

IEEE754 Floating point representation

Bias.

Biased Exponent.

Decimal to Floating point.

Floating point to Decimal.

Floating point Addition, sub.

Floating point mul.

Overflow - Underflow.

Chapter 3 (video 1)

$$12 \rightarrow X_{21}.$$

$$\text{addi } X_{21}, X_0, 12$$

$$12.0 \times 10^3 \rightarrow X_{21} \quad (\text{How?})$$

$\rightarrow (f_0 - f_{31})$

$$\begin{array}{c} 4.5 \times 10^6 \\ \downarrow \text{co-eff} \\ 4.5 \times 10^0 \end{array} \quad \begin{array}{c} \uparrow \text{Base} \\ \downarrow \text{Exponent} \end{array}$$

single precision (32-bit)

double " (64 bit)

IEEE
754

Lulinox[®]
Luliconazole 1% INN

Azonox[®]
Itraconazole USP

Die

Normalized

only one digit before the decimal/bin point.
Digit must be a non-zero number.

(1) $1.2 \underline{123.456}$ (X)
↓ ↑
3 digit

(2) 1.2345 (✓)
↓
1 digit, $\neq 0$

(3) 0.123 (X)
↓
1 digit, it is zero

(4) 9.123 (✓)
↓
1 digit, it's not zero

(i) 1.0111 (✓)

(ii) $\underline{1101.1000}$ (X)
↓
not one

(iii) 0.11101 (X)
↓
It's zero

Steps to normalize.

To normalize a number you need to shift the decimal point left or right, until there is only one single digit

$$123.45$$

$$123$$

$$1234.5678$$

$$1.2345678$$

$$11010.00110$$

$$1.101000110$$

if you left-shift, the number of times you left shift will be added to the exponent

$$112.54 \times 10^{35}$$

$$\rightarrow 1.1254 \times 10^{(35+2)} = 1.125 \times 10^{37}$$

$$456.123.$$

$$\rightarrow 4.56123 \times 10^2.$$

If you shift right, the number of
you right shift will be subtracted
the exponent

$$0.0065 \\ \Rightarrow 6.5 \times 10^{-3}$$

$$\begin{array}{l|l} * 110110.1101_2 & (0.000010101)_2 \\ \Rightarrow 110110.1101 \times 2^0 & \Rightarrow 1.0101 \times 2^{-5} \\ \Rightarrow 1.101101101 \times 2^5 & \end{array}$$

IEEE Floating point format.

(32) single p.	Sig. Bit	Expon.	Fraction.
	1 bit	8 bit	23 bit
Double p. (64)	1 bit	11 bit	52 bit.

single precision. (32 bit)

Example :

$$\frac{01011}{\downarrow} \times 2^{30}$$

6 bit

$$\frac{1101011101011}{\downarrow} \times 2^{30}$$

12 bit

Exponent bit will be unsigned \rightarrow Using Bias.

$\hookrightarrow 0 - 255$.

but ~~0 and~~ 0 and 255 is reserved.

\therefore range for bias $\Rightarrow 1 - 254$.

Double precision.

11 bit unsigned.

$$0 - 2^{11} - 1$$

$$\Rightarrow 0 - 2047$$

but 0 and 2047 are reserved.

so, the range,

$$1 \text{ to } 2047$$

$$\text{Bias} = (2^8 - 1) - 1$$

$$\Rightarrow 127$$

Actual $(1 - 127)$ to $(2047 - 127)$

(-126) to ~~(-126)~~ (127)

Bias & Bias exponent.

→ size of the exponent field.

$$\text{Bias} = 2^{n-1} - 1$$

$$\text{Biased exponent} = \text{Actual expo} + \text{Bias}$$

Ex: 1.1011×2^{34} → Actual exponent

→ Find the biased expo in IEEE 754 single precision format.

↳ Expo → 8 bits.

$$\begin{aligned} \text{Bias} &= 2^{8-1} - 1 \\ &= 127 \end{aligned}$$

① Check if the num is normalized.

$$\begin{aligned} \therefore \text{Biased exponent} &= (34 + 127) \\ &= 161. \end{aligned}$$

Ex: S. P. f → 11.1011×2^{-8} .

① The number is not normalized.

$$\rightarrow 11.1011 \times 2^{-8}$$

$$\hookrightarrow 1.11011 \times 2^{-8+1} =$$

$$= 1.11011 \times 2^{-7}$$

$$\therefore \text{Biased exponent} = (-7 + 127) = 120$$

11.1011 $\times 2^{-8}$, Find the biased exponent of
given number in 15-bit IEEE 754 format
where size of fraction is 8 bits.

Fraction = 8 bit

sign-bit = 1 bit

9 bit \therefore Exponent = $(15 - 9) = 6$ bits.

$$\therefore \text{Bias} = 2^{6-1} - 1 = 31$$

~~\therefore Biased~~

$$\text{number} = 11.1011 \times 2^{-8}$$

$$\hookrightarrow 1.11011 \times 2^{-7}$$

$$\therefore \text{Biased exponent} = (31 - 7) = 24 \text{ bits.}$$

Decimal to floating point.

- (i) Convert to binary
- (ii) Normalize
- (iii) Biased exponent. ✓
- (iv) sign-bit.
- (v) Find the fraction.

Ex

Convert 50.6749_{10} to 32 bit IEEE-754 floating point representation. 6486

$$(50.6749)_{10} = 110010.1010110011$$

Normalize.

$$1.100101010110011 \times 2^5$$

$$\text{Bias} = 127$$

$$\text{Biased exponent} = 127 + 5 = (132)_{10} = 1000\ 0100$$

$$\text{sign-bit} = 0 \checkmark$$

$$\text{fraction} = \underbrace{100101010110011}_{15\text{ bits}} \underbrace{0000\ 0000}_{8\text{ bits}}$$

0	1000 0100	100101010110011 00000000
---	-----------	--------------------------

Convert -0.0232 to 12 bit IEEE-754 floating point where biased exp is 4 bits.

14848

$$\therefore \text{sign} = 1$$

$$\text{bias} = 4$$

$$\frac{5}{5} \therefore \text{fraction} = 7 \text{ bit}$$

$$(0.0232)_{10} = \frac{1(10111)_2}{2^{10}}$$

$$= (-0.0000010)_2$$

$$\text{normalize, } 1.0 \times 2^{-6}$$

$$\text{bias} = 2^3 - 1$$

$$= 7$$

$$\therefore \text{biased exponent} = 7 - 6 = 1, = 0001$$

$$\text{sign} = 1$$

1	0000100	0000000
---	---------	---------

Floating point to Decimal.

① Hex to Binary

(i) Arrange the binary acc. to format.

(ii) Determine the sign.

(iv) Find out actual exponent.

(v) Conv. fraction to decimal.

(vi) Decimal number = $(-1)^{\text{sign bit}} \times (1 + \text{Fraction}) \times 2^{\text{Exponent}}$

$0x F2400120 \rightarrow$ Single Pre. Float point.

\rightarrow

1	1	1	1	0	0	1	0	0	1	0	0	0	0	0	0
0	0	0	0	0	0	0	1	0	0	1	0	0	0	0	0

fraction.

4194592

sign bit = 1 = negative number.

Exponent = 111000100 = 228.

bias = 127

\therefore actual exponent = $228 - 127$
= 101

fraction = 100 0000 0000 0001 0010 0000

= 0.100 0000 0000 0001 0010 0000

= 0.5000343323

Dec = $(-1)^1 \times (1 + 0.5000343323) \times (2)^{101}$

= - 3.803038843 $\times 10^{30}$

and it up to 6 dec. point.

$$- 3.80303\overset{87}{884} \times 10^{30}.$$

$$\Rightarrow -3.80304 \times 10^{30}$$

Floating point addition.

- Numbers are in binary.
- Normalize A & B.
- Align the bin point so that lower exp match the higher.
- Now add / sub.
- Normalize the result.
- Round if necessary.
- Overflow / underflow.

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

$$\Rightarrow (99.99) + (0.1610)$$

$$= \underline{1100611111} \ 0101 + 0.\underline{001010} \ 0100 \rightarrow \text{2nd exp}$$

$$\Rightarrow (1.1000111111110101 \times 2^6) + (1.0100100 \times 2^{-3})$$

$$\rightarrow \frac{(1.10001111110101 \times 2^6) + 0.000000001010}{100 \times 2^6}$$

$$\Rightarrow 2^6 (\quad + \quad)$$

$$\begin{array}{r}
 1.1000 \ 1111 \ 1111 \ 0101 \\
 0.0000 \ 0000 \ 1010 \ 0100 \\
 \hline
 1.1001 \ 0000 \ 1001 \ 1001
 \end{array}$$

$$\begin{array}{r}
 1.1000 \ 1111 \ 1111 \ 0101 \\
 0.0000 \ 0000 \ 1010 \ 0100 \\
 \hline
 1.1001 \ 0000 \ 1001 \ 1001
 \end{array}$$

Floating point mul.

→ Make sure that the values are in binary.

→ Normalize. → Add the expo.

→ multiply

→ Normalize. → If it was negative → same.

$$(1.110 \times 2^5) \times (1111 \times 2^{-7})$$

$$= (1.110 \times 2^5) \times (1.111 \times 2^{-5})$$

$$= (1.110 \times 1.111) \times (2^{5-5})$$

$$\hookrightarrow 1110 \times 1111 \text{ (each)}$$

$$= 110110010$$

$$\Rightarrow 11.0110010 \times 2^0$$

$$= 1.10110010 \times 2^1$$

$$\Rightarrow -1.10110010 \times 2^1$$

Overflow / underflow

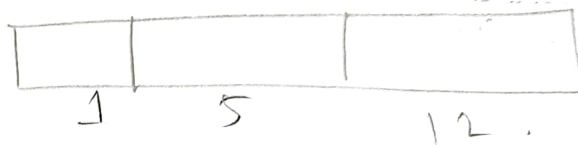
range = (x to y).
 ↑ ↓
 lower upper.

a is a number.

if $(x \leq a \leq y)$:
 perfect

if $(a < x)$:
 under.

if $(y < a)$:
 overflow.



$$1.01 \times 2^{-89} \times 1.01 \times 2^{-89}$$

$$= (1.01 \times 1.01) \times 2^{-178}$$

$$= 1.1001 \times 2^{-178} \rightarrow \text{actual.}$$

$$\text{bias} = 2^{5-1} - 1 = 15$$

$$\text{biased exponent} = -178 + 15 = -163$$

range of biased exponent,

$$0 - 2^5 - 1$$

$$\Rightarrow 0 - 31$$

$$\text{range} \rightarrow (1 - 30)$$

overflow underflow

Exponent
field.

Long multiplication

1000

1001

1000

0000x

0000xx

1000xxx

1001000

Optimize multiplier.

Iteration = $\lceil \text{multiplier bits} \rceil$

① multiplier last bit 1

→ multiplicand + product half MSB. (A)

→ 1 bit right shift (Product) (B)

or, multiplier last bit = 0.

→ 1 bit right shift (product) (C)

product + multiplier.

1000 X 1001

Ite	Multiplicand	Product
0	1000	0000 1001
1	1000	<div> <div>1000 1001 (A)</div> <div>0100 0100 (B)</div> </div>
2	1000	0010 0010 (C)
3	1000	0001 0001 (C)
	1000	<div> <div>1001 0001 (A)</div> <div>0100 1000 (B)</div> </div>

