

# COMPUTER SCIENCE



Computer Organization  
and Architecture

Floating Point  
Representation

Lecture\_01

Vijay Agarwal sir







TOPICS  
TO BE  
COVERED

o1

Floating Point Representation

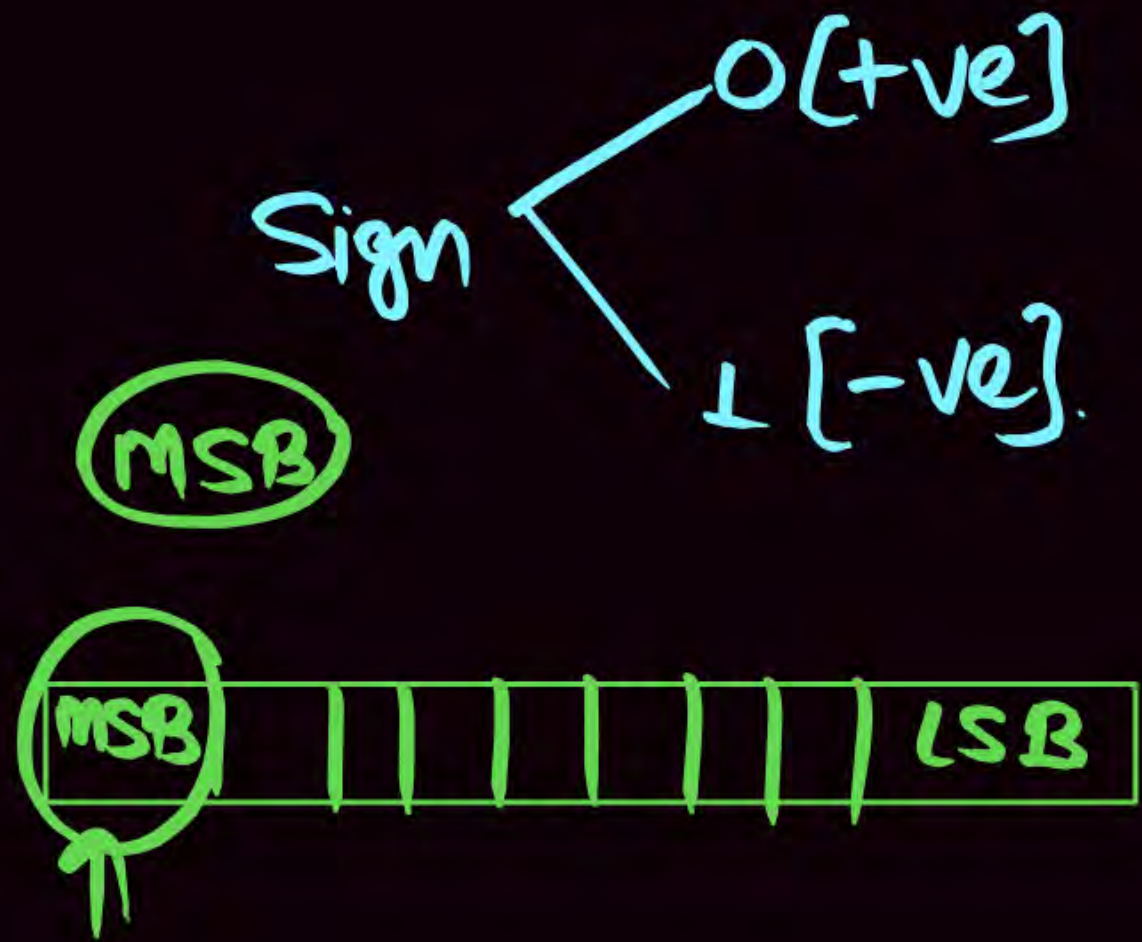
- ✓ ① Introduction of COA.
- ✓ ② Machine Instruction & Addressing Mode.
- ✓ ③ Floating Point Representation.

✓ Number System

✓ Floating Point Representation



## Number System:



- ① Magnitude Format
  - [Signed & Unsigned]
    - (+ve & -ve)
    - (only +ve)
- ② Complement Format
  - [1's Complement, 2's Complement]

①

Magnitude  
Format



unsigned

$(0 \text{ to } 2^n - 1)$

Signed

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

②

Complement  
Format



1's Complement

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

2's Complement

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



n bit number. ① Magnitude Format  
eg 4 bit

① Unsigned Range = 0 to  $2^n - 1$  0000  $\rightarrow$  0

eg 4 bit Unsigned Range = 0 to  $2^4 - 1 \Rightarrow$  0 to 15 1111  $\rightarrow$  15

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

4 bit Signed Range =  $-2^{4-1} - 1$  to  $+2^{4-1} - 1 \Rightarrow -7$  to  $+7$



n bit number. ①. Complement format.  
eg 4 bit

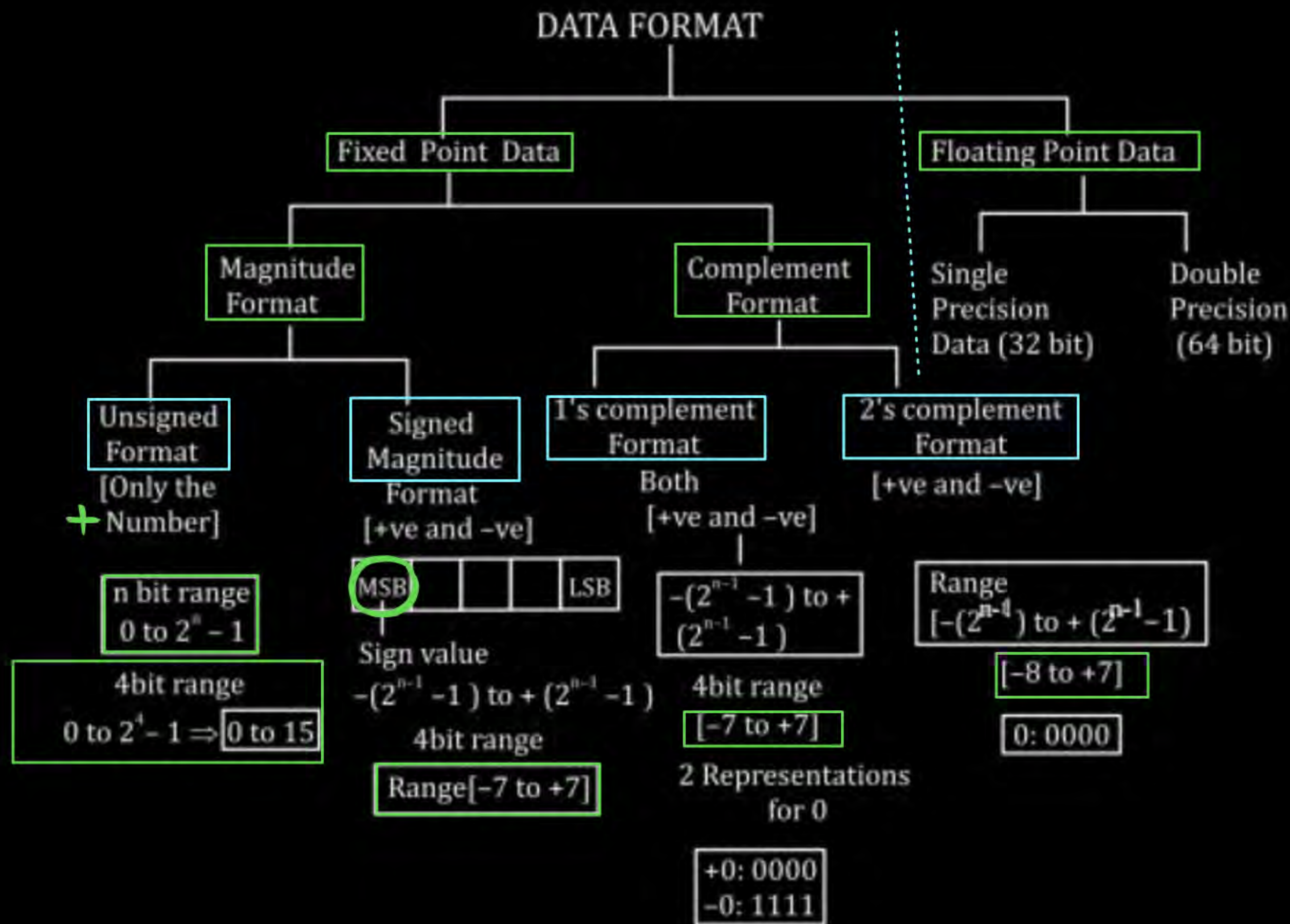
① 1's Complement Range =  $-(2^{n-1}-1)$  to  $+(2^{n-1}-1)$

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

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

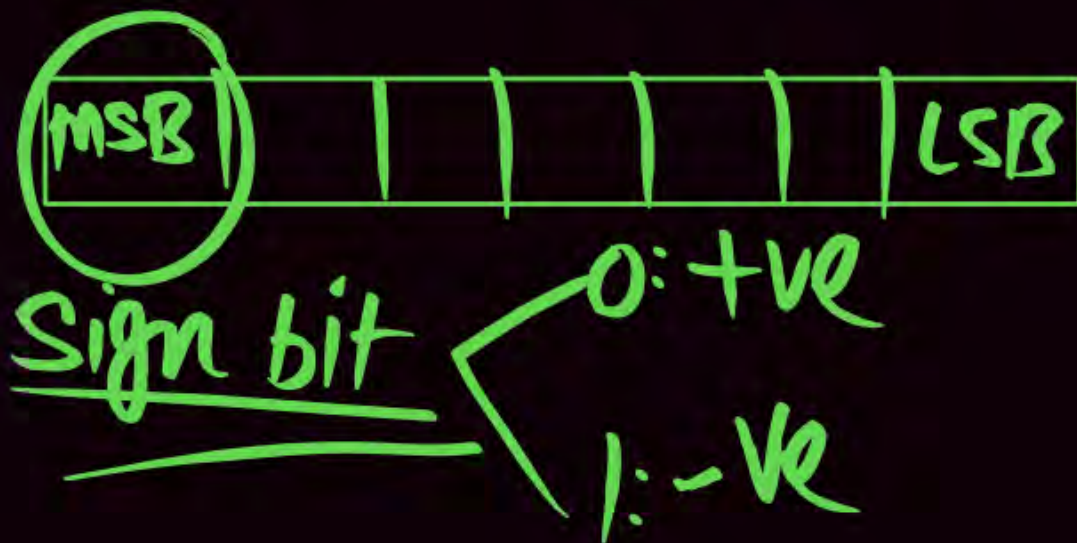
eg 4 bit 2's Complement =  $-2^{4-1}$  to  $+(2^{4-1}-1) \Rightarrow -8$  to  $+7$





Q How to Find Number is negative?

Sol<sup>n</sup> If MSB bit is 0 : +ve Number  
1 : -ve Number.





4 bit Binary	<u>Unsigned Data</u>	<u>Signed Magn. Data</u>	<u>1's Complement Data</u>	<u>2's Complement Data</u>
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'



# How to Find 1's Complement & 2's Complement:

1's Complement:

0  $\xrightarrow{\text{Converted into}}$  1  
1  $\longrightarrow$  0

1's Complement

(a)  $\underline{1000} \Rightarrow 0111 [-7]$

(b)  $1001 \Rightarrow 0110 [-6]$

(c)  $1101 \Rightarrow 0010 [-2]$

(d)  $\underline{1011} \Rightarrow 0100 [-4]$

2's Complement

1's Complement  
+ 1

2's Complement

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

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

$$\begin{array}{r} 0010 \\ + 1 \\ \hline 0011 \end{array} \Rightarrow 0011 [-3]$$

$$\begin{array}{r} 0100 \\ + 1 \\ \hline 0101 \end{array} \Rightarrow 0101 [-5]$$



Q) In Computer System Why 2's Complement  
are Used to Represent Negative Number, even  
we have Signed & 1's Complement also?

Sol<sup>n</sup>

Signed Magnitude.

+0 :	0000
-0 :	1000

1's Complement

+0 :	0000
-0 :	1111

Redundant Representation of '0' in Signed & 1's Complement  
So 2's Complement are used in Computer System.



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

③ 4 bit 2's Complement Range =  $-\left(2^{4-1}\right) \text{ to } +\left(2^{4-1}-1\right)$   
= -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



# Floating-Point Representation

16 bit fixed point data format then

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

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

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

Because range  $[-32k$  to  $+32k - 1]$

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

Floating point  
used to represent  $\left\{ \begin{array}{l} + 8.5641000000000000... [\Rightarrow \infty] \\ + 0.0000000000007892 \Rightarrow [\Rightarrow 0] \end{array} \right.$

# Floating Point Representation



S[sign]  $\begin{cases} 0 \text{ [+ve]} \\ 1 \text{ [-ve]} \end{cases}$

E/BE: Biased Exponent. | Exponent

M: Mantissa.

$$\begin{array}{c} \textcircled{\pm} \\ \downarrow \text{sign} \end{array} 0.\underbrace{\text{xxxxxxxx}}_{\downarrow \text{mantissa}} \times 2^e$$

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

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

$e \leftarrow AE$   
 $e$ : exponent.  
(Actual exponent)



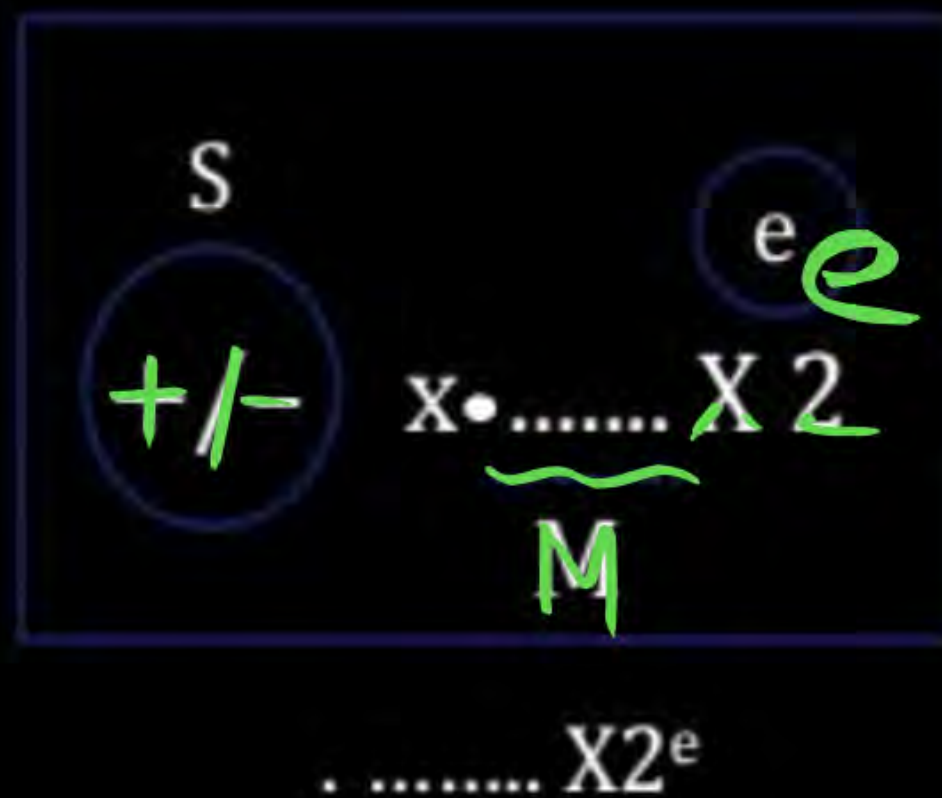
# Floating-Point Representation

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

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

E: exponent

M: Mantissa



$e$ : exponent  
Actual exponent [AE]

6.5 in Binary  $\Rightarrow 110.1$

## Number System

$$\begin{array}{cccccccccccccccc} 2^7 & 2^6 & 2^5 & 2^4 & 2^3 & 2^2 & 2^1 & 2^0 & \bullet & 2^{-1} & 2^{-2} & 2^{-3} & 2^{-4} & 2^{-5} & 2^{-6} \\ 128 & 64 & 32 & 16 & 8 & 4 & 2 & 1 & \bullet & \frac{1}{2} & \frac{1}{4} & \frac{1}{8} & \frac{1}{16} & \frac{1}{32} & \frac{1}{64} \\ & & & & & & & & & (.5) & (.25) & (.125) & (0.0625) & . & . & . \end{array}$$

$$\textcircled{1} \quad 6.5 \Rightarrow 110.1$$

$$\textcircled{2} \quad 5.0 \Rightarrow 101.0$$



16 8 4 2 1 • (.5) (.25) (.125) (0.0625)

③ 11.5  $\Rightarrow$  1011.1

④ 13.25  $\Rightarrow$  1101.01

⑤ 16.625  $\Rightarrow$  10000.101

⑥ 6.375  $\Rightarrow$  110.011

⑦ 19.25  $\Rightarrow$  10011.01  
16 8 4 2 1

$$\textcircled{8} \quad 4.75 \Rightarrow 100.11$$

$$\textcircled{9} \quad 11.625 \Rightarrow 1011.101$$

$$\textcircled{10} \quad 37 \Rightarrow \overset{32}{1} \overset{16}{0} \overset{8}{0} \overset{4}{1} \overset{2}{0} \overset{1}{1}$$

$$\textcircled{11} \quad .37 \Rightarrow 0.011\dots$$

$$\textcircled{12} \quad 13.75 \Rightarrow 1101.11$$

$$\textcircled{13} \quad 29.625 \Rightarrow 11101.101$$

$$.625 \Rightarrow 0.101$$

$$\begin{array}{r} 0.50 \\ 0.00 \\ \hline 0.125 \end{array}$$



+ (6.5)

110.1

110.1 ✓

$2^3$	$2^2$	$2^1$	$2^0$	.	$2^{-1}$	$2^{-2}$
8	4	2	1	.	(.5)	(.25)
(110.1)						

0110.1

00110.1

000110.1

0000110.1

Not

Q

110.1

+ (6.5)

+ 110.1

+/- 0.xxxx  $\times 2^e$

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

S	E	M
1bit	2bit	4bit

v.v.v.v.v.v.v.Imp.

6.5  $\Rightarrow 110.1$   
 $\Rightarrow 0.1101 \times 2^{+3}$

CROSS  
CHECK

$$\begin{aligned}
 & \Rightarrow [0.\overset{2^{-1}}{1}\overset{2^{-2}}{1}\overset{2^{-3}}{0}\overset{2^{-4}}{1}] \times 2^{+3} \\
 & [0.2^{-1} + 2^{-2} + 2^{-4}] \times 2^{+3} \\
 & \Rightarrow (2^{-1} \times 2^{+3}) + (2^{-2} \times 2^{+3}) + (2^{-4} \times 2^{+3}) \\
 & \Rightarrow 2^{+2} + 2^{+1} + 2^{-1} \\
 & \Rightarrow 4 + 2 + 0.5 = (6.5)
 \end{aligned}$$

6.5



Another method  
to check  $2^{+ve}$

⑧

$2^{-ve}$ .

110.1

Q

+ (6.5)

+ 110.1

+/- 0.xxxx  $\times 2^e$

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

2nd Technique

Databits

side

Right Alignment :

$2^{+ve}$

Left Alignment :  $2^{-ve}$

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

1bit 2bit 4bit

6.5

$\Rightarrow 110.1$

$\Rightarrow 0.1101 \times 2^{+3}$

v.v.v.v.v.v.v.Imp.

$[0.\overset{-1}{1}\overset{-2}{1}\overset{-3}{1}\overset{-4}{1}] \times 2^{+3}$

$[0.2^{-1} + 2^{-2} + 2^{-4}] \times 2^{+3}$

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

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

$\Rightarrow 4 + 2 + 0.5 = (6.5)$

6.5

CROSS CHECK



Q1  $+(6.5)$

$+110.1$

$+0.\underbrace{1101}_{\text{Mantissa}} \times 2^{+3}$

Sign = 0

M = 1101

$e = +3 \Rightarrow (11)$

S	E	M
1bit	2bit	4bit

S	e	M
0	11	1101

Ans

$$\textcircled{Q.2} + (4.5)$$

$$+ (100.1)$$

$$\Rightarrow + 0.1001 \times 2^{+3}$$

$$S: 0$$

$$e = 3 (11)$$

$$m: 1001$$

1bit	2bit	4bit
S	E	M

S	E(2bit)	M(4bit)	
0	11	1001	<u>Ans</u>





Q.3  $+ (4.75)$

$+ (100.11)$

$\Rightarrow + 0.10011 \times 2^{+3}$

$S = 0$

$e = 3 \text{ (11)}$

$M = 10011$

S	e (2bit)	Mantissa
1bit	2bit	5bit

Sign (1bit)	e (2bit)	Mantissa (5bit)
0	11	10011

Q. 1

+6.5

$$6.5 = (110.1)_2$$

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

$$S = 0 (+)$$

$$M = 1101$$

$$e = 3 = (11)_2$$

S	e	M
0	11	1101

Very. Imp

$$6.5 = 110.1$$

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

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

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

$$= 6.5$$



Q. 2

+ 4.5

100.1

$0.1001 \times 2^3$

$S = 0$  (+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.

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

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

### Right Alignment

6.5

110.1

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

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

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

$$\Rightarrow 4 + 2 + 0.5$$

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

### Left Alignment

Data: 0.0000000101  $\times 2^{+5}$

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

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

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

(Align to use upto 8 times)

Left

② WHY 'E' Required ?

'E' : Bias Exponent



⑧  $+0.001010$   
 $\Rightarrow +0.\underbrace{1010}_m \times 2^{-2}$

1bit	4bit	4bit
S	e	m

S: 0 [+ve]

m: 1010

~~e = -2~~

e = -ve

→ No Provision to tell that exponent is Negative.

Number is positive ( $\because S=0$ ) but exponent is Negative.

→ Sign bit (0) How to Deal with this Negative exponent.

if exponent is negative .

then 2's complement  
taking



⑧  $+0.001010$

$\Rightarrow +0.\underbrace{1010}_m \times 2^{-2}$

S: 0 [+ve]

M: 1010

$e = -2$

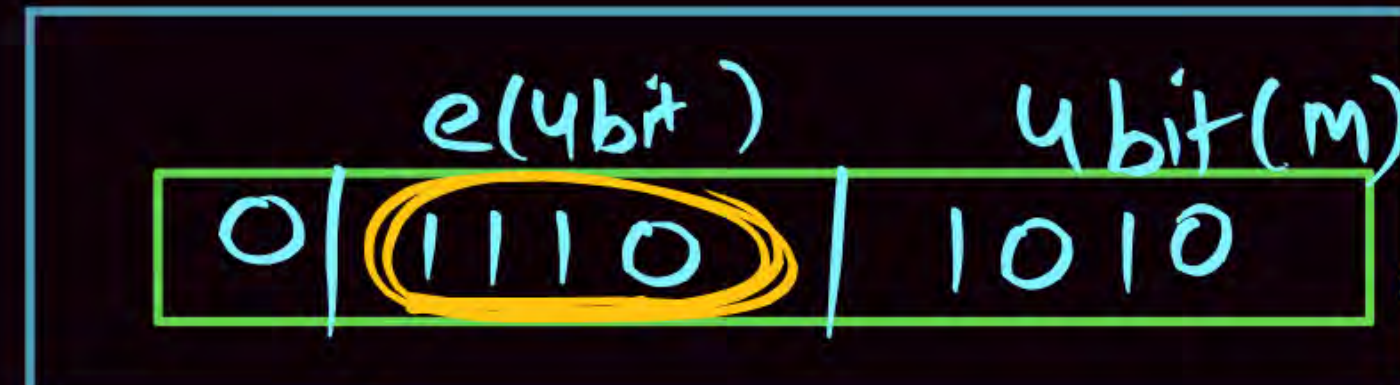
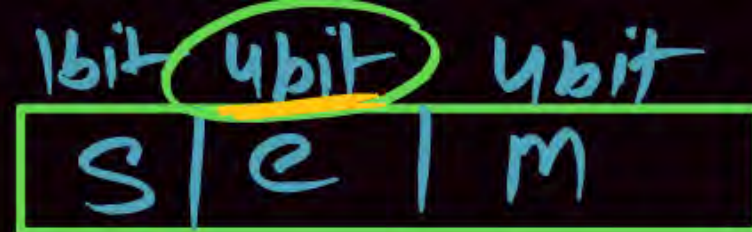
$e = 1110$

$e = -2$

2's complement So

2's complement

$0010 \Rightarrow \begin{array}{r} 1101 \\ +1 \\ \hline 1110 \end{array}$



Ans

⑧ But My Question is

How we can say

$e = 1110 \Rightarrow e = 14$

or

$e = 1110$  2's Complement of  $-2$

Here 'e' create confusion

if another example  
with having  $e=14$



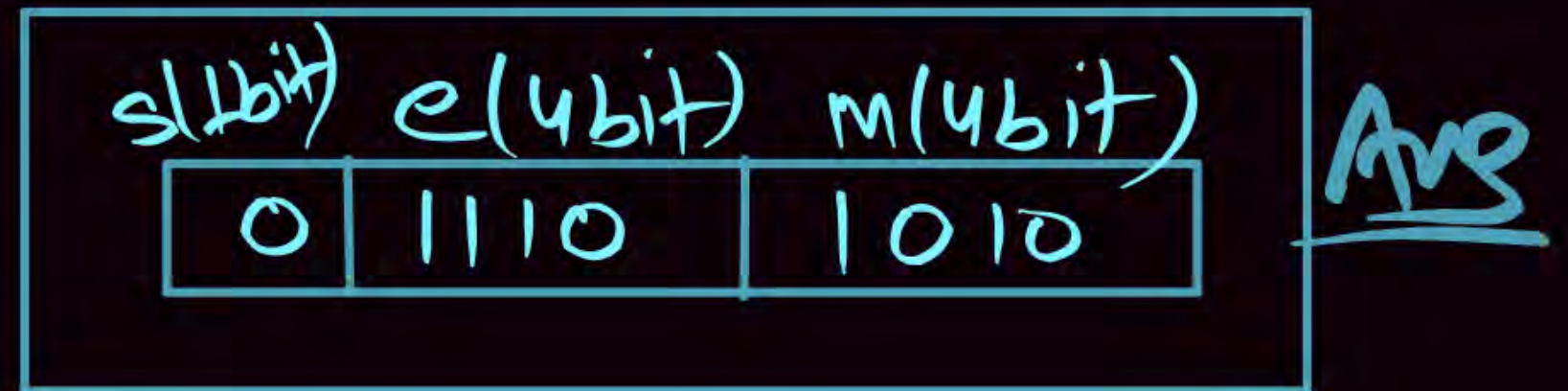
⑧  $+0.1010 \times 2^{+14}$

$S = 0$

$M : 1010$

$E = +14 \Rightarrow 1110$

$E = +14$



Q.1 Why Biasing (Biased Exponent(BE)/E) is Required?

Q.2 How bias value is decided/selected?



Q  $+0.001010$   
 $\Rightarrow +0.\underbrace{1010}_m \times 2^{-2}$

Soln)

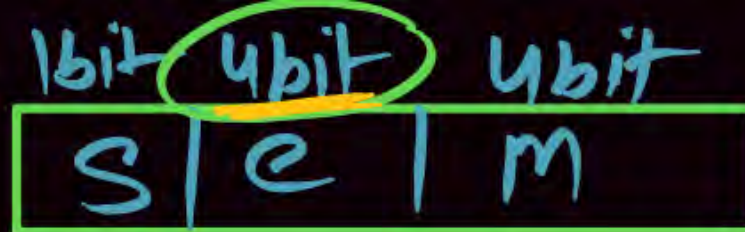
S: 0 [+ve]

M: 1010

$e = -2$

$e = -2$

$e = 1110$



Ans

Q But My Question is

2's complement.

How we can say

$e = 1110 \Rightarrow e = 14$   
 or

$e = 1110$  2's Complement of  $-2$

So here 'e' create confusion

$0010 \Rightarrow$   
 $\begin{array}{r} 1101 \\ +1 \\ \hline 1110 \end{array}$



Q

$$+0.1010 \times 2^{+14}$$

Soln

$$S = 0$$

$$M: 1010$$

$$E = +14 \Rightarrow 1110$$

$$E = +14$$

S	E	M
1bit	4bit	4bit

s(1bit)	e(4bit)	m(4bit)
0	1110	1010



Sol<sup>n</sup> 1

If we take Actual exponent  $[e]$  then there is No Provision to Represent the Negative exponent. Bcz Sign bit is telling Number +ve [0] or -ve [if  $S=1$ ].

if we take 2's complement of Negative exponent then it Creates Ambiguity. Like in Previous 2 example

$-2 : 1110$   
 $+14 : 1110$  } create Ambiguity.



So the solution is instead of writing in 2's Complement

used Bias exponent  $(E/BE)$

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

OR

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

$$0.\underbrace{\hspace{1cm}}_{\text{mantiss}} \times 2^{e(AE)}$$



Soln 2

All Number (-ve @ +ve)

biasing: To Convert a Number into '0'  
@ +ve Number.

How bias value Select: ?

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

4 bit 2's Complement =  $-(2^{4-1})$  to  $+(2^{4-1}-1)$   
=  $\boxed{-8 \text{ to } +7}$



How bias is select?

if Exponent = k bit

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

if exponent is k bit  
then bias =  $2^{k-1}$

if exponent k bit then 2's Complement =  $\left(-2^{k-1}\right)$  to  $\left(+2^{k-1}-1\right)$

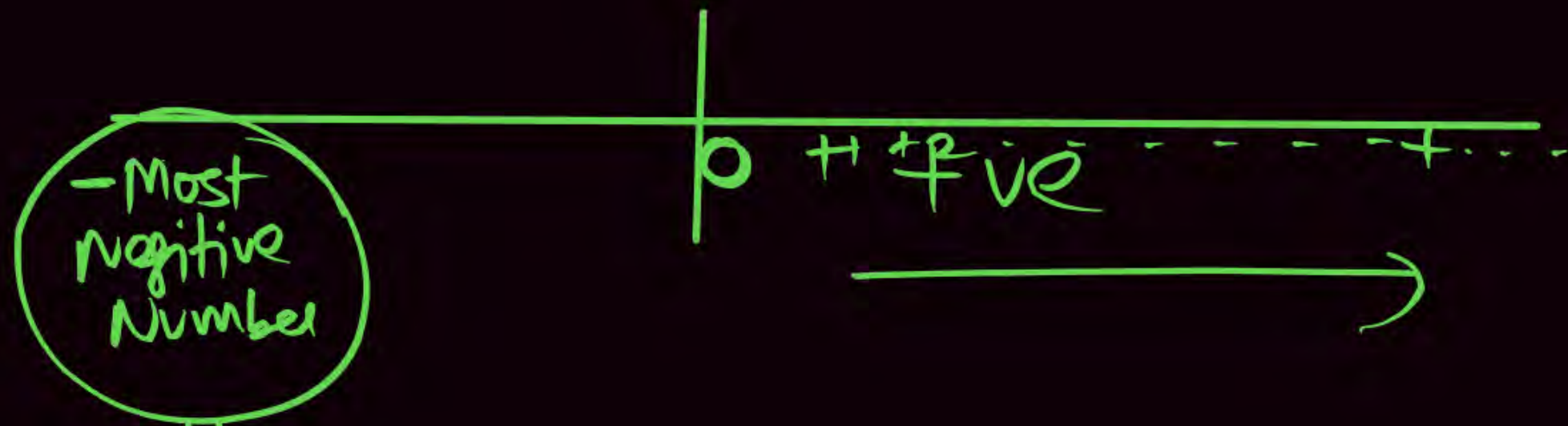
exponent 4 bit 2's Complement Range =  $-2^{4-1}$  to  $+2^{4-1}-1$

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

Note

In order to Convert ALL Number (Negative & Positive) into Positive Number take Most Highest Negative Number & ADD as a bias.





Add as bias so we get '0' @ +ve Number

Exponent = k bit

Exponent = 4 bit

$$\text{bias} = 2^{k-1} \\ \Rightarrow 2^{4-1}$$

$$\text{bias} = 8$$

$$2\text{'s Complement} = \left(-2^{\frac{n}{2}-1}\right) \text{ to } +\left(2^{\frac{n}{2}-1}\right) \\ = -8 \text{ to } +7$$

AE(e)	[E] RE [e + bias]
-8	$-8 + 8 = 0$
-7	$-7 + 8 = 1$
-6	$= 2$
-5	$= 3$
-4	$= 4$
...	
0	$0 + 8 = 8$
...	$1 + 8 = 9$
...	
+7	$7 + 8 = 15$



Q. 4

$+0.00101$

$0.101 \times 2^{-2}$

$M = 101$

$E = -2$

$S = 0$

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

$E = -2 = (1110)_2$  2's complement

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

Bit	Bit	Bit
S	E	M

S = Sign

E/BE = Exponent or

BE = bias exponent

M = Mantissa

$E = e + \text{bias}$

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

where K is exponent bits

Example

If  $K = 4$  bits

Exponent = 4 bit then

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

Exponent(K) = 4 bit

$4-1 = 8$   
bias = 2

$E = e + \text{bias}$   
 $E = e + 8$





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

$\text{bias} = 8$

$E = e + \text{bias}$

$E = e + 8$

$E = 4 \text{ bit}$

or

Excess 8 code

$2^{K-1} = 8$

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

$K - 1 = 3$

$K = 4$

$E = 4 \text{ bit}$

Excess 8 Code

↓  
 $\text{bias} = 8$

e [original exponent]		Stored exponent [BE] E $E = e + \text{bias} \Rightarrow e + 8$
-8	$-8 + 8$	0
-7	$-7 + 8$	1
-6	$-6 + 8$	2
-5	$-5 + 8$	3
-4	$-4 + 8$	4
-3	$-3 + 8$	5
-2	$-2 + 8$	6
-1	$-1 + 8$	7
0	$0 + 8$	8
1	$1 + 8$	9
2	$2 + 8$	10
3	$3 + 8$	11
4	$4 + 8$	12
5	$5 + 8$	13
6	$6 + 8$	14
7	$7 + 8$	15

Exponent =  $k$  bit

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

Excess - 8 : bias = 8

Excess - 16 : bias = 16

Excess - 32 : bias = 32

Excess - 64 : bias = 64

Excess : 8

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

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

$$k-1 = 3$$

$$\boxed{k=4} \text{ Exponent} = 4 \text{ bit}$$

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



Excess - 8 : bias = 8  
Excess - 16 : bias = 16  
Excess - 32 : bias = 32  
Excess - 64 : bias = 64.

Exponent = k bit

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

$$\text{Excess} : 16$$

$$2^{k-1} = 2^4$$

$$2^{k-1} = 2^4$$

$$k-1 = 4$$

$$k = 5$$

Exponent = 4 bit

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

Excess - 8 : bias = 8  
Excess - 16 : bias = 16  
Excess - 32 : bias = 32  
Excess - 64 : bias = 64.

Exponent =  $k$  bit

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

$$\text{Excess} : 32$$

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

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

$$k-1 = 5$$

$$k = 6 \quad \text{Exponent} = 6 \text{ bit}$$

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



Q.

From previous question

0.00101

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

$$M = 101$$

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

$$\text{Bias} = 16$$

$$e = -2$$

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

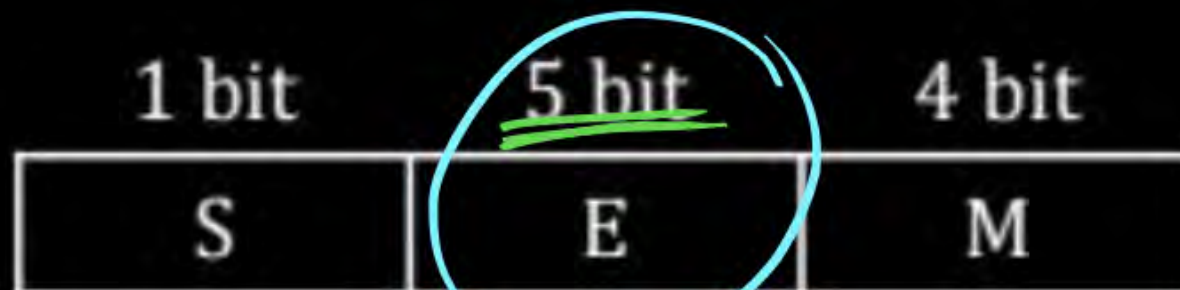
$$E = -2 + 16$$

$$E = 14$$

$$E = (01110)_2$$

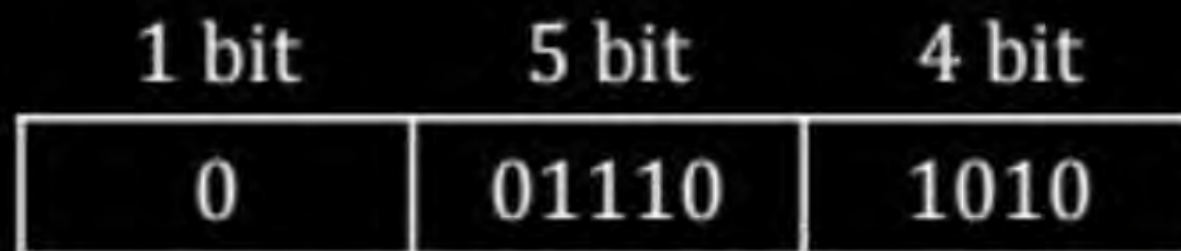
$$\text{Formula: } (-1)^S \times 0.M \times 2^E$$

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



$$\text{bias} = 2^{5-1} \Rightarrow 2^4 \Rightarrow \text{bias} = 16$$

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



Ans

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

$$0.000101 \text{ Ans}$$



1 bit

x bit

y bit

Normalized Mantissa

### Explicit Normalized Syntax

$$\frac{0.1..... \times 2^E}{M}$$

Formula to get number  
[value formula]

$$(-1)^S \times 0.M \times 2^E$$

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

### Implicit Normalized Syntax

$$\frac{1..... \times 2^E}{M}$$

Formula to get number  
[value formula]

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

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





## Explicit

0.1 After the point,

Immediate first bit should be 1

### Example

(101.11)

$0.10111 \times 2^3$

$M = 10111,$

$e = 3$

$E = e + \text{bias}$

## Implicit

Before the point 1 means 1. ....

### Example

(101.11)

$1.0111 \times 2^2$

$M = 0111,$

$e = 2$

$E = e + \text{bias}$

**Q. 1**  $+(4.875)$  format

Then do explicit and implicit normalization

**Explicit**

$(+4.875)$

100.111

$0.100111 \times 2^3$

$M = 100111$

$e = 3, \text{ bias} = 2^{4-1}$

$E = 3 + 8$

$E = 11$

$E = 1011$

1 bit	4 bit	5 bit
0	1011	10011

**Value Formula:**  $(-1)^s \times 0.M \times 2^e$

$(-1)^0 \times 0.10011 \times 2^{11-8}$

$0.10011 \times 2^3$

100.11

4.75

(Not getting very accurate)





Implicit

(+4.875)

100.111

$1.00111 \times 2^2$

$M = 00111$

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

$E = 2 + 8$

$E = 10$

$E = 1010$

1 bit	4 bit	5 bit
0	1010	00111

Value Formula:  $(-1)^S \times 1.M \times 2^e$

$(-1)^0 \times 1.00111 \times 2^{10-8}$

$1.00111 \times 2^2$

100.111

4.875

(Getting very accurate)



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?





# Solution

-29.75

-11101.11

$0.1110111 \times 2^5$

M: 1110111

$e = 5$

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

$\text{bias} = 32$

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

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

Q.

+21.75

Implicit?

10101.11

$1.010111 \times 2^4$

$M = 010111$

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

$E = 4 + 64$

$E = 68 = (1000100)_2$

1 bit

7 bit

8 bit

S

E

M

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





**THANK  
YOU!**

