# CS & IT ENGINEERING

## COMPUTER ORGANIZATION AND ARCHITECTURE

Floating Point Representation

**Lecture No.-02**

By- Vishvadeep Gothi sir

# Recap of Previous Lecture

**Topic** Floating-Point Numbers

**Topic** Biased Exponent

**Topic** Normalization: Implicit & Explicit

# Topics to be Covered

**Topic** Floating-Point Numbers

**Topic** Biased Exponent

**Topic** Number Range

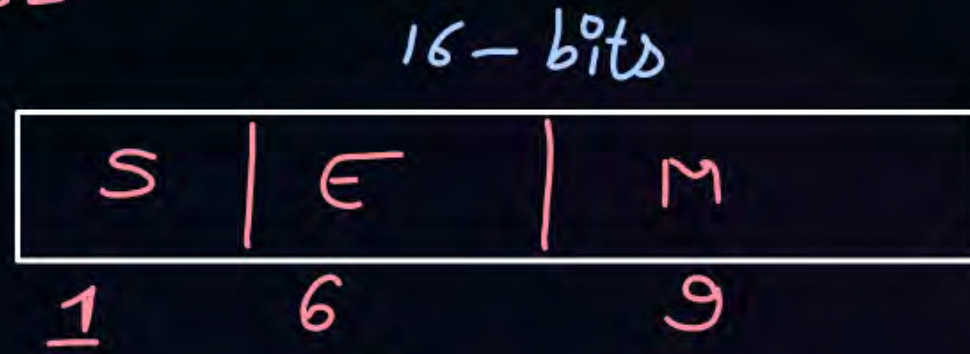**Topic** IEEE-754 Floating Point Representation

**Topic** Denormalized Number

- The number is represented in format:

| S | E | M |
|---|---|---|

- Mantissa is signed normalized (implicit/explicit) fraction number

- Exponent is stored in biased form.

*explicitly*

**#Q.** Consider a 16-bit register used to store floating point numbers. The mantissa is normalized signed fraction number. Exponent is represented in excess-32 form. What is the 16-bit value for $+(11.5)_{10}$ in this register?

$bias = 32$

16 - bits

| S | E | M |
|---|---|---|
| 1 | 6 | 9 |

$2^{k-1} = 32$

$2^{k-1} = 2^5$

$k - 1 = 5$

$k = 6$

$+(11.5)_{10} \Rightarrow$ positive $\Rightarrow S = 0$

$(11.5)_{10} = (1011.1)_2$

$\downarrow$

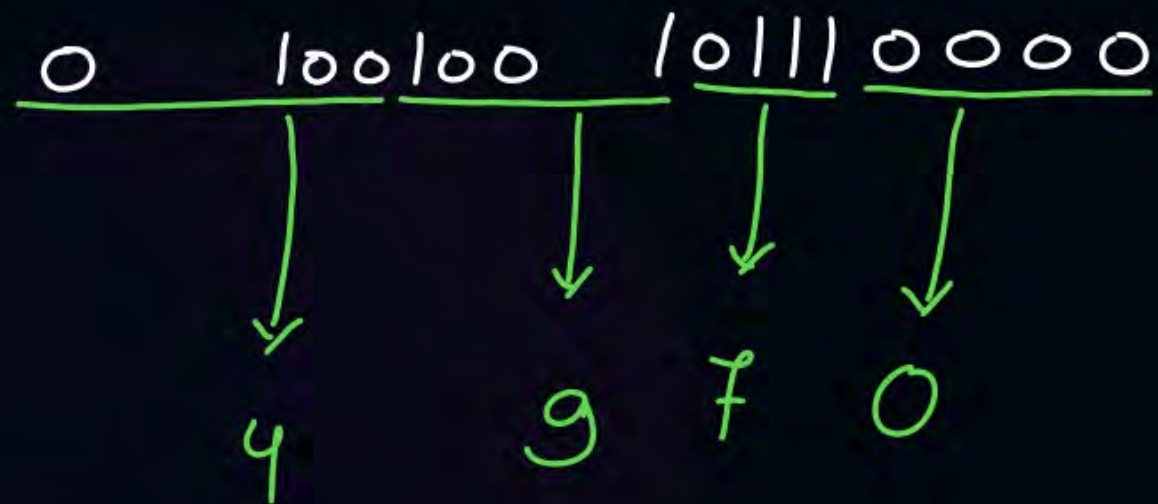explicit Normalizat$^n$

$\Downarrow$

$0.10111 * 2^4$

$M = 10111 0000$

$e = 4$

$E = 4 + 32 = 36 = \left(100100\right)_2$

| S | E | M |
|---|---|---|

0     100100     10111 0000

#Q. What is the 4-digit hexadecimal value for $+(11.5)_{10}$ in above question's register?

$$0 \quad \underline{100100 \quad 10111\,0000}$$

$$4 \qquad 9 \quad 7 \quad 0$$

$$(4970)_{16}$$

$$0x\,4970$$

$$4970\,H$$

$Ans = 4D2E$

#Q. What is the 4-digit hexadecimal value for $+(37.75)_{10}$ in above question's register?
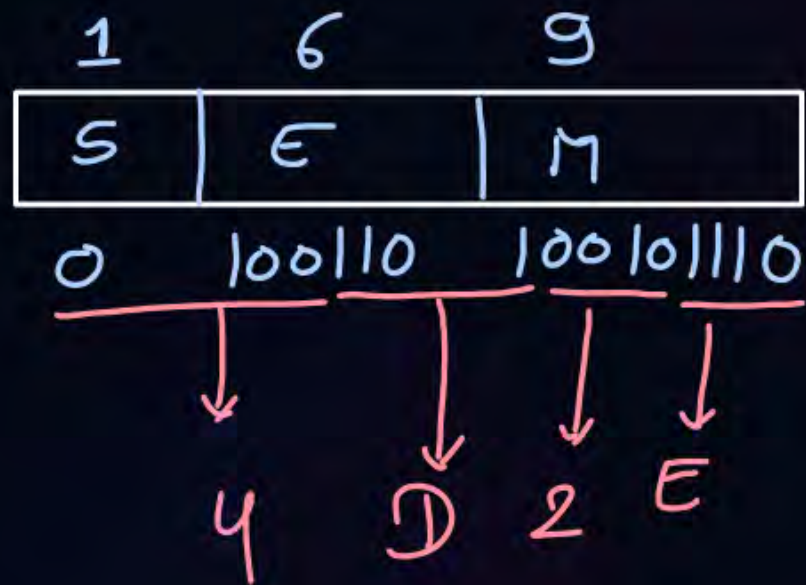
$$(37.75)_{10} = (100101.11)_2$$

$\Downarrow$

Expliat normalizat$^n$

$\Downarrow$

$0.100101111 * 2^6$

$M = 100101110$

$e = 6$

$E = 6 + 32 = 38 = (100110)_2$

| 1 | 6 | 9 |
|---|---|---|
| S | E | M |

| O | 100110 | 100101110 |
|---|---|---|

$\downarrow$ $\downarrow$ $\downarrow$ $\downarrow$

4   D   2   E

can not store (underflow)

overflow

overflow

overflow



min                                    mmax

example :-

$$E_{min} = 0$$

6-bits

$$e_{min} = 0 - 32$$

$$= -32$$

if any value $\Rightarrow$ 0.0000...... 11

$\Downarrow$

explicit normalizat$^n$

$\Downarrow$

$0.11 * 2^{-33}$ $\Leftarrow$ Cannot

$e = -33$ $\Leftarrow$ store

More no. of bits in $E$ $\Rightarrow$ larger range of numbers

_____ || _____ $M$ $\Rightarrow$ Better precision or Accuracy

$\longrightarrow$ Can not represent zero.

$\longrightarrow$ it can not store very small numbers around zero.

(has underflow)

IEEE-754
Representation

Single Precision

32 – bits

Double
Precision

64 – bits

$\xleftarrow{\quad 32 \quad}$

| S | E | M |
|---|---|---|

1    8     23

bias = 127

$\xleftarrow{\quad 64 \quad}$

| S | E | M |
|---|---|---|

1    11     52

bias = 1023

$$E = 00 \cdots \cdots - 0$$

or

$$E = 11 \cdots \cdots - - 1$$

$\left.\begin{array}{c} \\ \\ \\ \end{array}\right\}$ special number (exceptions)

---

$$E \neq 00 \cdots \cdots 0$$

and

$$E \neq 11 \cdots \cdots 1$$

$\left.\begin{array}{c} \\ \\ \\ \end{array}\right\}$ normal number (Implicitly normalized)

| S | E | M | Number |
|---|---|---|---|
| 0 | 00......0 | 0.......0 | $+0$ |
| 1 | 00......0 | 0.......0 | $-0$ |
| 0 | 11......1 | 0.......0 | $+\infty$ |
| 1 | 11......1 | 0.......0 | $-\infty$ |
| 0 or 1 | 11......1 | $M \neq 0$......0 | N.A.N. (Not A Number) |
| 0 or 1 | 00......0 | $M \neq 0$......0 | Denormalized number |
| 0 or 1 | $E \neq 0$......0 and $E \neq 11$......1 | $xxxxxx$....$x$ | Implicitly normalize |

A very-very small number which can not be implicitly normalized.

single precision

| S | E | M |
|---|---|---|
| 1 | 8 | 23 |

$bias = 127$

$E_{min}$ for normalized number = 1

$e_{min} = 1 - 127$
$\quad\quad = -126$

ex:-

$Value = 0.00000\ldots\ldots 11$

$\Downarrow$

Implicit normalizat$^n$

$\Downarrow$

$\left.\begin{array}{l} 1.1 * 2^{-128} \end{array}\right\}$ not allowed beyond $2^{-126}$

$\Downarrow$

if can not be normalized till $2^{-126}$ then store
the number as denormalized number

$$0.000 \cdots \cdots 11$$

$\downarrow$

try to normalized till $2^{-126}$

$\downarrow$

$$0.\underbrace{011}_{M} * 2^{-126}$$

$$\text{value} = 0.M * 2^{-126}$$

| S | 0 $\cdots\cdots$ 0 | 01100$\cdots$0 |
|---|---|---|
| | E | M |

$$\text{Value}_{(\text{Implicit})} = (-1)^{S} * 1.M * 2^{E-\text{bias}}$$

$$\text{Value}_{(\text{denormalized})} = (-1)^{S} * 0.M * 2^{-126 \text{ or } -1022}$$

$$\boxed{\quad E \quad}$$

$$8$$

max value $e = +127$

max value $E = 127 + bias$

$\qquad = 127 + 127$

$\qquad = (254)_{10}$

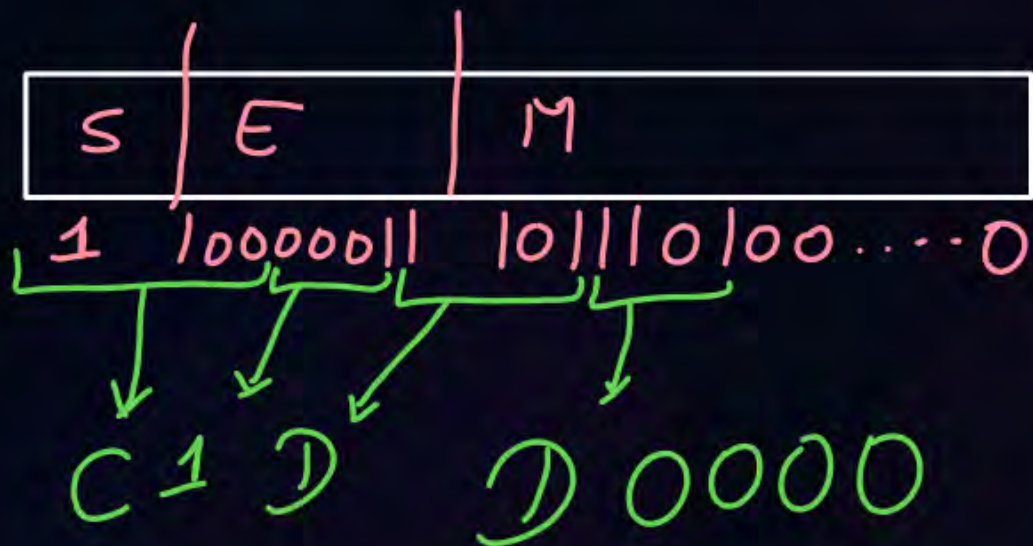$\qquad = (11111110)_2$

that's why $E = 111\ldots 1$

is preserved.

$Ans = (C1DD0000)_{16}$

#Q. The value of a float type variable is represented using the single- precision 32-bit floating point format IEEE-754 standard that uses 1bit for sign, 8 bits for biased exponent and 23 bits for mantissa. A float type variable X is assigned the decimal value of $-27.625$. The representation of X in hexadecimal notation is?

$(27.625)_{10} = (11011.101)_2 \Rightarrow$ Implicit normalizat$^n$

$$1.1011101 * 2^4$$

$$S = 1$$

| S | E | M |
|---|---|---|
| 1 | 10000011 | 10110100.....0 |

C1 D D 0000

$M = 10110100....0$
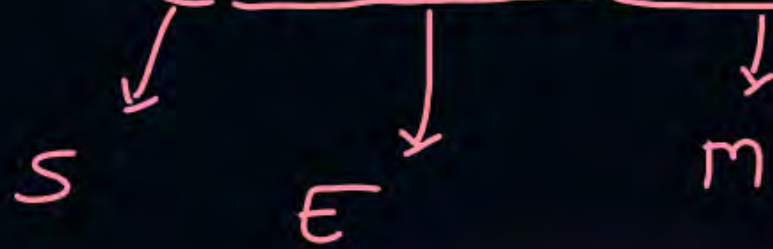
$e = 4$

$E = 4 + 127 = (131)_{10} = (10000011)_2$

$$Ans = +(28)_{10}$$

(P)(W)

#Q.    The value represented by the following 32-bits in IEEE-754 representation is?

01000001111100000...00

S    E    m

$S = 0 \Rightarrow +ve$

$E = 10000011 = (131)_{10}$

$M = 1100\cdots0$

$E \neq 0\cdots-0$
and
$E \neq 11.\frown.1$
$\left. \right\}$ Implicit normalized

$value = 1.1100\cdots0 * 2^{131-127}$

$= 1.110\cdots0 * 2^{4}$

$= (11100.0)_2$

$= +(28)_{10}$

#Q. The value represented by the following 32-bits in IEEE-754 representation is?

$$\underline{00000000\underline{01100000...00}}$$

$S = 0$

$E = 00000000$

$M = 1100 \cdots 0$

$\left.\begin{array}{l} E = 0 \cdots 0 \\ \text{and} \\ M \neq 0 \cdots 0 \end{array}\right\}$ denormalized

$$\text{value} = 0.11 * 2^{-126}$$
$$= 11.0 * 2^{-2} * 2^{-126}$$
$$= + 3 * 2^{-128}$$

#Q.    Maximum value represented in IEEE-754 single precision?
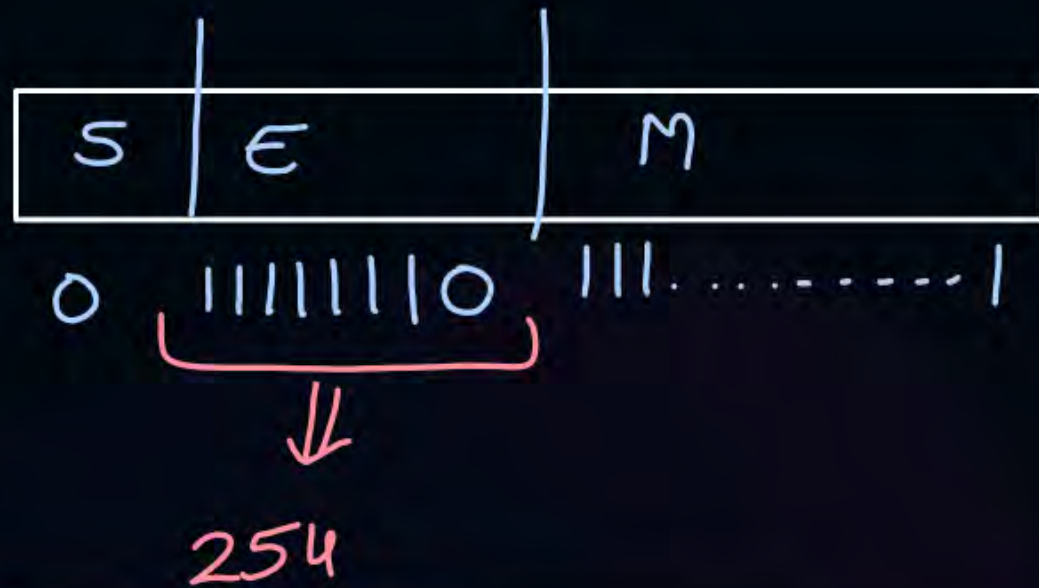
$$+\infty$$

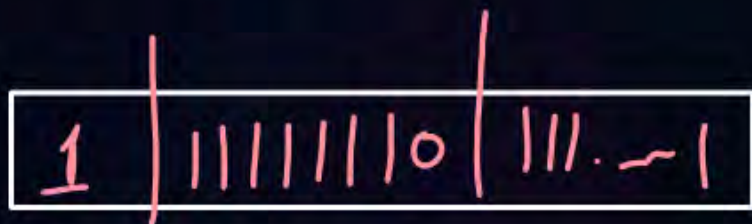#Q.    Minimum value represented in IEEE-754 single precision?

$$-\infty$$

**#Q.** ~~Minimum~~ Max positive normalized value represented in IEEE-754 single precision?

| S | $\epsilon$ | M |
|---|---|---|

0     11111110    111.........1

$$\underbrace{11111110}$$

254

$$Value = +\ 1.11\cdots1 * 2^{254-127}$$

$$= 111\cdots1.0 * 2^{-23} * 2^{127}$$

$$= +(2^{24}-1) * 2^{104}$$
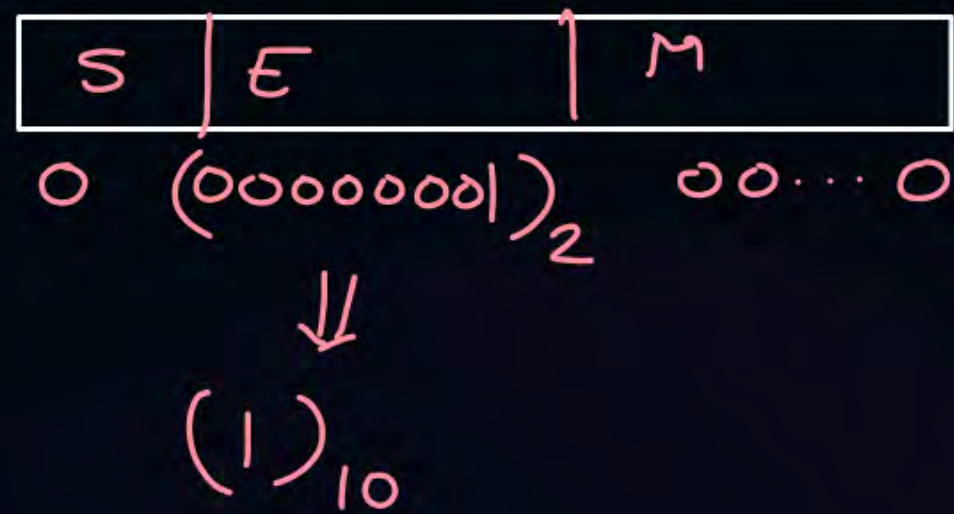
---

min possible normalized value.

| 1 | 11111110 | 111.~1 |
|---|---|---|

$$Value = -(2^{24}-1) * 2^{104}$$

#Q. Minimum positive normalized value represented in IEEE-754 single precision?

| S | E | M |
|---|---|---|

0  $(0000001)_2$  00...0

$\Downarrow$

$(1)_{10}$

$$\text{Value} = +\ 1.00\cdots0 * 2^{1-127}$$

$$= +\ 2^{-126}$$

#Q.    ~~Minimum~~ **Maximum** positive denormalized value represented in IEEE-754 single precision?

| S | E | M |
|---|---|---|
| 0 | 00000000 | 11......1 |

$$\text{Value} = 0.11.....1 * 2^{-126}$$

$$= 111....1.0 * 2^{-23} * 2^{-126}$$

$$= +(2^{23} - 1) * 2^{-149}$$

#Q. ~~Maximum~~ Minimum positive denormalized value represented in IEEE-754 single precision?

| S | E | M |
|---|---|---|
| 0 | 00000000 | 0000....0 <u>1</u> |

$$Value = + \; 0.000\cdots01 * 2^{-126}$$

$$= +1.0 * 2^{-23} * 2^{-126}$$

$$= 2^{-149}$$

H. ω.

#Q.   How to represent +1 and -1 in IEEE-754 single precision floating point number?

H. W.

#Q. How to represent +0.0000101 in IEEE-754 single precision floating point number?

H.W.

#Q. The value of a float type variable is represented using the single- precision 32-bit floating point format IEEE-754 standard that uses 1bit for sign, 8 bits for biased exponent and 23 bits for mantissa. A float type variable X is assigned the decimal value of −14.25. The representation of X in hexadecimal notation is

**A** C1640000H

**B** 416C0000H

**C** 41640000H

**D** C16C0000H

#Q.   Consider the following representation of a number in IEEE 754 single-precision floating point format with a bias of 127.

S : 1   E : 10000001 = $(129)_{10}$   F: 11110000000000000000000

Here S, E and F denote the sign, exponent and fraction components of the floating-point representation.

The decimal value corresponding to the above representation (rounded to 2 decimal places) is _____

$E \neq 0 \cdots \sim 0$
and
$E \neq 11 \cdots \sim 1$
} Implicit normalized

$$Value = -1.1111 * 2^{129-127}$$
$$= -1.1111 * 2^{2}$$
$$= -(111.11)_{2}$$
$$= -(7.75)_{10}$$

#Q.    The format of the single-precision floating-point representation of a real number as per the IEEE 754 standard is as follows:

| Sign | Exponent | mantissa |
|------|----------|----------|

Which one of the following choices is correct with respect to the _smallest normalized positive number_ represented using the standard?

A.    exponent = 00000001 and mantissa = 00000000000000000000001

B. ✓  exponent = 00000001 and mantissa = 00000000000000000000000

C.    exponent = 00000000 and mantissa = 00000000000000000000000

D.    exponent = 00000000 and mantissa = 00000000000000000000001

**Topic** Biased Exponent

**Topic** Normalized Mantissa

**Topic** Explicit vs Implicit Normalization

**Topic** IEEE-754 Floating Point Representation