



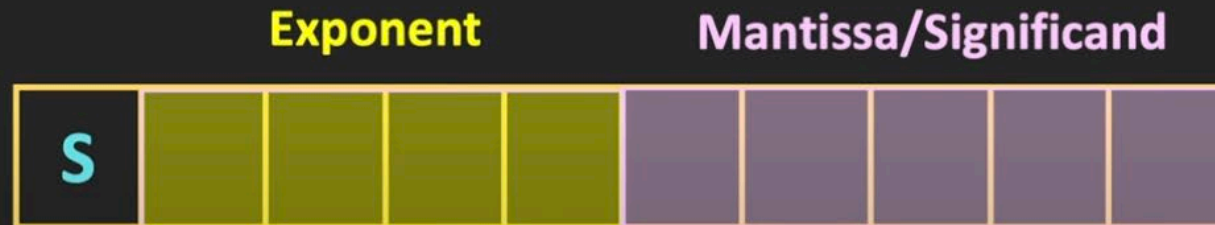
IEEE 754 Standard Single Precision Format



$$\pm 1.BBB \times 2^{\pm \text{exp}}$$

↑ ↑ ↑
Sign Fraction Exponent

No. of bits used to store the Floating Point Number
No. of reserved bits for Exponent and Mantissa
Format used for storing Exponent and Mantissa



- **IEEE 754 Standard**
 - **Half Precision (16 bits)**
 - **Single Precision (32 bits)**
 - **Double Precision (64 bits)**
 - **Quadruple Precision (128 bits)**
 - **Octuple Precision (256 bits)**

- **IEEE 754 Standard**
 - Half Precision (16 bits)
 - Single Precision (32 bits)
 - Double Precision (64 bits)
 - Quadruple Precision (128 bits)
 - Octuple Precision (256 bits)

IEEE 754 – Single Precision Format

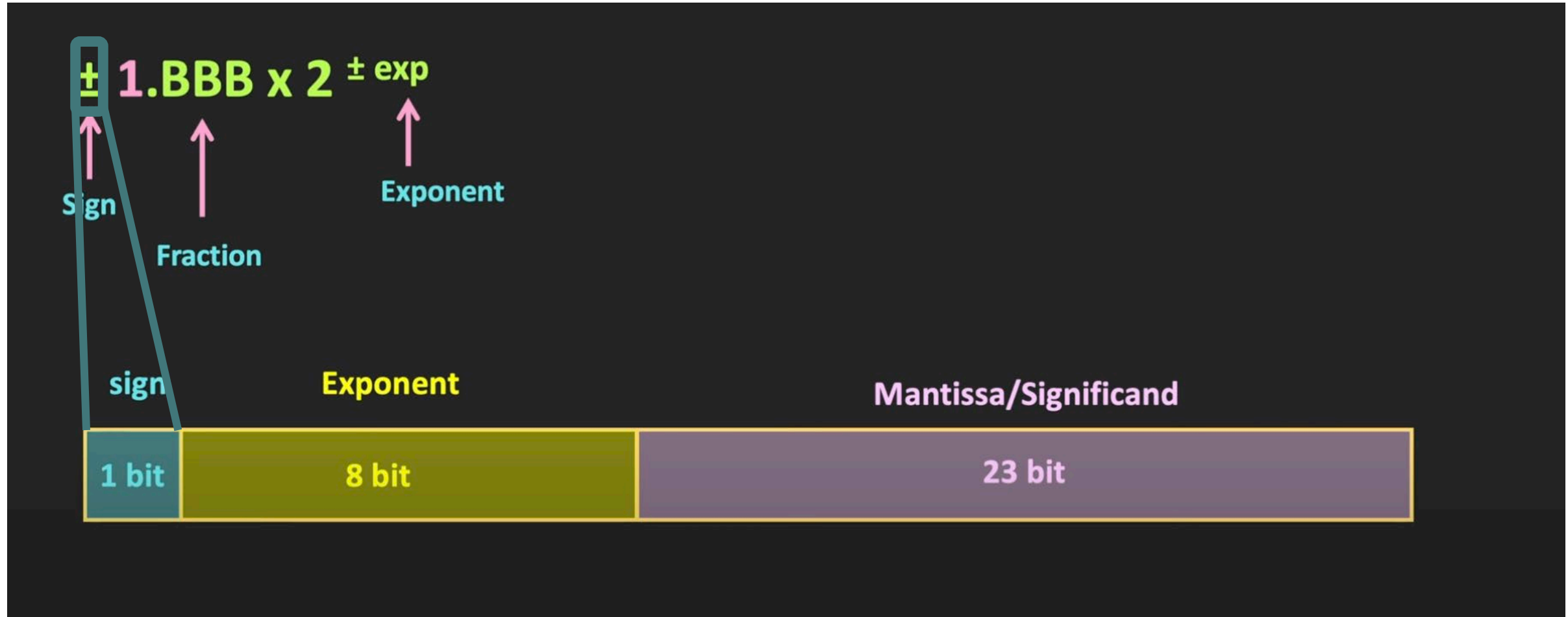


$$\pm 1.BBB \times 2^{\pm \text{exp}}$$

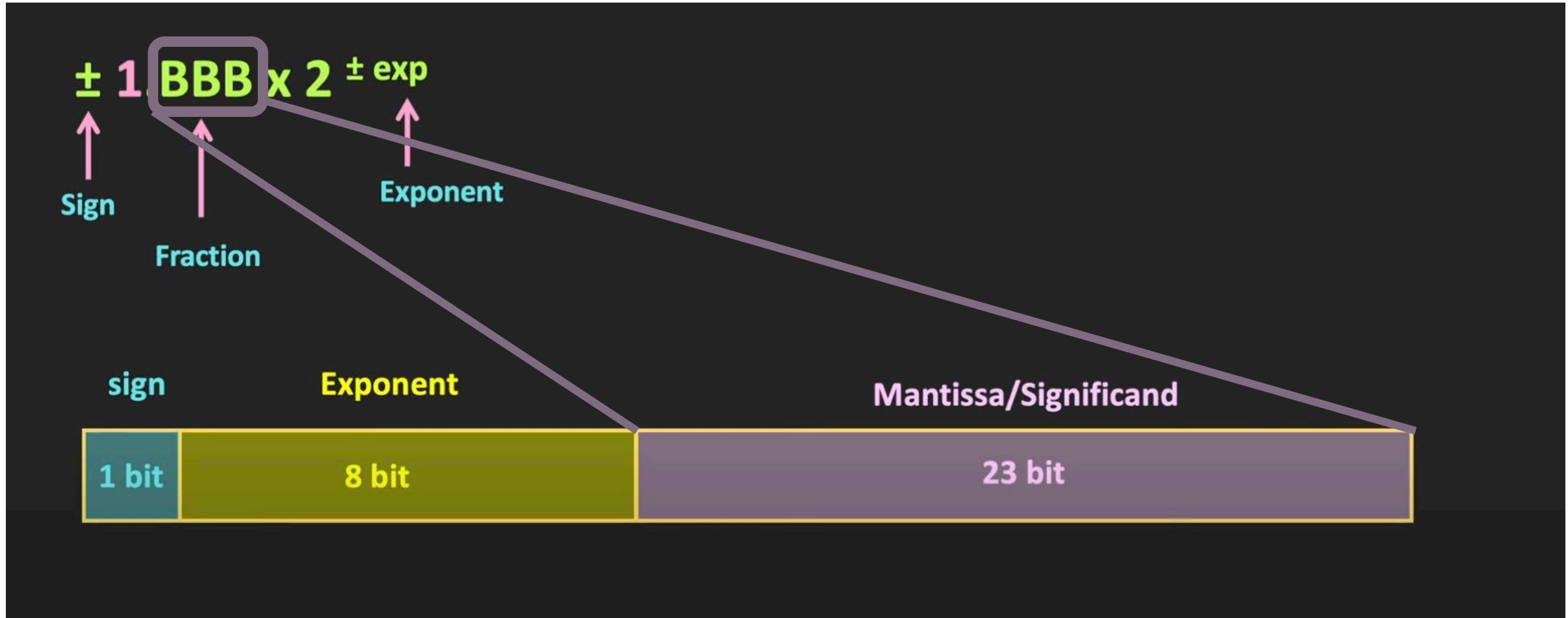
↑ ↑ ↑
Sign Fraction Exponent



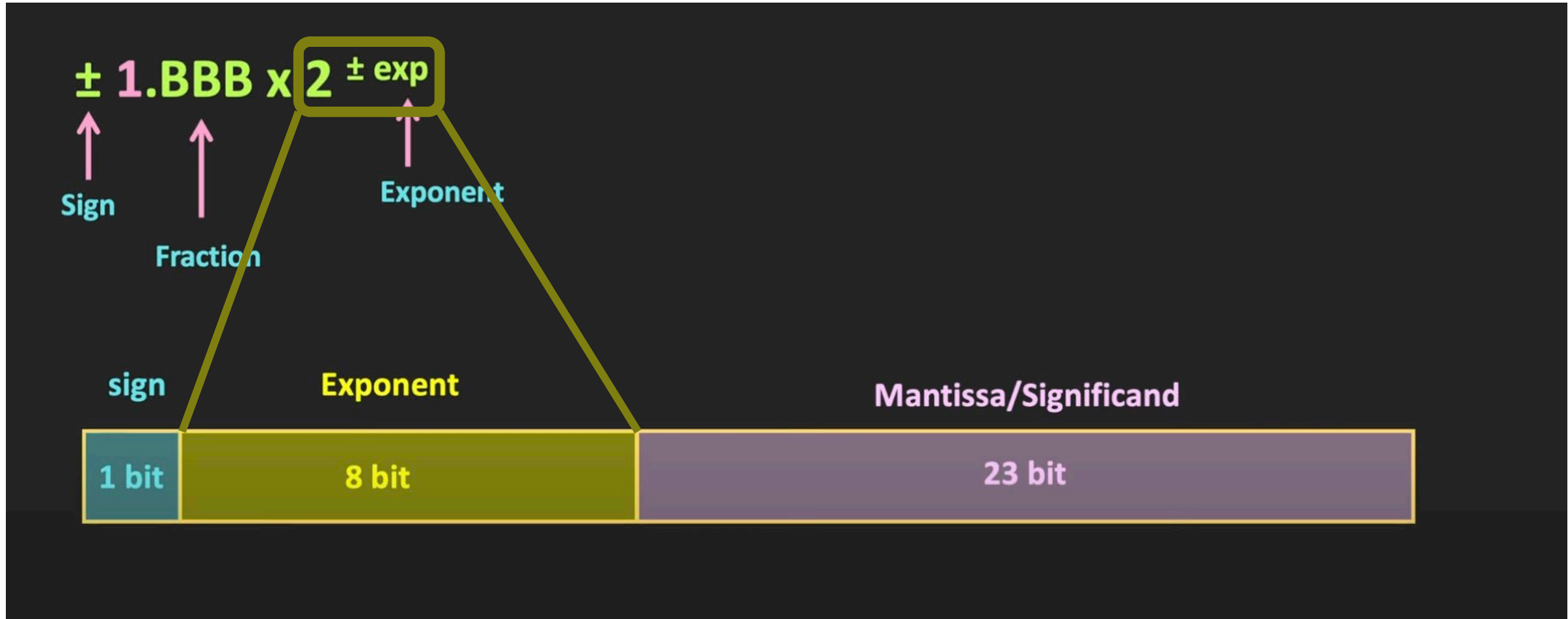
IEEE 754 – Single Precision Format



IEEE 754 – Single Precision Format



IEEE 754 – Single Precision Format



IEEE 754 – Single Precision Format



$$\pm 1.BBB \times 2^{\pm \text{exp}}$$

↑ Sign ↑ Fraction ↑ Exponent

Exponent
8 bits → 0 to 255 (unsigned)
How to represent negative numbers?



$$\pm 1.BBB \times 2^{\pm \text{exp}}$$

↑ Sign ↑ Fraction ↑ Exponent

Exponent

8-bits → 0 to 255 (unsigned)

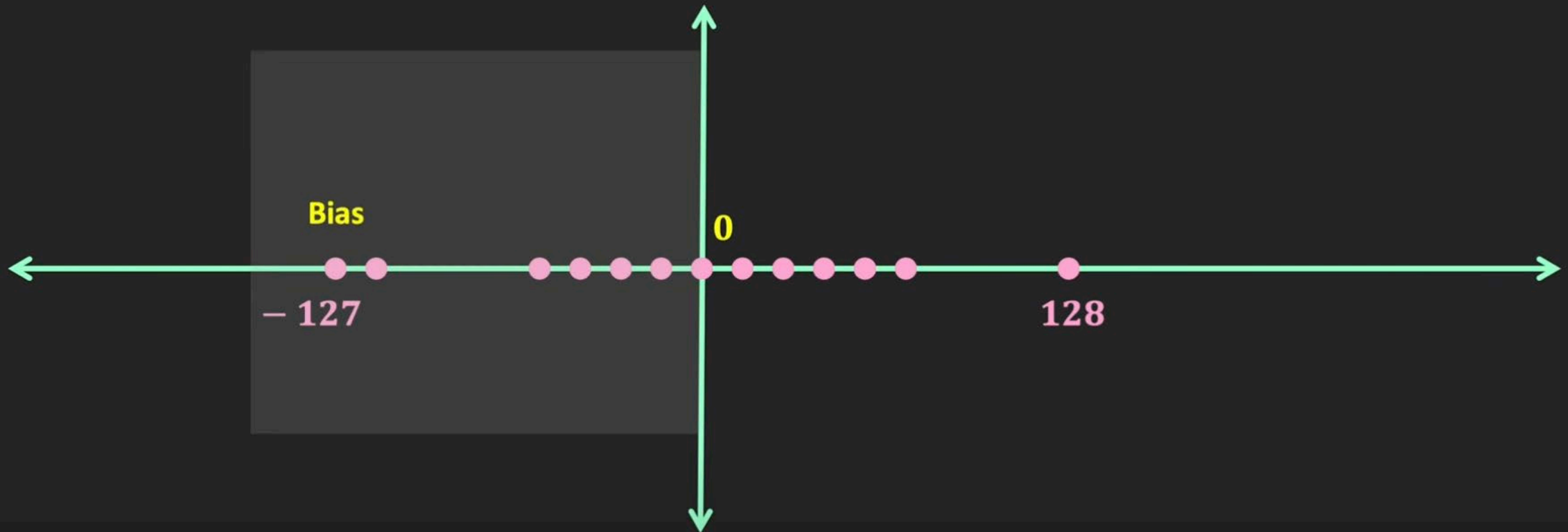
How to represent Negative exponent values ?

2's Complement

Signed Magnitude

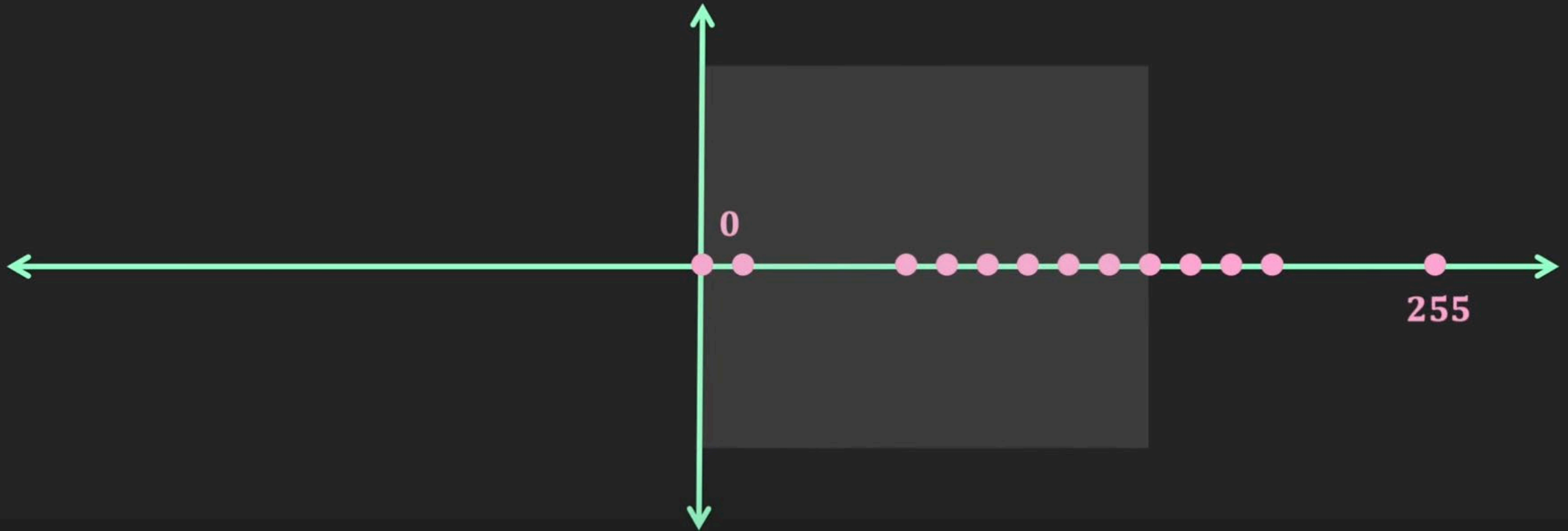
Biased Representation

Biased Representation



In Biased Representation, the bias or the fixed offset is added to the number in a such a way that the negative numbers get shifted to the positive side

Biased Representation



In Biased Representation, the bias or the fixed offset is added to the number in a such a way that the negative numbers get shifted to the positive side

Biased Representation



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

n - no of bits

8 bits  Bias = 127

Biased Representation



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

n - no of bits

8 bits \rightarrow Bias = 127

In biased Representation

0

+255

Biased Representation

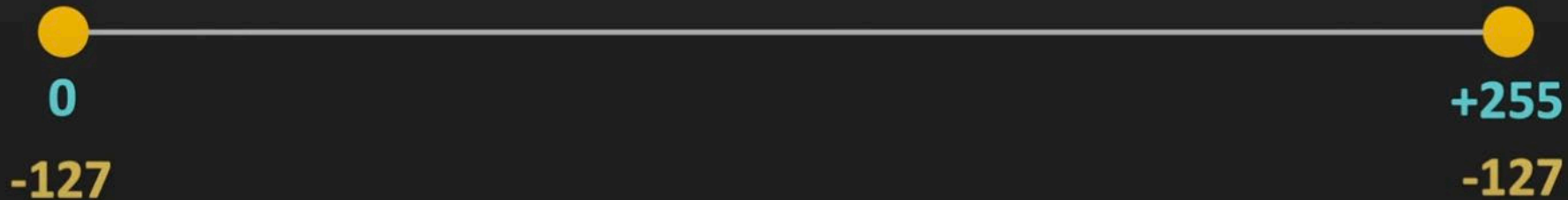


$$\text{Bias} = 2^{n-1} - 1$$

n - no of bits

8 bits \rightarrow Bias = 127

In biased Representation



Biased Representation



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

n - no of bits

8 bits \rightarrow Bias = 127

Actual Range

-127

+128

Biased Representation



Actual Number	Biased Number	Biased Representation
-127	0	0000 0000
-126	1	0000 0001
.....
-1	126	0111 1110
0	127	0111 1111
1	128	1000 0000
.....		
127	254	1111 1110
128	255	1111 1111

Biased Representation



Actual Number	Biased Number	Biased Representation
-127	0	0000 0000
-126	1	0000 0001
.....
-1	126	0111 1110
0	127	0111 1111
1	128	1000 0000
.....		
127	254	1111 1110
128	255	1111 1111

Special Values

Exponent Range

-126 to +127

Actual Number	Biased Number	Biased Representation
-127	0	0000 0000
-126	1	0000 0001
.....
-1	126	0111 1110
0	127	0111 1111
1	128	1000 0000
.....		
127	254	1111 1110
128	255	1111 1111

Special Values

Biased Representation



Exponent Range
-126 to +127

Actual Number	Biased Number	Biased Representation
-127	0	0000 0000
-126	1	0000 0001
.....
-1	126	0111 1110
0	127	0111 1111
1	128	1000 0000
.....		
127	254	1111 1110
128	255	1111 1111

Continuity



Biased Representation



Discontinuity

2's Complement

Number	2's Complement
-128	1000 0000
-127	1000 0001
-126	1000 0010
.....	
-1	1111 1111
0	0000 0000
1	0000 0001
.....	
126	0111 1110
127	0111 1111



Sign Magnitude

Number	Sign Magnitude
-127	1111 1111
-126	1111 1110
-125	1111 1101
.....	
-1	1000 0001
-0	1000 0000
+0	0000 0000
1	0000 0001
.....	
126	0111 1110
127	0111 1111



$$\pm 1.BBB \times 2^{\pm \text{exp}}$$

↑ ↑ ↑
Sign Fraction Exponent

What is the problem if the number representation is discontinue?

How to compare two numbers?



How to compare two numbers?

1. First compare the sign bit

$\pm 1.BBB \times 2^{\pm \text{exp}}$

↑ ↑ ↑

Sign Fraction Exponent



$$\pm 1.BBB \times 2^{\pm \text{exp}}$$

↑ ↑ ↑
Sign Fraction Exponent

How to compare two numbers?

1. First compare the sign bit
2. Then compare the Exponent

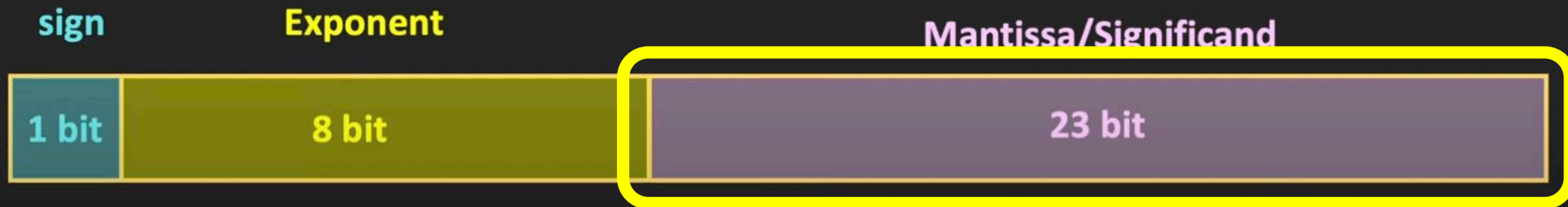


$$\pm 1.BBB \times 2^{\pm \text{exp}}$$

↑ ↑ ↑
Sign Fraction Exponent

How to compare two numbers?

1. First compare the sign bit
2. Then compare the Exponent
3. Finally, compare the Mantissa



Biased Representation



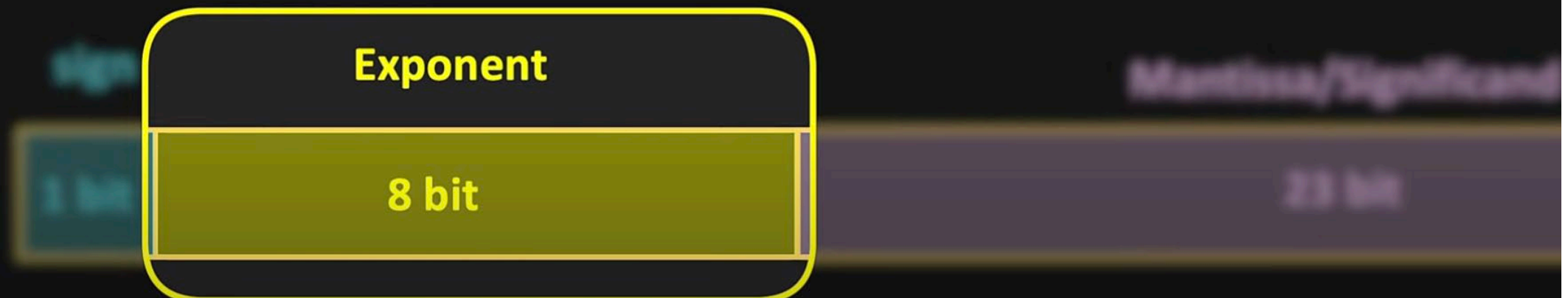
If there is a continuity in the representation of exponent value, it will be much easier to compare

$\pm 1.888 \times 2^{+exp}$

↑
Sign

↑
Fraction

↑
Exponent



Let's compare these two floating point numbers

Number 1

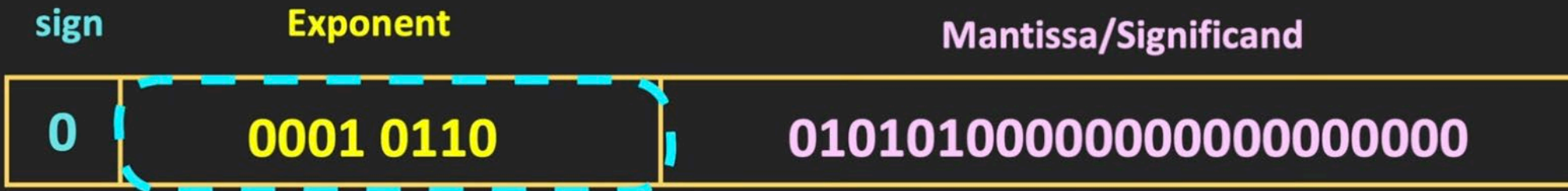
sign	Exponent	Mantissa/Significand
0	0001 0110	010101000000000000000000

Number 2

sign	Exponent	Mantissa/Significand
0	1011 0110	111101000000000000000000

By comparing the exponent parts, we can easily say that number 2 is greater than number 1

Number 1



Number 2



So, biased representation for is very useful

IEEE 754 uses this for floating point representation

**Let's see how to get the actual number from a IEEE 754
Single Precision format**

(Example 1)

Let's say this is a 32 bit number stored in the Single Precision format

sign	Exponent	Mantissa/Significand
0	1000 0101	001111000000000000000000

(Example 1)

Let's say this is a 32 bit number stored in the Single Precision format

sign	Exponent	Mantissa/Significand
0	1000 0101	001111000000000000000000

The MSB is 0. So this is a positive number

(Example 1)

Let's say this is a 32 bit number stored in the Single Precision format

sign	Exponent	Mantissa/Significand
0	1000 0101	001111000000000000000000

Actual value of the exponent

(Example 1)

Let's say this is a 32 bit number stored in the Single Precision format

sign	Exponent	Mantissa/Significand
0	1000 0101	001111000000000000000000

Actual value of the exponent

1000 0101  133

(Example 1)

Let's say this is a 32 bit number stored in the Single Precision format

sign	Exponent	Mantissa/Significand
0	1000 0101	001111000000000000000000

Actual value of the exponent

1000 0101 \longrightarrow 133

Actual Exponent = $133 - 127 = 6$

Since the number is stored using biased format, we need to subtract the bias to get the actual value

(Example 1)

Let's say this is a 32 bit number stored in the Single Precision format

sign	Exponent	Mantissa/Significand
0	1000 0101	001111000000000000000000

Exponent: 2^6

(Example 1)

Let's say this is a 32 bit number stored in the Single Precision format

sign	Exponent	Mantissa/Significand
0	1000 0101	001111000000000000000000

Exponent: 2^6

In normalized binary form, there is a 1 before the Mantissa:

1. 001111000000000000000000

(Example 1)

Let's say this is a 32 bit number stored in the Single Precision format

sign	Exponent	Mantissa/Significand
0	1000 0101	001111000000000000000000

Exponent: 2^6

In this fractional part, we can remove all the zeros from the right side

1. 001111

Significand

(Example 1)

Let's say this is a 32 bit number stored in the Single Precision format

sign	Exponent	Mantissa/Significand
0	1000 0101	001111000000000000000000

Exponent: 2^6

In this fractional part, we can remove all the zeros from the right side

Actual normalized binary number:

$$1.001111 \times 2^6$$

1.001111

Significand

(Example 1)

Let's say this is a 32 bit number stored in the Single Precision format

sign	Exponent	Mantissa/Significand
0	1000 0101	001111000000000000000000

Exponent: 2^6

Significand: 1.001111

Normalized binary number: 1.001111×2^6

Actual binary number: 1001111

(Example 1)

Let's say this is a 32 bit number stored in the Single Precision format

sign	Exponent	Mantissa/Significand
0	1000 0101	001111000000000000000000

Exponent: 2^6

Significand: 1.001111

Normalized binary number: 1.001111×2^6

Actual binary number: 1001111 \longrightarrow $(79)_{10}$

(Example 2)

sign	Exponent	Mantissa/Significand
1	1000 0011	110011000000000000000000

?

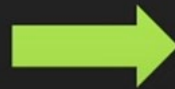
(Example 2)

sign	Exponent	Mantissa/Significand
1	1000 0011	110011000000000000000000

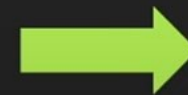
2^4

1.110011

1.110011×2^4



11100.11



$(-28.75)_{10}$

Let's try to represent a number

(Example 3)

$(12.625)_{10}$

(Example 3)

$(12.625)_{10}$ $(12)_{10} \rightarrow (1100)_2$ And $(.625)_{10} \rightarrow (101)_2$

(Example 3)

$(12.625)_{10}$ $(12)_{10} \rightarrow (1100)_2$ And $(.625)_{10} \rightarrow (101)_2$

$(1100.101)_2$

$\pm 1.BBB \times 2^{\pm \text{exp}}$

↑ ↑ ↑

Sign Fraction Exponent

(Example 3)

$(12.625)_{10}$ $(12)_{10} \rightarrow (1100)_2$ And $(.625)_{10} \rightarrow (101)_2$

$(1100.101)_2$

$\pm 1.BBB \times 2^{\pm \text{exp}}$

↑ ↑ ↑

Sign Fraction Exponent

$(1100.101)_2 = 1.100101 \times 2^3$

(Example 3)

1.100101×2^3

sign

Exponent

Mantissa/Significand

--	--	--

(Example 3)

+1.100101 x 2³



Sign

sign

Exponent

Mantissa/Significand

sign	Exponent	Mantissa/Significand
0		

(Example 3)

~~+~~**X.100101** x 2³

↑
Sign

sign

Exponent

Mantissa/Significand

0		
---	--	--

(Example 3)

~~+~~**X.100101** x 2³



Sign



Fraction

sign

Exponent

Mantissa/Significand

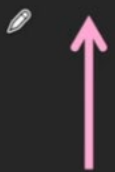
0		100101
---	--	--------

(Example 3)

~~+~~**X.100101** x 2³



Sign



Fraction

sign

Exponent

Mantissa/Significand

0		100101
---	--	--------

(Example 3)

~~+~~**X.100101** x 2³

↑
Sign

↑
Fraction

sign

Exponent

Mantissa/Significand

0

100101000000000000000000

(Example 3)

~~+~~**X.100101** x **2³**

↑
Sign

↑
Fraction

↑
Exponent

$$\begin{aligned}\text{Stored Exponent} &= \text{Actual Exponent} + \text{Bias} \\ &= 3 + 127 \\ &= 130\end{aligned}$$

sign

Exponent

Mantissa/Significand

0	1000 0010	100101000000000000000000
---	-----------	--------------------------

IEEE 754 – Single Precision Format Range



Largest Value of Exponent : + 127

Smallest Value of Exponent : - 126

IEEE 754 – Single Precision Format Range



1. bbbbbb.....bbb $\times 2^{127}$

1. bbbbbb.....bbb $\times 2^{-126}$

IEEE 754 – Single Precision Format Range



$$1.111111\dots111 \times 2^{127}$$

$$1.000000\dots000 \times 2^{-126}$$

Largest Value

* all the bits should be 1

Smallest Value

* all the bits should be 0

IEEE 754 – Single Precision Format Range



$$1.111111\dots111 \times 2^{127} \longrightarrow (2 - 2^{-23}) \times 2^{127} \approx 3.4 \times 10^{38}$$

$$1.000000\dots000 \times 2^{-126} \longrightarrow (1) \times 2^{-126} \approx 1.1 \times 10^{-38}$$

Single Precision Format (32 bit)

Largest Number $\approx 3.4 \times 10^{38}$

Smallest Number $\approx 1.1 \times 10^{-38}$

32 bit Fixed Point Representation (Signed Integer)

Largest Positive Number $\approx 2.1 \times 10^9$

Single Precision Format (32 bit)

Largest Number $\approx 3.4 \times 10^{38}$

Why it covers greater range ?

Smallest Number $\approx 1.1 \times 10^{-38}$

Floating point allows range at the cost of precision



Single Precision Format (32 bit)

Largest Number $\approx 3.4 \times 10^{38}$

Smallest Number $\approx 1.1 \times 10^{-38}$

7 Significant Digits in Decimal



IEEE 754 – Single Precision Format vs Fixed Point



1235646.895 **X**

7 Significant Digits in Decimal



IEEE 754 – Single Precision Format vs Fixed Point



1235646.895 **X** 1.235646895×10^6

7 Significant Digits in Decimal



IEEE 754 – Single Precision Format vs Fixed Point



1235646.895 **X** 1.2356469×10^6

7 Significant Digits in Decimal



IEEE 754 – Double Precision Format



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

n - no of bits

11 bits → Bias = 1023

111111111111

000000000000

Reserved

Max. value of Exponent = 2046
Min. value of Exponent = 1



Actual max. value of Exponent = $2046 - 1023 = 1023$

Actual min. value of Exponent = $1 - 1023 = -1022$

Max. value of Exponent = 2046

Min. value of Exponent = 1



Actual max. value of Exponent = $2046 - 1023 = 1023$ 2^{1023}
Actual min. value of Exponent = $1 - 1023 = -1022$ 2^{-1022}

Smallest Positive Number $\approx 2.2 \times 10^{-308}$

Largest Positive Number $\approx 1.797 \times 10^{308}$



Thank you

Any Question?

