

# Floating Point

부동소수점

## ■ Representing non-integer numbers

- $\pi = 3.141592\dots$
- $e = 2.71828\dots$
- $0.000000001 = 1.0 \times 10^{-9}$
- $3155760000 = 3.15576 \times 10^9$

} scientific notations

## ■ Scientific notation

- Single digit to the left of the decimal point
- Normalized number:  $1.0 \times 10^{-9}$ 
  - no leading zero
  - $0.1 \times 10^{-8}$ ,  $10.0 \times 10^{-10}$  are not normalized

- Binary number in scientific notation:  $1.0_2 \times 2^{-1}$

## ■ Floating point numbers

- Numbers in which binary point is not fixed
  - No fixed number of digits before and after the point

$\Rightarrow 2.34 \times 10^{56}$  (normalized)  
 $= \underline{\underline{2.34e^{56}}}$

$0.002 \times 10^{(-4)}$  (not normalized)  
 $= \underline{\underline{2.0e^{-7}}}$

action (mantissa)

23 bits

mantissa

1.bbbbbb  $\times 2^{\text{bbb}}$

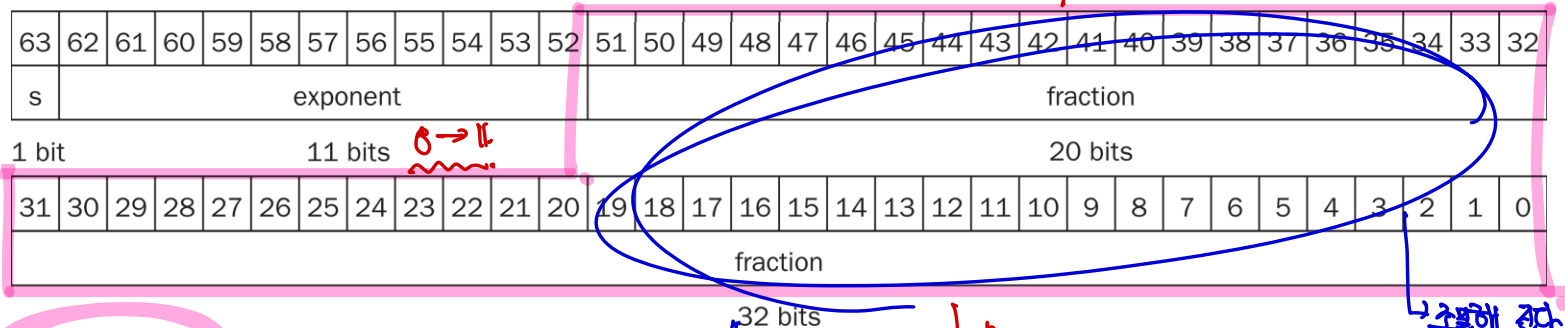
b  $\rightarrow$  decimal digit

Binary Floating point representation

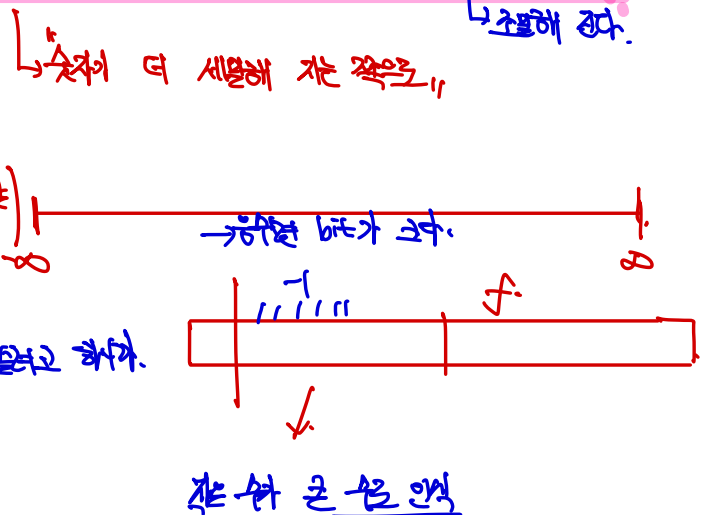
1.110110101  $\times 2^6$

# Floating Point Representation

- Single precision: 32 bits
- Double precision: 64 bits
  - 11 bits for exponent, 52 bits for fraction
  - Range (in decimal):  $2.0 \times 10^{-308} \sim 2.0 \times 10^{308}$
  - benefit: Increased precision from larger fraction bits



- IEEE 754 floating point standard  $1.2 \times 2^9$ 
  - Exponent 00000000, Fraction 0  $\rightarrow 0$
  - E: 11111111, F: 0  $\rightarrow$  infinity
  - E: 1-254, F: anything  $\rightarrow$  normal FP number
  - Consideration for sorting
    - MSB is used as a sign bit
    - Exponent comes before fraction part



# Floating Point Representation

## Biased notation

•  $X = 1.0 \times 2^{-1}$

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	...
0	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	...

•  $Y = 1.0 \times 2^{+1}$

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	...
0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	...

- In sorting, X looks larger than Y.

## Rearranging the exponent value range

• 00000000 ~ 11111111

most negative ~ most positive

## IEEE uses a bias of 127 for single precision (1023 for double P)

• -1:  $-1 + 127 = 126 = 0111\ 1110_2$

• 1:  $1 + 127 = 128 = 1000\ 0000_2$

single precision

후에 설명

Precision Range

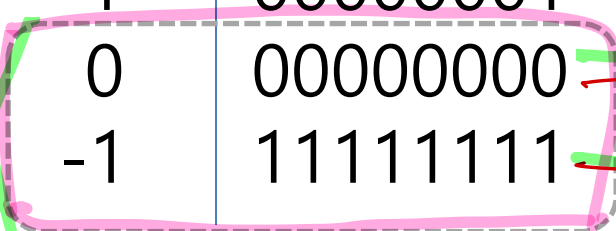
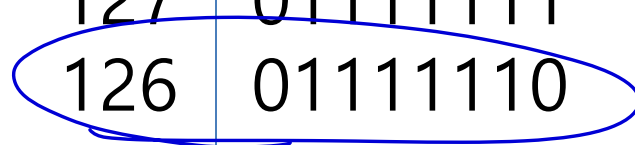
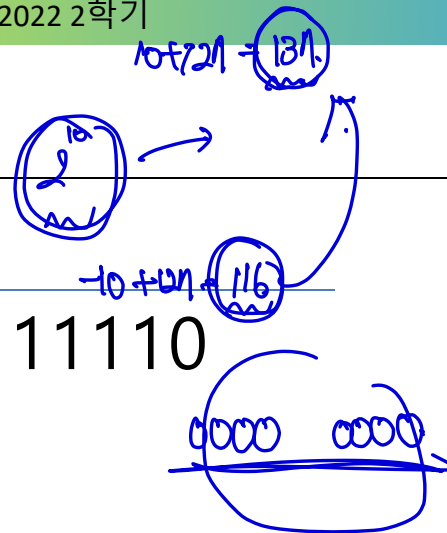
$\pm 1.000000000000000000000000 \times 2^{-126}$

$\pm 1.111111111111111111111111 \times 2^{+127}$

# Floating Point Bias

Exponent	Adjusted Exponent
127	254
.	.
.	.
.	.
1	128
0	127
-1	126
.	.
.	.
.	.
-126	1
-127	0
-128	-1

2's complement



예외

# IEEE 754 Floating-Point Format

## ■ Encoding of the single precision floating-point numbers

- Exponent : 0, Fraction 0 → 0
- Exponent : 255, Fraction : 0 → infinity
- Exponent : 1-254, Fraction : anything → normal FP number
- Exponent : 255, Fraction : Nonzero → NaN *not a number*
  - NaN is a symbol for the result of invalid operations (e.g. 0/0).

Single precision		Double precision		Object represented
Exponent	Fraction	Exponent	Fraction	
0	0	0	0	0
0	Nonzero	0	Nonzero	± denormalized number
1-254	Anything	1-2046	Anything	± floating-point number
255	0	2047	0	± infinity
255	Nonzero	2047	Nonzero	NaN (Not a Number)

# Floating Point Example → 실수도 문젠 많이 풀어야 하는.

## Represent -0.75

- $= -3/4_{10}$  or  $-3/2^2_{10}$
- $= -11_2/2^2_{10} = -0.11_2$
- $= -0.11_2 \times 2^0$
- $= -1.1_2 \times 2^{-1} =$

$$(-1)^s \times (1 + \text{Fraction}) \times 2^{(\text{Exponent}-127)}$$

$$(-1)^1 \times (1 + 0.100...000) \times 2^{126}$$

Exponent

0.10000 ...

$$-0.75 = -\frac{3}{4} = -11 \div 2^2 = -0.11$$

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
1	0	1	1	1	1	1	1	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	

↑ bit

8 bits

23 bits

1 bit

8 bits

23 bits

## What number does this represent?

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
1	1	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

$$(-1)^s \times (1 + \text{Fraction}) \times 2^{(\text{Exponent}-\text{Bias})} = (-1)^1 \times (1 + 0.25) \times 2^{(129-127)}$$

$$\begin{aligned}
 &= -1 \times 1.25 \times 2^2 \\
 &= -1.25 \times 4 \\
 &= -5.0
 \end{aligned}$$

$$\begin{aligned}
 &(-1)^1 \times (1 + 0.01) \times 2^2 \\
 &= -1 \times 1.01 \times 4 \\
 &= -4.04
 \end{aligned}$$

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

Exponent - 121.

$$(-1)^e \times (1 + \text{fraction}) \times 2.$$

$$\begin{aligned} &\underline{-0.15} \\ &-3/4 \end{aligned}$$

$$\begin{aligned} &-11. \\ &= -0.11 \quad \leftarrow 284. \\ &\approx -1.1 \times 2^{-1} \\ &\quad \text{~~~~~} \end{aligned}$$

$$\frac{1+121}{1} = 126$$



1000 0001  
0111 1111  
1000 0001  
1000 0001

# Floating Point Addition

## ■ Consider:

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

- Assume we can store only 4 digits of significand

→ 좌와 exponent를 분석해서 이동

① 둘 다 같은 차수로  
맞아야겠다.

## ■ Steps

- ① Align decimal points – shift smaller exponent number

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

→ 좌를 오른쪽 쪽에 맞춘다.

- ② Add significands → 가산 숫자를 더한다.

- $9.999 + 0.016 = 10.015$

- Result:  $10.015 \times 10^1$

- ③ Normalize, check over/underflow

- $1.0015 \times 10^2$  → "좌를 맞춘다."

- ④ Round it to 4 digits

Round off

- $1.002 \times 10^2$

반올림

# Binary FP Addition

■  $0.5_{10} + (-0.4375_{10})$

•  $0.5_{10} = 0.1_2 = 0.1 \times 2^0 = 1.000 \times 2^{-1}$

•  $0.4375_{10} = 7/16 = 7/2^4 = 111 \times 2^{-4} = 1.110 \times 2^{-2}$

■ Align

•  $1.000 \times 2^{-1} - 1.110 \times 2^{-2} = 1.000 \times 2^{-1} - 0.110 \times 2^{-1}$

■ Add significands

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

■ Normalize, check over/underflow

•  $1.0 \times 2^{-4}$

■ Round

• No need

$0.5 + (-0.4375)$

$= 0.5 = 1.000 \times 2^{-1}$

$0.4375 = \frac{7}{16} = \frac{111}{1000} = \frac{1.110}{1000} = 1.110 \times 2^{-2}$

$1.110 \times 2^{-2} = 0.111$

$0.5 \rightarrow \frac{1}{2} = 1 \times 2^{-1} = 1.000 \times 2^{-1}$   
 $-0.4375 = \frac{7}{16} = 111 \times 2^{-4} = 1.110 \times 2^{-2}$   
 $= 0.111 \times 2^{-1}$

$\begin{array}{r} 1.000 \\ - 0.111 \\ \hline 1001 \end{array}$

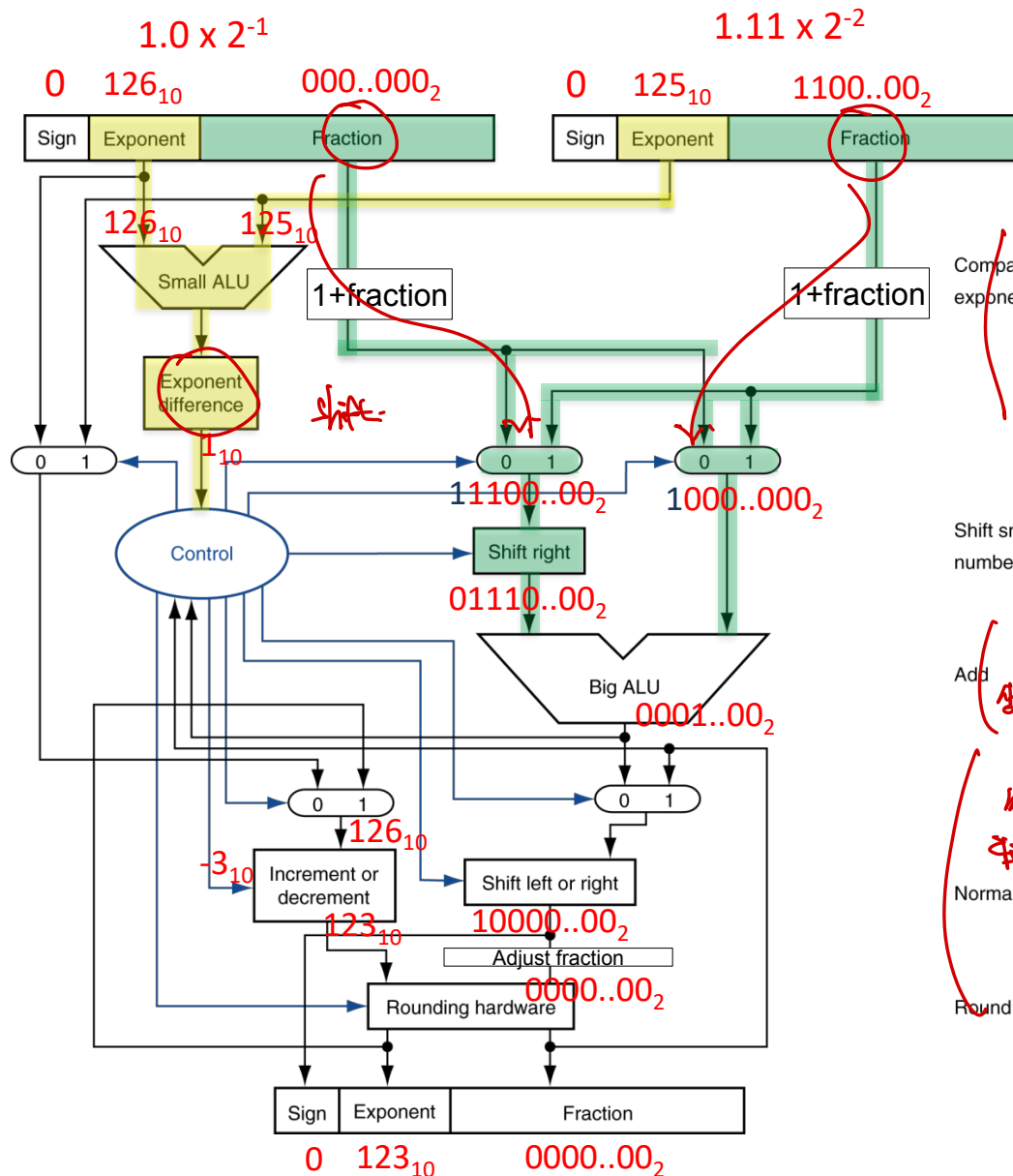
$0.001$

1000

$\frac{0.001 \times 2^{-1}}{1 \times 2^{-1}} = \frac{1}{8}$

$\frac{1}{8} = \frac{1}{8}$

# FP Adder Hardware



Compare exponents  
→ Exponent  
Fraction  
여로작도 취

Shift smaller number right

Add  
하드웨어의 상용화물 취  
↓  
mega fcs  
화 화 화  
Normalize  
화 화 화  
Round  
화 화 화



Floating Point Unit

