  $-1 + 1023 = 1022 = 011111111110_2$

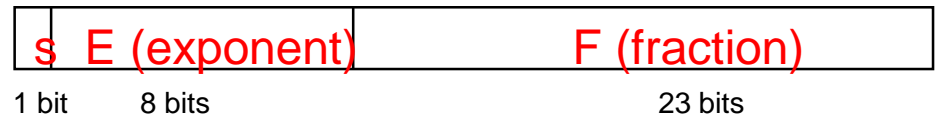
- Single: **1011111110** $1000...00$

- Double: **10111111111110** $1000...00$



# EXAMPLE : What number is represented by the single-precision float

**110000001**01000...00



- $S = 1$
- Fraction =  $01000...00_2$
- Exponent =  $10000001_2 = 129$
- $X = (-1)^1 \times (1.01_2) \times 2^{(129 - 127)} \leftarrow \text{Exponent-Bias}$   
 $= (-1) \times 1.01_2 \times 2^2$   
 $= (-1) \times 101_2$   
 $= -5.0$

# FLOATING-POINT ADDITION

Consider a 4-digit decimal example

$$9.999 \times 10^1 + 1.610 \times 10^{-1}$$

1. Align decimal points

Shift number with **smaller exponent**

$$9.999 \times 10^1 + 0.016 \times 10^1$$

2. Add significands

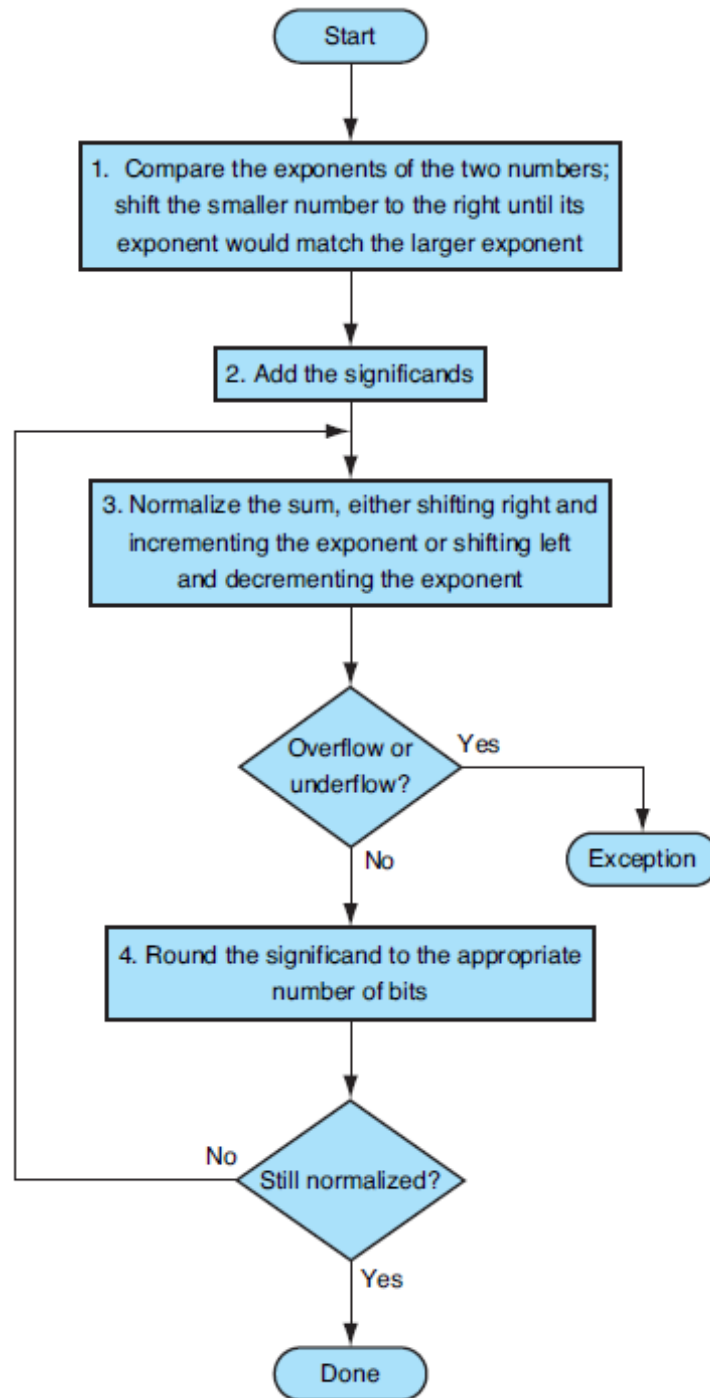
$$9.999 \times 10^1 + 0.016 \times 10^1 = 10.015 \times 10^1$$

3. Normalize result & check for over/underflow

$$1.0015 \times 10^2$$

4. Round and renormalize if necessary

$$1.002 \times 10^2$$



# FLOATING-POINT ADDITION

Now consider a 4-digit binary example

$$1.000_2 \times 2^{-1} + -1.110_2 \times 2^{-2} \quad (0.5 + -0.4375)$$

1. Align binary points

Shift number with smaller exponent

$$1.000_2 \times 2^{-1} + -0.111_2 \times 2^{-1}$$

2. Add significands

$$1.000_2 \times 2^{-1} + -0.111_2 \times 2^{-1} = 0.001_2 \times 2^{-1}$$

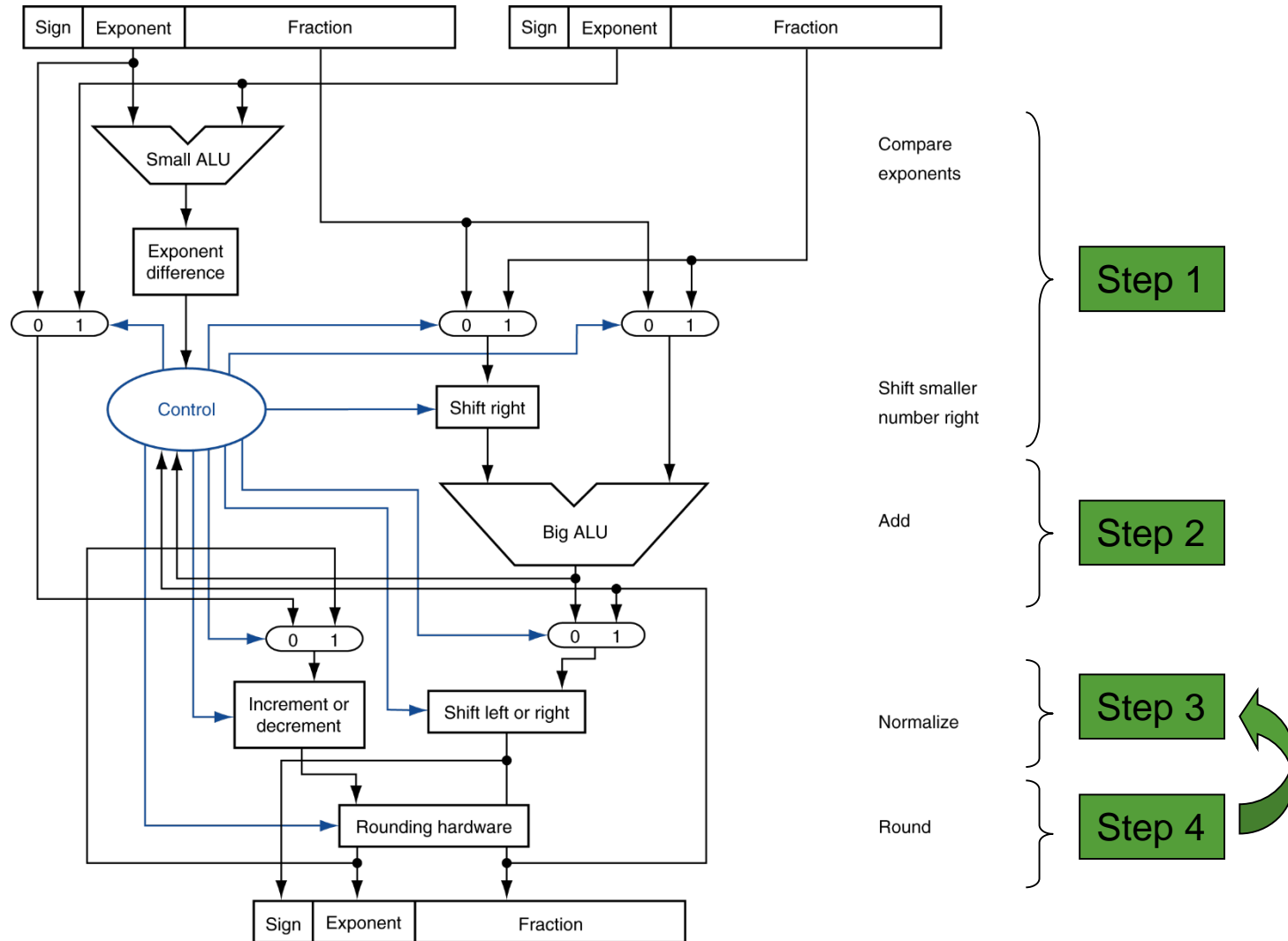
3. Normalize result & check for over/underflow

$$1.000_2 \times 2^{-4}, \text{ with no over/underflow}$$

4. Round and renormalize if necessary

$$1.000_2 \times 2^{-4} \text{ (no change)} = 0.0625$$

# FP ADDER HARDWARE



# FP ADDER HARDWARE

- Much more complex than integer adder
- Doing it in one clock cycle would take too long
  - Much longer than integer operations
  - Slower clock would penalize all instructions
- FP adder usually takes several cycles
  - Can be pipelined