

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.

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

S: Sign
0 (+ve)
1 (-ve)

$$E @ BE = e + bias$$

$$BE = AE + bias$$

S	E		m.
---	---	--	----

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

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

Explicit

$$0.1 \dots \times 2^e$$

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

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

Implicit

$$1. nnnnxx2^e$$

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

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

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

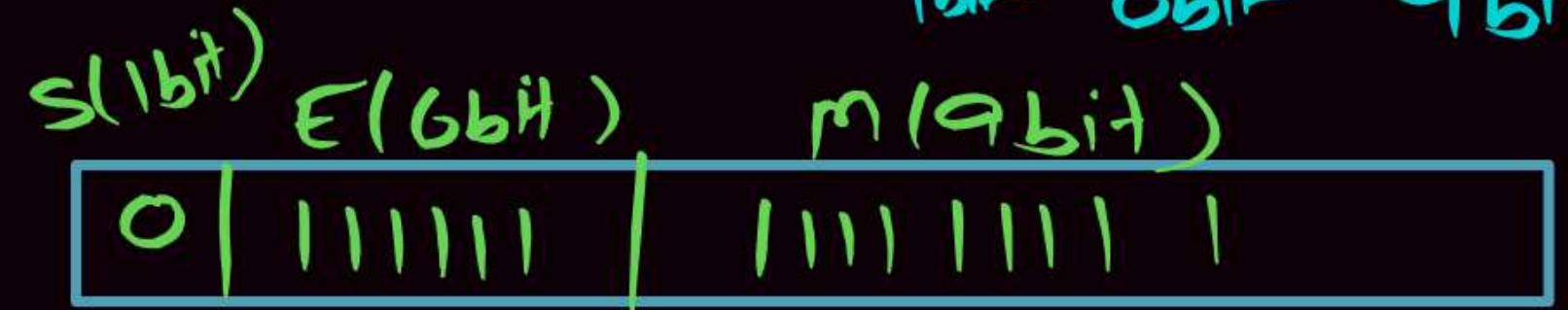
$$(-1)^s 1.m \times 2^{e-\text{bias}}$$

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

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

GP series
yesterday class

⑥ Max +ve Value
Using Explicit



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

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

$$(-1)^S \cdot 0.11111111 \times 2^{63-32}$$

$$+ 0 \cdot \underbrace{11111111}_{9 \text{ times}} \times 2^{31}$$

$$11111111.0 \times 2^{-9} \times 2^{31}$$

Left Alignment



Conventional &
Explicit form

Max +ve Value

$$\text{bias} = 2^{6-1} = 2^5$$

$$\text{bias} = 32$$

$$E = 111111 \Rightarrow E = 63$$

$$\begin{aligned}
 & M: 111111111 \\
 & \underline{111111111} \times 2^{22} + 2^{31} - 2^{22} \\
 & (511) (2-1) \times 2^{22} \\
 & 2^{9+22} - 2^{22} = \underline{\underline{+2^{31}} \text{ Any}}
 \end{aligned}$$

Alternate

$$0.1111111 \times 2^{31}$$



$$\left(1 - \frac{1}{2^9}\right) \times 2^{31}$$

$$2^{31} - 2^{31-9} = +2^{31} \text{ Amp}$$



Normalized Mantissa

1 bit x bit y bit

Explicit Normalized Syntax

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

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

$$\frac{1.1 \dots \times 2^e}{M}$$

Formula to get number
[value formula]

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

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

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

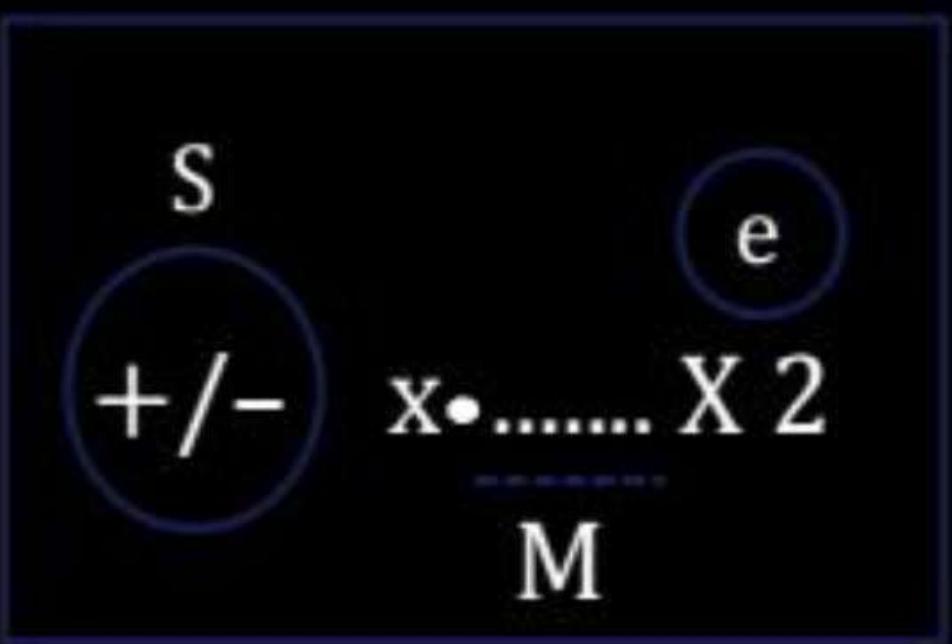
E: Biased exponent

M: Mantissa

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

or

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



$$\dots \times 2^e$$

Q.

+21.75

1 bit	7 bit	8 bit
S	E	M

Implicit?

10101.11

$$1.010111 \times 2^4$$

$$M = 010111$$

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

$$E = 4 + 64 \quad \text{bias} = 64$$

$$E = \underline{68} = (1000100)_2$$

Value Formula:

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

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

$$1.010111 \times 2^4$$

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

Ans

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

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

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?

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

-29.75

$$-1.11011 \times 2^5$$

Sign = 1

$$M: 1110111 \quad 0.1110111 \times 2^{+5}$$

e=5

bias = 32

$$E = 5 + 32 - 37 \Rightarrow 100101$$

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

$$k-1 = 5$$

k=6

Solution

-29.75

-11101.11

0.1110111×2^5

M: 1110111

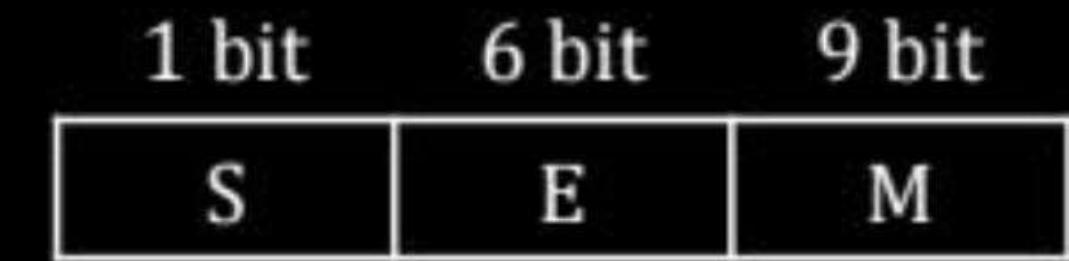
e = 5

bias = 2^{6-1}

bias = 32

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

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



bias = 32

$$(-1)^S \cdot M \times 2^{\frac{E - bias}{2}} \Rightarrow (-1)^S \cdot M \times 2^{\frac{37 - 32}{2}}$$

$$(-1)^L \cdot 0.111011100 \times 2$$

$$-0.111011100 \times 2^{+5}$$

-11101.1100

(-29.75) Ans

Q.

+21.75

Implicit?

1 bit	7 bit	8 bit
S	E	M

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:

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

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

$$1.010111 \times 2^4$$

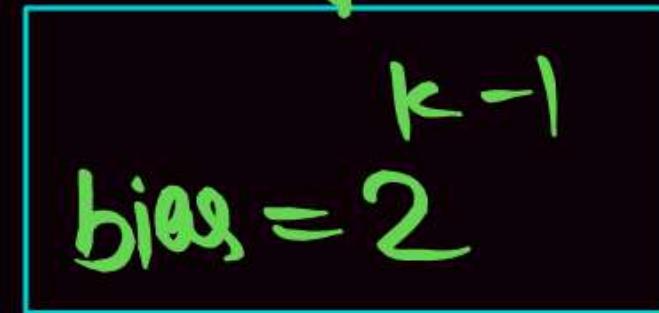
$$10101.11 = (21.75)_{10}$$

Ans

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

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

Conventional Representation.



Can not Represent '0' & ∞ .

Implicit \Rightarrow 1. Something

Represent
1. $n \Rightarrow$ Not '0'.

Explicit \Rightarrow 0.1....

Represent.
0.1 \Rightarrow Not '0'

In CN

IEEE 802.3 Ethernet

IEEE 802.5 Token Ring

IEEE 802.11 WiFi

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

Default is Implicit.

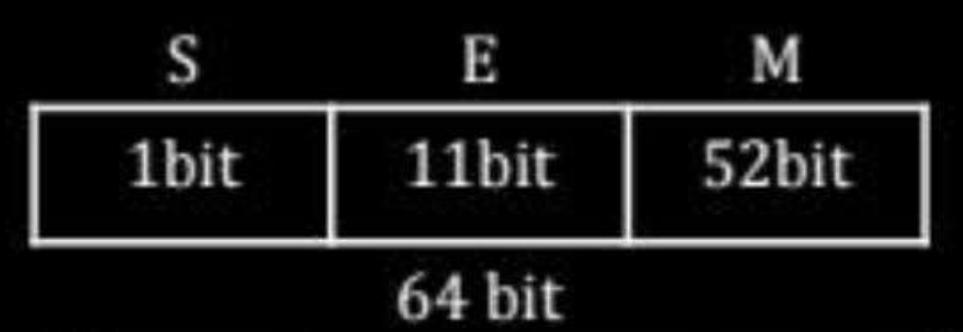
IEEE 754 Floating Point Representation

Single Precision
(32 bit)
Excess 127

Double Precision
(64 bit)
Excess 1023



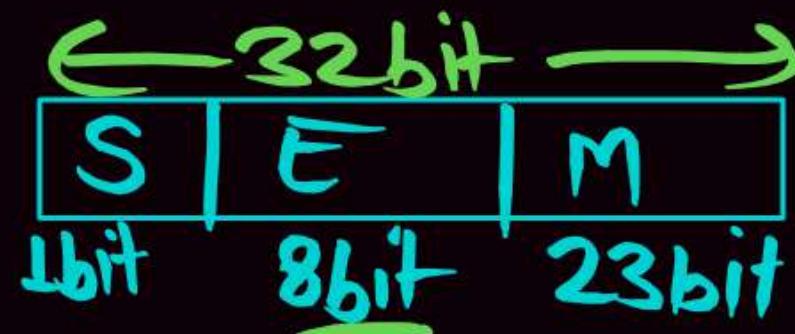
$$\begin{aligned} \text{bias} &= 2^{K-1} - 1 \\ &= 2^8 - 1 \\ \text{Bias} &= 127 \end{aligned}$$



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

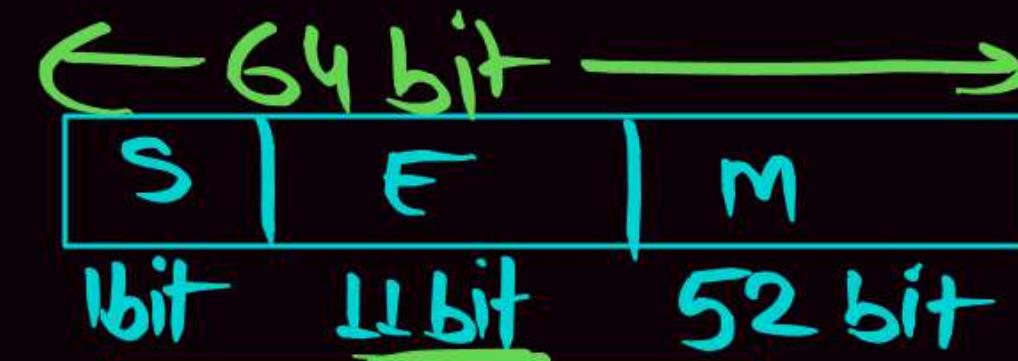
IEEE 754

Single
Precision
(32bit)



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

Double
Precision
(64 bit)



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

S	E	m
1bit	8bit	23bit 8-1

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

$$\text{bias} = 127$$

Excess - 127

Single Precision

S	E	m
1bit	11bit	52bit

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

$$\text{bias} = 1023$$

Excess - 1023.



Double Precision

111

$$\begin{array}{ccccccc} 2^7 & 2^6 & 2^5 & 2^4 & 2^3 & 2^2 & 2^1 & 2^0 \\ 128 & 64 & 32 & 16 & 8 & 4 & 2 & 1 \end{array}$$

$$A = 10$$

$$B = 11$$

$$C = 12$$

$$D = 13$$

$$E = 14$$

$$F: 15$$

133:

10000101

Next Term is
Add of all
Previous term + 1

$2^7 \quad 2^6 \quad 2^5 \quad 2^4 \quad 2^3 \quad 2^2 \quad 2^1 \quad 2^0$
128 64 32 16 8 4 2 1

$\frac{32 \quad 16 \quad 8 \quad 4 \quad 2 \quad 1}{100} = 3$
 $100 = 4$
 $111 = 7$
 $1000 = 8$
 $1111 = 15$
 $10000 = 16$
 $11111 \rightarrow 31$
 $10.0000 \rightarrow 32$

133: 1 0 0 0 0 1 0 1

Q. 1

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

P
WInfotechSingle Precision. $(+119)$ $+110111 \times 2^0$ $+1.110111 \times 2^{+6}$ $e = +6$ $S = 0$ $M = 110111$ $bias = 127$ $E = e + bias = 6 + 127$ $E = 133$ 10000101

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

1bit 8bit 23bit

S(1bit)	E(8bit)	Mantissa (23bit)
0	10000101	11011100000000000000000

$(42EE0000)_F$ Ans

Value formula

 $(-1)^S \underline{1 \cdot M} \times 2^E \rightarrow (-1)^S 1 \cdot M \times 2^{E - bias}$
 $E = 133$ $bias = 127$
 $(-1)^0 1 \cdot 11011100000000000000000 \times 2^{+6}$
 $1 \cdot 11011100000000000000000 \times 2^{+6}$
 $110111 \cdot 00000000000000000000000 \dots$

(119) Ans

1029

E \Rightarrow 11 bit

$$\begin{array}{cccccccccc} 2^{10} & 2^9 & 2^8 & 2^7 & 2^6 & 2^5 & 2^4 & 2^3 & 2^2 & 2^1 & 2^0 \\ 1024 & 512 & 256 & 128 & 64 & 32 & 16 & 8 & 4 & 2 & 1 \end{array}$$

$$1029 \Rightarrow 10000000101$$

A: 10

B: 11

C: 12

D: 13

E: 14

F: 15

Q.

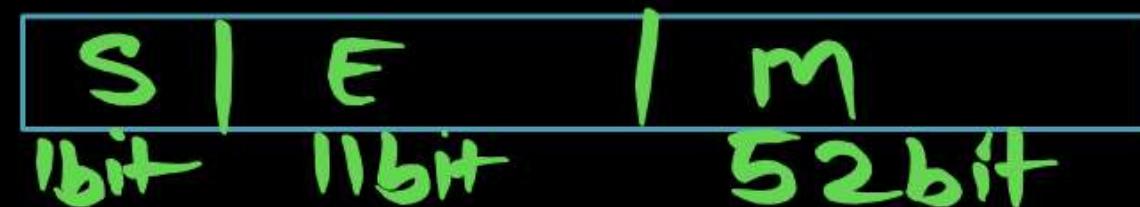
How to represent +(119) into IEEE 754 single precision &

P
W

Double precision floating Point Representation?

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

Double Precision



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

+ 119

$$+ 111011 \times 2^0$$

L.LLO111x2⁺⁶

$$S=0$$

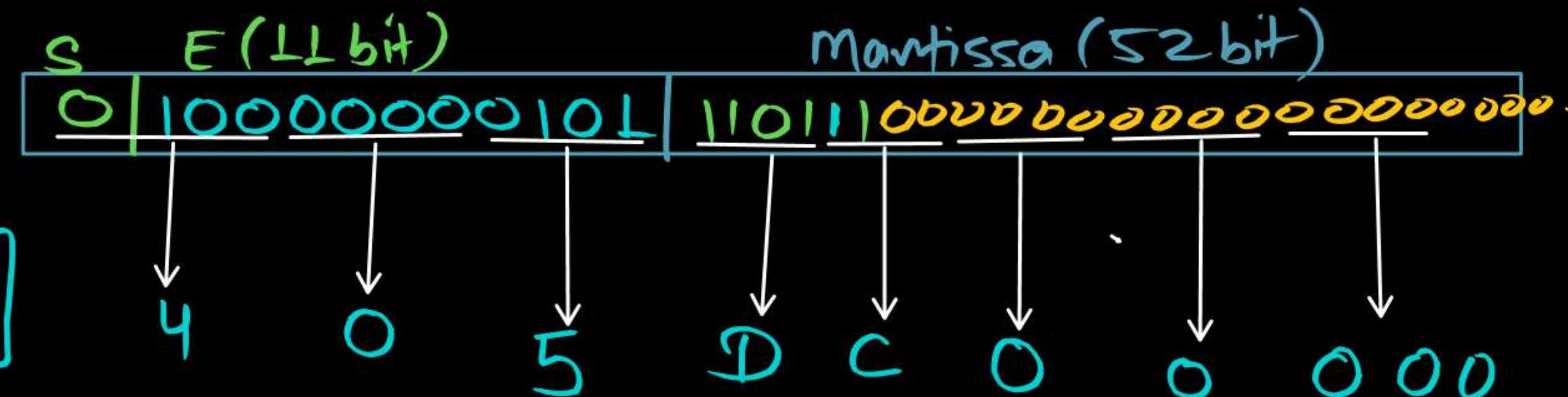
$$m = 11011$$

$$e = +6$$

bias = 1023

$$E = 6 + 1023 = 1029$$

$$E = 100000000101$$



$(405D\underset{(16)}{000000000000})_{H_j}$, Aug

Q.1 How to represent +(119) into IEEE 754 single precision & Double precision floating Point Representation? [Ans. d/c]

PW

Double Precision

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

$$(-1)^S \cdot m \times 2^{E_{\text{bias}}}$$

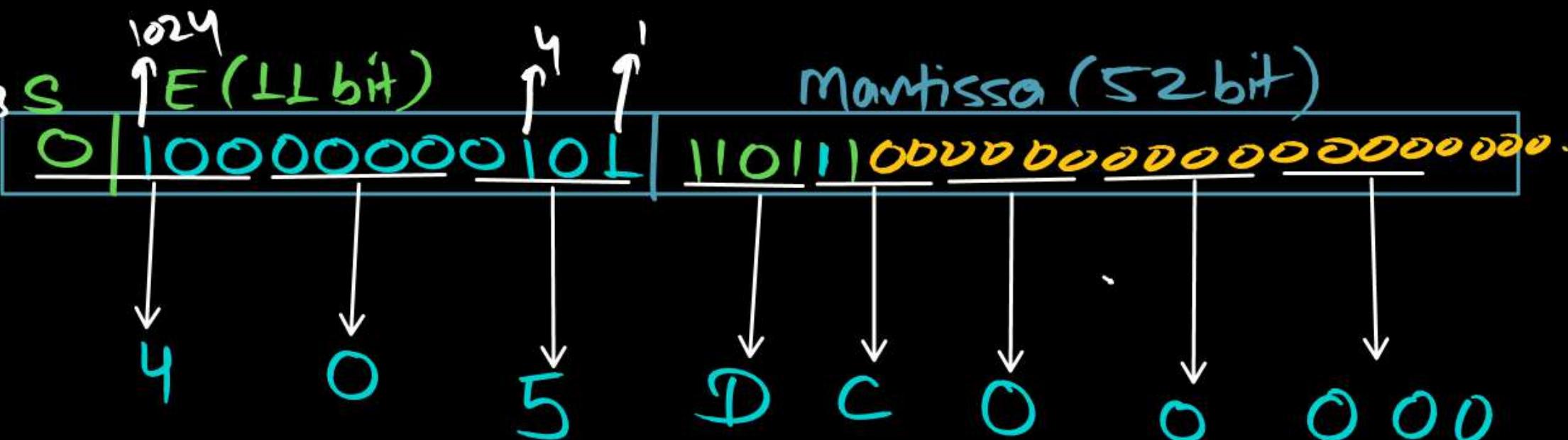
1110111.0000000000X2

(119) Aug

E-1029

The diagram illustrates a 64-bit memory address format. It is divided into three fields: S (1bit), E (11bit), and M (52bit). The fields are represented by boxes of increasing width from left to right, separated by vertical lines.

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



Q.

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

$$E=132$$

$$\text{bias} = 127$$

Sign(1bit)	Exponent(8bit)	Mantissa(23bit)
0	10000100	11000000000000000000000 0000

- A 48
- B 56
- C 64
- D 60

Ans (B)

$$\begin{aligned}
 & (-1)^{\text{Ex}} \cdot 1 \cdot M \times 2^{E-\text{bias}} \\
 & (-1)^0 \cdot 1 \cdot 11000000000000000000000 \times 2^{132-127} \\
 & \Rightarrow +1 \cdot 11000000000000000000000 \times 2^{+5} \\
 & \Rightarrow 11100000000000000000000 \\
 & \quad -756
 \end{aligned}$$

Solution

Sign(1bit)	Exponent(8bit)	Mantissa(23bit)
0	10000100	1100000000000000000 0000

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

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

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

$$= 111000$$

$$= 56 \text{ Ans}$$

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

$$\text{bias} = 127$$

$$E - \text{bias}$$

$$E = 10000100$$

$$E = 132 \leftarrow$$

Q.

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?

s	e (8bit)	m (23bit)
-----	------------	-------------

A

416C0000H

B

41640000H

-14.25

C

C16C0000H

D

C1640000H

-1110.01

 1.11001×2^3
 $e = +3$
 $bias = 127$, $\epsilon = 1/30$

Ans(D)

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

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

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

$$= 3 + 127$$

$$E = 130 \Rightarrow 10000010$$

S	E	M
1	10000010	11001 0000000000000

↓ ↓ ↓ ↓ ↓ ↓
 C E 6 4 0 D

$$\text{Hexadecimal} = (C1640000)_H$$

IEEE 754 Floating Point Representation



Represent

① '0'



② ∞



③ NaN.

④ etc

IEEE 754 Floating Point Representation

Single Precision

(32 bit)

Excess 127



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

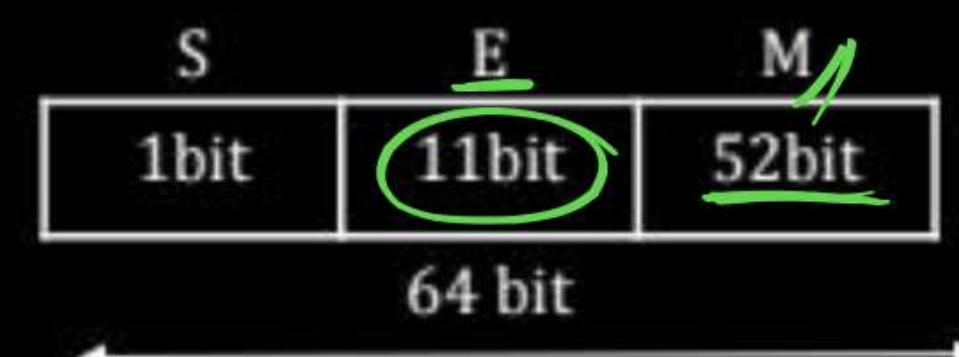
$$= 2^8 - 1$$

$$\text{Bias} = 127$$

Double Precision

(64 bit)

Excess 1023



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

$$\text{Bias} = 1023$$

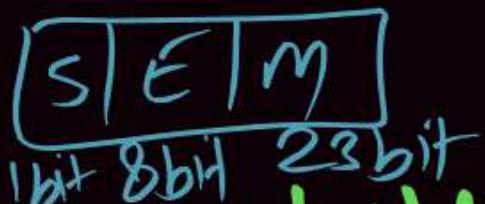
In Conventional Representation

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

But In IEEE 754 Representation

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

WHY? (P.T.O)



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

In Single Precision

Because If we take bias = $128(2^{8-1})$ Excess 128, then

there is a chance of getting 'E' is 255. (when $E=255$)
 then either ∞ or Nan

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

Assume if we have $e=127$

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

$$E = e + \text{bias} \quad (\text{if } \text{bias} = 128)$$

$$= 127 + 128$$

$$= 255$$

So bias will be 2^{k-1}

③ In Single Precision = $2^{8-1} - 1 = 127$

Excess = 127

11111111 → 255.

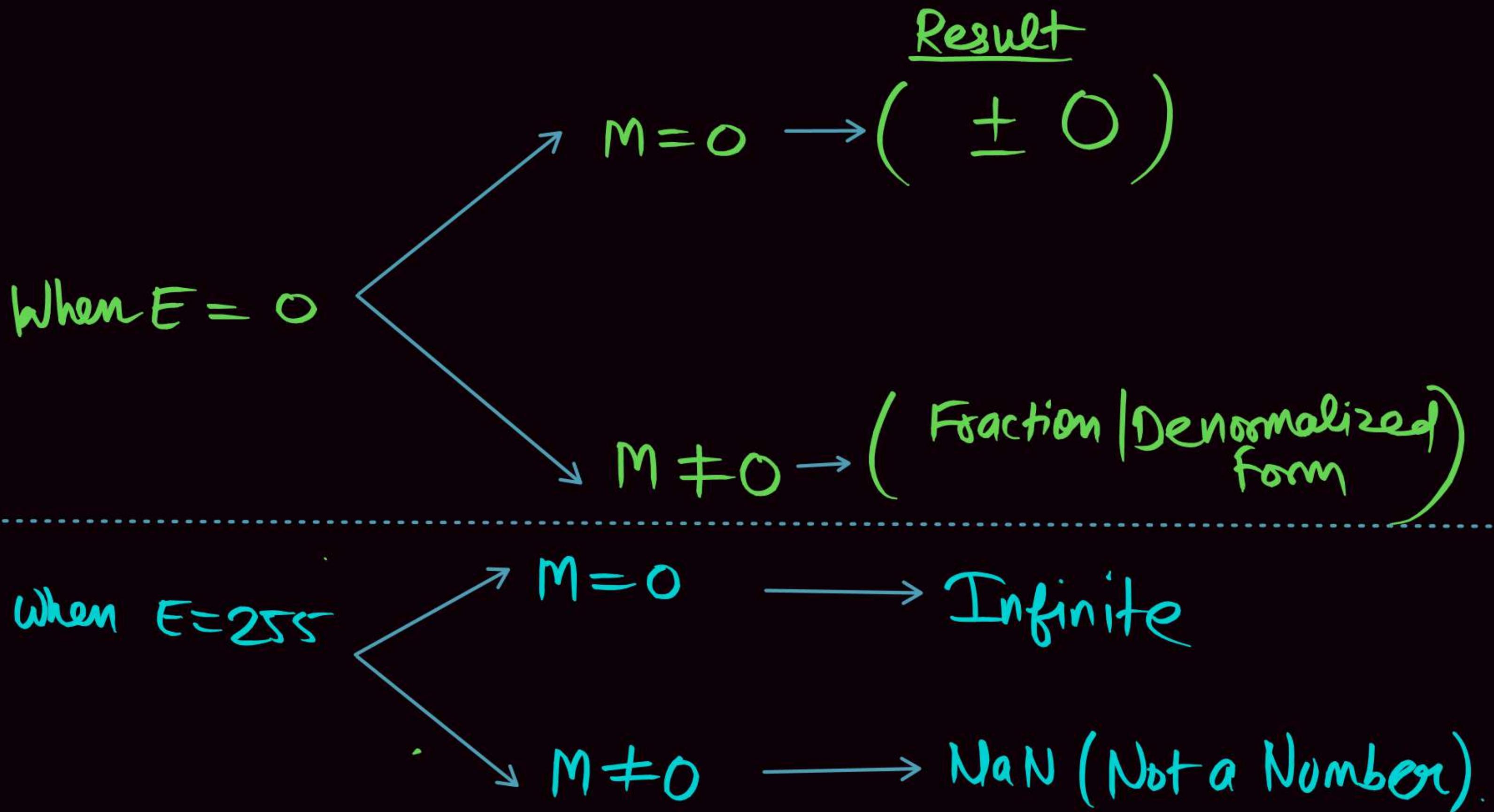
Single Precision (32 bit)

S	E	M
1 bit	8 bit	23 bit

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

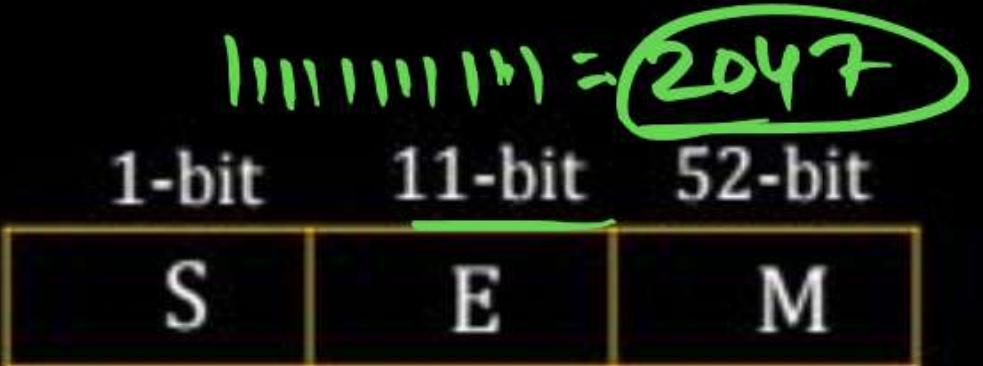
L.L.V.V. Impl

	Sign(1 bit)	E(1 bit)	M(23 bit)	Value
①	0 or 1	00000000 <u>E = 0</u>	0000000000000000 <u>M = 0</u>	<u>±0</u>
②	0 or 1	11111111 <u>E = 255</u>	0000000000000000 <u>M = 0</u>	<u>±∞</u> Infinity
③	0 or 1	1 ≤ E ≤ 254 1 to 254	M = <u>Any</u> <u>0</u>	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	<u>E = 0</u>	<u>M ≠ 0</u>	Denormalized number/Fractional form $(-1)^S \times 0.M \times 2^{E-127 \text{ bias}}$
	0 or 1	<u>E = 255</u>	<u>M ≠ 0</u>	Not a Number (NAN)



Double Precision

Excess - 1023



$$1111111111 = 2047$$

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

$$\text{bias} \Rightarrow 1023$$

P
W

Sign (1 bit)	E(11 bit)	M(52 bit)	Value
0 or 1	0000 0000 000 <u>E = 0</u>	000000000000.. <u>M = 0</u>	<u>± 0</u>
0 or 1	<u>1111 1111 111</u> <u>E = 2047</u>	0000000000.. M = 0	$\pm \infty$ (<u>infinity</u>)
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 $(-1)^S 0.M \times 2^{E-1023}$
	E = 2047	M $\neq 0$	Not a number

NOTE: When E = 0 then Value 0

when M = 0

or

fractional form

when M ≠ 0

NOTE: When E = 2047

then Value ∞

when M = 0

or

Not a Number(NAN)

when M ≠ 0

Note

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



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.

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

$$1.0 \times 2^0$$

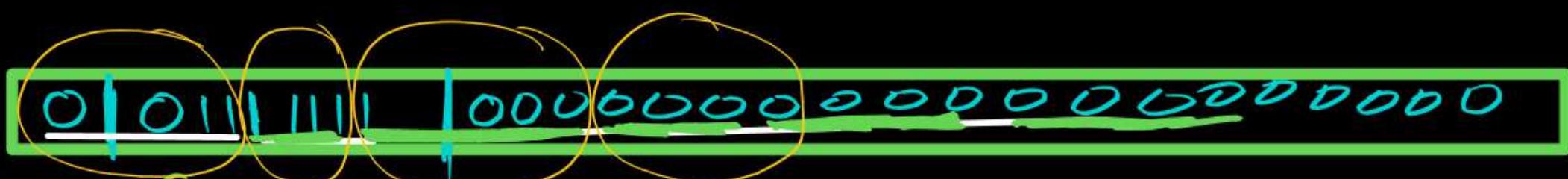
$$e=0$$

$$\text{bias} = 127$$

$$E = e + \text{bias} = 0 + 127$$

$$E = 127$$

$$E = 0111111$$



$$(3F800000)_{16} \text{ Ans}$$

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

$$(-1)^0 1.0000000 \times 2^{127-127} \Rightarrow +1.0000000 \times 2^0 \\ = (1.0) \text{ Ans}$$

$$M = 0$$

$$S = 0$$

Sol.

To represent +(1.0) into IEEE 754 single precision floating Point P
Repartition.. W



$$+ 1.0 \times 2^0$$

$$M = 0$$

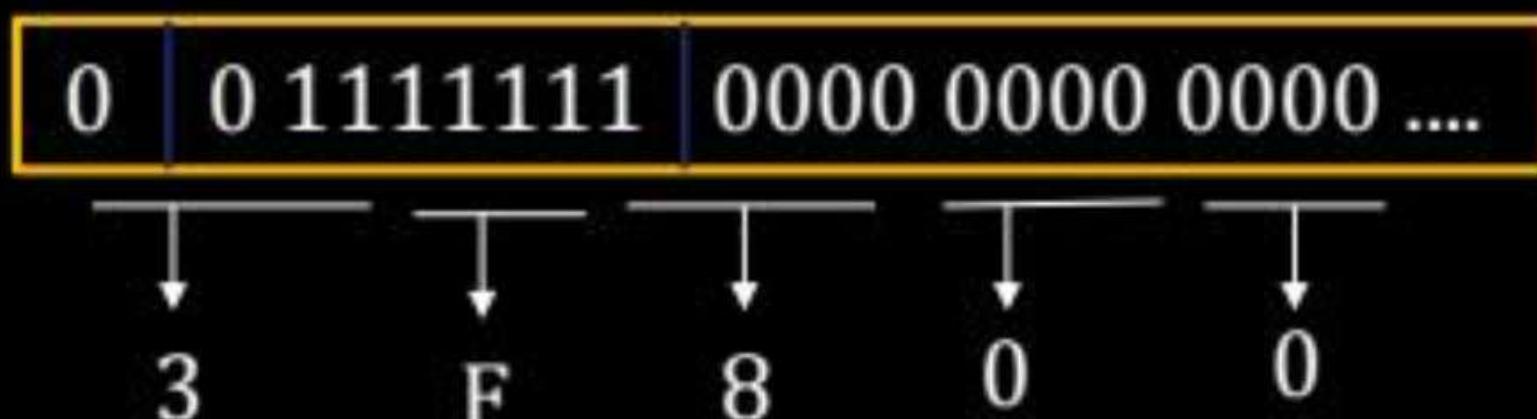
$$e = 0$$

$$E = 127$$

S	E	M
1	8	23

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

$$\begin{aligned}E &= e + \text{bias} \\E &= 0 + 127\end{aligned}$$



(3F80 0000)

$$2^3 - 1$$

↑

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

'7'

$\overline{1111}$ then $2^4 - 1$

4bit

Is 24 bit then $2^{24} - 1$

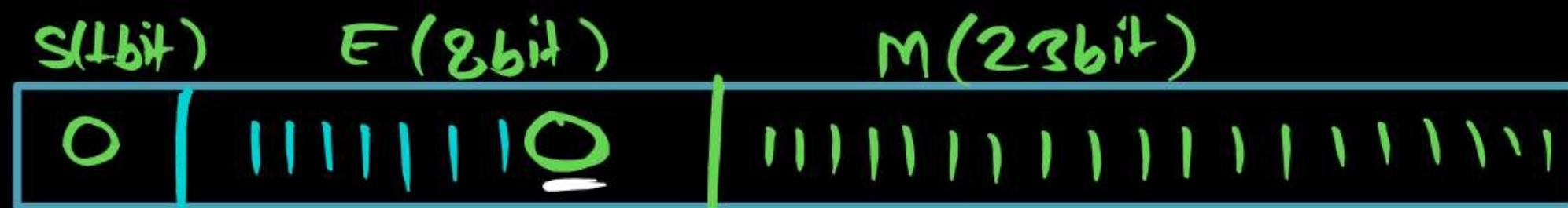
$$\overbrace{\overline{1111111111111111111111}}^{24 \text{ times } 1's} \times 2^{104}$$

$$(2^{24} - 1) \times 2^{104}$$

$$+ 2^{128} \text{ Ans}$$

Q.

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



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

$$S=0$$

$$E=1111110 = 254$$

$$\text{bias} = 127$$

$$\text{Mantissa} = 1111111111111111111$$

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

$$(-1)^0 \cdot \overbrace{1 \cdot 1111111111111111111}^{23 \text{ times}} \times 2^{254-127}$$

$$1 \cdot \overbrace{1111111111111111111}^{23 \text{ times}} \times 2^{127}$$

$$\frac{2^{24}+1}{2} \times 2^{104} + 2^{128} \text{ Ans}$$

$$(2^{24}-1) \times 2^{104} \times 2^{-23} \times 2^{127}$$

$$\frac{127}{2^{-23}} \times 104$$

Q. L

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

1100 0000 0011 1101 1000 0000000000....

P
W

Q.2

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	0000000000000

Here all exponent bits is 1

E = 2047

→ Infinity (when M = 0)

→ Not a number (when M ≠ 0)

Here M = 0 so its +∞

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

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.

1. L

1 + (-L)

④

1. LL

(L · 75)

⑤

10 · 11

(2 · 75)

①

1. LL \Rightarrow 1 + (-11)

1 + (0.5 0.25)

(1.75)

②

10 · 11

10 + 2^1 2^{-2}

2 + ·5 ·25
(2.75)

Floating-Point Representation

Example:

Show the IEEE 754 binary representation of the number -0.75_{ten} in single and double precision.

Answer:

The number -0.75_{ten} is also.

$$\begin{aligned} & -0.75 \\ & -(0.\underline{11}) \end{aligned}$$

$$-3/4_{\text{ten}} \text{ or } -3/2^2_{\text{ten}}$$

Single Precision

It is also represented by the binary fraction.

$$-0.11_{\text{two}}/2^2_{\text{ten}} \text{ or } -0.11_{\text{two}}$$

In scientific notation, the value is.

$$-0.11_{\text{two}} \times 2^0$$

And in normalized scientific notation, it is.

$$-1.1_{\text{two}} \times 2^{-1}$$

The general representation for a single precision number is.

$$(-1)^S \times (1 + \text{Fraction}) \times 2^{(\text{Exponent} - 127)}$$

When we subtract the bias 127 from the exponent of -1.1_{two} × 2⁻¹, the result.

$$(-1)^S \times (1 + .1000\ 0000\ 0000\ 0000\\ 0000\ 000_{\text{two}}) \times 2^{(126 - 127)} \quad |+ \cdot 1$$

The single precision binary representation of -0.75_{ten} is then



31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1
1	0	1	1	1	1	1	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

1 bit 8 bits 23 bits

Floating-Point Addition

Let's add numbers in scientific notation by hand to illustrate the problems in floating-point addition:

$$9.999_{\text{ten}} \times 10^1 + 1.610_{\text{ten}} \times 10^{-1}.$$

Assume that we can store only four decimal digits of the significand and two decimal digits of the exponent.

Step 1. To be able to add these numbers properly, we must align the decimal point of the number that has the smaller exponent. Hence, we need a form of the smaller number, $1.610_{\text{ten}} \times 10^{-1}$, that matches the larger exponent.

We obtain this by observing that there are multiple representations of an unnormalized floating-point number in scientific notation:

$$1.610_{\text{ten}} \times 10^{-1} = 0.1610_{\text{ten}} \times 10^0 = 0.01610_{\text{ten}} \times 10^1$$

The number on the right is the version we desire, since its exponent matches the exponent of the larger number, $9.999_{\text{ten}} \times 10^1$. Thus, the first step shifts the significand of the smaller number to the right until its corrected exponent matches that of the larger number. But we can represent only four decimal digits so, after shifting, the number is really:

$$0.016_{\text{ten}} \times 10^1$$

Step 2. Next comes the addition of the significands:

$$\begin{array}{r} 9.999_{\text{ten}} \\ + 0.016_{\text{ten}} \\ \hline 10.015_{\text{ten}} \end{array}$$

The sum is $10.015_{\text{ten}} \times 10^1$.

Step 3.

:This sum is not in normalized scientific notation, so we need to adjust it

$$10.015_{\text{ten}} \times 10^1 = 1.0015_{\text{ten}} \times 10^2$$

Thus, after the addition we may have to shift the sum to put it into normalized form, adjusting the exponent appropriately. This example shows shifting to the right, but if one number were positive and the other were negative, it would be possible for the sum to have many leading 0s, requiring left shifts. Whenever the exponent is increased or decreased, we must check for overflow or underflow—that is, we must make sure that the exponent still fits in its field.

Step 4.

Since we assumed that the significand can be only four digits long (excluding the sign), we must round the number. In our grammar school algorithm, the rules truncate the number if the digit to the right of the desired point is between 0 and 4 and add 1 to the digit if the number to the right is between 5 and 9. The number

$$1.0015_{\text{ten}} \times 10^2$$

Is rounded to four digits in the significant to

$$1.002_{\text{ten}} \times 10^2$$

since the fourth digit to the right of the decimal point was between 5 and 9. Notice that if we have bad luck on rounding, such as adding 1 to a string of 9s, the sum may no longer be normalized, and we would need to perform step 3 again.

Steps 1 and 2 are similar to the example just discussed: adjust the significand of the number with the smaller exponent and then add the two significands. Step 3 normalizes the results, forcing a check for overflow or underflow. The test for overflow and underflow in step 3 depends on the precision of the operands. Recall that the pattern of all zero bits in the exponent is reserved and used for the floating-point representation of zero. Also, the pattern of all one bits in the exponent is reserved for indicating values and situations outside the scope of normal floating-point numbers. Thus, for single precision, the maximum exponent is 127, and the minimum exponent is -126. The limits for double precision are 1023 and -1022.

Decimal Floating-Point Addition

Example:

Try adding the numbers 0.5_{ten} and -0.4375_{ten} in binary using the algorithm in Figure:

Answer:

Let's first look at the binary version of the two numbers in normalized scientific notation, assuming that we keep 4 bits of precision:

$$\begin{array}{lll} 0.5_{\text{ten}} & = 1/2_{\text{ten}} & = 1/2^1_{\text{ten}} \\ & = 0.1_{\text{two}} & = 0.1_{\text{two}} \times 2^0 & = 1.000_{\text{two}} \times 2^{-1} \\ -0.4375_{\text{ten}} & = -7/16_{\text{ten}} & = -7/2^4_{\text{ten}} \\ & = -0.0111_{\text{two}} & = -0.0111_{\text{two}} \times 2^0 & = -1.110_{\text{two}} \times 2^{-2} \end{array}$$

Now we follow the algorithm:

Step 1. The significand of the number with the lesser exponent ($-1.11_{\text{two}} \times 2^{-2}$) is shifted right until its exponent matches the larger number:

$$-1.110_{\text{two}} \times 2^{-2} = -0.111_{\text{two}} \times 2^{-1}$$

Step 2. Add the significands:

$$1.000_{\text{two}} \times 2^{-1} + (-0.111_{\text{two}} \times 2^{-1}) = 0.001_{\text{two}} \times 2^{-1}$$

Step 3. Normalization the sum, checking for overflow or underflow:

$$\begin{aligned} 0.001_{\text{two}} \times 2^{-1} &= 0.010_{\text{two}} \times 2^{-2} = 0.100_{\text{two}} \times 2^{-3} \\ &= 1.000_{\text{two}} \times 2^{-4} \end{aligned}$$

Since $127 \geq -4 \geq -126$, there is no overflow or underflow. (The biased exponent would be $-4 + 127$, or 123, which is between 1 and 254, the smallest and largest unreserved biased exponents.)

Step 4. Round the sum:

$$1.000_{\text{two}} \times 2^{-4}$$

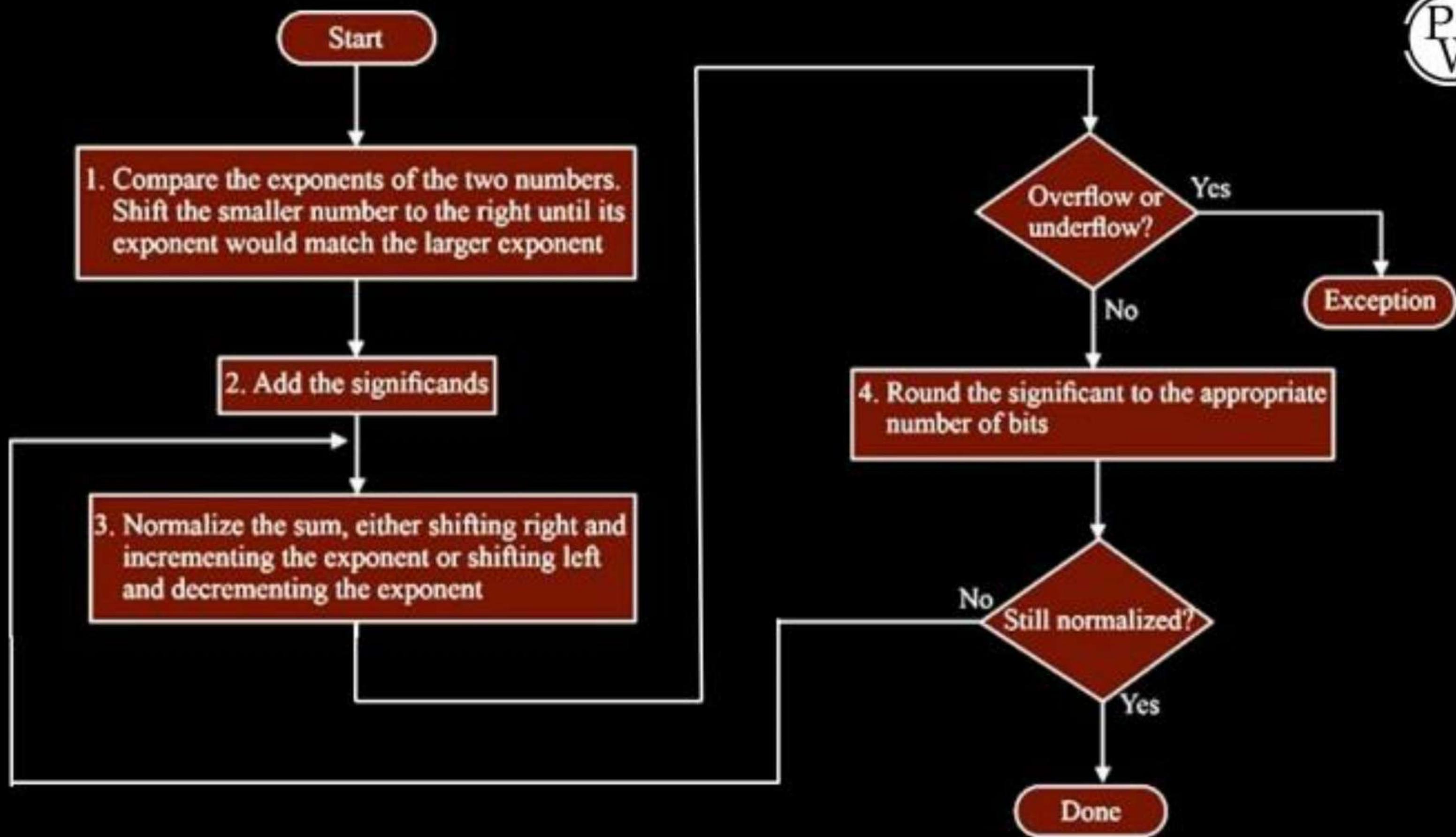
The sum already fits exactly in 4 bits, so there is no change to the bits due to rounding.

This sum is then

$$\begin{aligned}1.000_{\text{two}} \times 2^{-4} &= 0.0001000_{\text{two}} = 0.0001_{\text{two}} \\&= 1 / 2_{\text{ten}}^4 &= 1/16_{\text{ten}} &= 0.0625_{\text{ten}}\end{aligned}$$

This sum is what we would expect from adding 0.5_{ten} to -0.4375_{ten} .

Many computers dedicate hardware to run floating-point operations as fast as possible. Figure sketches the basic organization of hardware for floating-point addition.



Floating-point addition. The normal path is to execute steps 3 and 4 once, but if rounding causes the sum to be abnormalized, we must repeat step 3.

Floating-Point Multiplication

Now that we have explained floating-point addition, let's try floating-point multiplication. We start by multiplying decimal numbers in scientific notation by hand: $1.110_{\text{ten}} \times 10^{10} \times 9.200_{\text{ten}} \times 10^{-5}$. Assume that we can store only four digits of the significand and two digits of the exponent. Step 1. Unlike addition, we calculate the exponent of the product by simply adding the exponents of the operands together:

Step 1. Unlike addition, we calculate the exponent of the product by simply adding the exponents of the operands together:

$$\text{New exponent} = 10 + (-5) = 5$$

Let's do this with the biased exponents as well to make sure we obtain the same result:

$$10 + 127 = 137, \text{ and } -5 + 127 = 122, \text{ so}$$

$$\text{New exponent} = 137 + 122 = 259$$

This result is too large for the 8-bit exponent field, so something is amiss! The problem is with the bias because we are adding the biases as well as the exponents:

$$\text{New exponent} = (10 + 127) + (-5 + 127) = (5 + 2 \times 127) = 259$$

Accordingly, to get the correct biased sum when we add biased numbers, we must subtract the bias from the sum:

$$\text{New exponent} = 137 + 122 - 127 = 259 - 127 = 132 = (5 + 127)$$
 and 5 is indeed the exponent we calculated initially.

Step 2. Next comes the multiplication of the significands:

$$\begin{array}{r} 1.110_{\text{ten}} \\ \times \underline{9.200}_{\text{ten}} \\ 0000 \\ 0000 \\ 2220 \\ \underline{9990} \\ 10212000_{\text{ten}} \end{array}$$

There are three digits to the right of the decimal for each operand, so the decimal point is placed six digits from the right in the product significand:

$$10.212000_{\text{ten}}$$

Assuming that we can keep only three digits to the right of the decimal point, the product is 10.212×10^5 .

Step 3. This product is unnormalized, so we need to normalize it:

$$10.212_{\text{ten}} \times 10^5 = 1.0212_{\text{ten}} \times 10^6$$

Thus, after multiplication, the product can be shifted right one digit to put it in normalized form, adding 1 to the exponent. At this point, we can check for overflow and underflow. Underflow may occur if both operands are small—that is, if both have large negative exponents.

Step 4. We assumed that the significand is only four digits long (excluding the sign), so we must round the number.

The number

$$1.0212_{\text{ten}} \times 10^6$$

is rounded to four digits in the significand to

$$1.021_{\text{ten}} \times 10^6$$

Step 5. The sign of the product depends on the signs of the original operands. If they are both the same, the sign is positive; otherwise, it's negative. Hence the product is

$$+1.021_{\text{ten}} \times 10^6$$

The sign of the sum in the addition algorithm was determined by addition of the significands, but in multiplication the sign of the product is determined by the signs of the operands.

Once again, as Figure shows, multiplication of binary floating-point numbers is quite similar to the steps we have just completed. We start with calculating the new exponent of the product by adding the biased exponents, being sure to subtract one bias to get the proper result. Next is multiplication of significands, followed by an optional normalization step. The size of the exponent is checked for overflow or underflow, and then the product is rounded. If rounding leads to further normalization, we once again check for exponent size. Finally, set the sign bit to 1 if the signs of the operands were different (negative product) or to 0 if they were the same (positive product).

Decimal Floating-Point Multiplication

Example:

Let's try multiplying the numbers 0.5_{ten} and -0.4375_{ten} , using the steps in Figure.

Answer:

In binary, the task is multiplying $1.000_{\text{two}} \times 2^{-1}$ by $-1.110_{\text{two}} \times 2^{-2}$.

Step 1. Adding the exponents without bias:

$$-1 + (-2) = -3$$

or, using the biased representation:

$$\begin{aligned}(-1 + 127) + (-2 + 127) - 127 &= (-1 - 2) + (127 + 127 - 127) \\&= -3 + 127 = 124\end{aligned}$$

Step 2. Multiplying the significands:

$$\begin{array}{r} 1.000_{\text{two}} \\ \times \underline{1.110}_{\text{two}} \\ 0000 \\ 1000 \\ 1000 \\ \underline{1000} \\ 1110000_{\text{two}} \end{array}$$

The product is $1.110000_{\text{two}} \times 2^{-3}$, but we need to keep it to 4 bits, so it is $1.110_{\text{two}} \times 2^{-3}$.

Step 3. Now we check the product to make sure it is normalized, and then check the exponent for overflow or underflow. The product is already normalized and, since $127 \geq -3 \geq -126$, there is no overflow or underflow. (Using the biased representation, $254 \geq 124 \geq 1$, so the exponent fits.)

Step 4. Rounding the product makes no change:

$$1.110_{\text{two}} \times 2^{-3}$$

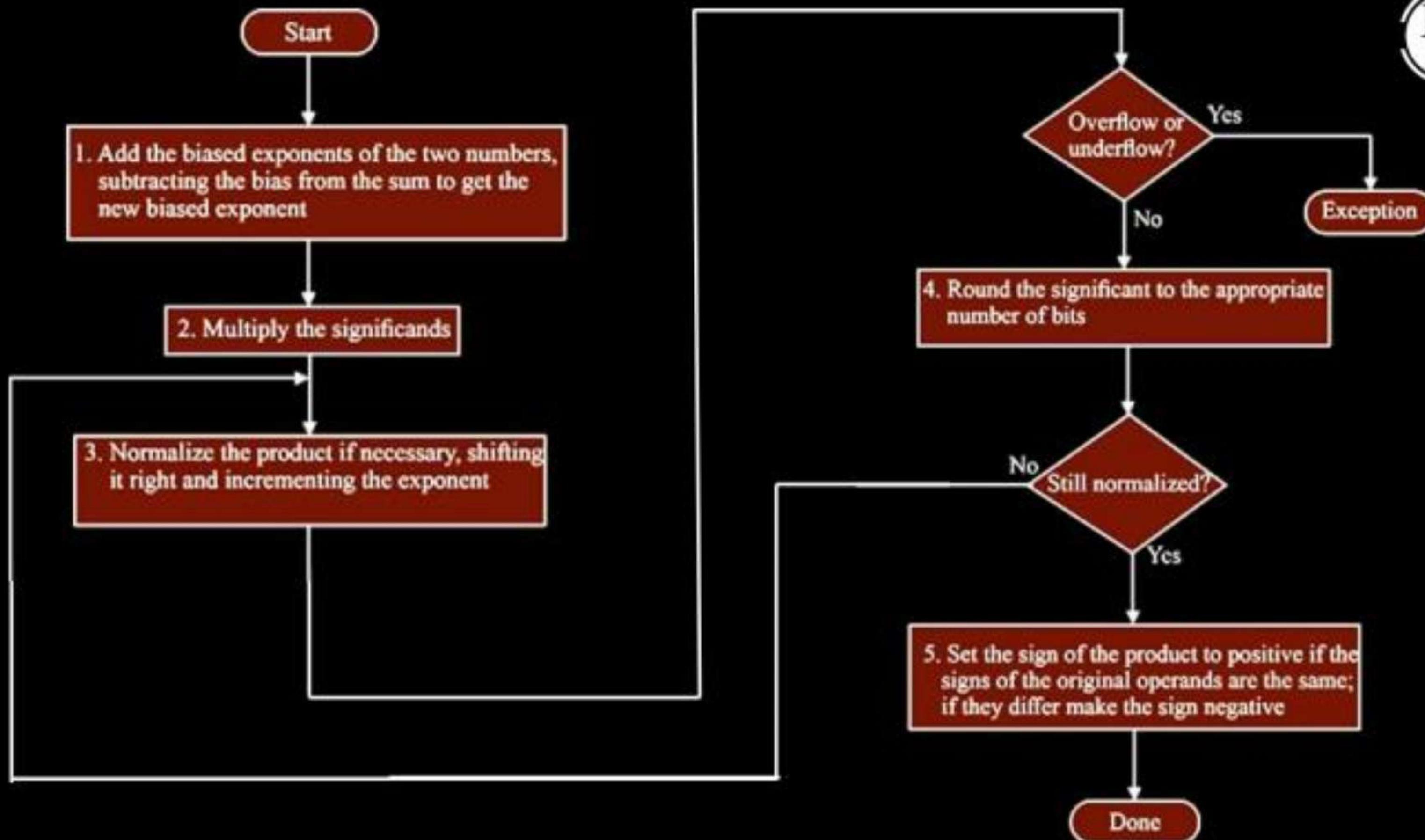
Step 5. Since the signs of the original operands differ, make the sign of the product negative. Hence the product is

$$-1.110_{\text{two}} \times 2^{-3}$$

Converting to decimal to check our results:

$$\begin{aligned} -1.110_{\text{two}} \times 2^{-3} &= -0.001110_{\text{two}} = -0.00111_{\text{two}} \\ &= -7/2^5_{\text{ten}} = -7/32_{\text{ten}} = -0.21875_{\text{ten}} \end{aligned}$$

The product of 0.5_{ten} and -0.4375_{ten} is indeed -0.21875_{ten} .



Floating multiplication:

The normal path is to execute steps 3 and 4 once, but if rounding causes the sum to be unnormalized, we must repeat step 3.

Arithmetic Operations on Floating-Point Number

In this section, we outline the general procedure for addition, subtraction, multiplication, and division of floating-point numbers. The rules we give apply to the single-precision IEEE standard format. These rules specify only the major steps needed to perform the four operations; for example, the possibility that overflow or underflow might occur is not discussed. Furthermore, intermediate results for both mantissas and exponents might require more than 24 and 8 bits, respectively. These and other aspects of the operations must be carefully considered in designing an arithmetic unit that meets the standard. Although we do not provide full details in specifying the rules, we consider some aspects of implementation, including rounding, in later sections.

If their exponents differ, the mantissas of floating-point numbers must be shifted with respect to each other before they are added or subtracted. Consider a decimal example in which we wish to add 2.9400×10^2 to 4.3100×10^4 . We rewrite 2.9400×10^2 as 0.0294×10^4 and then perform addition of the mantissas to get 4.3394×10^4 . The rule for addition and subtraction can be stated as follows:

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!**

