

COMPUTER SCIENCE

Computer Organization and Architecture

Floating Point Representation

Floating Point Representation

Lecture_01



Vijay Agarwal sir

A graphic of a construction barrier with orange and white diagonal stripes and two yellow bollards. A yellow diamond-shaped sign is mounted on a post next to the barrier, containing the text 'TOPICS TO BE COVERED'.

**TOPICS
TO BE
COVERED**

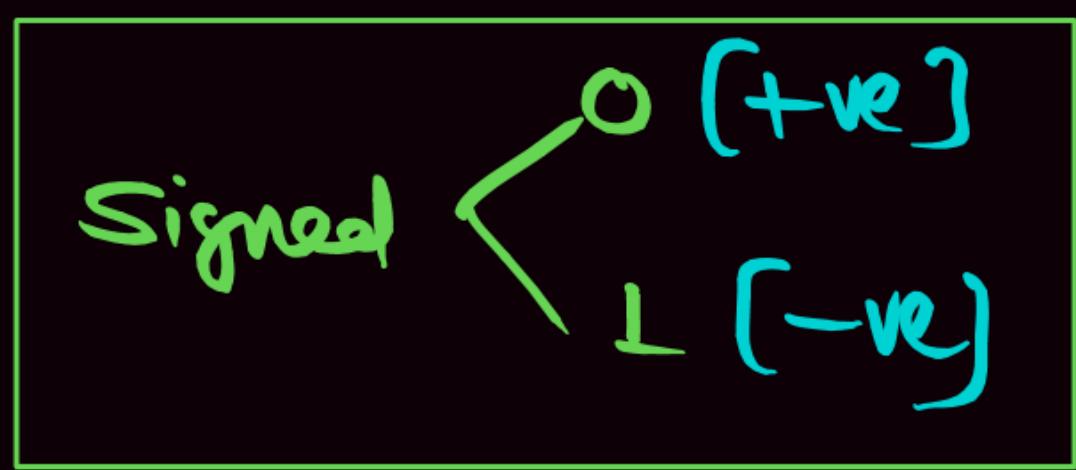
o1

Floating Point Representation

- ① Introduction of COA
- ② Instrⁿ format & Am.

Number System

Floating Point Representation



Magnitude → Unsigned (Only +ve Number)

Magnitude → Signed (+ve & -ve Number)

Why 2's
Complement:

In Signed magnitude &
1's Complement
0 has Redundant Representation

+0
-0

DATA FORMAT

Fixed Point Data

Magnitude Format

ubit
 $0000 \rightarrow 0$

Unsigned Format
[Only the Number]

Signed Magnitude Format
[+ve and -ve]

n bit range
 $0 \text{ to } 2^n - 1$



Sign value

$-(2^{n-1} - 1) \text{ to } + (2^{n-1} - 1)$

4bit range

$0 \text{ to } 2^4 - 1 \Rightarrow 0 \text{ to } 15$

$0 \oplus 2^4 - 1$

0⊕15

ubit

$0000 \rightarrow 0$

$1111 \rightarrow 15$

Complement Format

1's complement Format

Both
[+ve and -ve]

$-(2^{n-1} - 1) \text{ to } + (2^{n-1} - 1)$

4bit range

$[-7 \text{ to } +7]$

2 Representations for 0

+0: 0000
-0: 1111

Floating Point Data

Single Precision Data (32 bit)

Double Precision (64 bit)

$n-1$
 -2^{n-1} to $+2^{n-1}$

4bit Range
 $4-1$
 $4-1$

-2^{4-1} to $+2^{4-1}$
 $\Leftrightarrow -8 \text{ to } +7$

1's Complement

1's \rightarrow 0

0 \rightarrow 1

1011

↓↓

1's
Complement: 0100

DE | DUD.

~~Flags~~

2's Complement

1's
Complement + 1.

Carry
Parity
Sign
zero
Auxilliary
Overflow

Conditional
Flag

2's Compl.

$$\begin{array}{r} 0111 \\ +1 \\ \hline 1000 \end{array} \Rightarrow -8$$

$$\begin{array}{r} 0110 \\ +1 \\ \hline 0111 \end{array} \Rightarrow -7$$

$$\begin{array}{r} 0101 \\ +1 \\ \hline 0110 \end{array} = -6$$

$$\begin{array}{r} 0100 \\ +1 \\ \hline 0101 \end{array} \xrightarrow{\text{9} \sim 1} -5$$

$$\begin{array}{r} 0000 \\ +1 \\ \hline 0001 \end{array}$$

1's Complement

$$\begin{array}{r} 1000 \\ 0000 \\ \hline 0111 \end{array} \Rightarrow 0111 = -7$$

$$1001 \Rightarrow \underline{0110} = -6$$

$$1010 \Rightarrow 0101 = -5$$

$$1011 \Rightarrow 0100 = -4$$

$$1111 \Rightarrow 0000 = -0$$

4 bit Binary	<u>Unsigned Data</u>	<u>Signed Magn. Data</u>	1's Complement Data	2's Complement Data
0000	0	+0	+0	+0
0001	1	+1	+1	+1
0010	2	+2	+2	+2
0011	3	+3	+3	+3
0100	4	+4	+4	+4
0101	5	+5	+5	+5
0110	6	+6	+6	+6
0111	7	+7	+7	+7
1000	8	-0	-7	-8
1001	9	-1	-6	-7
1010	10	-2	-5	-6
1011	11	-3	-4	-5
1100	12	-4	-3	-4
1101	13	-5	-2	-3
1110	14	-6	-1	-2
1111	15	-7	-0	-1

Redundant Representations for '0'

Note

2's Complement are used in
the Computer System to Represent
the Negative Number.

(e) 4 bit 2's Complement Range = $-(2^{4-1})$ to $+(2^{4-1} - 1)$ Ans
 $\Rightarrow -8$ to $+7$

Floating-Point Representation

Principles

- With a fixed-point notation it is possible to represent a range of positive and negative integers centered on or near 0.
- By assuming a fixed binary or radix point, this format allows the representation of numbers with a fractional component as well
- Limitations:
 - ❖ Very large numbers cannot be represented nor can very small fractions
 - ❖ The fractional part of the quotient in a division of two large numbers could be lost

Why Floating Point Representation: ?

Very - Very Large Number \Rightarrow [∞] nearly 9764.1234.....

Very - very Small Number \Rightarrow [0] 00. -12

16 bit Number \Rightarrow $-(2^{16-1})$ to $+(2^{16-1})$

$\Rightarrow -2^{15}$ to $+2^{15} - 1$

\Rightarrow -32k to +32k - 1

As we want
51000 Not Possible with
~~(51k)~~ 16bit Data

Floating-Point Representation

16 bit fixed point data format then

Range = -2^{16-1} to $+ (2^{16-1} - 1)$

$\Rightarrow \underline{-(2^{15}) \text{ to } + (2^{15} - 1)}$

If we want to store 61,000 then we cannot store

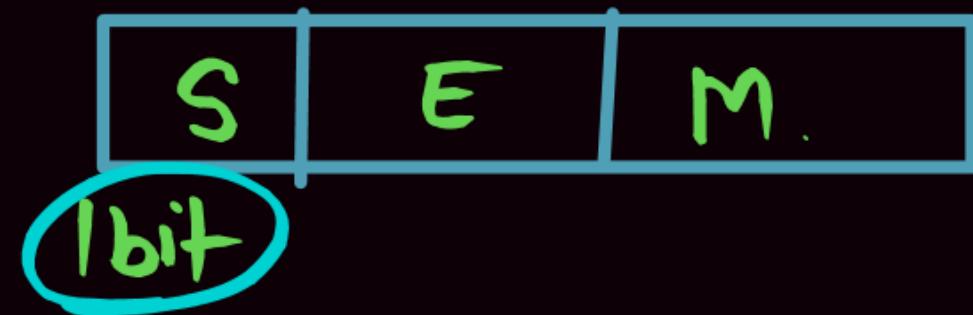
Because range $[-32k \text{ to } + 32 k - 1]$

So floating point representation is to represent very large data and very small fraction and consume less memory

Floating point
used to represent

$$\left\{ \begin{array}{l} + \underline{8564100000000000\dots} [\Rightarrow \infty] \\ + \underline{0.000000000007892} \Rightarrow [0] \end{array} \right.$$

FLOATING POINT Rep.



S[sign] \rightarrow 0 (+ve)
 \rightarrow 1 (-ve)

E | BE : Exponent | Biased Exponent

M : Mantissa .



③ Sign 1bit
 $0.\underline{nnnn} \times 2^e$
 ↓
 Mantissa .

e : exponent

$$E = e + bias$$

(OR)

$$BE = AE + bias$$

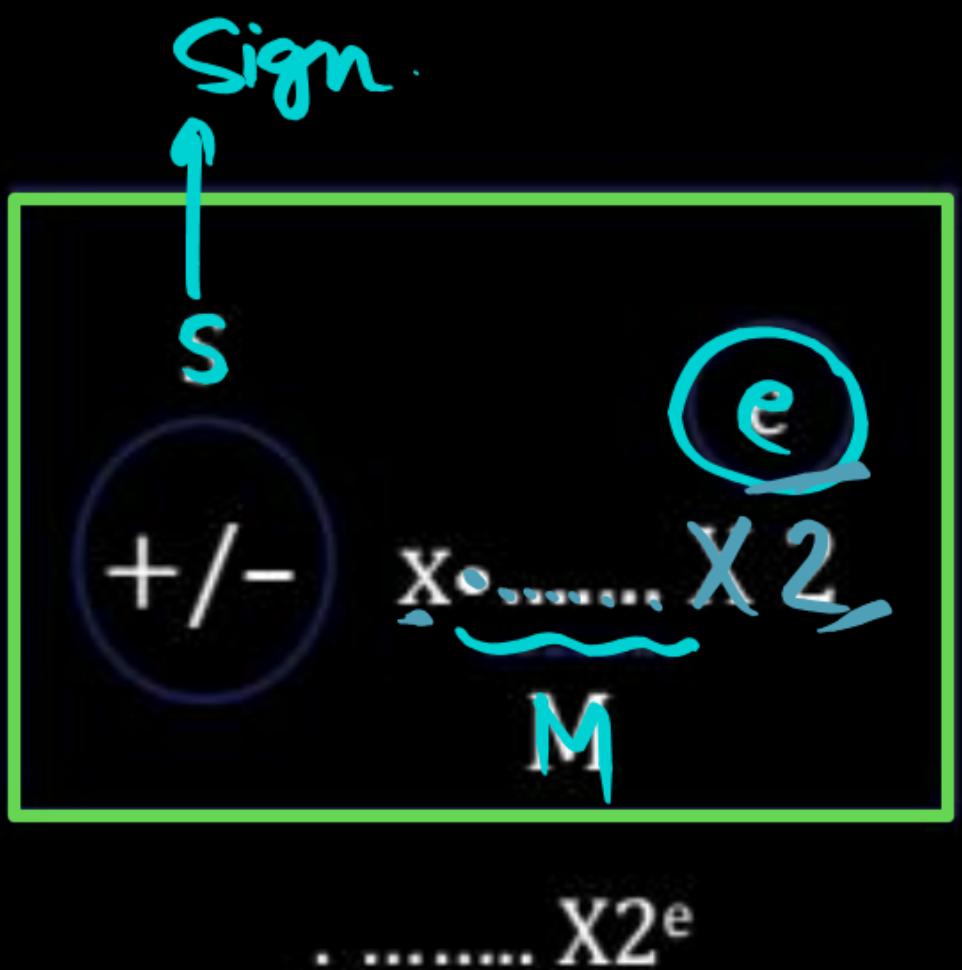
Floating-Point Representation



S: sign bit $\begin{cases} 0 & +\text{ve} \\ 1 & -\text{ve} \end{cases}$

E: exponent

M: Mantissa

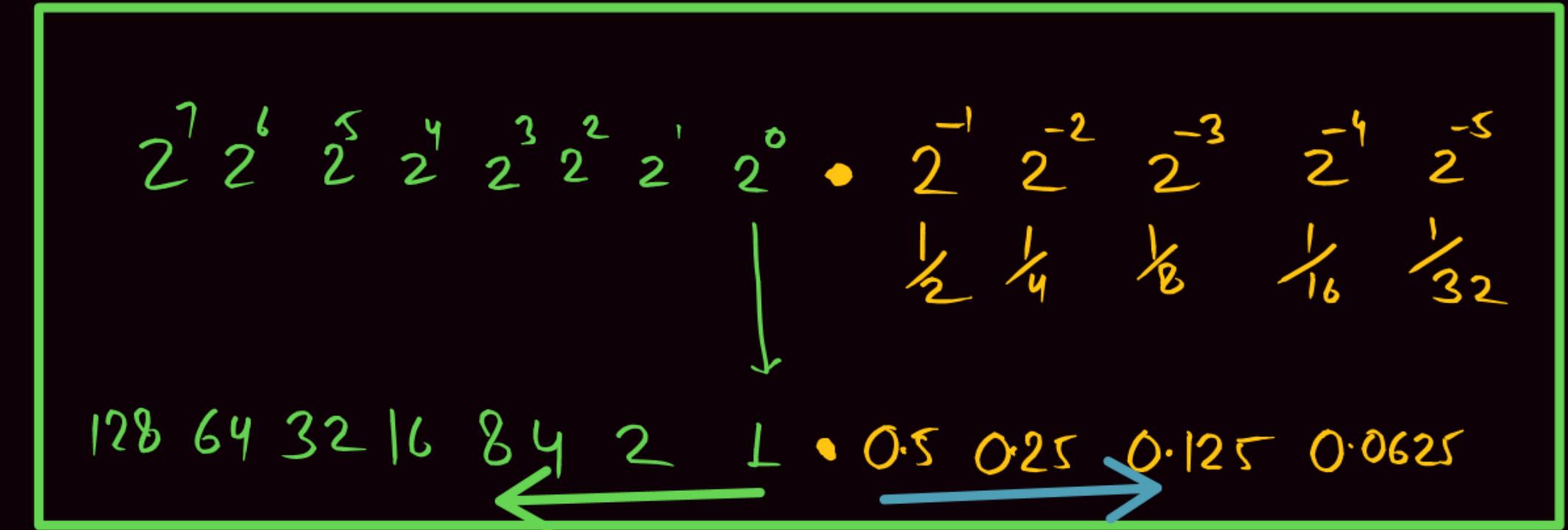


6.5 in Binary \Rightarrow 110.1

(Q.L)

6.5

110.1



$+/- x \text{ Mantissa} \times 2^e$

G

• 4bit 0110
• 5bit 00110 } X
• 6bit 000110 }

6
110

Q.1

6.5

110.1

$$[0.1101] \times 2^{+3}$$

2nd Technique

Data bits
Side Direction
Right Alignment $\Rightarrow 2^0$

Left Alignment $\Rightarrow 2^{-1}$

1st
Technique
to check

e^{+3}

S	e	M
1bit	2bit	4bit

V.V.V. Disp.

Q. $e = +3$

is Correct or Not?

$$[0.1101] \times 2^{+3}$$

$$\left(2^{-1} + 2^{-2} + 2^{-4}\right) \times 2^{+3}$$

$$\left(\frac{-1}{2} \times 2^0\right) + \left(\frac{-2}{2} \times 2^0\right) + \frac{-4}{2} \times 2^{-3}$$

$$2^2 + 2^1 + 2^{-1}$$

$$4 + 2 + 0.5 = 6.5$$

Q.1

6.5

110.1
...

$$+ \cdot 1101 \times 2^{+3}$$

S = 0 (+ve)

e = 3 [LL]

M = 1101

S	e (2bit)	m (4bit)
0	LL	1101

S	e	M
1bit	2bit	4bit

V.V.V. Disp.

Q) e = +3 is correct or not?

$$[0.1101] \times 2^{+3}$$

$$\left(2^{-1} + 2^{-2} + 2^{-4}\right) \times 2^{+3}$$

$$\left(\frac{1}{2} \times 2^3\right) + \left(\frac{1}{4} \times 2^3\right) + \left(\frac{1}{16} \times 2^3\right)$$

$$4 + 2 + 0.5 = 6.5$$

$$2^4 2^3 2^2 2^1 2^0 \cdot 2^{-1} 2^{-2} 2^{-3} 2^{-4}$$

$$0.\underline{1101} \times 2^{+3} \Rightarrow (6.5)$$

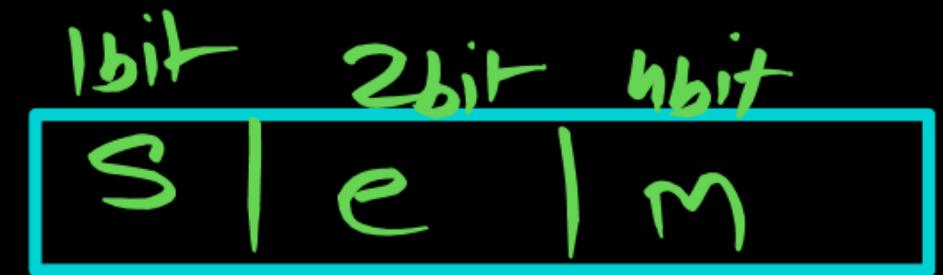
$$\left[2^{-1} + 2^{-2} + 2^{-4}\right] \times 2^{+3}$$

$$2^{+2} + 2^{+1} + 2^{-1}$$

$$4 + 2 + \frac{1}{2}$$

6.5

(Q.2) + (4.5)



(Q.3) + 4.75

S	e	m
0	11	1001

S	e	m
0	11	10011

Q. 1

+6.5

$$6.5 = (110.1)_2$$

$$\frac{0.1101}{S} \frac{2^3}{M} \times \frac{2^e}{2^e}$$

S = 0 (+)

M = 1101

e = 3 = (11)₂

S(1bit)	e(2bit)	M(4bit)
0	11	1101

Very. Imp

$$6.5 = 110.1$$

$$= .1101 \times 2^3$$

$$= [2^{-1} + 2^{-2} + 2^{-4}] \times 2^3$$

$$= [2^2 + 2^1 + 2^{-1}]$$

$$= 6.5$$

Right Alignment = $\frac{+ve}{2^e} \Rightarrow +ve$

Left Alignment = $\frac{-ve}{2^{-e}} \Rightarrow -ve$

P
W

Q. 2

+ 4.5

100.1

0.1001×2^3

S = 0 (+ve)

M = 1001

e = 3 [11]

S	e	M
0	11	1101

✓ Q. 3

+ 4.75

100.11

.10011 $\times 2^3$

S = 0

M: 10011

e = 3 $\Rightarrow (11)_2$

S	e	M
0	11	10011

NOTE:

Mantissa alignment process is used to adjust the decimal point; in this process right alignment increments the exponent and left alignment decrements the exponent.

2^{+shift} power (+) = Right alignment \Rightarrow Increment the exponent

2^{-shift} power (-) = Left alignment \Rightarrow Decrease the exponent

Right Alignment

6.5

110.1

$$\Rightarrow .1101 \times 2^3$$

$$\Rightarrow [2^{-1} + 2^{-2} + 2^{-4}] \times 2^3$$

$$\Rightarrow 2^2 + 2^1 + 2^{-1}$$

$$\Rightarrow 4 + 2 + 0.5$$

$$\Rightarrow 6.5 \text{ Ans}$$

$$(2^{+3})$$

Left AlignmentData: 0.0000000101 $\times 2^{+5}$

$$[1.01 \times 2^{+5-8}]$$

$$+ 1.01 \times 2^{-3}$$

(Align to use upto 8 times)

$$1.01 \times 2^{-8} \times (2^{+5})$$

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

Q) $+0.00101$

$$+0.\underline{101} \times 2^{-2}$$

Left Alignment

1bit	4bit	5bit
S	e	m

S = 0 (+ve)

M = 101

e = -2

e = -ve

No Provision to Represent (+ve)
exponent is Negative.

Here exponent is Negative but Sign bit [0] Bcz Number is +ve.

So How to Deal this Negative Exponent.

(Q.L)

Why Biasing is Required ?

(Q.2)

How bias Value Decide>Select ?

Q + 0.00101
e ← Left Alignment

$$+ 0.\underline{101} \times 2^{-2}$$

S = 0 (+ve)

M = 101

e = -2

Here e is -ve

So taking 2's Complement

1bit	4bit	5bit
S	e	m

1bit	4bit	M (\leq 5bit)
0	1110	10100

Ans

Need of e ?

$$E = e + bias$$

..

Why biasing Required ?

(Note) if e = -ve then taking 2's complement.

$$2\text{'s Complement} = -(2^{n-1}) \text{ to } +(2^{n-1} - 1)$$

$$4\text{bit } 2\text{'s Complement} = -2^{4-1} \text{ to } +2^{4-1} - 1$$

5 bit then

$$-2^4 \text{ to } +2^4 - 1$$

$$\Rightarrow -16 \text{ to } +15$$

$$\Rightarrow -2^3 \text{ to } +2^3 - 1$$

$$= \boxed{-8 \text{ to } +7}$$

$\xrightarrow{-ve.}$ $\xrightarrow{0}$ $\xrightarrow{+ve}$

if we take 2's Complement = 5 bit \Rightarrow -16 to +15

So Solution

Note

Convert the Number into '0' ~~or~~ Any tve Number.
So biasing will be Used.



Q. 4

$$+ 0.\underline{00}101\ldots$$

$$0.\underline{101} \times 2^{-2}$$

$$M = 101$$

$$E = -2$$

$$S = 0$$

$$E = -2 = (1110)_2 \text{ 2's complement}$$

So taking
Complement

S	E(4bit)	M(5 bit)
0	1110	10100

E

M

$$e = 14 \Rightarrow 1110$$

$$e = -2 \Rightarrow 1110$$

$$e = 14 \Rightarrow 1110$$

Biasing: is method in which we convert the negative number into

the positive number [exponent].

(Any Number)
Biasing = to Convert -ve Number into 0 or Any +ve Number.

Q2) How bias is select ?

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

$$e=k\text{bit}$$

Sol) If exponent is k bit then $2^{\text{exponent}} = \frac{k-1}{-2} \text{ to } +2^{-1}$

e3) 4 bit Range = $-2^{4-1} \text{ to } +2^{4-1} \Rightarrow -8 \text{ to } +7$

Note In Order to Convert ALL Number into +ve (Positive Number)
then take the Most (Highest) Negative Number & Add as Bias.
 (2^{k-1}) :

if exponent K bit

③ $e = 4 \text{ bit} \Rightarrow -2^{4-1} \text{ to } +2^{4-1} = -8 \text{ to } +7$

bias = -8

~~so~~

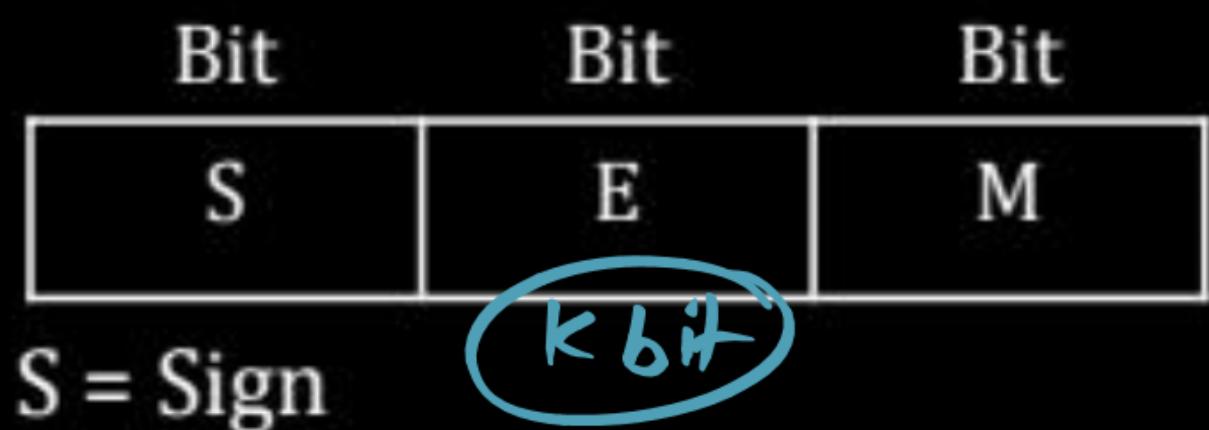
$$E = e + \text{bias}$$

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

$$\begin{array}{r} -8 + 8 = 0 \\ -7 \\ -6 \\ -5 \\ -4 \\ -3 \\ -2 \\ -1 \\ -0 \\ +1 \end{array}$$

4 bits

$$+7 + 8 < 15$$



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

(OR)

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

E/BE = Exponent or

BE = bias exponent

M = Mantissa

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

$$\text{Bias} = 2^{K-1}$$

where K is exponent bits

$$\text{Bias} = 2^{K-1}$$

K : # bits in exponent

Example

If K = 4 bits

Exponent = 4 bit then

$$\text{bias} = 2^{K-1} = 2^{4-1} = 8$$

1 Bit

4 Bit



x bit

$$\text{Bias} = 2^{K-1} = 2^{4-1}$$

$$\text{bias} = 8$$

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

$$E = e + 8$$

$$E = 4 \text{ bit}$$

or

Excess 8 code

$$2^{K-1} = 8$$

$$2^{K-1} = 2^3$$

$$K - 1 = 3$$

$$K = 4$$

$$E = 4 \text{ bit}$$

$$2^{K-1} = 8$$

$$2^{K-1} = 2^3$$

$$K - 1 = 3$$

$$K = 4$$

*If K=4
then bias
 $2^{4-1} = 8$*

e [original exponent]	Stored exponent [BE] E
-8	+8 → 0
-7	+8 1
-6	+8 2
-5	+8 3
-4	.
-3	.
-2	.
-1	1
0	8
1	9
2	10
3	11
4	12
5	13
6	+8 → 14
7	+8 15

Q.

From previous question

0.00101

$$0.101 \times 2^{-2}$$

$$M = 101$$

$$\text{Bias} = 2^{5-1}$$

$$\text{Bias} = 16$$

$$e = -2$$

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

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

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

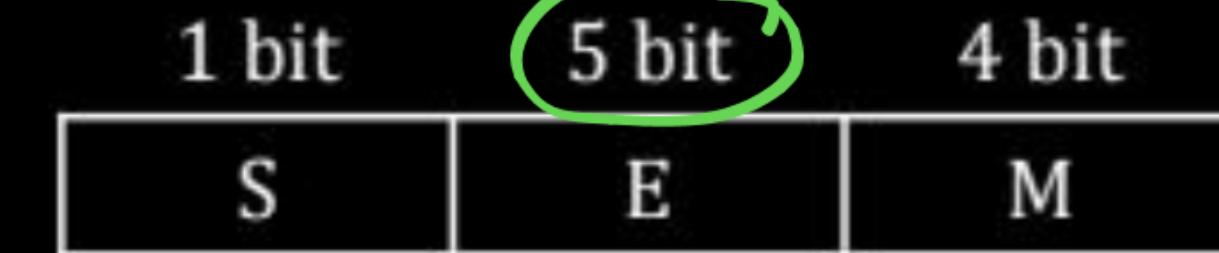
$$E = -2 + 16$$

$$E = 14$$

$$E = (01110)_2$$

$$\text{Formula: } (-1)^S \times 0.M \times 2^E$$

$$(-1)^0 \times 0.101 \times 2^{E-\text{bias}}$$



$$E = 5 \text{ bit}$$

$$\text{bias} = 2^{5-1} = 16$$

$$E = -2 + 16$$

$$E = 14$$



$$E = 14$$

$$\text{bias} = 16$$

M.
Ans

$$0.101 \times 2^{14-16} = 0.101 \times 2^{-2}$$

$$0.00101 \text{ Ans}$$

P
W

$$+ \underbrace{x_0}_{M} \cdot \underbrace{x_2}_{M}^e$$

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

$$+ \underbrace{x_0}_{M} \cdot \underbrace{x_2}_{M}^{E-\text{bias}}$$

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



Normalized Mantissa

1 bit x bit y bit

Explicit Normalized Syntax

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

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

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

Formula to get number
[value formula]

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

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

**THANK
YOU!**

