

COMPUTER SCIENCE

Computer Organization and Architecture

Floating Point Representation

Lecture_03



Vijay Agarwal sir



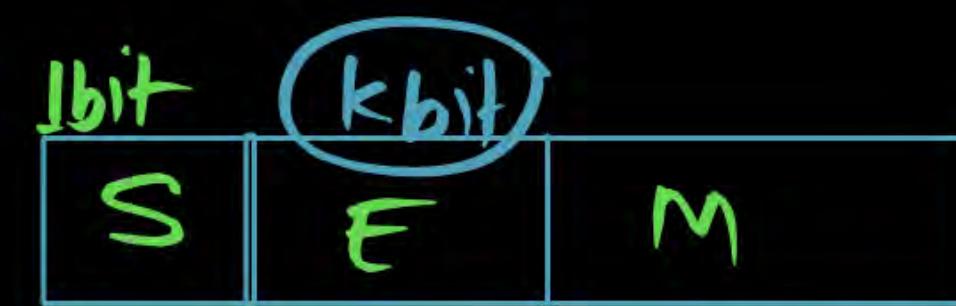
A graphic of a construction barrier with orange and white diagonal stripes and two yellow bollards at the top.

**TOPICS
TO BE
COVERED**

o1 Floating Point Representation

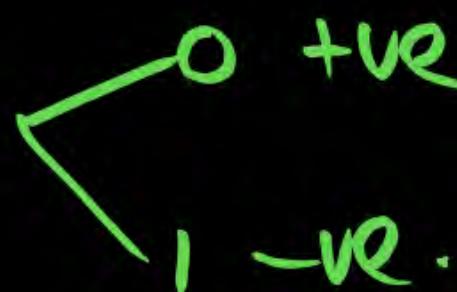
**o2 IEEE 754 Floating Point
Representation**

Floating Point Representation.



M : Mantissa .

S : Sign bit



BE | E : Exponent
Biased Exponent

$$E = e + \text{bias}$$

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

Exponent bit

$$\text{bias} = 2^{\lfloor k-1 \rfloor}$$

$E = e + \text{bias}$

Implicit
Normalized

$\pm \cdot \text{Something} \times 2^e$

$(-1)^s \pm \cdot M \times 2^e$

$(-1)^s \pm \cdot M \times 2^{E-\text{bias}}$

Mantissa

S	E	M
		^{kbit}

Explicit
Normalized

$0 \cdot 1xxx \times 2^e$

$(-1)^s 0 \cdot m \times 2^e$

$(-1)^s 0 \cdot m \times 2^{E-\text{bias}}$

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

Normalized Mantissa

1 bit x bit y bit

Explicit Normalized Syntax

$$\frac{0.1\ldots \times 2^e}{M}$$

Formula to get number
[value formula]

$$(-1)^s \times 0.M \times 2^e$$

$$(-1)^s \times 0.M \times 2^{E-\text{bias}}$$

Implicit Normalized Syntax

$$\frac{1.\ldots \times 2^e}{M}$$

Formula to get number
[value formula]

$$(-1)^s \times 1.M \times 2^e$$

$$(-1)^s \times 1.M \times 2^{E-\text{bias}}$$

Explicit

0.1 After the point,

Immediate first bit should be 1

Example

(101.11)

0.10111×2^3

$M = 10111,$

$e = 3$

$E = e + \text{bias}$

Implicit

Before the point 1 means 1.

Example

(101.11)

1.0111×2^2

$M = 0111,$

$e = 2$

$E = e + \text{bias}$

Floating-Point Representation



S: sign bit
0 +ve
1 -ve

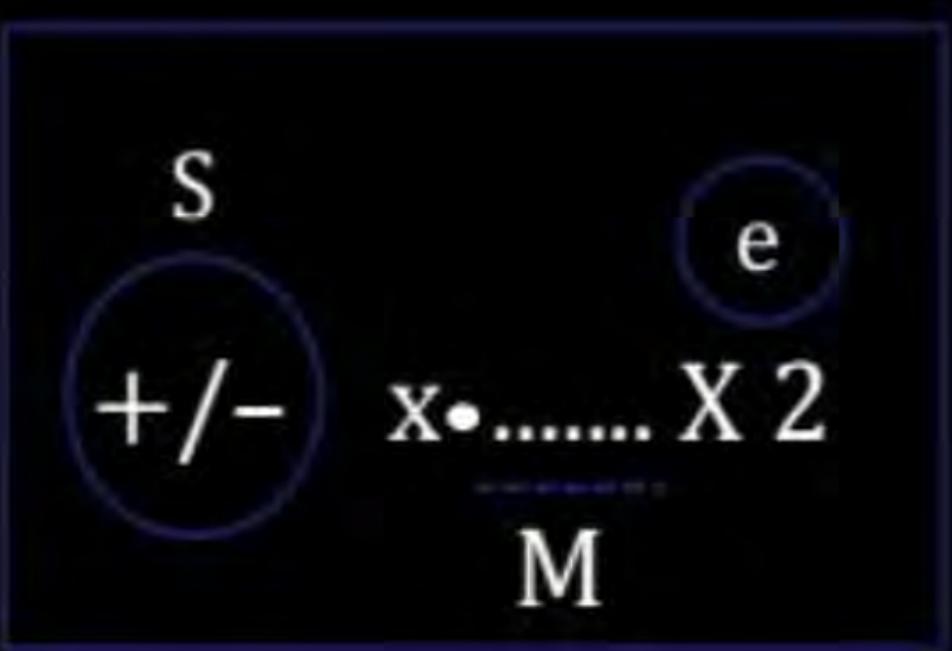
E: Biased exponent

M: Mantissa

$$E = e + \text{bias}$$

or

$$BE = AE + \text{bias}$$



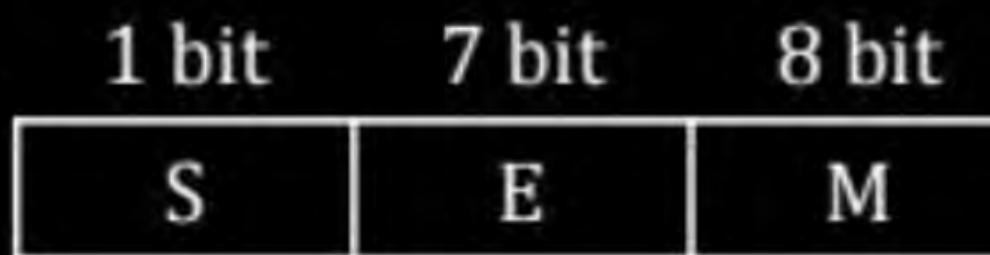
$$\dots \cdot X 2^e$$

Q.

+21.75

Implicit ?

10101.11

 1.010111×2^4 $M = 010111$ $e = 4, \text{ bias} = 2^{7-1}$ $E = 4 + 64$ $E = 68 = (1000100)_2$ 

Value Formula:

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

$$(-1)^0 \times 1.010111 \times 2^{68-64}$$

$$1.010111 \times 2^4$$

$$10101.11 = (21.75)_{10}$$

Ans

S(1bit)	E(7bit)	M(8 bit)
0	1000100	01011100

Hexadecimal = $(445C)_{16}$

Q.

Consider a 16 bit register used to store floating point number.
Mantissa is Explicit normalized signed fraction number.
Exponent is in Excess-32 form then what is 16-bit for
 $-(29.75)_{10}$ in the register?

P
W

Solution

-29.75

-11101.11

0:1110111 $\times 2^5$

M: 1110111

e = 5

bias = 2^{6-1}

bias = 32

$$E = 5 + 32 = 37 = (100101)_2$$

1 bit	6 bit	9 bit
S	E	M

Implicit.

$$1.110111 \times 2^4$$

S(1 bit)	E(6 bit)	M(9 bit)
1	100101	111011100

Q.

+21.75

Implicit ?

10101.11

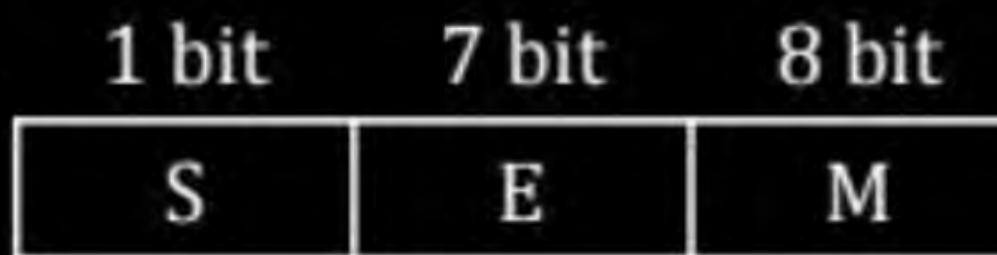
$$1.010111 \times 2^4$$

$$M = 010111$$

$$e = 4, \text{ bias} = 2^{7-1}$$

$$E = 4 + 64$$

$$E = 68 = (1000100)_2$$



Value Formula:

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

$$(-1)^0 \times 1.010111 \times 2^{68-64}$$

$$1.010111 \times 2^4$$

$$10101.11 = (21.75)_{10}$$

Ans

S(1bit)	E(7bit)	M(8 bit)
0	1000100	01011100

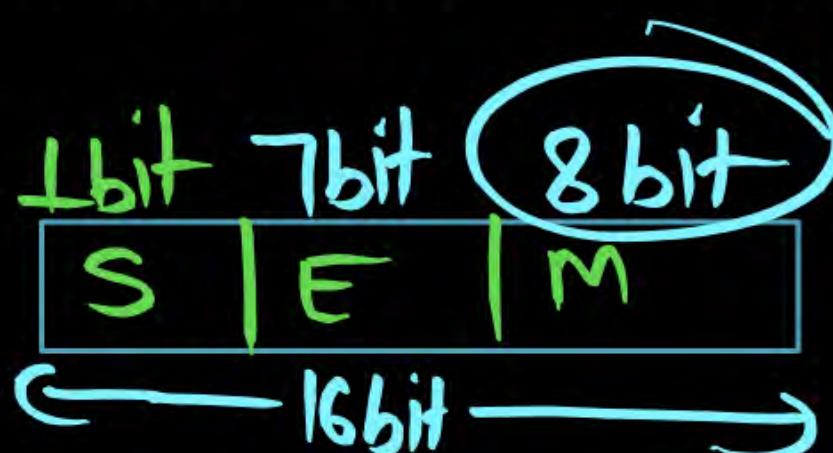
Hexadecimal = (445C)₁₆

Q.

Consider a 16 bit register used to store floating point number.
 Mantissa is **Implicit** normalized signed fraction number.

Exponent is in **Excess-64** form then

- (i) what is the First Smallest Positive number?
- (ii) what is the Second Smallest Positive number?
- (iii) what is the Difference between First Smallest & Second Smallest Positive number?



$$\text{Excess - 64} \Rightarrow \text{bias} = 64.$$

$$2^{k-1} = 64 \Rightarrow 2^{k-1} = 2^6$$

$$\begin{array}{l} k-1=6 \\ \boxed{k=7\text{bit}} \end{array}$$

We Can not Represent '0':

Implicit $\Rightarrow 1 \cdot \text{Something} \neq 0$ (Not zero)

OR
Explicit $\Rightarrow 0 \cdot 1 \text{nnn} \neq 0$ (Not zero)

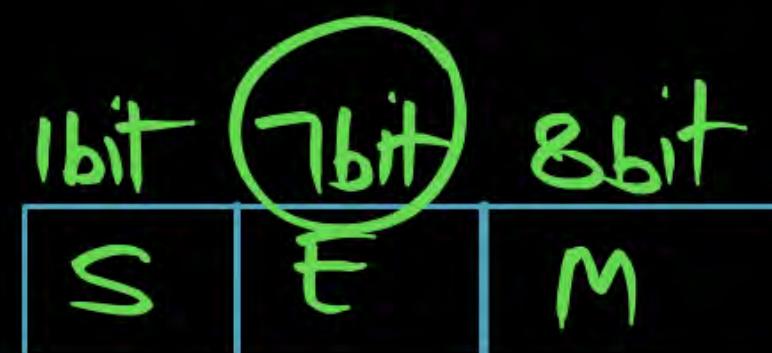
Note In the **Implicit** & **Explicit** Representation we can not Represent '0'. In IEEE 754 floating point we can represent '0'.

Q.

Consider a 16 bit register used to store floating point number.
Mantissa is Implicit normalized signed fraction number.

Exponent is in Excess-64 form then

- (i) what is the First Smallest Positive number?
- (ii) what is the Second Smallest Positive number?
- (iii) what is the Difference between First Smallest & Second Smallest Positive number?



$$\text{bias} = 2^{7-1} - 2^6$$

bias = 64

I Smallest
Implicit

S(1bit)	E(7bit)	M(8bit)
0	0000000	00000000

$$\begin{array}{l} \text{bias} = 64 \\ E = 0 \end{array}$$

$$\begin{aligned} (-1)^S & L \cdot M \times 2^{E-\text{bias}} & 0-64 \\ (-1)^0 & 1 \cdot 00000000 \times 2^{-64} \end{aligned}$$

$$I_{\text{Smallest}} = 1 \cdot 00000000 \times 2^{-64}$$

II Smallest
x 1e Number

S(1bit)	E(7bit)	M(8bit)
0	0000000	00000001

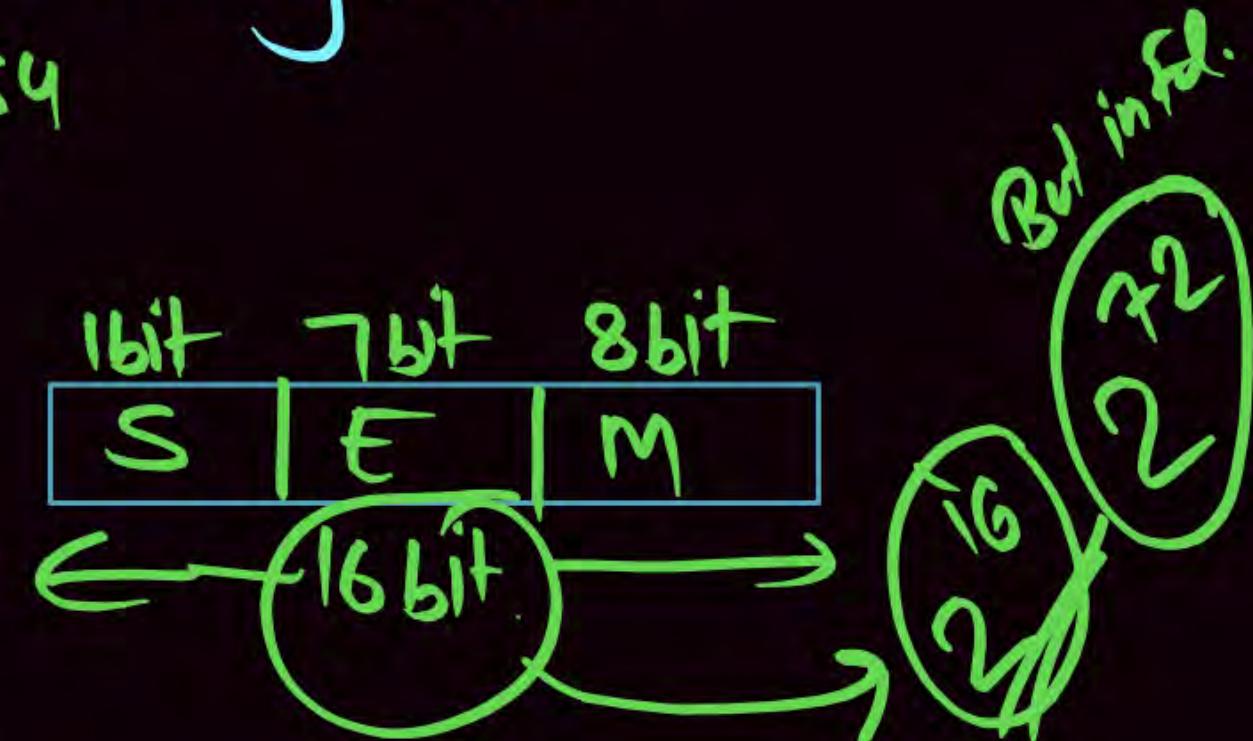
$$\begin{aligned} (-1)^S & L \cdot M \times 2^{E-\text{bias}} \\ (-1)^0 & 1 \cdot 00000001 \times 2^{-64} \end{aligned}$$

$$II_{\text{Smallest}} = 1 \cdot 00000001 \times 2^{-64}$$

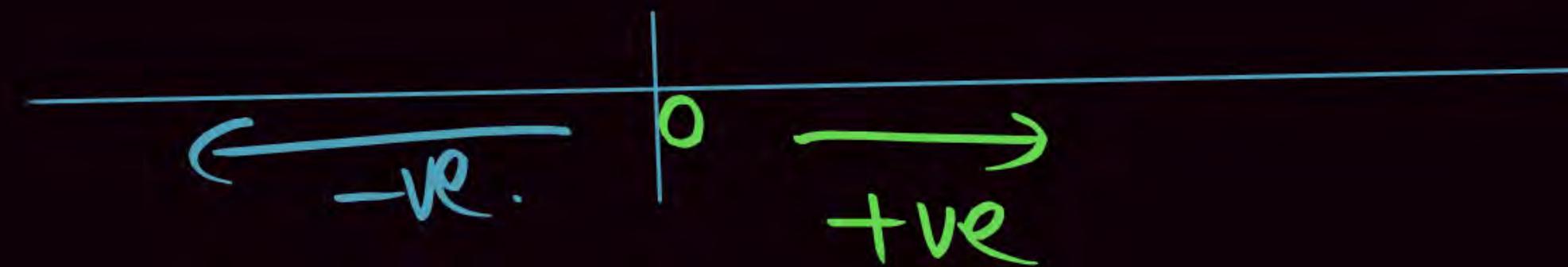
Difference b/w 1st & 2nd smallest.

$$\begin{aligned}\text{Difference} &= 1.0000000L \times 2^{-64} - 1.0000000 \times 2^{-64} \\ &= [1.0000000L - 1.0000000] \times 2^{-64} \\ &\Rightarrow (0.0000000001L) \times 2^{-64}\end{aligned}$$

Difference b/w 2 smallest +ve Number. $\Rightarrow \boxed{\frac{1}{2^{72}}} \underline{\text{Ans}}$



If we want the same thing for Negative.



Q.2

Consider a 16 bit register used to store floating point number.
 Mantissa is **Implicit** normalized signed fraction number.

Exponent is in **Excess-64** form then

- (i) what is the First Highest Positive number?
- (ii) what is the Second Highest Positive number?
- (iii) what is the Difference between First Highest & Second Highest Positive number?

Excess - 64

bias = 64

1bit	7bit	8bit
S	E	M

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

$$k-1 = 6$$

$k = 7$

(i) 1st Highest
+ve Number.

(ii) 2nd Highest
+ve Number :

S	E(7bit)	Mantissa (8bit)
0	1111111	11111111

S(1bit)	E(7bit)	M18bit
0	1111111	111111110

Difference : **Home Work**

$$\textcircled{1} \quad 0.\overline{111} \quad [1 - 2^{-3}] @ 1 - \frac{1}{2^3}$$

$$\textcircled{2} \quad 0.\overline{1111} \quad [1 - 2^{-4}] @ 1 - \frac{1}{2^4}$$

$$\textcircled{3} \quad 0.\overline{11111} \quad [1 - 2^{-5}] @ 1 - \frac{1}{2^5}$$

$$\textcircled{4} \quad 0.\overline{111111} \quad [1 - 2^{-6}] @ 1 - \frac{1}{2^6}$$

Proof in Next Pages.

V.V. Dubb.

(i)

$$0.\overline{111} \Rightarrow 1 - \frac{1}{2^3}$$

(ii)

$$0.\overline{1111} \Rightarrow 1 - \frac{1}{2^4}$$

(iii)

$$0.\overline{11111} = 1 - \frac{1}{2^5}$$

(iv)

$$0.\overline{111111} = 1 - \frac{1}{2^6}$$

$$1 - \frac{1}{2^3}$$

OR

~~proof~~ Left Alignment

$$0.\overline{111} \\ \Rightarrow 1\overline{111} \times 2^{-3} \\ \textcircled{7} [2^3 - 1]$$

$$\Rightarrow [2^3 - 1] \times 2^{-3}$$

$$\Rightarrow 2^3 \times 2^{-3} - 1 \times 2^{-3}$$

$$1 - \frac{1}{2^3}$$

$$= 1 - \frac{1}{2^3}$$

OR GP Series

~~Prob~~

$$0 \cdot 1 \mid 1$$

$$0 \cdot \frac{1}{2} \mid \frac{1}{2^2} \mid \frac{1}{2^3}$$

GP Series.

$$a = \frac{1}{2}$$

$$\delta = \frac{1}{2}$$

$$\frac{a(1-\delta^n)}{1-\delta}$$

$$\Rightarrow \frac{\frac{1}{2}(1 - \frac{1}{2^3})}{(1 - \frac{1}{2})}$$

= $\frac{1}{2} \left(1 - \frac{1}{8}\right)$ Ans

left align Q2 $0 \cdot 11111 \Rightarrow$

$0 \cdot 11111$

$\Rightarrow \frac{11111}{31} \cdot 0 \times 2^{-5}$

$\Rightarrow (\sum_{i=1}^5 1) \times 2^{-5}$ OR

$\Rightarrow 2^5 \times 2^{-5} = 2^0$

$\boxed{1 - 2^{-5}}$ $\textcircled{\times} \quad \boxed{1 - \frac{1}{2^5}}$

$0 \cdot 11111 \Rightarrow$

$1 - \frac{1}{2^5}$ Ans

Gp series $\frac{1}{2} \frac{1}{2^2} \frac{1}{2^3} \frac{1}{2^4} \frac{1}{2^5}$

$a = \frac{1}{2}, r = \frac{1}{2}$

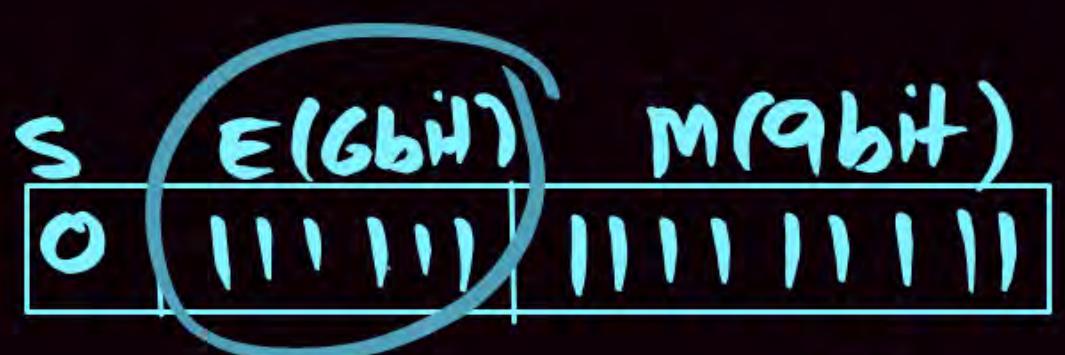
$\frac{a(1-r^n)}{1-r} \Rightarrow \frac{\frac{1}{2}(1-\frac{1}{2^5})}{1-\frac{1}{2}}$

$\Rightarrow \boxed{1 - \frac{1}{2^5}}$ Ans

64 32 16 8 4 2 1

11 → 3
100 → 4
111 → 7
1000 → 8
1111 → 15
10000 → 16
11111 → 31
100000 → 32
111111 → 63
1000000 → 64.

(Q)



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

$$(-1)^S \cdot M \times 2^{E - \text{bias}}$$

$$(-1)^0 \cdot 111111111 \times 2^{63 - 32}$$

maximum +ve
using Explicit ?

$$\text{bias} = 2^{6-1} = 32$$

$$\boxed{\text{bias} = 32}$$

$$\boxed{E = 111111}$$

$$\boxed{E = 63}$$

$$\begin{aligned}
 & + 0 \cdot \underbrace{11111111}_{\text{Ans}} \times 2^{31} \\
 & + (-2^{-9}) \times 2^{31} \\
 & + 2^{31} - 2^{31-9} \\
 & + 2^{31} - 2^{22} \quad \alpha \\
 & + 2^{31} \quad \text{Ans} \\
 & + 2^{31} \quad \cancel{\text{Ans}}
 \end{aligned}$$

$0 \cdot \underbrace{11111111}_{\text{Ans}} \times 2^{+31}$
 $\left[\underbrace{11111111}_{\text{Ans}} \cdot 0 \times 2^{-9} \right] \times 2^{31}$
 $\left[(2^9 - 1) \times 2^{-9} \right] \times 2^{31}$
 $(-2^{-9}) \times 2^{31}$
 $\Rightarrow 2^{31} - 2^{31-9} \quad \Rightarrow +2^{31} - 2^{22}$
 $\Rightarrow +\left(2^{31}\right) \underline{\text{Ans}}$

1bit	6bit	9bit
S	E	M

16bit →

Conventional Representation.

↳ Can not Represent '0' @ ' ∞ '

Implicit : L. Something (Not zero)

Explicit : 0.L... (Not zero)

Disadvantage of Conventional Representation

It can not store '0'

Explicit = 0.1

Implicit 1.0

It cannot Represent ' ∞ '

[Infinity]

It can not represent
number which is not
normalized

IEEE 754 Floating Point Representation

Single Precision
(32 bit)
Excess 127

Double Precision
(64 bit)
Excess 1023

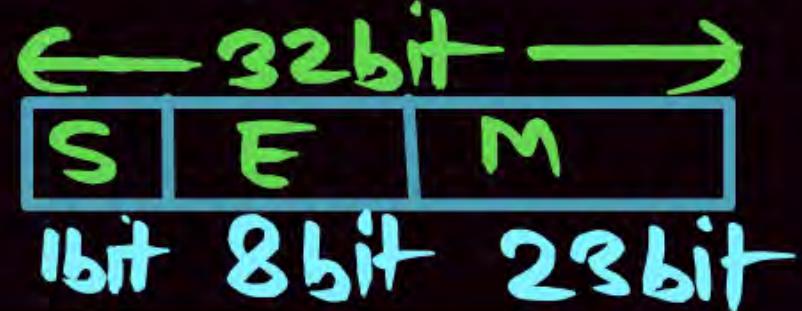


$$\begin{aligned}\text{bias} &= 2^{K-1} - 1 \\ &= 2^8 - 1 \\ \text{Bias} &= 127\end{aligned}$$

$$\begin{aligned}\text{bias} &= 2^{11-1} - 1 \\ \text{Bias} &= 1023\end{aligned}$$

IEEE-754 Floating Point Representation.

Single Precision Format (32 bit)

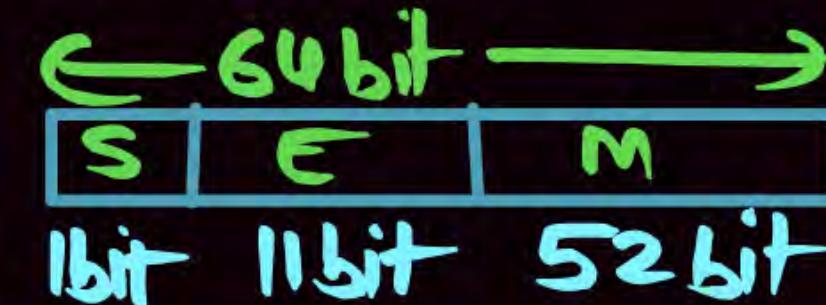


$$\text{bias} = 2^{k-1} - 1$$

Excess-127

$$\text{bias} = 2^8 - 1 \Rightarrow \text{bias} = 127$$

Double Precision Format (64 bit)



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

$$\text{bias} = 2^{11-1} - 1$$
$$\text{bias} = 1023$$

Excess-1023.

IEEE 754 Floating Point Representation.

Single Precision

S	E	M
1bit	8bit	23bit

$$\text{bias} = 127$$

OR

$$\text{Excess} - 127$$

S	E	M
1bit	11bit	52bit

$$\text{bias} = 1023$$

OR

$$\text{Excess} = 1023$$

IEEE 754 Floating Point Representation.

Single Precision

S	E	M
1bit	8bit	23bit

$$\text{bias} = 127$$

(OR)

$$\text{Excess} - 127$$

Note

In IEEE 754 Default is

S	E	M
1bit	11bit	52bit

$$\text{bias} = 1023$$

(OR)

$$\text{Excess} = 1023$$

Implicit Normalization

Q. 1

How to represent +(119) into IEEE 754 single precision & Double precision floating Point Representation?

P
W

Single Precision.

$$119: 1110111 \times 2^6$$

$$+ 1.110111 \times 2^{+6}$$

$$S=0$$

$$m = 1101110000000000000\dots$$

$$e=6$$

$$\text{bias} = 127$$

$$E = e + \text{bias} = 6 + 127 = 133 = E$$

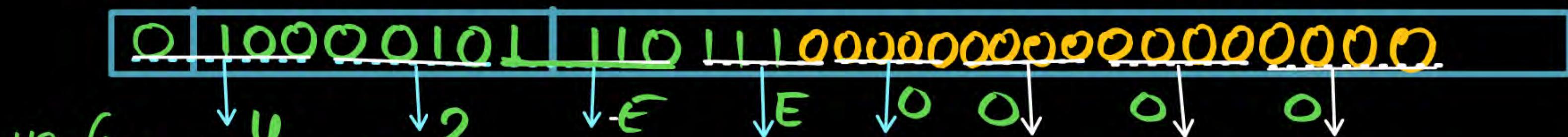
$$E = 133$$

$$E = 10000101$$

S	E	M
1bit	8bit	23bit

$$\text{bias} = 2^{8-1} - 1$$

$$\text{bias} = 127$$



$$119: (42EE0000)_{16}$$

+ (119) Single Precision.

(42EE 0000)₁₆ Ans

Now get the value from

(42EE 0000)₁₆ (P.T.O)



119 Ans

Single Precision.

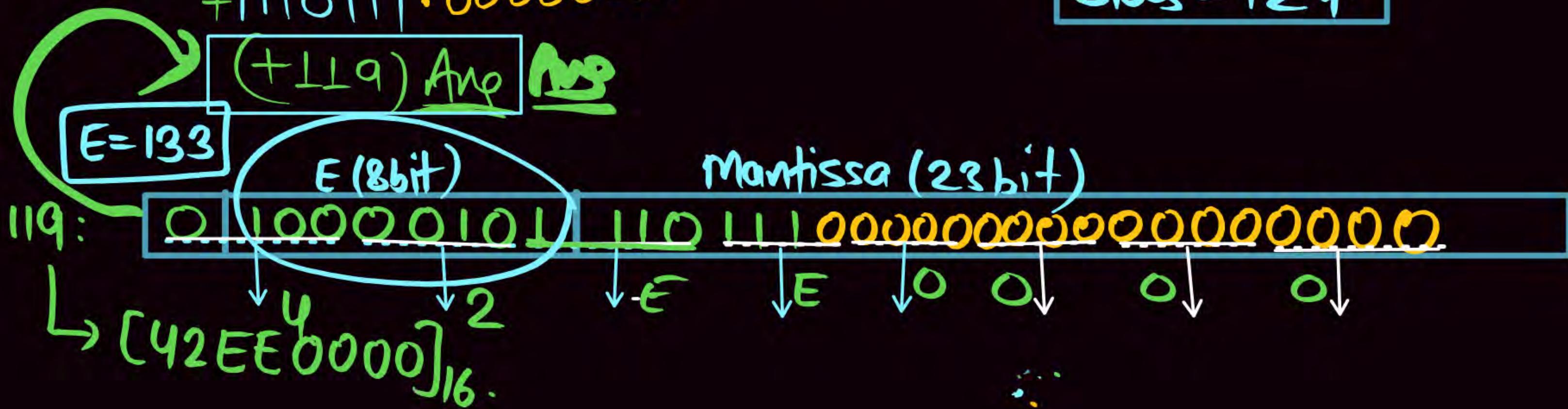
$$(-1)^S \cdot M \times 2^{E - bias}$$

$$\begin{aligned} & (-1)^0 | 0 \cdot 11011100000000000000000000000000 \times 2 \\ & + 1 \cdot 11011100000000000000000000000000 \cdot \cdot \cdot \times 2^{+6} \\ & + 1110111 \cdot 000000000 \end{aligned}$$

S	E	M
133-127 bit	8bit	23bit

$$bias = 2^8 - 1$$

$$bias = 127$$



+119)

Now Double Precision.

Q.2

How to represent +(119) into IEEE 754 single precision & Double precision floating Point Representation?

P
W

Double Precision.

119: 1110111 x²

+ 1.110111 x 2⁺⁶

S=O

$$e = +6$$

$$\text{bias} = 1023$$

$$E = 1023 + 6 = 1029$$

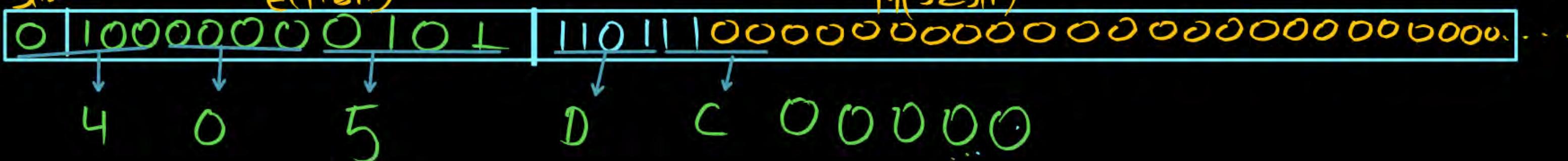
$$E = 1029$$

$$E = 100000000010$$

४१८

E(1161)

M(525t)



Q.

How to represent +(119) into IEEE 754 single precision & Double precision floating Point Representation?

P
W

+119 Double Precision

[405 D C 0000 0000 0000] 16. Ave

Now get the Value.

Q

Double Precision

E-bias

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

1029-1023

$(-1)^S | \cdot 1| 1011100000000000... \times 2^E$

$\Rightarrow + 1 \cdot 11011100000000 \times 2^{+6}$

$\Rightarrow 111011100000000$

$\rightarrow (+119) \text{ Ans}$

$E = 1029$

$E = 100000000101$

all 1's

E(11bit)

$M(52\text{bit})$

The diagram illustrates the bit layout of a double precision floating-point number. It consists of three main fields: Sign (S), Exponent (E), and Mantissa (M).

- Sign (S):** A single bit representing the sign of the number.
- Exponent (E):** An 11-bit field representing the exponent. In double precision, the exponent is biased by 1023, meaning the value stored is the actual exponent plus 1023.
- Mantissa (M):** A 52-bit field representing the fraction part of the number. It includes an implicit leading 1, so the total precision is 53 bits.

The diagram also labels specific bits:

- S:** The 4th bit from the left.
- E:** The 0th bit from the left.
- D:** The 5th bit from the left.
- C:** The 0th bit from the right.

S	E	M
1bit	11bit	52bit

$$\text{bias} = 2^{11-1} \rightarrow 2^{10}-1$$

$$\text{bias} = 1023.$$

Q.

2 Marks Consider a 32 bit register which stores floating numbers in IEEE single precision format(Implicit). The value of the number, f₃₂ bit are given below is _____

P
W

Sign(1bit)	Exponent(8bit)	Mantissa(23bit)
0	10000100 $E = 132$ bias = 127	11000000000000000000000 0000

- A 48
 B 56
 C 64

- (B) 56
 D 60

$$\begin{aligned}
 & (-1)^S \cdot 1.M \times 2^{E - \text{bias}} \\
 & (-1)^0 \cdot 1 \cdot 1100000000 \times 2^{132 - 127} \\
 & + 1 \cdot 1100000000 \times 2^5 \\
 & \Rightarrow 1.11000 \cdot 00000 \times 2^0 \\
 & = \underline{\underline{+56}}
 \end{aligned}$$

Solution

Sign(1bit)	Exponent(8bit)	Mantissa(23bit)
0	10000100	11000000000000000 0000

$$(-1)^S (1.M) \times 2^{E-127}$$

$$(-1)^0 (1.110....) \times 2^{132 - 127}$$

$$E = e + \text{bias}$$

$$(1.11) \times 2^5$$

$$E - \text{bias}$$

$$= 111000$$

$$E = 10000100$$

$$= 56 \text{ Ans}$$

$$E = 132$$

Q.

P
W

The value of float type variable is represented using the single precision 32 bit floating point IEEE 754 standard use 1 bit sign, 8bit for E and 23 bit mantissa. A float type variable X is assigned the decimal value (-14.25). Then representation of X in Hexa decimal notation is?

A 416C0000H

B 41640000H

C C16C0000H

C1640000H

()₁₆ ()H

Solution

-14.25

1110.01

$$1.11001 \times 2^3$$

$$S = 1$$

$$M: 11001$$

$$e = 3$$

1 bit	8 bit	23 bit
S	E	M

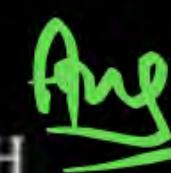
$$\text{Bias} = 2^{8-1} - 1 = 127$$

$$E = e + \text{bias}$$

$$= 3 + 127$$

$$E = 130$$

S	E	M
1	10000010	11001 0000000000000

Hexadecimal = $(C1640000)_H$ 

IEEE 754 Floating Point Representation

IEEE 754 Floating Point Representation

Single Precision
(32 bit)
Excess 127



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

Double Precision
(64 bit)
Excess 1023



$$\begin{aligned}\text{bias} &= 2^{11-1} - 1 \\ \text{Bias} &= 1023\end{aligned}$$

IEEE 754 Floating Point Representation

IEEE 754 Floating Point Representation

Conventional Representation

Drawback
Can not Represent

Single Precision
(32 bit)
Excess 127

Double Precision
(64 bit)
Excess 1023

(i) 0



(ii) ∞



(iii) (a) NaN (Not a Number)

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

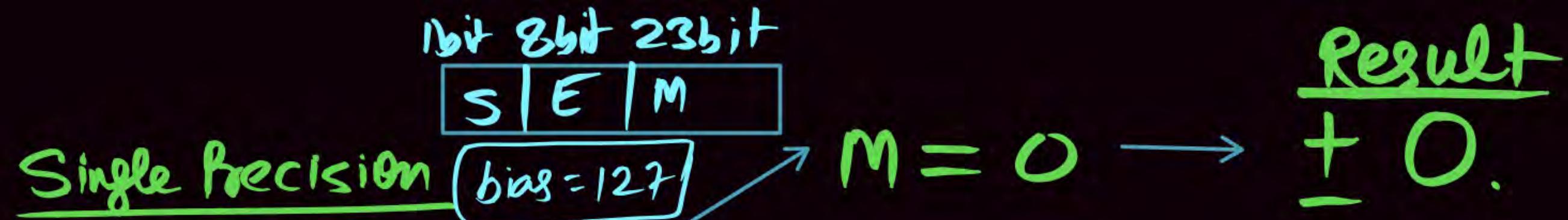
$$\begin{aligned} \text{bias} &= 2^{11-1} - 1 \\ \text{Bias} &= 1023 \end{aligned}$$

(iv) (a) Denormalized Number

Single Precision (32 bit)

S	E	M
1 bit	8 bit	23 bit

Sign(1 bit)	E(1 bit)	M(23 bit)	Value
0 or 1	00000000 E = 0	0000000000000000 M = 0	± 0
0 or 1	11111111 E = 255	0000000000000000 M = 0	$\pm \infty$
0 or 1	$1 \leq E \leq 254$	M =	Implicit Normalized form $(-1)^S \times 1.M \times 2^E$ $(-1)^S \times 1.M \times 2^{E-127 \text{ bias}}$
0 or 1	E = 0	M ≠ 0	Denormalized number/Fractional form $(-1)^S \times 0.M \times 2^{E-127 \text{ bias}}$
0 or 1	E = 255	M ≠ 0	Not a Number (NAN)



$E = 0.$

$E = 0000000000$

$M \neq 0.$

Denormalized /
Fractional Number.

$E = 11111111$

$M = 0 \rightarrow \pm \infty$

$E = 255$

$M \neq 0$

NaN [Not a Number].

Q Why $\text{bias} = 2^{k-1} - 1$? Single Precision

	1bit	8bit	23bit
S	E	M	

If in Single Precision, if $\text{bias} = 2^7 \Rightarrow 2^7$ ($\text{bias} = 128$) then

there is a possibility to getting $E = 255$

$$E = e + \text{bias}$$

$$(e=+127) \text{ then } E = 127 + 128$$

8 bit 2's complement Range = -2^{8-1} to $+2^{8-1}$

$$\Rightarrow -128 \text{ to } +127$$

$$\text{So bias} = 2^{k-1} - 1$$

\cong
Avg

$$E = 255$$

$$\pm \infty$$

$$if E=255$$

Not a Number.

Double precision

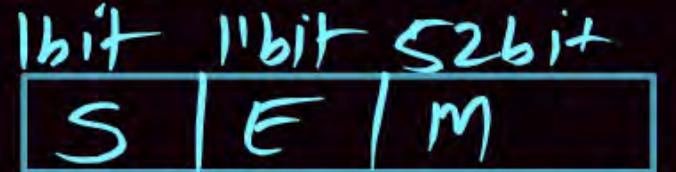
Double Precision

1-bit 11-bit 52-bit

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

Excess - 1023

Sign (1 bit)	E(11 bit)	M(52 bit)	Value
0 or 1	0000 0000 000 E = 0	000000000000.. M = 0	± 0
0 or 1	1111 1111 111 E = 2047	000000000000.. M = 0	$\pm \infty$
0 or 1	$1 \leq E \leq 2046$	M = -----	Implicit Normalization $(-1)^S \cdot M \times 2^E$ $(-1)^S \times 1 \cdot M \times 2^{E-1023}$
0 or 1	E = 0	M \neq 0	Denormalized number/ Fractional Form $(-1)^S 0 \cdot M \times 2^{E-1023}$
	E = 2047	M \neq 0	Not a number



Double Precision

bias=1023

Result

$M = 0 \rightarrow \pm 0.$

$E = 0.$

$E = 000000000000$

$M \neq 0. \rightarrow$

Denormalized /
Fractional Number.

$E = 1111111111$

$\rightarrow E = 2047$

$M = 0 \rightarrow \pm \infty$

$M \neq 0$

$\rightarrow \text{NaN} [\text{Not a Number}]$

NOTE: When $E = 0$ then Value 0 or fractional form
when $M = 0$ when $M \neq 0$

NOTE: When $E = 2047$
then Value ∞ or Not a Number(NAN)
when $M = 0$ when $M \neq 0$

TEEE - 754

Features.

Note

Features IEEE 754 are special symbols to represent unusual events....

- ① For example instead of interrupting on a divide by 0, software can set the result to a bit pattern representing $+\infty$, $-\infty$; The largest exponent is reserved for these special symbols.
- ② IEEE 754 has a symbols for the result of invalid operations, such 0/0, or subtract infinity from infinity. This symbols is NaN, for Not a Number.

The purpose of NaN is to allow programmer to postpone some test and decisions to a later time in this program. (when it is convenient)

Q.

How to represent +(1.0) into IEEE 754 single precision floating Point Repartition?

P
W

Sol.

To represent +(1.0) into IEEE 754 single precision floating Point

P
W

Repartition..



$$+ 1.0 \times 2^0$$

$$M = 0$$

$$e = 0$$

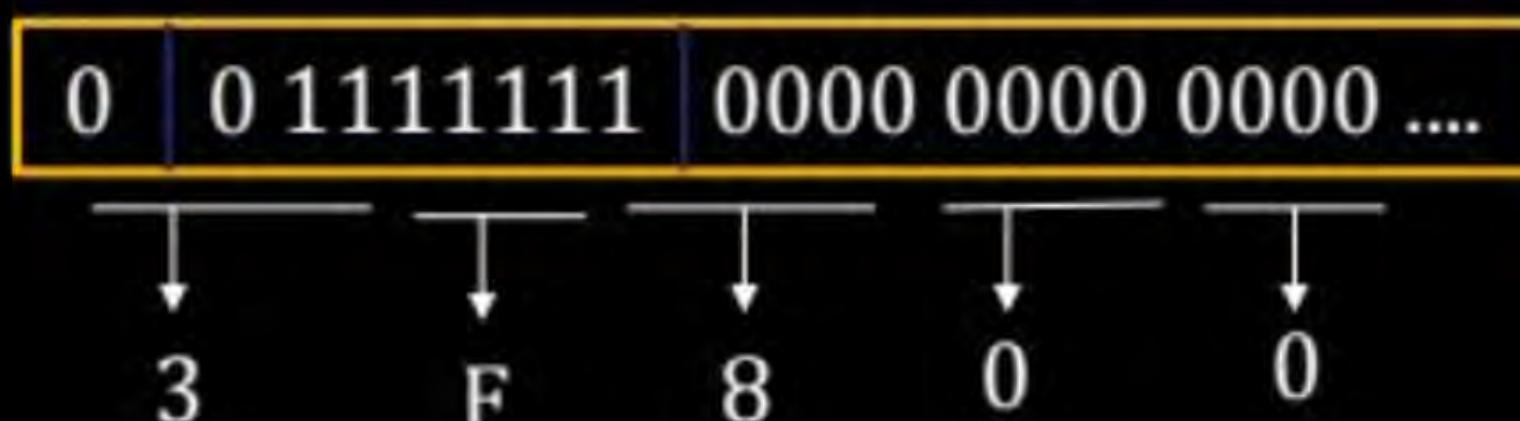
$$E = 127$$

S	E	M
1	8	23

$$\text{Bias} = 2^{8-1}-1 = 127$$

$$E = e + \text{bias}$$

$$E = 0 + 127$$



(3F80 0000)

Q.

What will be maximum positive value(+ve) in IEEE 754 single precision floating Point Representation?

P
W

Q.

Consider a 64 bit Register which store floating point Number in IEEE 754 Double Precision Format.

What is the value of the Number if 64 bit are given as
0111 1111 1111 0000 0000 0000000000....

P
W

Sol.

Consider a 64 bit Register which store floating point Number in IEEE 754 Double Precision Format.



What is the value of the Number if 64 bit are given as
0111 1111 1111 0000 0000 0000000000....

S(1 bit)	E(11 bit)	M(52 bit)
0	1111 1111 111	00000000000000

Here all exponent bits is 1

E = 2047

• Infinity (when $M = 0$)

• Not a number (when $M \neq 0$)

Here $M = 0$ so its $+\infty$

The value of a *float* type variable is represented using the single-precision 32-bit floating point format of IEEE-754 standard that uses 1 bit 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

[GATE-2014-Set2-CS: 2M]

- A C1640000H
- B 416C0000H
- C 41640000H
- D C16C0000H

Consider the IEEE-754 single precision floating point numbers
 $P = 0xC1800000$ and $Q = 0x3F5C2EF4$.

Which one of the following corresponds to the product of these numbers
(i.e., $P \times Q$), represented in the IEEE-754 single precision format?

[GATE-2023-CS: 2M]

- A 0x404C2EF4
- B 0x405C2EF4
- C 0xC15C2EF4
- D 0xC14C2EF4

Round off.

Add/Subtract Rule

1. Choose the number with the smaller exponent and shift its mantissa right a number of steps equal to the difference in exponents.
2. Set the exponent of the result equal to the larger exponent.
3. Perform addition/subtraction on the mantissas and determine the sign of the result.

Normalize the resulting value, if necessary.

Multiplication and division are somewhat easier than addition and subtraction, in that no alignment of mantissas is needed.

Multiply Rule

1. Add the exponents and subtract 127.
2. Multiply the mantissas and determine the sign of the result.
3. Normalize the resulting value, if necessary.

Divide Rule

1. Subtract the exponents and add 127.
2. Divide the mantissas and determine the sign of the result.
3. Normalize the resulting value, if necessary.

The addition or subtraction of 127 in the multiply and divide rules results from using the excess -127 notation for exponents.

**THANK
YOU!**

