

$$(120)_{10} \rightarrow (100000010)_2$$

$$\begin{array}{ccc} \boxed{0} & \boxed{100000010} & \boxed{11000100 \dots 0_{23}} \\ S & F & M \end{array}$$

$$-14.125 \quad 255$$

$$\begin{array}{ccc} \boxed{1} & \boxed{100000010} & \boxed{110001000 \rightarrow 0} \\ S & F & M \end{array}$$

Q: 14.125 \rightarrow single precision means 32 bit.

~~00~~

Step 1: convert binary.

2: Normalization.

3: take the bias.

4: convert bias exponent to binary.

5: substitute into required format.

$(14.125)_{10} \rightarrow (1110.001)_{2}$ \rightarrow positive value.

$1110.001 \rightarrow 1.110001 \times 2^3$

$$(-1)^S \times 1.M \times 2^E$$

$$= (-1)^0 \times 1.110001 \times 2^3 \rightarrow \text{True exponent}$$

single precision \rightarrow 127 (bias exponent).

$$2^{127+3} = 2^{130}$$

$$(-1)^0 \times 1.110001 \times$$

→ Power problem
(-) solution.

Exponent 10 (-) :

bias exponent use করে,

Bias exponent = True exponent + Bias.

$$BE = TE + Bias,$$

$\xrightarrow{(127)}$ → bit যা জন্য

⊛ Support equal number of positive and negative.

$$2^8 = 256 = 256 \text{ (০২ ০ পাঠ্য)} = \frac{256}{2} = 127.$$

১২৭ নিম্ন উপস্থাপন avoid করে, -১২৮ নিম্ন, ০ ১২৭ নিম্ন পাঠ্য না।

তারি $(-127 \text{ and } 127)$,
-১২৮ to ১২৮ নিম্ন space exceed করে,

$$2^{11} = 2048 \rightarrow 0 \text{ আর } 2047 = \frac{2047}{2}$$

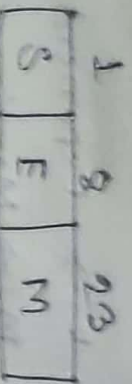
⊛ ১০২৩ → ১০২৩.৫
উপস্থাপন omit করে।

S.E.M (Sign, Exponent and Mantissa)

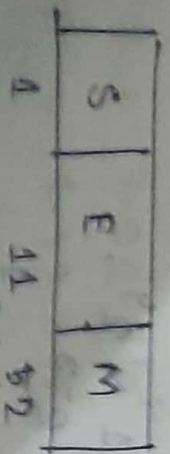
$$(-1)^S \times 1.M \times 2^E$$

→ Required normalization form.

IEEE - 754 - 32 bit format



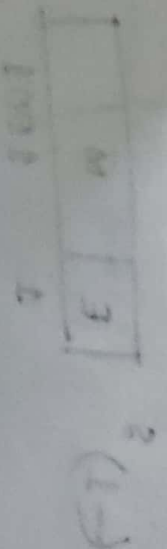
IEEE - 754 - 64 bit format



→ biased exponent

$$2^{OE} = 2^{E + \text{bias value}}$$

(bit)



101.101
0000, 000000

Exercise - 2

floating point number \rightarrow point store 247
 \rightarrow ignore zeros

$$70.0 = 7.00 \times 10^1$$

$$645 = 6.45 \times 10^2$$

$$0.0004 = 4.00 \times 10^{-4}$$

$$5.63 = 5.63 \times 10^0$$

position of \rightarrow (digit representation) like 1.23, 2.45.
the point



0 563

2 645

$$0101.0001 \rightarrow 01.01001 \times 2^2$$

$$11111.101 \rightarrow 1.111101 \times 2^4$$

$$0.000101 \rightarrow 1.01 \times 2^{-3}$$

$$-10.01 \rightarrow -1.001 \times 2^1$$

E \rightarrow Exponent (position of the point)
M \rightarrow Mantissa (digit in the number)



1 1001

neg \rightarrow sign = 1
pos \rightarrow sign = 0

001
1 omitted as bits same.

$$\begin{array}{r}
 9 \rightarrow +9 \quad (01001) \\
 9 \rightarrow -9 \quad (10110) \\
 \hline
 10111
 \end{array}$$

Fixed point numbers

56.7
 78.2
 42.5

} point at
 position fixed.

Floating point numbers

0.6152
 1.275
 23.5

} point at
 position fixed at

Any number where the position of the point is not fixed is called floating point number.

convert the floating to fixed point (cause

floating point number is processing the computer)

$$73.2 = 7.32 \times 10^1$$

$$649 = 6.49 \times 10^2$$

$$.247 = 2.47 \times 10^{-1}$$

$$5.63 = 5.63 \times 10^0$$

$$\begin{array}{ccccccc}
 16 & \leftarrow & 8 & \leftarrow & 4 & \leftarrow & 2 \\
 2^4 & & 2^3 & & 2^2 & & 2^1 \\
 (1) & & 0 & & 1 & & 0
 \end{array}
 \rightarrow (10)_2 \rightarrow (10)_{10}$$

$$\begin{array}{c}
 \boxed{\begin{array}{c} 3 \text{ bit} \\ \downarrow \\ 2^3 = 8 \end{array}}
 \end{array}$$

$$(-4 \text{ to } -1) \quad (0 \text{ to } 3) \quad (-8 \text{ to } -1) \quad (0 \text{ to } 7)$$

$$\begin{array}{c}
 \boxed{\begin{array}{c} 4 \text{ bit} \\ \downarrow \\ 2^4 = 16 \end{array}}
 \end{array}$$

$$\boxed{\begin{array}{c} 8 \text{ bit} \\ \downarrow \\ 2^8 = 256 \end{array}}$$

$$(-128 \text{ to } -1) \quad (0 \text{ to } 127)$$

$$\begin{array}{c}
 7 \text{ ---} \\
 \swarrow \quad \searrow \\
 +7 \quad -7 \\
 (0111) \quad (1000) \\
 +1 \\
 \hline
 1001
 \end{array}$$

101 \rightarrow 2's complement
because starts with 1.

2's complement 777
and 777 777 777

0101 \rightarrow 010
 $\frac{+1}{011}$
011 means (+3)
So, main number (-3)

Unique 2's complement

2's complement:

1's complement + 1

positive number same (0 to +255).

(-1 to -255) convert first 1's complement + 1

8 bit combination is 256 values

1's complement

100 missing 100 255 125 complement = 011

2's complement

$$\begin{array}{r} 011 \\ + 1 \\ \hline 100 \end{array}$$

100 (100)

So, (-255) problem solved so this method can be used.

Short cut:

From right side, copy the number as it is till you get first 1. After getting first 1, convert the remaining digits

एक 100 मिसिंग : 100 means (-0) is
 a represents neg sign and 00 represents 0.
 (-) 0 not possible because zero always positive.

Major failure of this method is (-0).
 So, this method can not be used.

— X —

1's complement

(+) का 1's complement : (+2) का 00000000.

(-2) का 1's complement : 11111111.

2's combination का 2's complement मिसिंग → सिग्न।

111 मिसिंग : 111 means 1's complement
 means 00000000 का 1's complement 11111111.

0 means sign and 00 means 0 का corresponding

binary में convert 00000000 का 1's complement

convert 00000000 (-0) है, Again (-0) problem

So 1's complement can not be used.

	1's complement	Sign magnitude	2's complement
+3	011	0 11	011
+2	010	0 10	010
+1	001	0 01	001
0	000	0 00	000
-1	110	1 01	111
-2	101	1 10	110
-3	100	1 11	101

Sign magnitude:

Method: 1
 $2^3 \ 2^2 \ 2^1 \ 2^0$

1st bit

sign represent value

(+ or -)

(+) or (-) 0 (-) or 1

zero or positive or negative

combination or any

Method: 1

0	0 0 0 0
1	0 0 1
2	0 1 0
3	0 1 1
4	1 0 0
5	1 0 1
6	1 1 0
7	1 1 1

missing check

Lecture - 4

$$(-1)^S \times 1.M \times 2^E$$

S	E	M
1	2	563

S	E	M
1	1	001

$$- 1.001 \times 2^1$$

$$= - 10.01$$

$$\frac{10 \times 10^2}{(+)\ 5 \times 10^2}$$

$$15 \times 10^2$$

* bias क्या होता है ?

To make the calculation faster. exponent pos and neg 2^n addition \rightarrow check and comparison time consuming and multiplication and division \rightarrow वास्तव में time consuming \rightarrow so bias imp.

* simple precision exponent 7 bits

मान , $2^8 = 256 \rightarrow 255$ (02 00 positive कावा).

$$\frac{255}{2} = 127.5 \text{ (fraction omit कावा)}.$$

so, 127 so that it can support a positive and negative. equal number of

* $0.00101.$

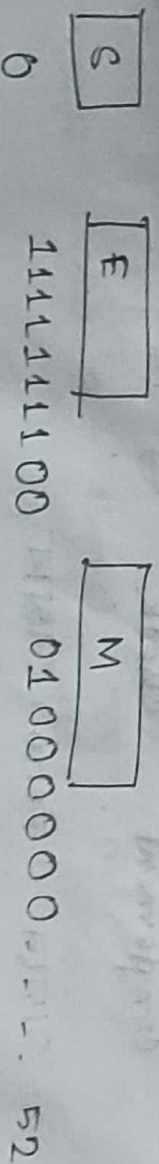
$$= 1.01 \times 10^{-3}$$

$$(-1)^S 1.M \times 2^E.$$

$$M = 01, E = -3.$$

$$D.P = 1023.$$

$$2^{-3} + 1023 = 2^{1020}$$



$$(-1)^0 (1.01) 1111011100.$$

Q: 0101.001 → single precision

$$(-1)^0 \times 1.01001 \times 2^2$$

$$(-1)^0 \times 1.001001 \times 2^{129}$$

$$(-1)^0 \times 1.01001 \times 2^{10000001}$$

Exponent 10 - 255 bias value always positive

Exponent 10 - 255 bias value always positive

Exponent 10 - 255 bias value always positive

Exponent 10 - 255 bias value always positive

Exponent 10 - 255 bias value always positive

Exponent 10 - 255 bias value always positive

Exponent 10 - 255 bias value always positive

Exponent 10 - 255 bias value always positive

Exponent 10 - 255 bias value always positive

Book 2

William

Stalling

→ (Writer)

Q: $(2A3B)_{16}$

0010 1010 0011

1011

$(-1)^0 \times 1.0101000111011 \times 2^{12}$

$(-1)^0 \times 1.0101$

$\boxed{0} \quad \boxed{10001100}$

⊛ Exponent 2 bit and double precision 2 11 bit

power . 210 210 2 00000010 210, 210

0 210 value 2 change 210

float32 means point 32 bits
 float16 means single precision 16 bits
 double precision 64 bits
 32 bit represent range 2³²
 1.01 2³² 1.0100 → 0
 range 2³² to 2³²
 0 to 2³²