

COMPUTER SCIENCE

Computer Organization and Architecture

Floating Point Representation

Lecture_04



Vijay Agarwal sir



TOPICS
TO BE
COVERED

o1

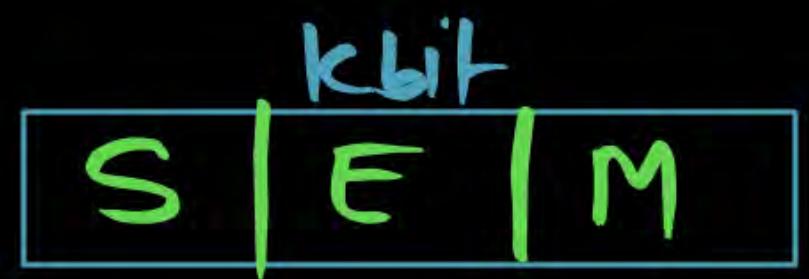
IEEE 754 Floating Point
Representation

o2

Expand Opcode Technique
Little Endian & Big Endian



Floating point Representation.



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

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

BE | E [Bias Exponent | Exponent]

$$BE = AE + bias$$

$$E = e + bias$$

IEEE 754 Floating Point Representation

Single Precision
Format (32bit)

S	E	m
1bit	8bit	23bit

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

$$\boxed{\text{bias} = 127}$$

$$\boxed{\text{Excess-127}}$$

Double Precision
Format (64bit)

S	E	m
1bit	11bit	52bit

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

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

~~Excess~~

$$\boxed{\text{Excess} = 1023.}$$

Introduction of C++ CHAPTER: 1

Mc Instn & AM CHAPTER: 2

Floating point Rep. CHAPTER: 3

Complete till 2nd Oct.

Consider the IEEE-754 single precision floating point numbers

P = 0xC1800000 and Q = 0x3F5C2EF4.

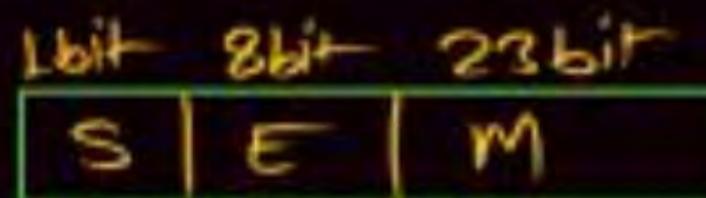
Which one of the following corresponds to the product of these numbers
(i.e., P × Q), represented in the IEEE-754 single precision format?

[GATE-2023-CS: 2M]

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

Ans (C)

$$P = 0x C1800000$$



$$Q = 0x 3F5C2EF4$$

$$bias = 2^8 - 1 \Rightarrow bias = 127$$

$$P = 0x C1800000$$

1	100 0001	1000 0000 0000 0000 0000 0000
---	----------	-------------------------------

Sign

E(8bit)

Mantissa (23bit)

$$BE = AE + bias$$

$$S = 1 (-ve)$$

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

$$131 - 127$$

$$E = 10000001 = 131$$

$$(-1)^1 \cdot 1.0000000 \times 2^{131}$$

$$E = e + bias$$

$$BE @ E = 131$$

$$M = 00000000$$

$$P = -(1.0000000) \times 2^{+4}$$

$$E = 131$$

$$bias = 127$$

$$Q = 3F5C2EF4$$

0011	$1LLL$	0011100000101110100
Sign 1 bit	$E(8\text{bit})$	$M(23\text{bit})$

Sign = 0 (+ve)

$$E = 0111110 \Rightarrow E = 126$$

$$BE @ E = 126.$$

$$M = 10111000010111011110100$$

$$\text{bias} = 127$$

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

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

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

$$(-1)^0 \cdot 1 \cdot 10111000010111011110100 \times 2^{126-127}$$

$$Q = (1.10111000010111011110100) \times 2^{-1}$$

$$\begin{aligned} E &= 126 & \text{bias} &= 127 \\ e &= E - \text{bias} & & \\ & & & 126 - 127 \end{aligned}$$

Sign = -ve

$$+ \times - = -ve$$

P × Q = exponent = (+4) + (-1) = +3

P × Q = Mantissa = (1.0000) × (1.101 1100 0010 1110 1111 0100)

$$- (1.101 1100 0010 1110 1111 0100) \times 2^{+3}$$

Sign = 1 (-ve)

e = +3

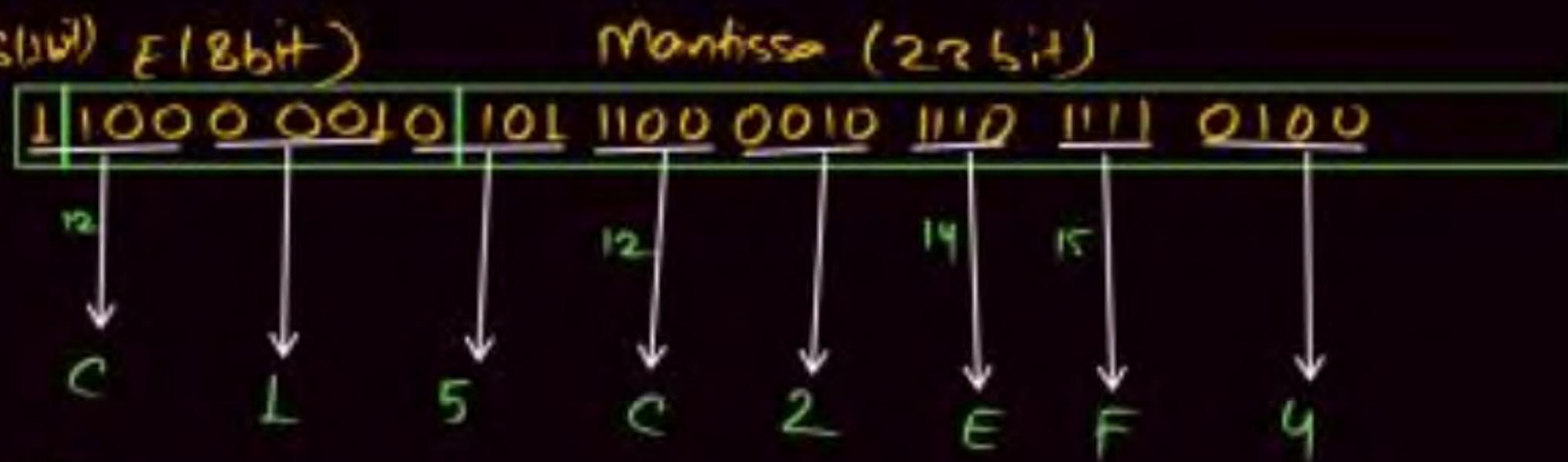
bias = 127

$$E = e + bias \text{ (11bit) } E/8\text{bit}$$

$$= 3 + 127$$

$$E = 130$$

$$E = 1000 0010$$



(C15C2EF4H) Amp

Alternate Approach

If we Don't Know
Floating Point Multiplication
then How to Deal ?

2nd Approach

Alternate Approach

$$P = C18\ 00000$$

S(1bit) E(8bit)

S(1bit)	E(8bit)	Mantissa (23bit)
1	100 0001	1000 0000 0000 0000 0000 0000

$$S = 1 \text{ (-ve)}$$

$$E = 10000001_2$$

$$\beta E @ E = 131$$

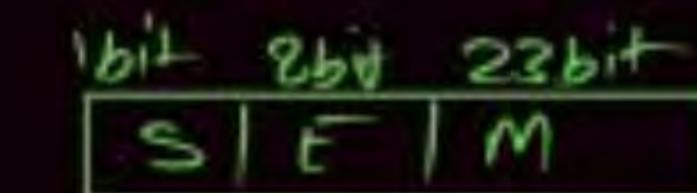
$$M = 000000000$$

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

$$\text{bias} = 127$$

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

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



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

$$(-1)^1 \cdot 1.0000000 \times 2$$

$$-1.0000000 \times 2^{+4}$$

$$-10000 \cdot 00$$

$$P = -16$$

$$P = -16.$$

$$Q = 3F5C2EF4$$

Sign	E(8bit)	Mantissa (23bit)
$b_9 b_8$	0011 1111 0101 1100 0010 1110 1111 0100	

$$e = E - b_{15}$$

$s = 0$ (+ve)

$E = 0111110$

$$RE \oplus E = 126$$

$m = 10111000010\ldots$

$b_{10}g \sim 12.9$

$$E = e + b_{10}g$$

$$e = E - b_{10}g$$

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

$$(-1)^0 \cdot 1 \cdot 101110000101110 \times 2$$

$$\begin{aligned} Q &= 1 \cdot 1011100 \times 2^{-1} \\ &= 0.11011100 \end{aligned}$$

$$Q = 0.8593$$

$$Q = 0.8593$$

$$126 - 12.7$$

$$P \times Q = -16 \times -8593 \\ = \underline{\underline{-(13.75)}}$$

$$P \times Q = -13.75$$

-1101-11

$$\Rightarrow -1.10111 \times 2^{+2}$$

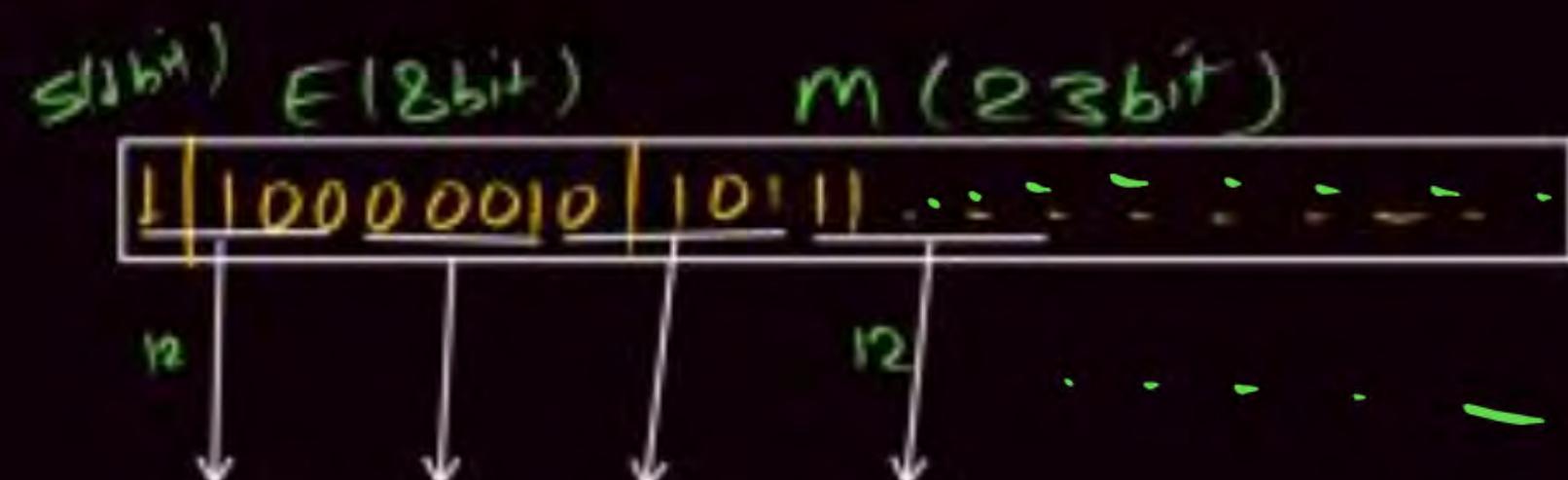
$$C = +3$$

$$\log = 127$$

$$E_{\text{at}} \beta E = e + b \cos \theta \Rightarrow 3 + 127 \\ = 130$$

m = 1011

$$S=1 \text{ -ve}$$



S 159

11

$$(C_1 5 \subset 2 E F Y)_k \text{ Avg}$$

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



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

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 00000000 M = 0	± 0
0 or 1	11111111 E = 255	0000000000000000 00000000 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)

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 <u>E = 0</u>	000000000000.. M = 0	± 0
0 or 1	1111 1111 111 <u>E = 2047</u>	0000000000.. 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	<u>E = 0</u>	M $\neq 0$	Denormalized number/ Fractional Form $(-1)^S 0 \cdot M \times 2^{E-1023}$
	<u>E = 2047</u>	M $\neq 0$	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 $M \neq 0$

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)

IEEE 754

✓ 0
✓ ∞

✓ Nan (Not a Number).
✓ Denormalized Number

Concept of Denormalized Number :

In a Single Precision

Denormalized Number ?

$$E=1$$

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

$$\boxed{e = -126}$$

$$-126 + 127$$

~~$$\boxed{E = 1}$$~~

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

$$E = e + 127$$

Single Precision.

$$\text{bias} = 127$$

$$E = 000000001$$

$$E = 1$$

$$E = e + bias$$

$$E = e + 127$$

$$e = E - 127$$

$$e = 1 - 126$$

$$e = -126$$

worst case .

$$e = -126$$

$$\text{if } e = -127$$

$$E = e + bias = -127 + 127 \leftarrow 0$$

$$\text{if } e = -128$$

$$E = e + bias = -128 + 127 \rightarrow E = -1$$

Denormalized Number



Minimum Possible value of $E=00000001$, $E=1$

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

$$e = E - \text{bias}, \quad 1 - 127 = 126, \quad \text{so } e = -126$$

That means in worst case $e = -126$, if value of e is smaller than -126, then number is not able to normalize & we store as Denormalize number.

eg 1.1001×2^{-127} $e = -127, E = e + \text{bias}, -127 + 127 = 0$ so $E = 0$, not able to normalize

eg 1.1001×2^{-128} $e = -128$, so $E = -1$, not able to normalize

or 1.1001×2^{-129}

Here $M = 1001$

$e = -129, E = e + \text{bias}, -129 + 127 = -2,$

So $E = -2$ not able to normalize because E must be 1 After biasing.

In a Single Precision.

for a Normalized Number in

$$\text{bias} = 127$$

Worst Case:

$$L \cdot M \times 2^{-126}$$

$$e = -126$$

then its Normalized

Otherwise Not Able to

Normalize the Number & Store as a
Denormalized Manner.

126 Time Shift
for Normalization

In a Double Precision.

$$E = e + b10s$$

$$E = e + 1023$$

Normalized
Number
Worst Case.

$$1.0 \times 2^{-1022}$$

$$e = -1022$$

$$E = -1022 + 1023$$

$$E = 1$$

$$e = -1024$$

$$E = -1024 + 1023 \Rightarrow$$

$$E = -1$$

Denormalized Number

Minimum Possible value of E=00000001, E=1

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

$$e = E - \text{bias}, \quad 1-127 = 126, \quad \text{so } e = -126$$

That means in worst case $e = -126$, if value of e is smaller than -126, then number is not able to normalize & we store as Denormalize number.

For eg. The Smallest Positive(+ve) single precision Normalized Number

is: $1.0000\ 0000\ 0000\ 0000\ 0000\ 000 \times 2^{-126}$ $1.000\ 0000\ 0000\ 0000\ 0000\ 000 \times 2^{-126}$

But the Smallest single precision denormalized number is:

$$0.0000\ 0000\ 0000\ 0000\ 0001 \times 2^{-126} \quad \text{Or} \quad 1.0 \times 2^{-149}$$

& Double precision Range is: 1.0×2^{-1022}

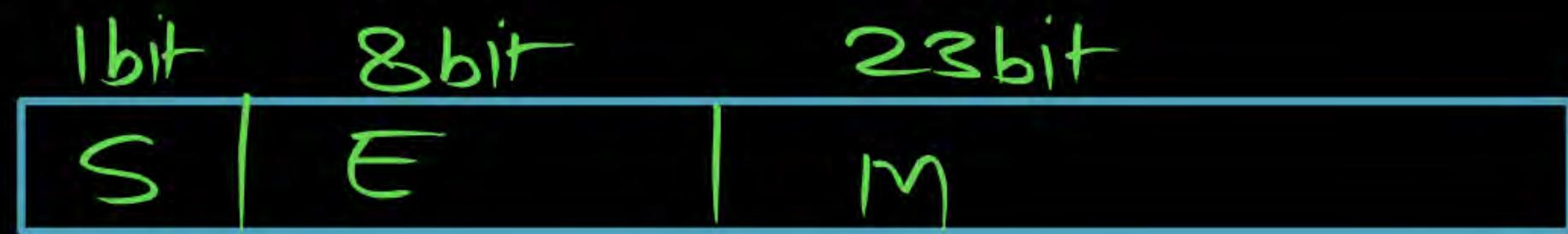
$$\text{To } 1.0 \times 2^{-1074}$$

$$\frac{-126}{-23} = -149$$

$$\frac{-1022}{-52} = -1074$$

Q.

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

P
W

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

$$\text{bias} = 127$$

Sol.

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

P
W

Repartition..



$$+ 1.0 \times 2^0$$

$$M = 0$$

$$e = 0$$

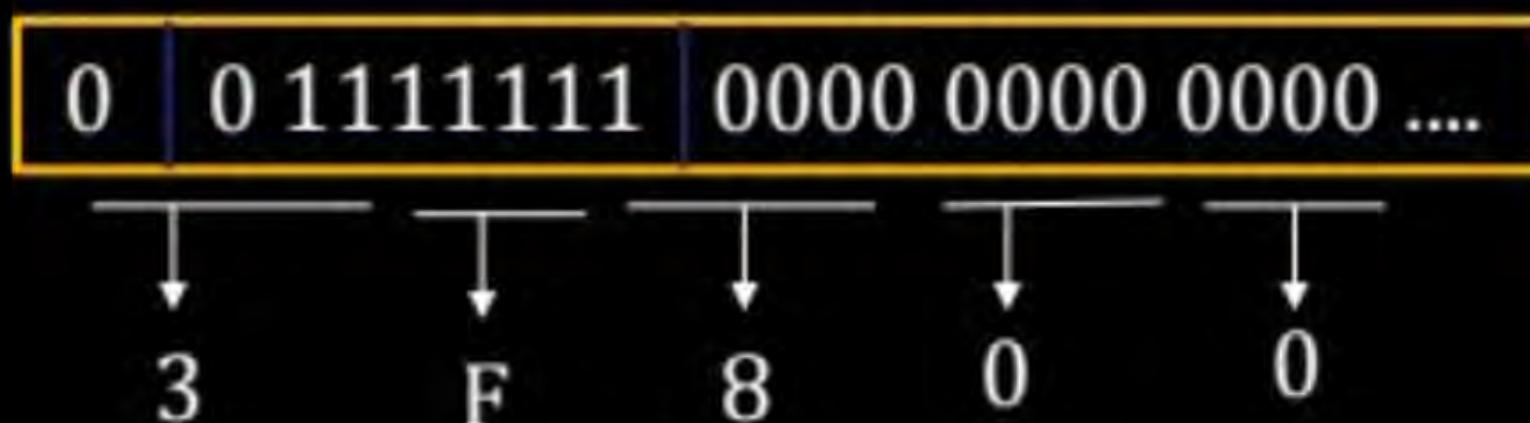
$$E = 127$$

S	E	M
1	8	23

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

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

$$E = 0 + 127$$



(3F80 0000)

Ans

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 000000000000....

↳ Sign

$$E = 2047$$

$$M = 0$$

too big

+ Infinite

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 0000000000000000...

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 +∞

MCQ

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

Additional [Extra] Concept of Floating Point Representation.

- ① Addition | Subtraction
- ② Multiplication | Division.

Number

①



Number

②

for Sub @ Addition: Power (exponent) must be equal

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.
4. Normalize the resulting value, if necessary. [1. Something X 2^e.]



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. ($1 \cdot \text{Something} \times 2^e$)

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.

A: 10

B: 11

C: 12

D: 13

E: 14

F: 15

Q) -0.75 in Single precision Format

$$\Rightarrow -0.11 \times 2^0.$$

$$\Rightarrow -1.1 \times 2^{-1}$$

$$e = -1$$

S	E	M
1bit	8bit	23bit

$$\text{bias} = 127$$

$$E = -1 + 127 = E = 126$$

$$E = 01111110$$

sign	E(8bit)	Mantissa (23bit)
$\boxed{1}$	$\boxed{0111110}$	$\boxed{10000000000000000000000}$
\Downarrow	$\downarrow 15(F)$	

∴ BF 40 0000

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.

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

It is also represented by the binary fraction.

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

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 $-1.1_{\text{two}} \times 2^{-1}$, the result.

$$(-1)^S \times (1 + .1000\ 0000\ 0000\ 0000)$$

$$0000\ 000_{\text{two}}) \times 2^{(126 - 127)}$$

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

Floating Point Addition:

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.
4. Normalize the resulting value, if necessary. [1. Something X 2^e]



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

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:

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

Addition

Decimal :

$$\textcircled{1} \quad 9.999 \times 10^{+1}$$

Now:

$$9.999$$

$$0.016$$

$$\frac{9.999}{0.016} = 10.015 \times 10^{+1}$$

Step 4: Normalize
1. Something $\times 2^e$

$$1.0015 \times 10^{+2}$$

Ans $\textcircled{2}$

Smallest

$$\textcircled{2} \quad 1.610 \times 10^{-1}$$

\Downarrow

$$0.1610 \times 10^0$$

\Downarrow

$$0.01610 \times 10^{+1}$$

New exponent
Equal
So Now
Perform
Addition

$$1.002 \times 10^{+2}$$

Rough off

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

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

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

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

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.

Now Addition for Binary Number

Decimal Floating-Point Addition

Example:

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

0.1

0.0111

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:

$$0.5_{\text{ten}} = 1/2_{\text{ten}} = 1/2^1_{\text{ten}}$$

$$= 0.1_{\text{two}} = 0.1_{\text{two}} \times 2^0$$

$$0.5 = 1.000_{\text{two}} \times 2^{-1} \quad \textcircled{1}$$

$$-0.4375_{\text{ten}} = -7/16_{\text{ten}} = -7/2^4_{\text{ten}}$$

$$= -0.0111_{\text{two}} = -0.0111_{\text{two}} \times 2^0$$

$$\begin{array}{r} -0.4375 \\ \hline -1.110 \end{array} \times 2^{-2} \quad \textcircled{2}$$

0.5

① 0.10

$$1.0 \times 2^{-1}$$

- 0.4375

② -0.0111

$$-1.11 \times 2^{-2}$$

$$-0.111 \times 2^{-1} - ②$$

Now

Addition

$$1.0 \times 2^{-1} + (-0.111 \times 2^{-1})$$

$$\Rightarrow (1.0 + (-0.111)) \times 2^{-1} = 0.001 \times 2^{-1}$$

Step Normalized
1.something

$$1.0 \times 2^{-4}$$

Avg

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:

$$-1.110_{\text{two}} \times 2^{-2} = -0.111_{\text{two}} \times 2^{-1} \quad | \cdot 0 \times 2^{-1}$$

Step 2. Add the significands:

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

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

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

V.DNP Since $\underline{127} \geq \underline{-4} \geq \underline{-126}$, 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:

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

Ans

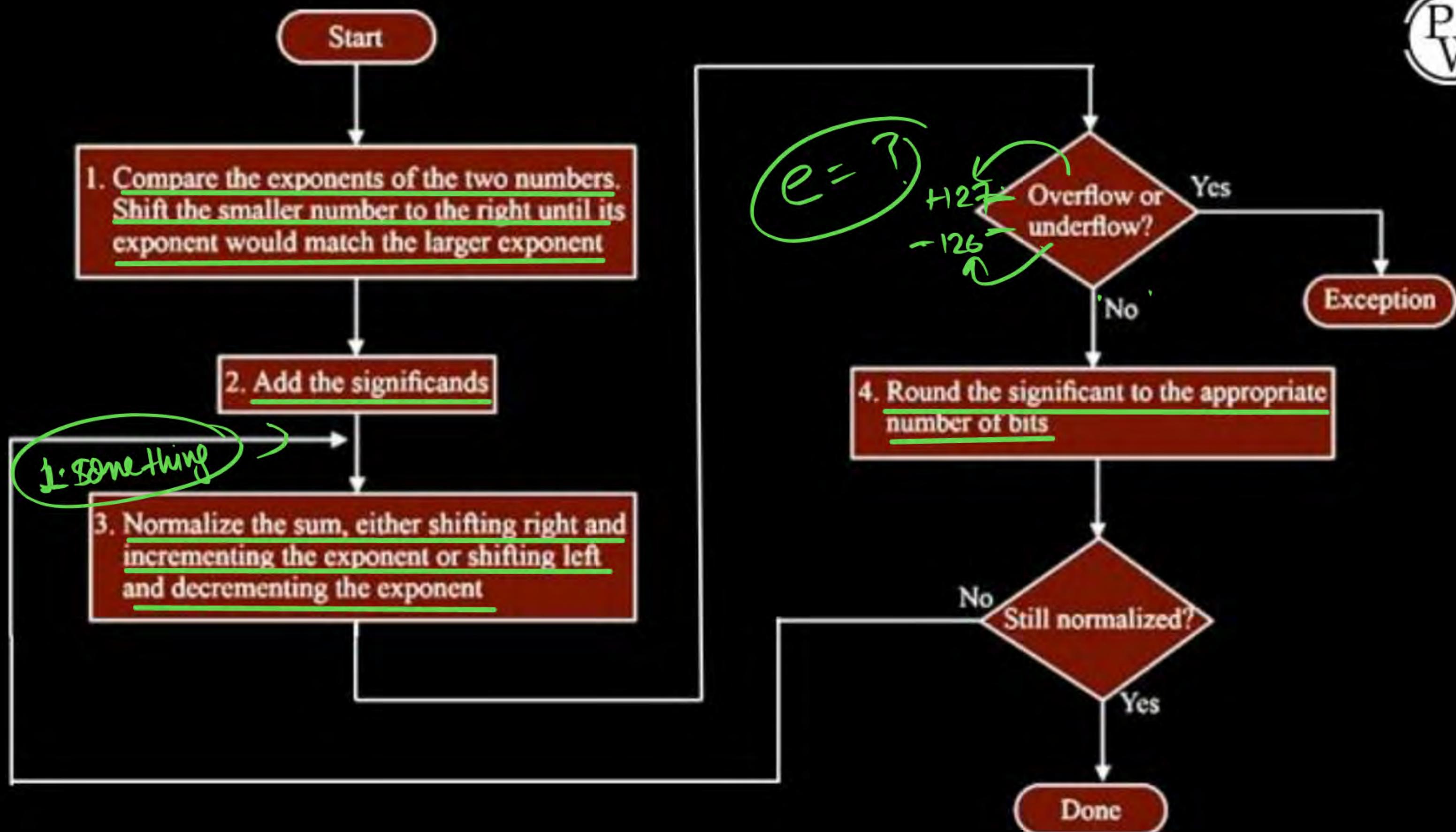
The sum already fits exactly in 4 bits, so there is no change to the bits due to rounding.

This sum is then

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

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.

Addition & Subtraction

- ① First step we do
 - Exponent (Power) must be equal
- ② Addition
- ③ Normalize 1.something

Floating Point multiplication:

~~easy~~ ↳ Direct Multiplication (Exponent (Power) must
Not be equal & same)

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. ($1 \cdot \text{Something} \times 2^e$).

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.

Decimal Multiplication.

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:

Number 1 : 1.110×10^{10}

Number 2 : 9.200×10^{-5}

Direct ~~e~~
 ① 1.110×10^{10} e
 ② 9.200×10^{-5} e

$$\begin{array}{r} 1.110 \\ \times 9.200 \\ \hline \end{array}$$

χ

$$e = 10 + (-5)$$

$$\boxed{e=5}$$

$$E = e + 127 \Rightarrow E = 5 + 127$$

$$\boxed{E=132} \text{ Ang}$$

Asteriate

OR

$$E = e + b \cos \theta$$

① $1.110 \times 10^{10} e$

$$E = e + 127$$

$$E = 10 + 127$$

$$E = 137 \checkmark$$

$$E = 137 + 122 - 127$$

$$E \Rightarrow 132$$

$$e = 132 - 127$$

② $9.200 \times 10^{-5} e$

$$E = e + 127$$

$$= -5 + 127$$

$$E = 122$$

$$\Rightarrow e = 5$$

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

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

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:

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

Step 2 Next comes the multiplication of the significands:

$$\begin{array}{r} 1.110_{\text{ten}} \\ \times 9.200_{\text{ten}} \\ \hline 0000 \\ 0000 \times \\ \hline 2220 \times \times \\ \hline 9990 \times \times \times \\ \hline 10212000_{\text{ten}} \end{array} \Rightarrow \boxed{10 \cdot 212 \times 10^{+5}}$$

Normalize : $1.0212 \times 10^{+6}$ Ans

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:

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

Ans

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).

Binary Floating Point Multiplication

Step 2

$$\begin{cases} \textcircled{1} \quad 0.5 \Rightarrow 0.1 \Rightarrow 1.0 \times 2^{-1} \\ \textcircled{2} \quad -0.4375 \Rightarrow -0.0111 \Rightarrow 1.11 \times 2^{-2} \end{cases}$$

Step 2: Direct Multiply

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

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:

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

e = -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:

$$\begin{array}{r} 1.000_{\text{two}} \\ \times 1.110_{\text{two}} \\ \hline 0000 \\ 1000X \\ \hline 1000XX \\ \hline 1000XXX \\ \hline 1110000_{\text{two}} \end{array}$$

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:

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

Ans

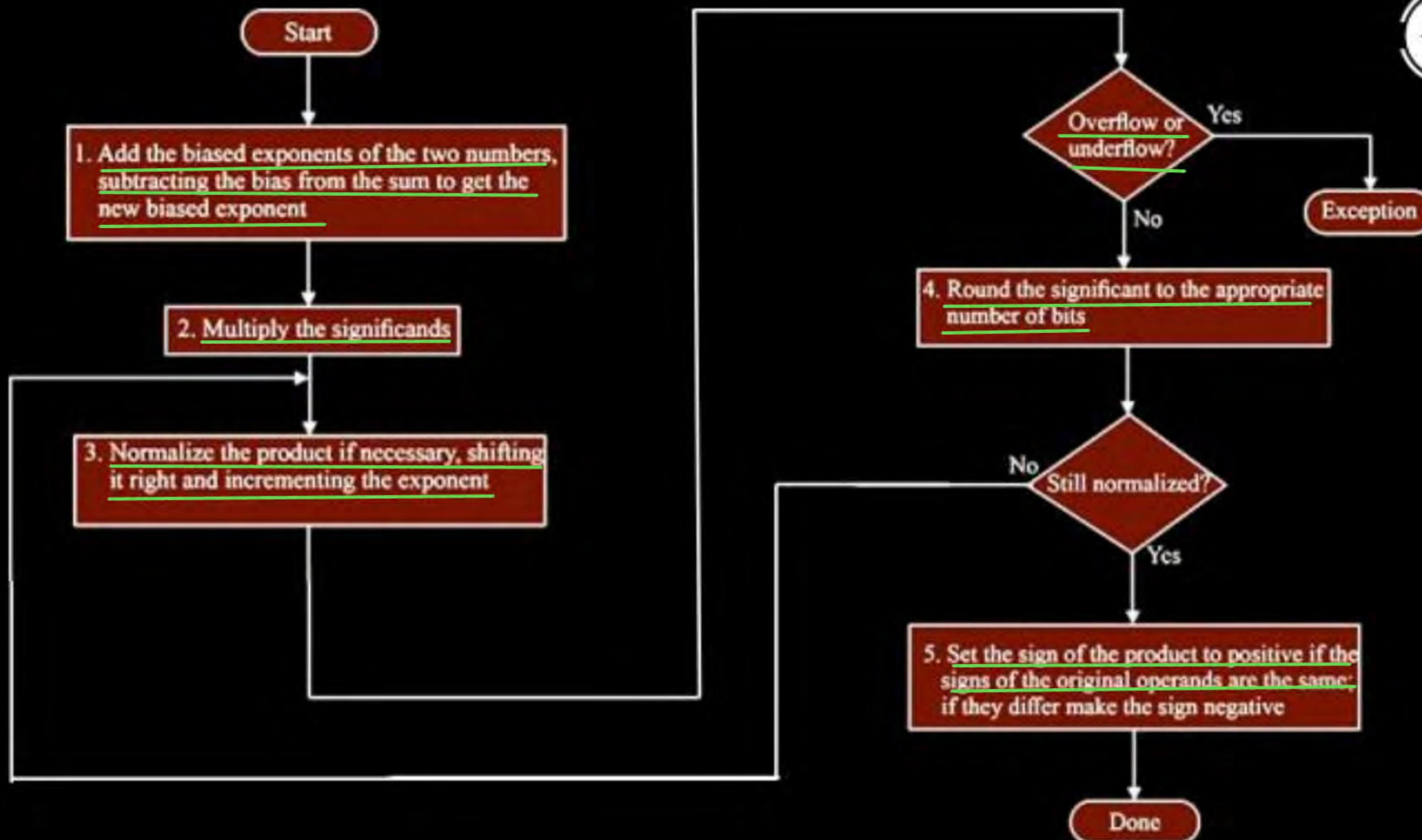
Step 5. Since the signs of the original operands differ, make the sign of the product negative. Hence the product is

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

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.

Little Endian & Big Endian.

Basic Terms and Notation

The alphabet of computers, more precisely digital computers, consists of 0 and 1.

Each is called a *bit*, which stands for the binary digit.

The term *byte* is used to represent a group of 8 bits.

The term *word* is used to refer to a group of bytes that is processed simultaneously.

The exact number of bytes that constitute a word depends on the system. For example, in the Pentium, a word refers to four bytes or 32 bits. On the other hand, eight bytes are grouped into a word in the Itanium processor.

We use the abbreviation “b” for bits, “B” for bytes, and “W” for words.

Sometimes we also use doubleword and quadword. A doubleword has twice the number of bits as the word and the quadword has four times the number of bits in a word.

$$\text{DoubleWord} = 2 \times \text{Word Size}$$
$$\text{QuadWord} = 4 \times \text{Word Size}.$$



Bits in a word are usually ordered from right to left, as you would write digits in a decimal number. The rightmost bit is called the least significant bit (LSB), and the leftmost bit is called the most significant bit (MSB).

1 Byte = 8 bit

Memory $\xrightarrow[\text{Byte}]{\text{Data in the form}}$ Byte

Word Size = Word Length
of the Processor.

CPU
(Processor) $\xrightarrow[\text{form of}]{\text{Data in the word}}$ Word

16 bit Processor

Word Size = 16 bit

8085 Processor

AD₀ - AD₇ & A_B - A₁₅

Word length = 8 bit

8086 Processor

AD₀ - AD₁₅ & A₁₆ - A₁₉

Word length = 16 bit

Operation performed
On 16bit Data format .

Note

Default Configuration in the Memory chip
is Byte Addressable. So in the Memory Data
is stored Byte wise.

Note

In the Processor operation are performed on a Word
format. So when the Word length of the Processor is
greater than 8bit (eg 16bit, 32bit) then Mutiple Cells
Accessing is Required to Access the Data from Memory Simultaneously.

16 bit Processor

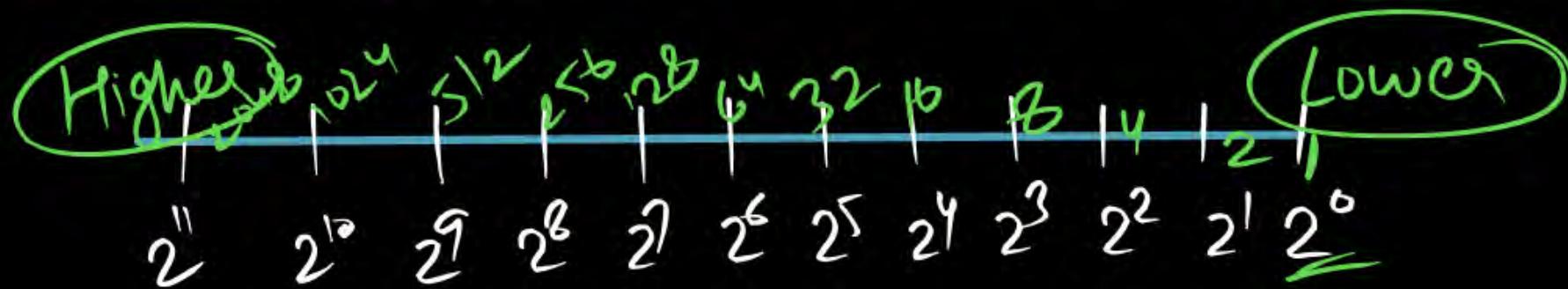
MOV Ax [4000]

$$Ax \leftarrow \begin{bmatrix} M[4000] \\ M[4001] \end{bmatrix}$$

32 bit Processor

MOV R0 [4000]

$$R0 \leftarrow \begin{bmatrix} M[4000] \\ M[4001] \\ M[4002] \\ M[4003] \end{bmatrix}$$



P
W

16 bit Processor

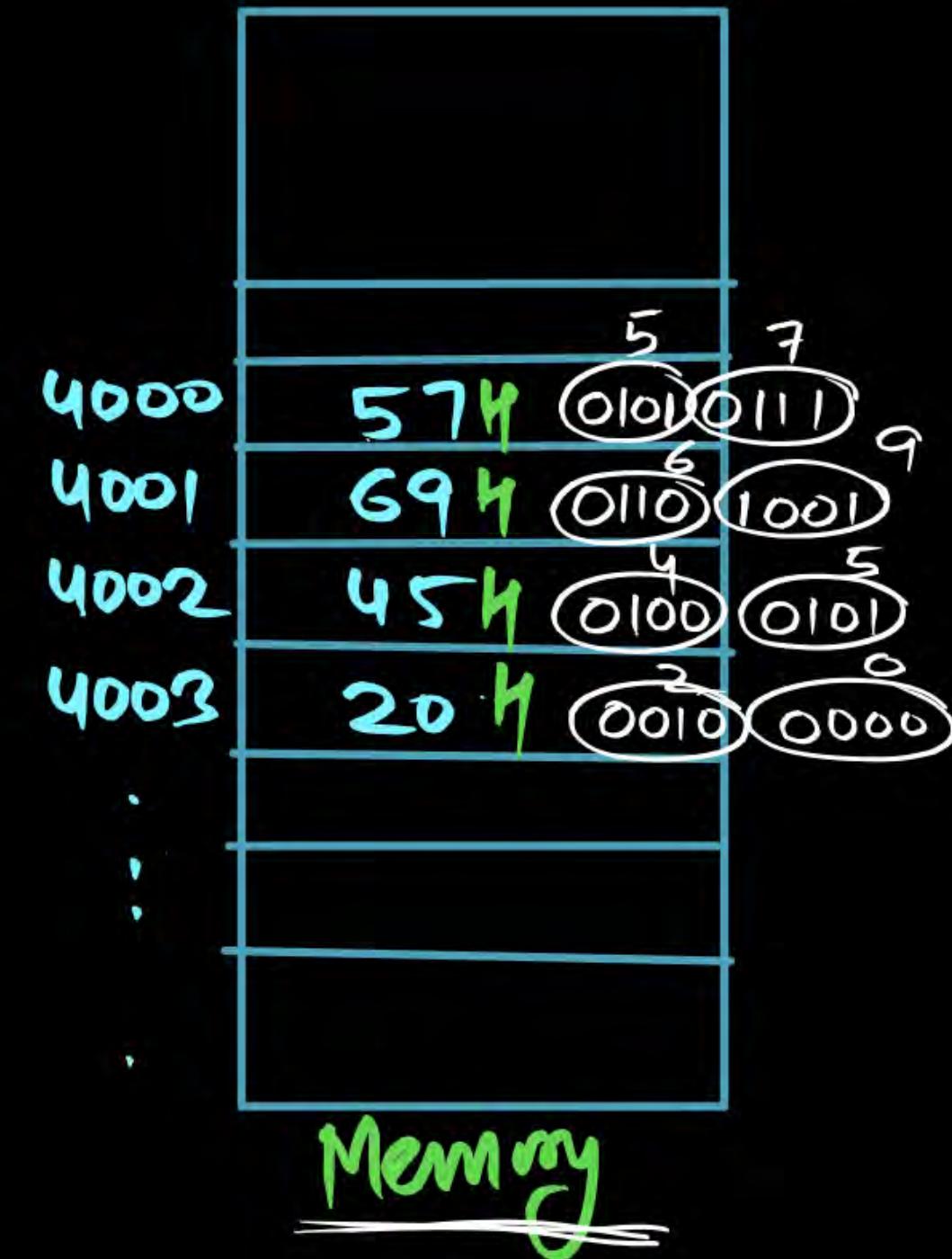
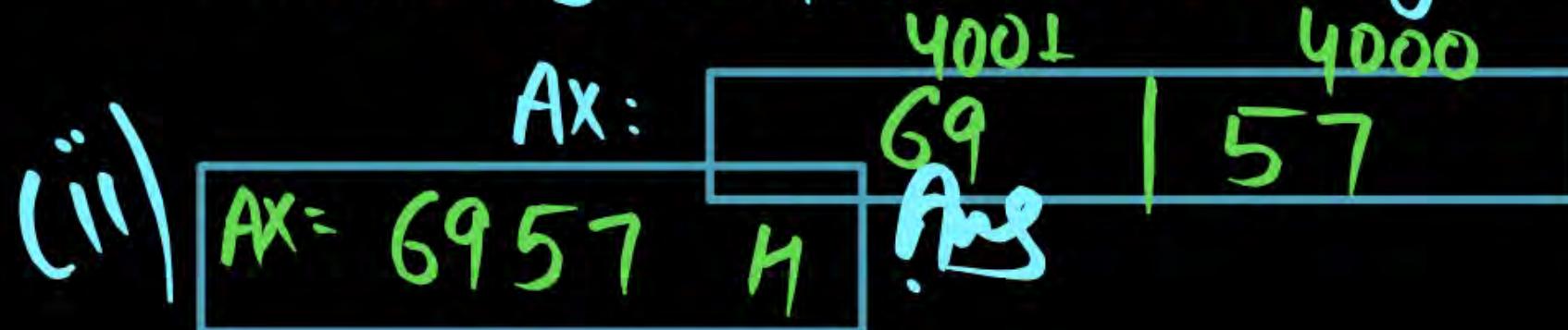
MOV Ax [4000]

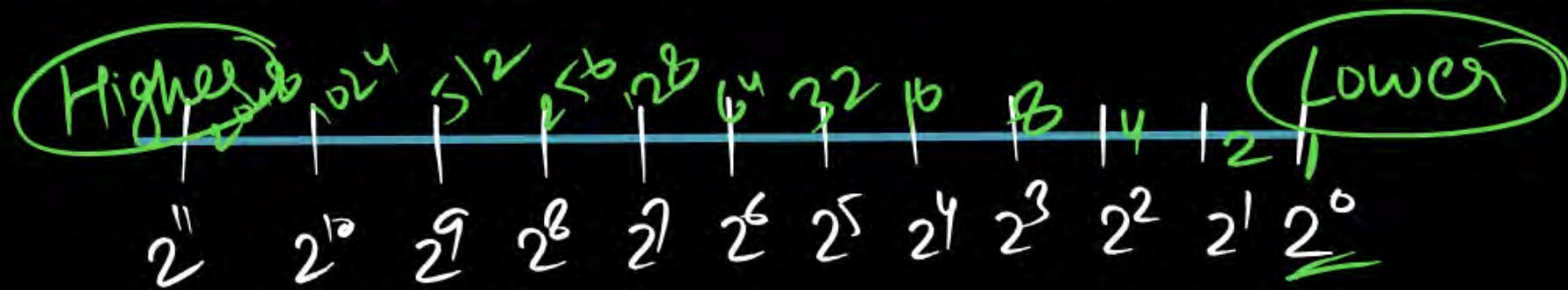
(i) If $M[4000]$ contain Higher Byte



(i) $AX = 57\ 69\ H$ Ans

(ii) If $M[4000]$ contain Lower Byte





P
W

32 bit Processor

MOV AX [4000]

(i) If $m[4000]$ contain Higher Byte.

AX:

4000	4001	4002	4003
57	69	45	20

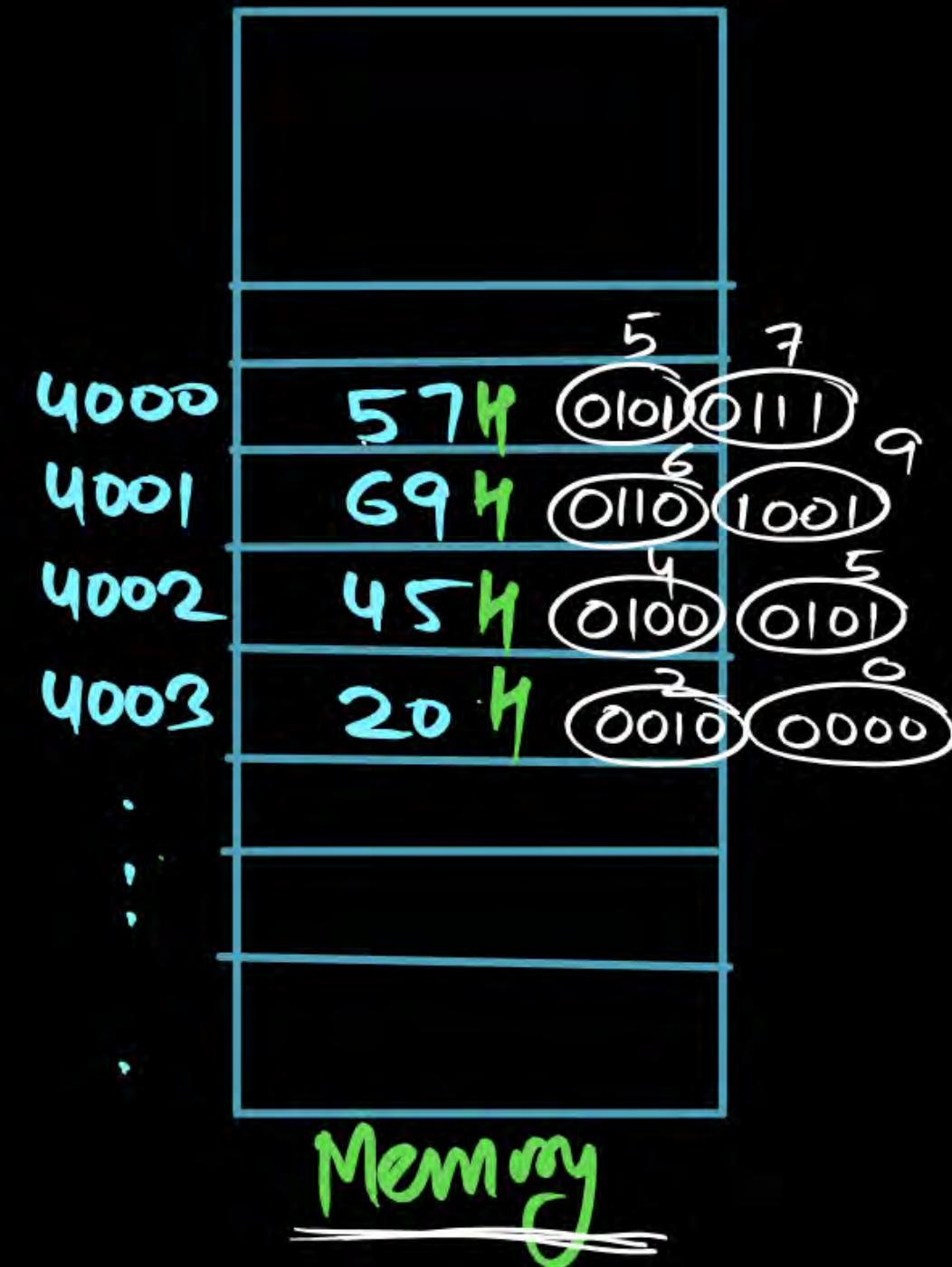
(i) $AX = 5769\ 4520H$ Ans

(ii) If $m[4000]$ contain Lower Byte.

AX:

4003	4002	4001	4000
20	45	69	57

$AX = 2045\ 6957H$ Ans



Here L Instruction Create 2 output
Create Ambiguity.

To solve the problem of ambiguity there is need of Memory Address interpretation Mechanism(MIC ism)
Called Endian MIC ism.

- ① Little Endian
- ② Big Endian

Byte ordering
or

Endian M/c ism : Show the order of Data Storage in Memory.

① Little Endian : Lower Addresses Contain Lower Byte

Address → 3 2 1 0 & Higher Addresses Contain Higher Byte.
Byte storage →

3	2	1	0
---	---	---	---

② Big Endian : Lower Addresses Contain Higher Byte & Higher Addresses Contain Lower Byte.

Address → 3 2 1 0
Byte storage →

0	1	2	3
---	---	---	---

Note

Default M_C ism is Little Endian

Used in Processor Design

But Some Processor use Big Endian also

So Little
Endian

16 bit Processor

Ax = 69 57 H

32 bit Processor

D₀ = 2045 6957 H

My Computer Properties.

✓ 4GB RAM

↳ 32bit Addressing line.

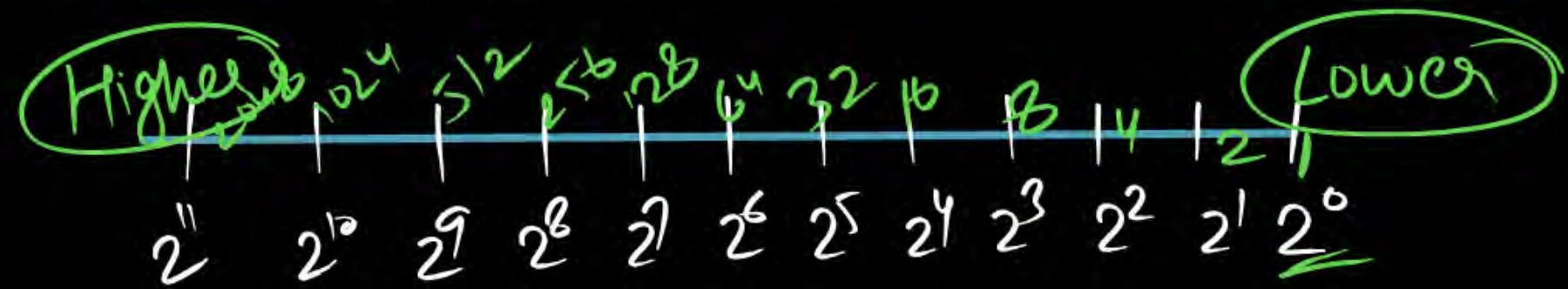
✓ 32/64 bit Processor

- 1TB (2^{40} Byte) HD.
- 4GHz Processor.

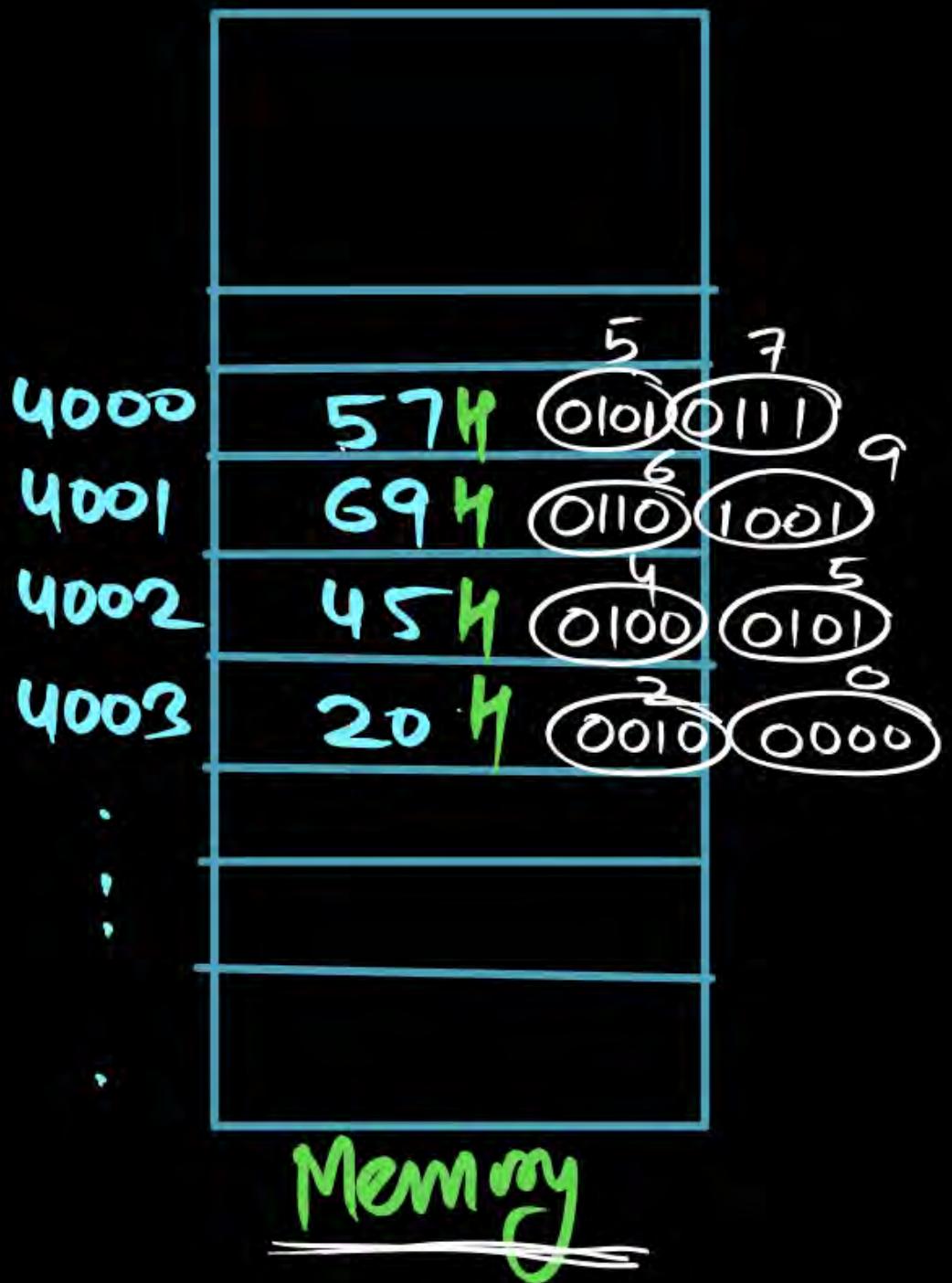
Word Length = 32/64 bit

Note

Little Endian : Right to Left
Big Endian : Left to Right



P
W



Byte Ordering

Storing data often requires more than a byte.

Suppose that we want to store these 4-byte data in memory at locations 100 through 103.

How do we store them?

① Figure Shows two possibilities: Least significant byte or Most significant byte is stored at location 100. These two byte ordering schemes are referred to as the little endian and big endian.

{ 2B, 4B, 8Byte
(16bit, 32bit, 64bit)

Note

Little Endian : Right to Left

Big Endian : Left to Right

0001 0001 0001 0010 0001 0011 0001 0100

(11, 12, 13, 14)H

32 bit Processor

Stored at Starting address
100.

x0 =

① Little Endian
(Right to Left)

$x_0 = \boxed{11} \boxed{12} \boxed{13} \boxed{14}$

$x_0 = 11\ 12\ 13\ 14\ H$

Ans

② Big Endian
(Left to Right)

$x_0 = \boxed{14} \boxed{13} \boxed{12} \boxed{11}$

$x_0 = 14\ 13\ 12\ 11\ H$

$(25, 27, 26, 29)_H$. Stored at Starting address 200
32 bit Processor.

① Little Endian
(Right to Left)

203	202	201	200
25	27	26	29

$$\gamma_1 = 25\ 27\ 26\ 29\ H \quad \text{Ans}$$

② Big Endian
(Left to Right)

203	202	201	200
29	26	27	25

$$\gamma_1 = 29\ 26\ 27\ 25\ H$$

$(51, 21, 15, 25)H$. Stored at Starting address 400
32 bit Processor.

① Little Endian
(Right to left)

	403	402	401	400
$\delta_2 =$	51	21	15	25

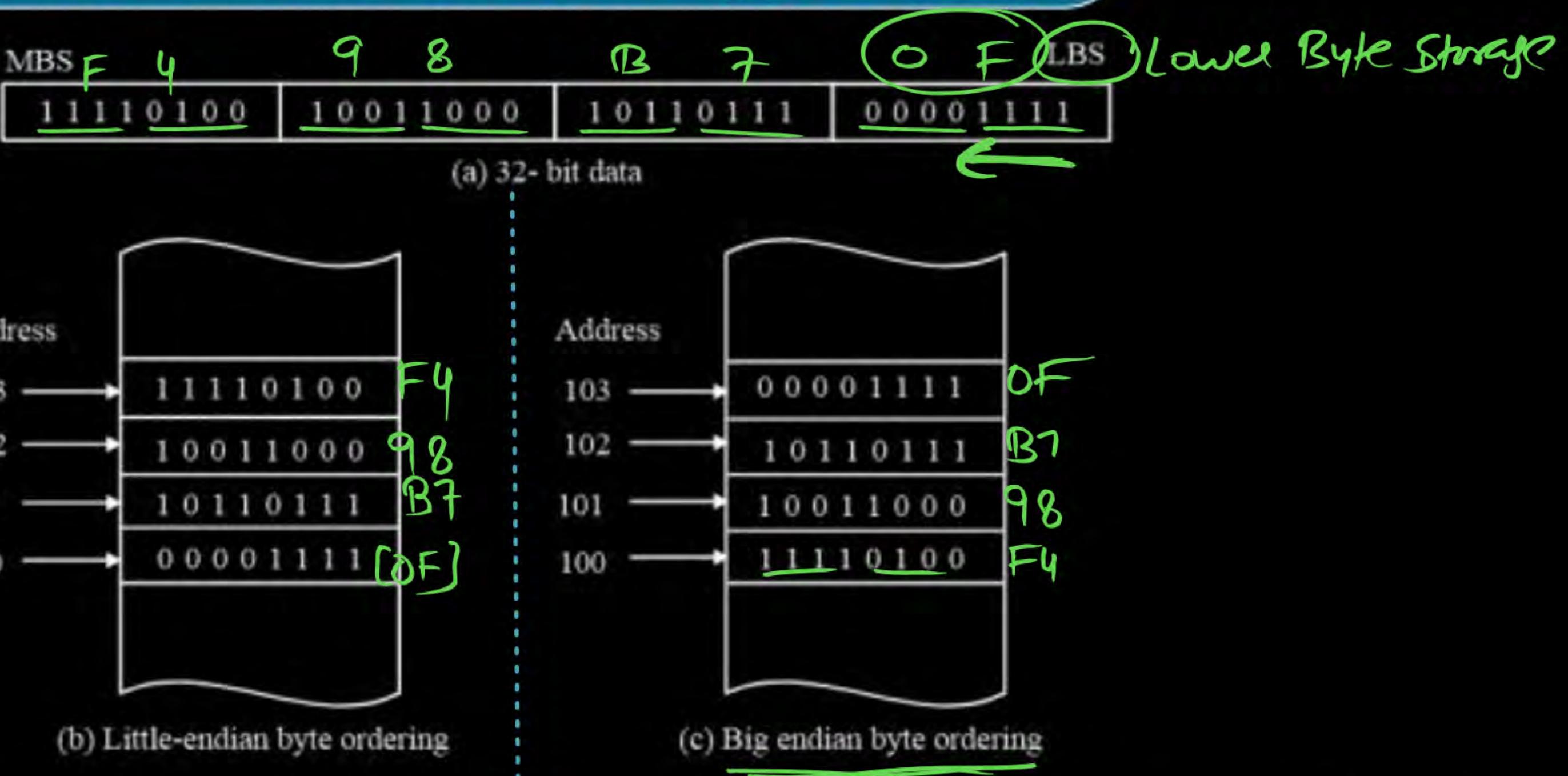
$$\boxed{\delta_2 = 51 21 15 25 H} \text{ Ans}$$

② Big Endian

	403	402	401	400
$\delta_2 =$	25	15	21	51

$$\boxed{\delta_2 = 25 15 21 51 H} \text{ Ans}$$

Two Important Memory Design Issues



Two byte ordering schemes commonly used by computer systems.

**THANK
YOU!**

