# CS & IT ENGINEERING

## COMPUTER ORGANIZATION AND ARCHITECTURE

Floating Point Representation

Lecture No.-01

By- Vishvadeep Gothi sir

# Recap of Previous Lecture

**Topic** — Control Unit ✓

**Topic** — Hardwired Control Unit

**Topic** — Micro-Programmed Control Unit

→ Horizontal
→ Vertical

**Topic** — RISC vs CISC

# Topics to be Covered

**Topic** Floating-Point Numbers

**Topic** Biased Exponent

**Topic** Number Range

**Topic** IEEE-754 Floating Point Representation

**Topic** Denormalized Number

## Need or requirement of Floating point representation

### Fixed point representation :-

Given 8-bits storage

$\rightarrow$ Range of number (unsigned) $\Rightarrow$ 0 to $2^8 - 1$ | 0 to 255

$\rightarrow$ ——— " ——— (signed) $\Rightarrow$ -128 to +127

2's comp.

Floating point representation provides a larger range of numbers as compared to fixed point representation.

- The number is represented in format:

| Sign | Exponent | Mantissa |
|------|----------|----------|
| S | E | M |

1 bit

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

- Exponent is stored in biased form.

$$S \begin{cases} 0 \Rightarrow +ve \\ 1 \Rightarrow -ve \end{cases}$$

$\rightarrow$ arithmetic operations on floating point numbers become easy.

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

ex:- $(1.11 * 2^2) + (1.01 * 2^{-3})$ vs $(1.11 * 2^4) + (1.01 * 2^9)$

original exponent + bias $\Rightarrow$ stored Exponent
(e)                              (E)

Assume, $E$ is represented $\Rightarrow$ 4-bits

originally number range with 4-bits $\Rightarrow$ −8 to +7 $\nearrow$ Transform into unsigned number range.
$\Downarrow$
0 to 15

| original exponent (e) | stored exponent (E) |
|---|---|
| $-8$ | 0 |
| $-7$ | 1 |
| $-6$ | 2 |
| $-5$ | 3 |
| $-4$ | 4 |
| $-3$ | 5 |
| $-2$ | 6 |
| $-1$ | 7 |
| 0 | 8 |
| 1 | 9 |
| 2 | 10 |
| 3 | 11 |
| $\vdots$ | $\vdots$ |
| 7 | 15 |

excess-8 Code

$$E = e + 8$$

Here in this example 8 is bias. ←

if $k$-bits are used to represent $E$, then

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

$101.11$

→ Explicit normalization $\Rightarrow 0.\underline{1}0111 * 2^3$

↑ should be one

$e = 3$
$E = 3 + bias$

$M = $ Number after point
$\quad = 10111$

→ Implicit Normalization $\Rightarrow 1.\underline{0}111 * 2^2$

↑ should be one

$e = 2$
$E = 2 + bias$

$M = 0111$

ex:-

Explicit normalization $\Rightarrow 0.10 * 2^1$

$e = 1$
$E = 1 + bias$
$M = 10$

1.0

Implicit Normalization $\Rightarrow 1.0 * 2^0$

$e = 0$
$E = 0 + bias$
$M = 0$

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

$$\text{Value}_{(\text{explicit})} = (-1)^S * 0.M * 2^{E-\text{bias}}$$

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

$$\begin{array}{|c|c|c|}
\hline
 & E & M \\
\hline
S & e_6 e_5 e_4 e_3 e_2 e_1 e_0 & \\
\hline
\end{array}$$
7 bits

#Q. A certain well-known computer family represents the exponents of its floating-point numbers as 'excess-64' integers; i.e., a typical exponent $e_6 e_5 e_4 e_3 e_2 e_1 e_0$ represents the number: → bias = 64

A $\quad e = -64 + \sum_{i=0}^{6} 2^i e_i$

B $\quad e = -64 + \sum_{i=0}^{6} 2 e_i$

C $\quad e = 64 - \sum_{i=0}^{6} 2^i e_i$

D $\quad e = 64 - \sum_{i=0}^{6} 2 e_i$

$$e_6 \; e_5 \; e_4 \; e_3 \; e_2 \; e_1 \; e_0$$

$$\Downarrow$$
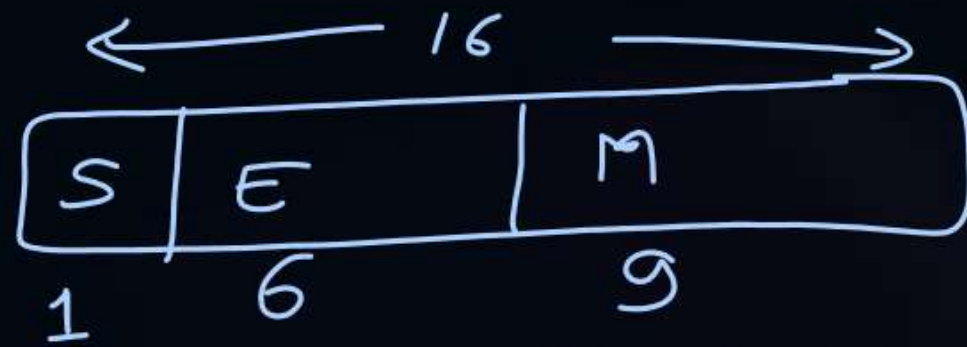
decimal

$$\Downarrow$$

$$(e_6 * 2^6) + (e_5 * 2^5) + \ldots + (e_0 * 2^0)$$

$$= \sum_{i=0}^{6} e_i * 2^i$$

_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?

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

1    6      9

0    100100    101110000

$bias = 32$

$2^{k-1} \Rightarrow k = 6$

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

$$(11.5)_{10} = (1011.1)_2$$
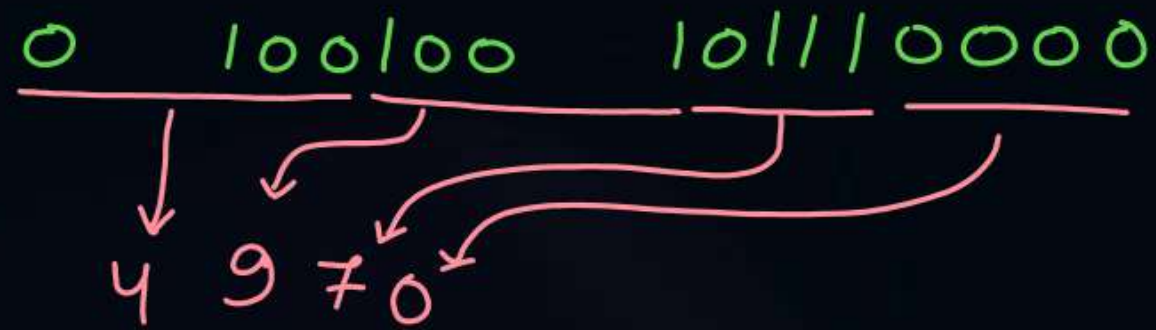
$(1011.1)_2 \Rightarrow$ explict normalize$^n$

$\Downarrow$

$0.10111 * 2^4$

$m = 10111$

$e = 4$

$E = 4 + 32 = (36)_{10} = (100100)_2$

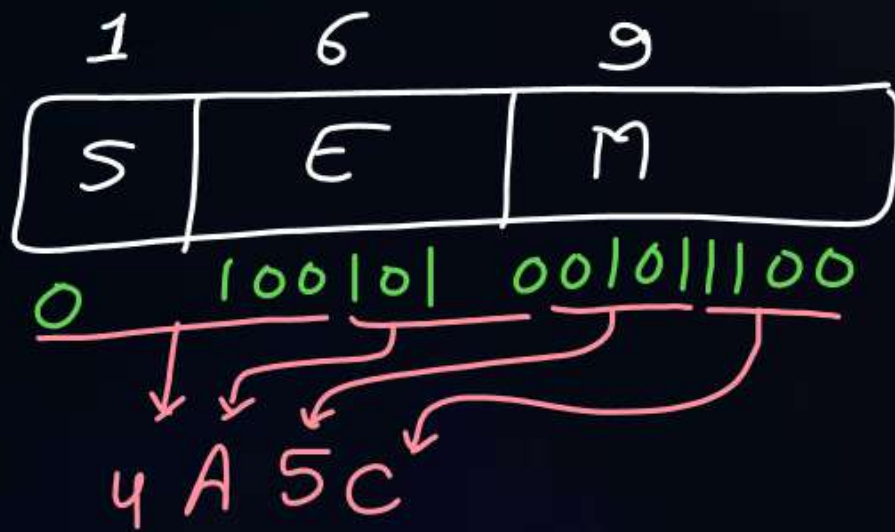#Q. What is the 4-digit hexadecimal value for $+(11.5)_{10}$ in above question's register?

$$0 \quad 100100 \quad 10111 \, 0000$$

$$4 \; 9 \; 7 \; 0$$

$$Ans \Rightarrow (4970)_{16} = 0 \times 4970 = 4970 \, H$$

$$Ans = (4A5C)_{16}$$

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

$$(37.75)_{10} = (100101.11)_2$$

$$\Downarrow$$

implicit normalization

$$\Downarrow$$

$$1.0010111 * 2^5$$

$$e = 5, \ E = 5 + 32 = 37 = (100101)_2$$

$$M = 00101111 00$$

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

0    100101    00101111 00

4 A 5 C

#Q.

what is 4-digit hexa-decimal representation of $+(0.000101)_2$ in prev. question's register with explicit normalization?
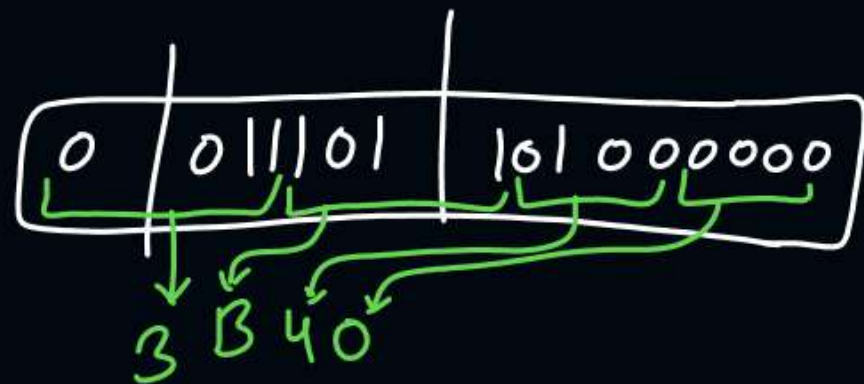
__Sol<sup>n</sup>__

$0.000101 \Rightarrow$ explicit normalization

$\Downarrow$

$0.101 * 2^{-3}$

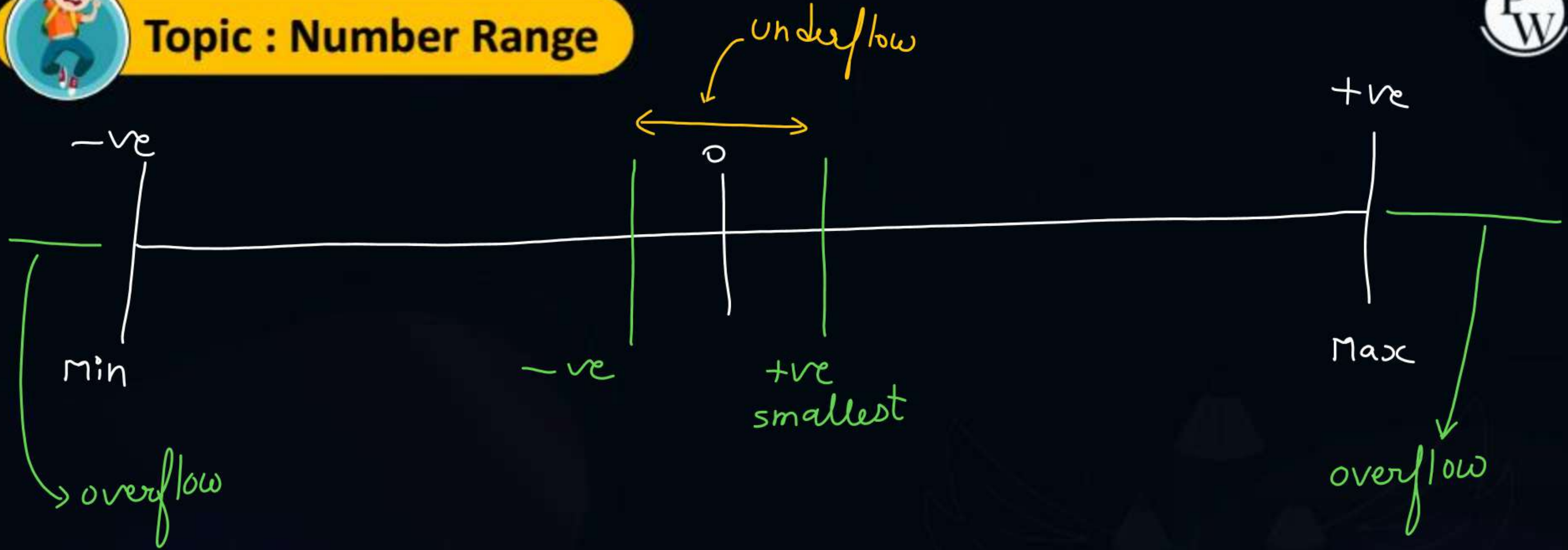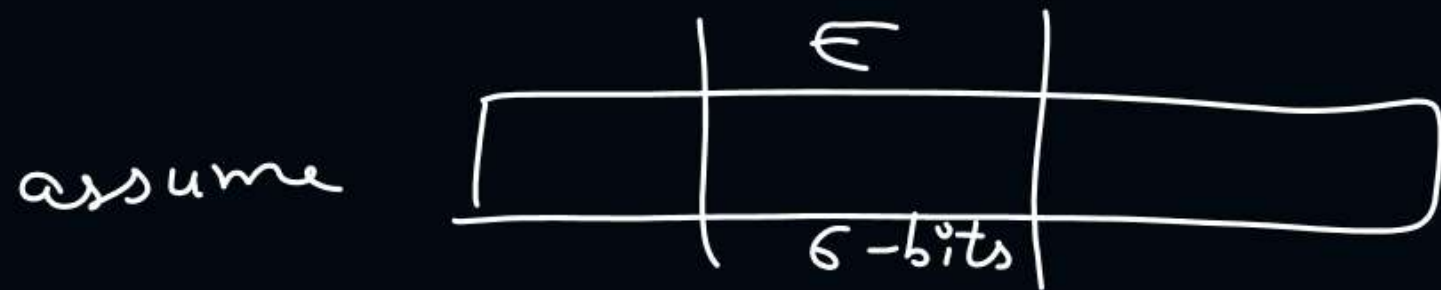$e = -3, \quad \epsilon = -3 + 32 = (29)_{10} = (011101)_2$

$M = 101$

| 0 | 011101 | 101 000000 |

$3\ B\ 40$

Ans = $(3B40)_{16}$

underflow

−ve

+ve

0

Min

−ve

+ve
smallest

Max

overflow

overflow

assume



$E$

6-bits

bias = 32

number after normalizat$^n$ $\Rightarrow$ $1.001 * 2^{32}$

$e = 32$

$E = 32 + 32 = \boxed{64}$

Can not be stored in 6-bit $E$.

$\Downarrow$

overflow

|  | min | max |
| --- | --- | --- |
| here $E$ $\Rightarrow$ | 0 | 63 |
| $e$ $\Rightarrow$ | $-32$ | 31 |

number after normalization $= 1.0101 * 2^{-33}$

$e = -33$

$E = -33 + 32 = -1$

underflow

More bits are used for $E \Rightarrow$ Range of numbers will be larger

— ।। ——————————— $M \Rightarrow$ More precision or accuracy

1. It can not store zero.

2. It suffers from underflow.

$\Rightarrow sol^n \Rightarrow$ IEEE-754

floating point

representation

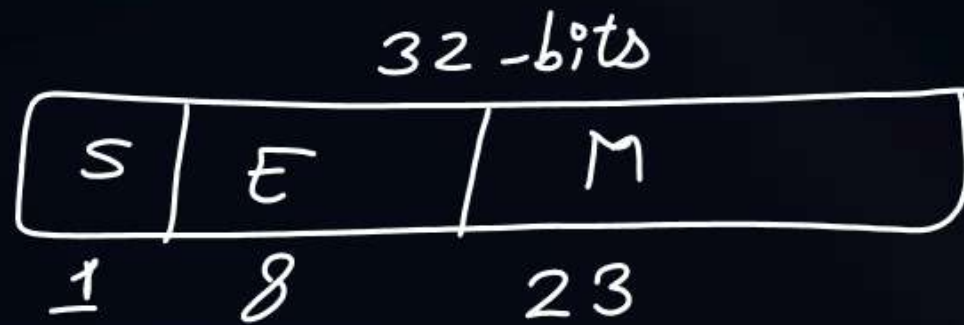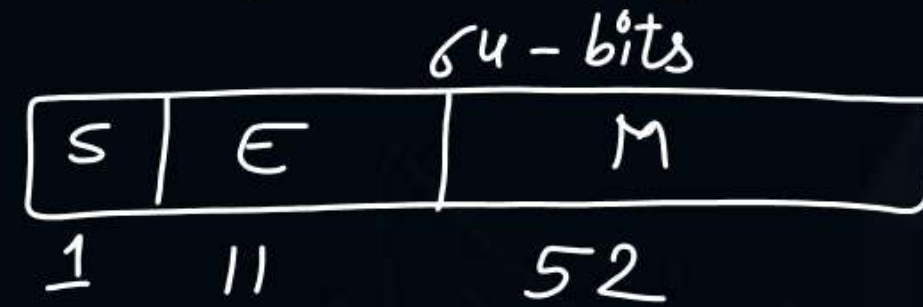IEEE-754
Representation

Single Precision

Double
Precision

32 -bits

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

64 - bits

| S | E | M |
|---|---|---|
| 1 | 11 | 52 |

bias = 127

bias = 1023

if $\quad E = 00 \cdots \cdots 0$

$$\text{or}$$

$$E = 111 \cdots \cdots 1$$

$\}$ special case

---

if $\quad E \neq 00 \cdots \cdots 0$

$$\text{and}$$

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

$\}$ normal number
(always implicitly normalized)

#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?

**2 mins Summary**

**Topic** — Biased Exponent

**Topic** — Normalized Mantissa

**Topic** — Explicit vs Implicit Normalization

**Topic** — IEEE-754 Floating Point Representation