# Floating Point Representation

## 32 Bit    IEEE 754

$(-10.75)_{10}$ $\longrightarrow$ Floating point representation $(?)$

Solution →

__step(I)__: Convert the number in binary form.

$(10.75)_{10}$ $\longrightarrow$ Binary

$$
\begin{array}{r|l|l}
2 & 10 & \\
2 & 5 & 0 \\
2 & 2 & 1 \\
2 & 1 & 0 \\
 & 0 & 1
\end{array}
$$

$0.75 \times 2 \rightarrow 1.50 \rightarrow 1$
$0.50 \times 2 \rightarrow 1.00 \rightarrow 1$

$(10.75)_{10}$ $\longrightarrow$ $(1010.11)_2$

__step II__:  Scientific form representation

$(1010.11)_2$ $\longrightarrow$ $(1.\underbrace{01011......}_{Mantissa}) \times 2^3$

__step III__:  Exponent calculation

$x = 3$

Exponent $= 127 + x$

$= 127 + 3 = 130$

$130 \xrightarrow{\text{Binary}} (10000010)_2$

__Step IV__

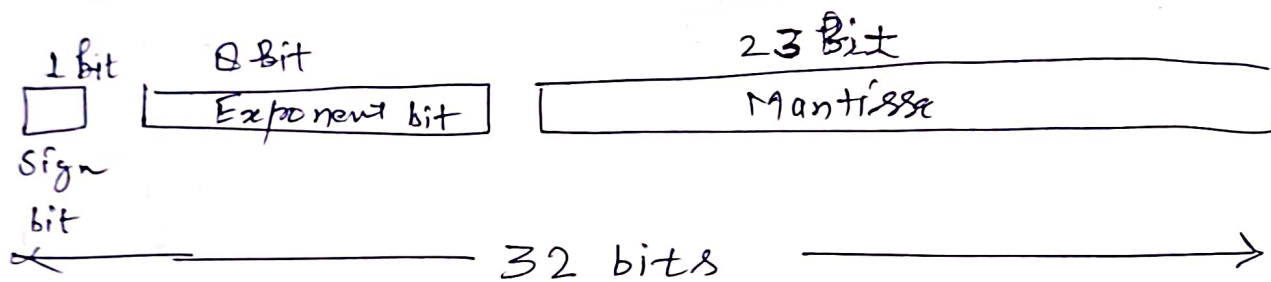| 1bit | 8bits | 23bits |
|---|---|---|
| 1 | 1 0 0 0 0 0 1 0 | 0 1 0 1 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 |
| Sign | Exponent | Mantissa |

$\longleftarrow$ 32 bits $\longrightarrow$
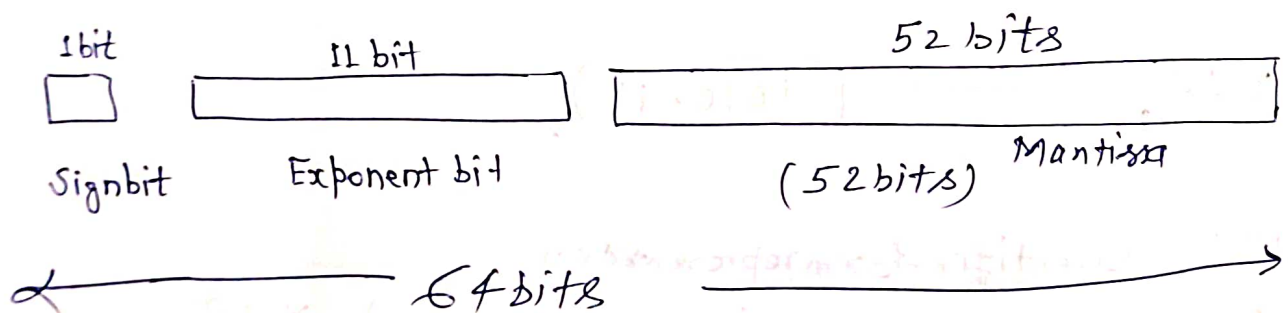
If given number is positive then Sign bit = 0

If given number is negative then sign bit = 1

* Combining sign bit, exponent bits and mantissa bits will be
the floating point representation [IEEE 754 format].
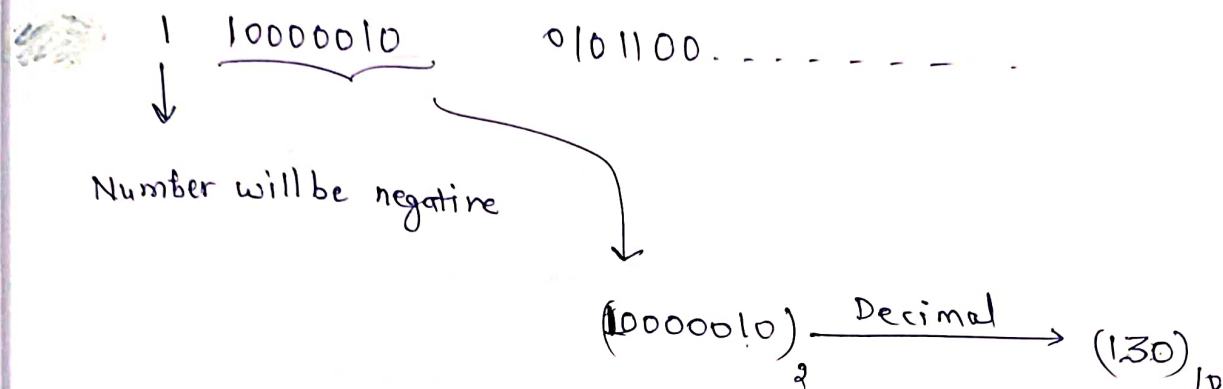
* Single precision [32bit IEEE 754]

| 1 bit | 8 Bit | | 23 Bit |
|---|---|---|---|
| | Exponent bit | | Mantissa |

Sign bit

$\xleftarrow{\hspace{4cm}}$ 32 bits $\xrightarrow{\hspace{4cm}}$

* Double precision [64bit IEEE 754]

| 1bit | 11 bit | | 52 bits |
|---|---|---|---|
| | | | |

Signbit       Exponent bit       (52bits) Mantissa

$\xleftarrow{\hspace{3cm}}$ 64bits $\xrightarrow{\hspace{3cm}}$

# Floating point ⟶ Decimal number
## (32 bit IEEE 754)

Q: 1 10000010 01011000000000000000000

Solution: Step(I)

1    10000010      0101100. . . . . — . . . .

↓

Number will be negative

$(10000010)_2 \xrightarrow{\text{Decimal}} (130)_{10}$

$$127 + x = 130$$

$$x = 3$$

Step (II)    ∴ Scientific representation ⟹

$$[1.xxxx \_\_\_\_ \quad] \times 2^x$$

$$[1.Mantissa \_\_\_\_\_ ] \times 2^x$$

$$[1.01011000000000 \_\_\_\_ \quad ] \times 2^3$$

$$= 1010.110000000000 \_\_\_\_$$

$$= (10.75)_{10} \quad \text{Decimal value}$$

Step III    final answer $= (-10.75)_{10}$

# Floating point addition :-  (32 bit IEEE 754)

$(100)_{10}$

$+(0 \cdot 25)_{10}$

?

$(100)_{10} \longrightarrow 1100100 \longrightarrow 1 \cdot 100100 \times 2^6$

$(0 \cdot 25)_{10} \longrightarrow 0 \cdot 0\overset{\frown}{1}00 \longrightarrow [1 \cdot 0000 \ldots ] \times 2^{-2}$

IEEE representation ⇒

| | | | |
|---|---|---|---|
| 0 | 1000 0101 | 1001 0000 0000 0000 0000 000 | ; $\underline{(100)}_{10}$ |
| 0 | 0111 1101 | 0000 0000 0000 0000 0000 000 | ; $(0 \cdot 25)_{10}$ |

Note

Exponent of $0.25$ is $(0111\,1101)_2$ and exponent of $100$ is $(1000\,0101)_2$

We have to make both exponent same. Because exponent of $0.25$ is lesser than exponent of $100$ so we will increase the exponent of $0.25$.

Hidden bit →

| | exponent | Hidden bit | mantissa |
|---|---|---|---|
| Initially | 0111 1101 | 1 | 0 000 0000 0000 0000 0000 000 |
| | 0111 1110 | 0 | 1000 0000 0000.0000 0000 000 |
| | 0111 1111 | 0 | 0100 0000 0000 0000 0000 000 |
| | 1000 0000 | 0 | 0010 0000 0000 0000 0000 000 |
| | 1000 0001 | 0 | 0001 0000 0000 0000 0000 000 |
| | 1000 0010 | 0 | 0000 1000 0000 0000 0000 000 |
| | 1000 0011 | 0 | 0000 0100 0000 0000 0000 000 |
| | 1000 0100 | 0 | 0000 0010 0000 0000 0000 000 |
| | 1000 0101 | 0 | 0000 0001 0000 0000 0000 000 |

**Hiddenbit**

| | | ↓ | | | | | |
|---|---|---|---|---|---|---|---|
| 0 | 1000 0101 | 1 | 1001 0000 | 0000 | 0000 | 0000 000 | |
| 0 | 1000 0101 | 0 | 0000 0001 | 0000 | 0000 | 0000 000 | |

---

| 0 | 1000 0101 | 1 | 1001 0001 | 0000 | 0000 | 0000 000 |
|---|---|---|---|---|---|---|

**( Add the mantissa)**

→ Remove this

\* Remove the hidden bit and final answer is given as :-

**Ans**    0    1000 0101     1001 0001 0000 0000 0000 000

**Note**

0    1000 0101     1001 0001 0000 0000 0000 000

$\underbrace{\qquad}$         $\underbrace{\qquad\qquad\qquad\qquad}$

↓              **Mantissa**

$(133)_{10}$

$127 + x = 133$

$\therefore x = 6$

$[1 \cdot \text{Mantissa}] \times 2^x$

$[1 \cdot 1001000 10000 \ldots] \times 2^6$

$= 1\,100100 \cdot 010000 \ldots$

$= (100 \cdot 25)_{10} \; (\text{Decimal}),$

**Ans**

# Floating point Subtraction (32bit IEEE754)

$$(100)_{10} - (0.25)_{10} = ?$$

solution

Hidden Bit

| 0 | 1000 0101 | 1 | 1001 0000 0000 0000 0000 000 | ← $(100)_{10}$ |
| 1 | 0111 1101 | 1 | 0000 0000 0000 0000 0000 000 | ← $-(0.25)_{10}$ |

| 0 | 1000 0101 | 1 | 1001 0000 0000 0000 0000 000 |
| 1 | 1000 0101 | 0 | 0000 0001 0000 0000 0000 000 |

| 0 | 1000 0101 | 1 | 1000 1111 0000 0000 0000 000 | Difference of both the Mantissa |

→ Remove this

\* Remove the hidden bit and final answer is given as:-

Ans   0  1000 0101      1000 1111 0000 0000 0000 000

$(133)_{10}$

Mantissa

$127 + x = 133$    ∴ $x = 6$

Ans = $[1. \text{Mantissa}] \times 2^{x}$

$[1. 1000 1111 0000 \_\_\_\_ ] \times 2^{6}$

$= 1100011. 11 0000 _____$

$= (99.75)_{10}$    [ Answer in decimal form ]

# Floating point multiplication (32 bit IEEE 754)

$$(100)_{10}$$
$$\times (0.25)_{10}$$
$$\overline{\phantom{xx}25\phantom{xx}}$$

```
  1100100
×     .01
---------
11001.00
```

IEEE 754 representation:-

$100 \Rightarrow 1.100100 \times 2^6$
$0.25 \Rightarrow 1.0000\text{---} \times 2^{-2}$ } Scientific form

$100 \Rightarrow$    0    | 1000 0101 |    1001 0000 0000 0000 0000 000

$0.25 \Rightarrow$    0    | 0111 1101 |    0000 0000 0000 0000 0000 000

---

10000 0010 (Add the exponent)

$$- \quad 0 0111\ 1111 \quad \Leftarrow (127)_{10}$$
$$\overline{\phantom{xxx}1 000\ 0001\phantom{xxx}}$$

| |
|---|
| 127 + 6 |
| 127 + (-2) |
| 127+127+4 |

\* Include hidden bit for both the mantissa and then multiply.

Hidden bit
↓
1 . 1001 0000   → Mantissa

1 . 0000 0000   → Mantissa

---

1 . 1001 0000

Multiply the above number

Remove this

\* Remove the hidden bit and your final answer will be as ⇒

0    1000 0001    1001 0000 - - - - - - -

1000 0001 ⟶ 131

$$127 + x = 131$$
$$x = 4$$

$$[ 1. \text{Mantissa} ] \times 2^x$$

$$[ 1.10010000 \ldots ] \times 2^4$$

$$= 11001.0000$$

$$= 11001$$

$$= (25)_{10} \qquad \text{(Decimal Value)}$$

**Note:**

* If multiplication output (result) is not in scientific form then we have to make multiplication result in scientific form. Accordingly we have to adjust the exponent.

# Flaating point division

(32bit IEEE 754)

$$\frac{(10.35)_{10}}{(2.25)_{10}} = \frac{1010.\ 0101100\ 11001100\ 1100.1100}{10.01}$$

2·25) 10·35 ( 4·6

$(10.35)_{10} \implies \left[1.010010\ 1100\ 1100\ 1100\ 1100\ 1100\right] \times 2^3$

$(2.25)_{10} \implies [1.001] \times 2^1$

23$^{rd}$ bit

**IEEE representation :-**

| | | | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| 10·35 ; | 0 | 10000010 | 0100 | 1011 | 0011 | 0011 | 0011 | 001 . |
| 2·25 ; | 0 | 10000000 | 0010 | 0000 | 0000 | 0000 | 0000 | 000 |

00000010   (Take the difference of exponent)

\* Include hidden bit for both mantissa and then divide.

$\underline{1.0100\ 1011\ \underbrace{0011\ 0011\ \ 0011\ 001}_{mantissa}}$

$1.0010\ 0000\ \ 0000\ 0000\ \underbrace{0000\ 000}_{Mantissa}$

$= \frac{1.0100\ 1011\ 0011\ 0011\ 0011\ 011}{1.001}$

1.00100110011001100011

1001) 1010 01011001100110011001001
      1001
      0001010
       1001
       0001110
        1001
        01010
        1001
        0001

Ans $\Rightarrow$ (IEEE format)

0    0000 0010    001001 ‒ ‒ ‒ ‒ ‒ ‒

Mantissa

Decimal value = 2

**Ans**

$$[1. \text{Mantissa}] \times 2^x$$

$$= [1.001001100110011] \times 2^2$$

$$= 100.10011 0011 0011 0011$$

$$= 4 + \frac{1}{2^1} + \frac{1}{2^4} + \frac{1}{2^5} + \cdots$$

$$= 4.59375 \underline{\ \ \ \ \ }$$

$$\frac{\begin{array}{r} 127+3 \\ 127+1 \end{array}}{2}$$

$$\boxed{\therefore x = 2}$$

**Answer**

Decimal

Decimal value

**Note**   If quotient is not in scientific format then we have to make quotient in scientific format. Accordingly we have to adjust the exponent.

$$1.\ 0010011\ 0011\ 0011\ \ 0011$$

```
1001 ⟞ 1010.010 1100  1100  1100  11001
          1001
         ─────────
         0001 010
          1001
       ──────────
         0 001 110
            1001
           ─────────
         x 101 0
            1001
          ─────────
         x x x 1 110
              1001
            ─────────
             1010
             1001
            ─────────
         x x x 1 110
              1001
            ─────────
             x 1010
              1001
            ─────────
             0001 110
              1001
            ─────────
              0 1010
               1001
             ─────────
               00011
```