



IEEE 754 Standard Single Precision Format



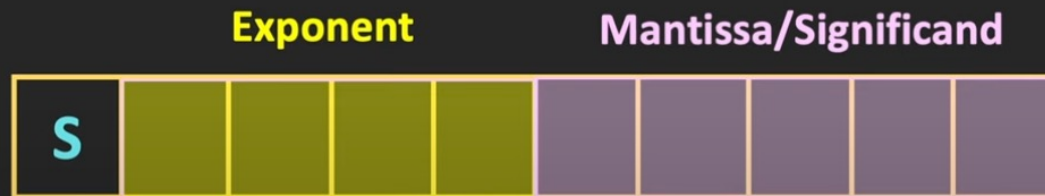
IEEE 754 Standard



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

↑ ↑ ↑
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



How Floating Point Numbers are Stored in Memory?



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

How Floating Point Numbers are Stored in Memory?



- **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

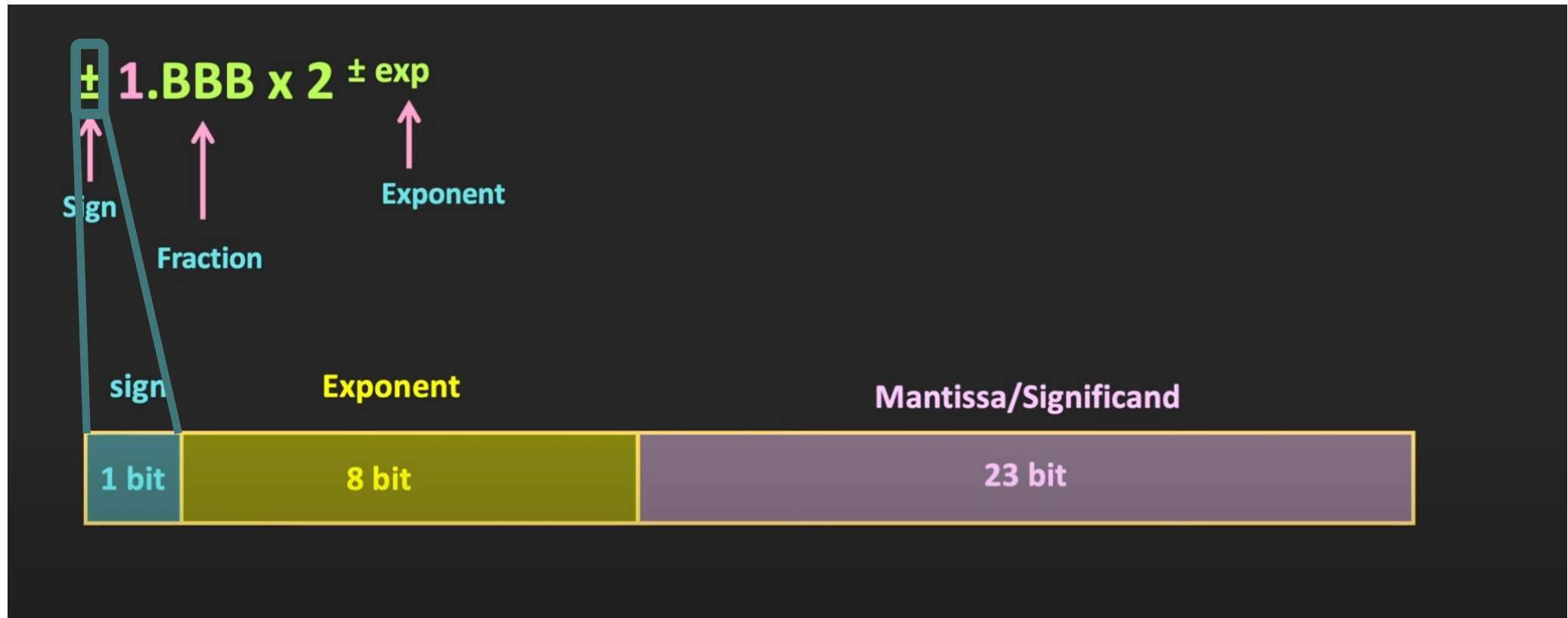


$$\pm 1.BBB \times 2^{\pm \text{exp}}$$

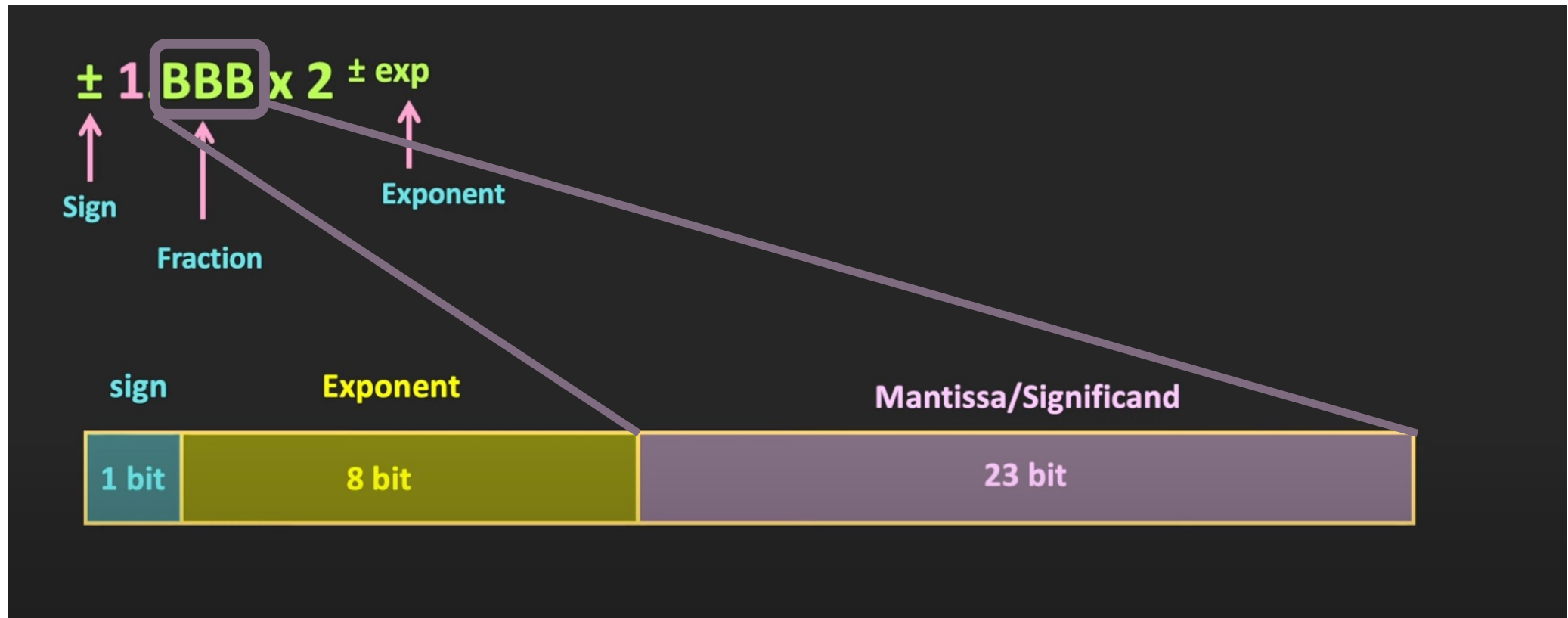
↑ ↑ ↑
Sign Fraction Exponent



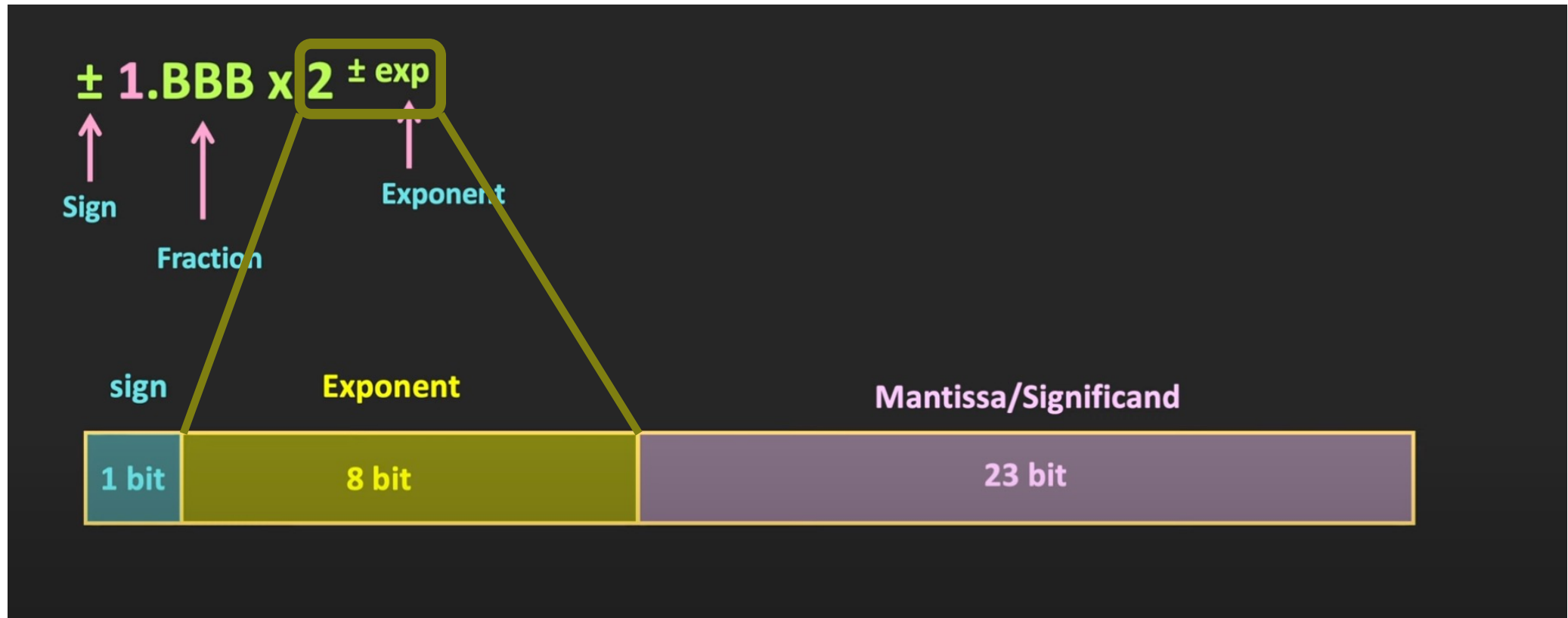
IEEE 754 – Single Precision Format



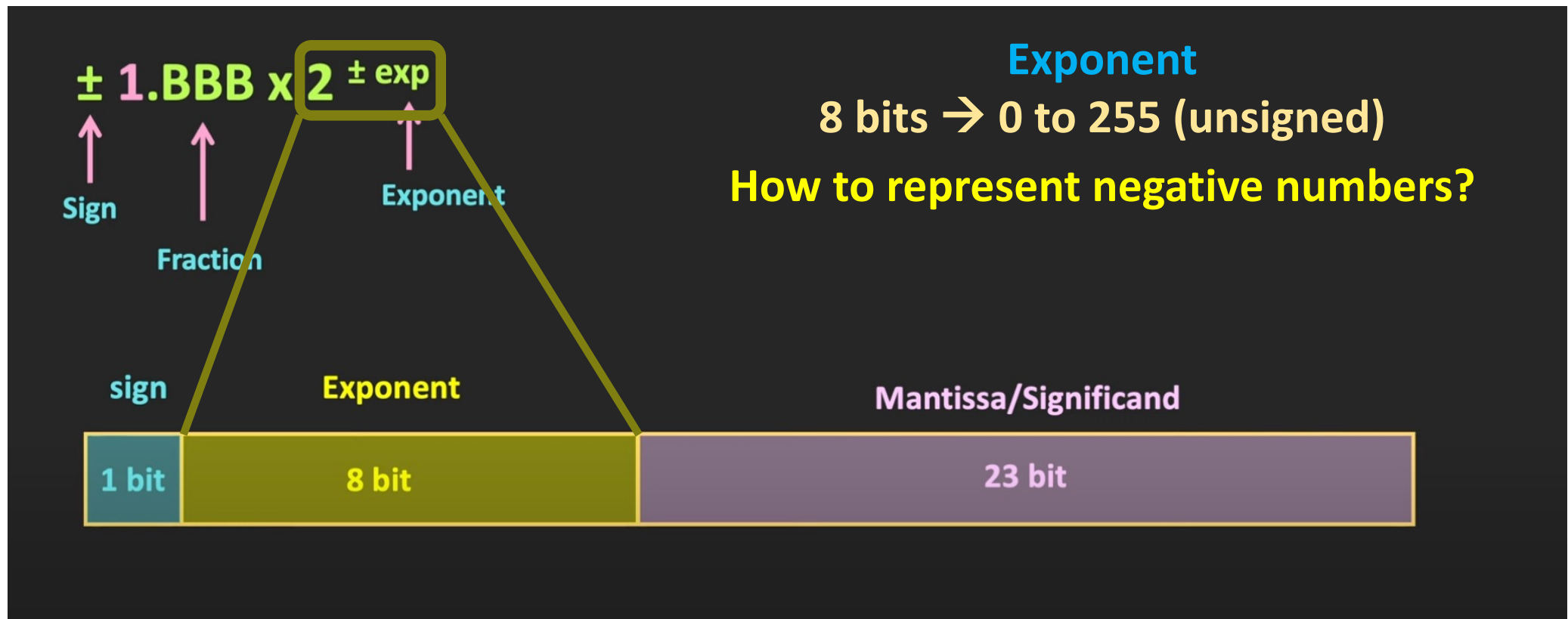
IEEE 754 – Single Precision Format



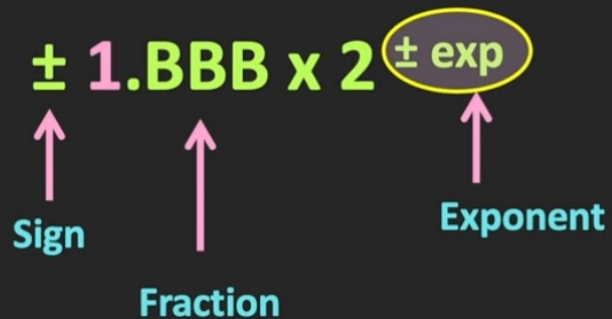
IEEE 754 – Single Precision Format



IEEE 754 – Single Precision Format



IEEE 754 – Single Precision Format



Exponent

8-bits → 0 to 255 (unsigned)

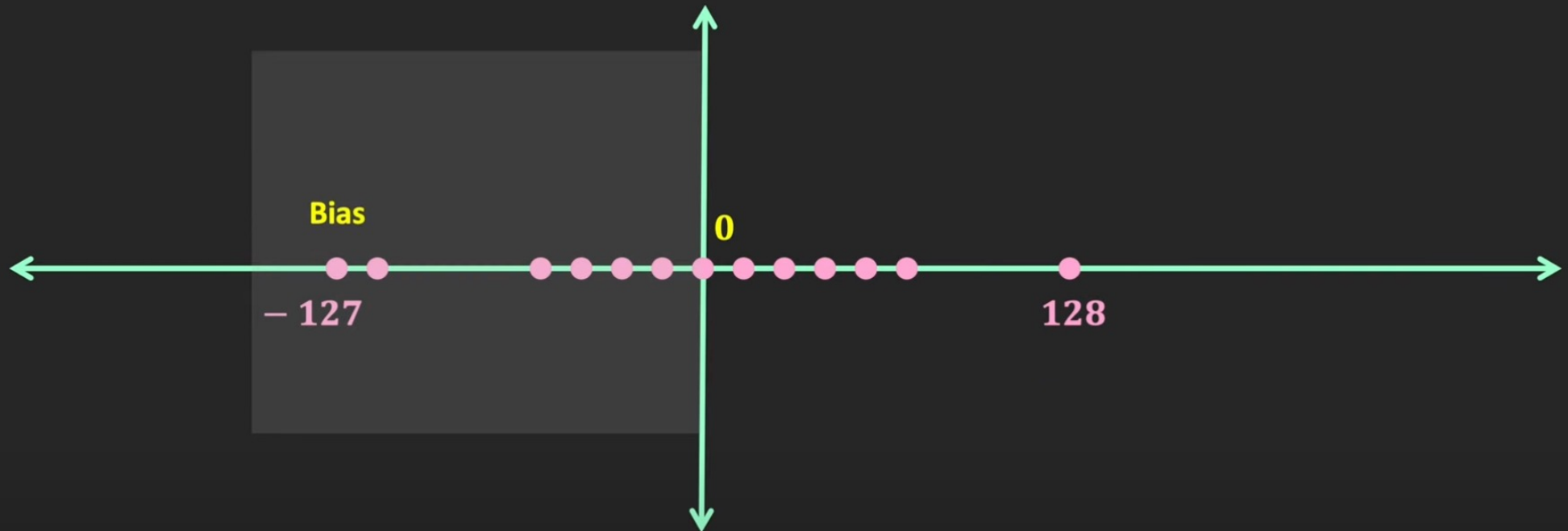
How to represent Negative exponent values ?

2's Complement

Signed Magnitude

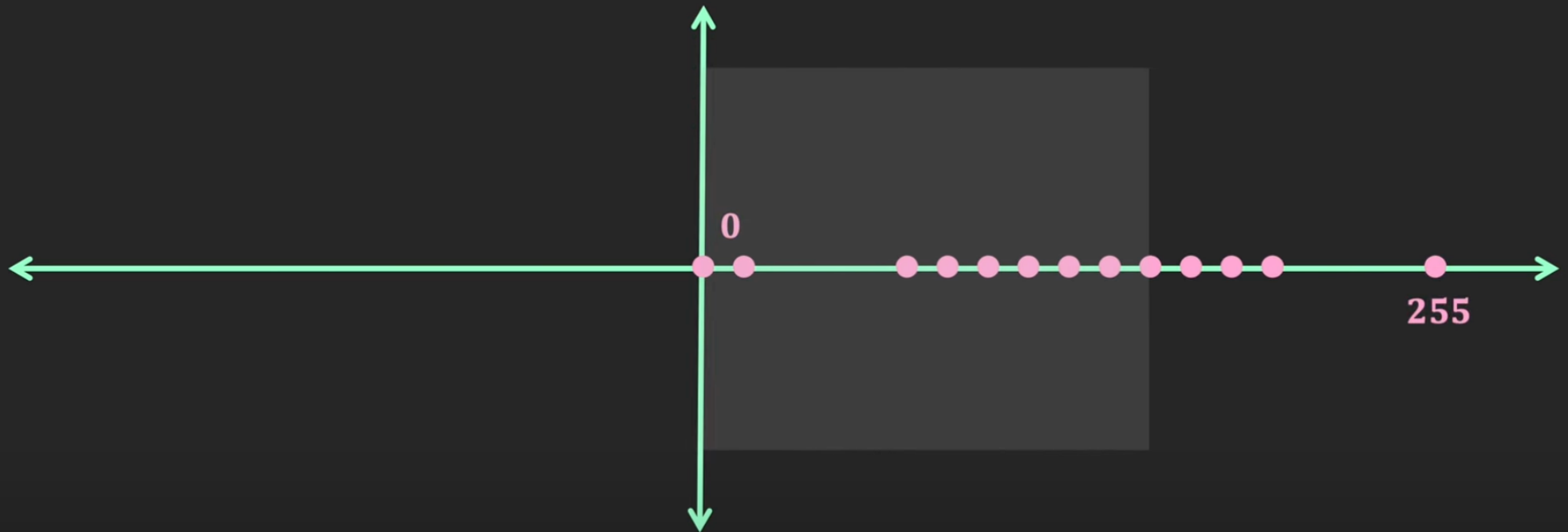
Biased Representation

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 \rightarrow 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

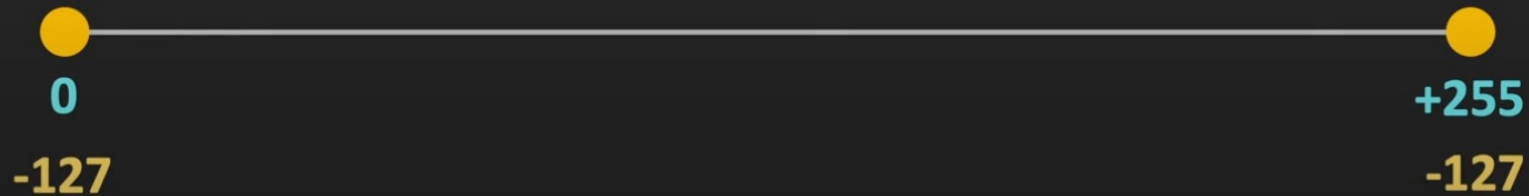


$$\text{Bias} = 2^{n-1} - 1$$

n - no of bits

8 bits \rightarrow Bias = 127

In biased Representation



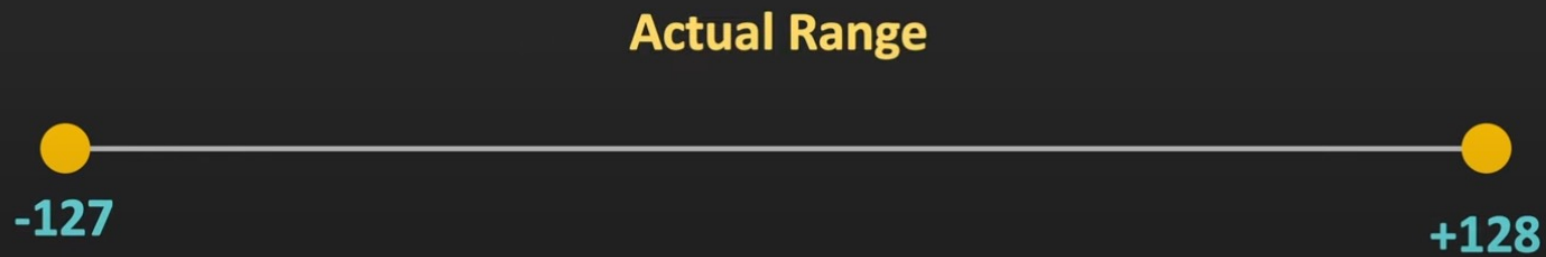
Biased Representation



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

n - no of bits

8 bits \rightarrow Bias = 127



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

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

Special Values

Biased Representation

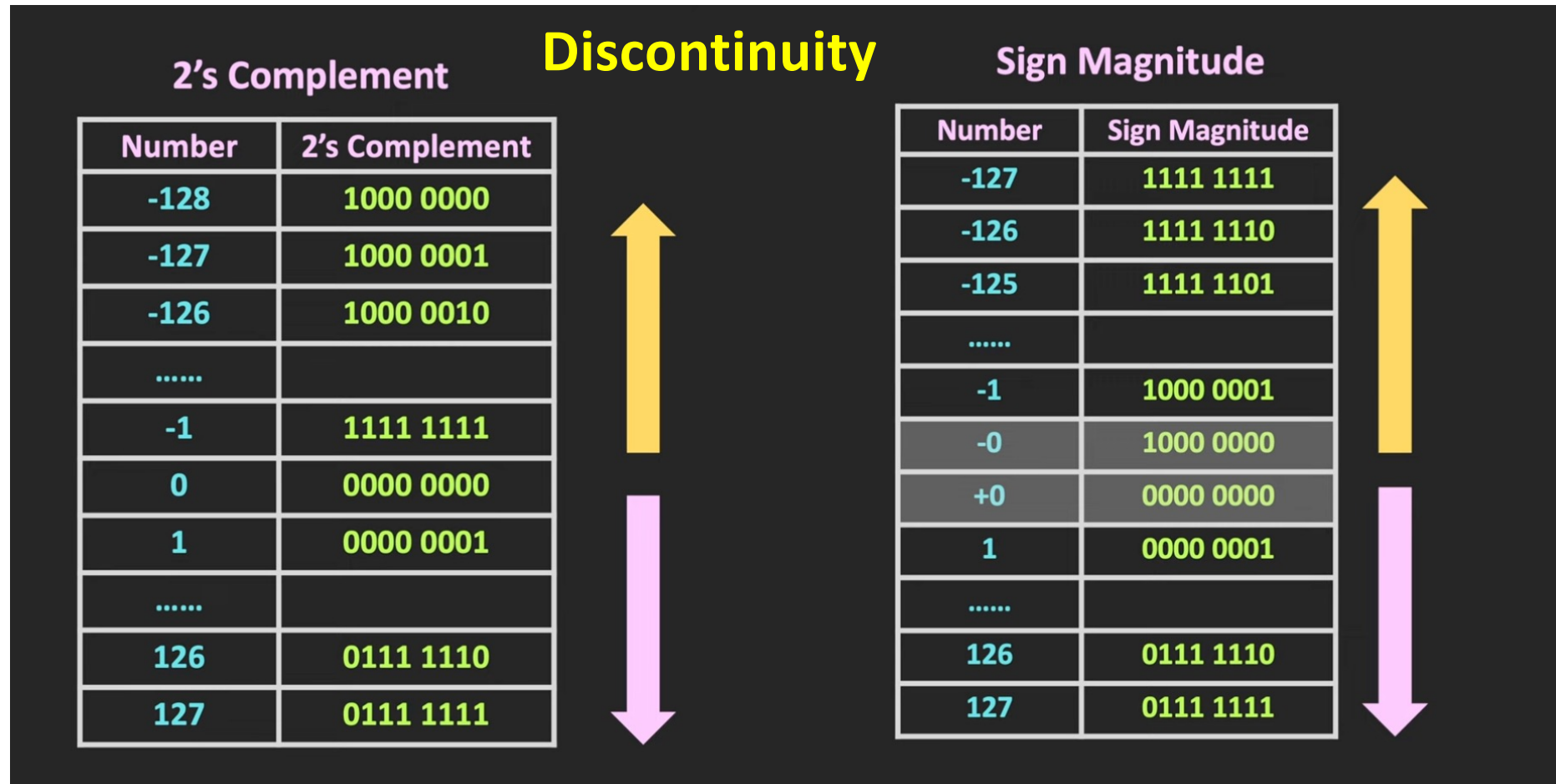


Exponent Range
-126 to +127

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



Biased Representation



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

↑ ↑ ↑
Sign Fraction Exponent

What is the problem if the number representation is discontinue?

How to compare two numbers?



Biased Representation



How to compare two numbers?

1. First compare the sign bit

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

↑ ↑ ↑
Sign Fraction Exponent



Biased Representation

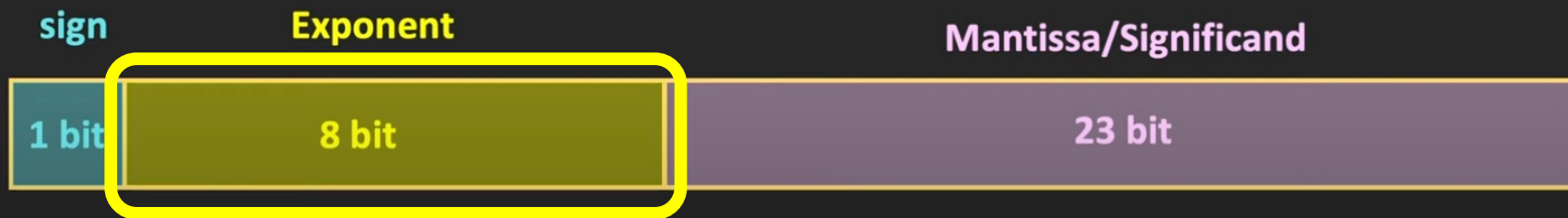


$$\pm 1.BBB \times 2^{\pm \text{exp}}$$

↑ ↑ ↑
Sign Fraction Exponent

How to compare two numbers?

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



Biased Representation



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

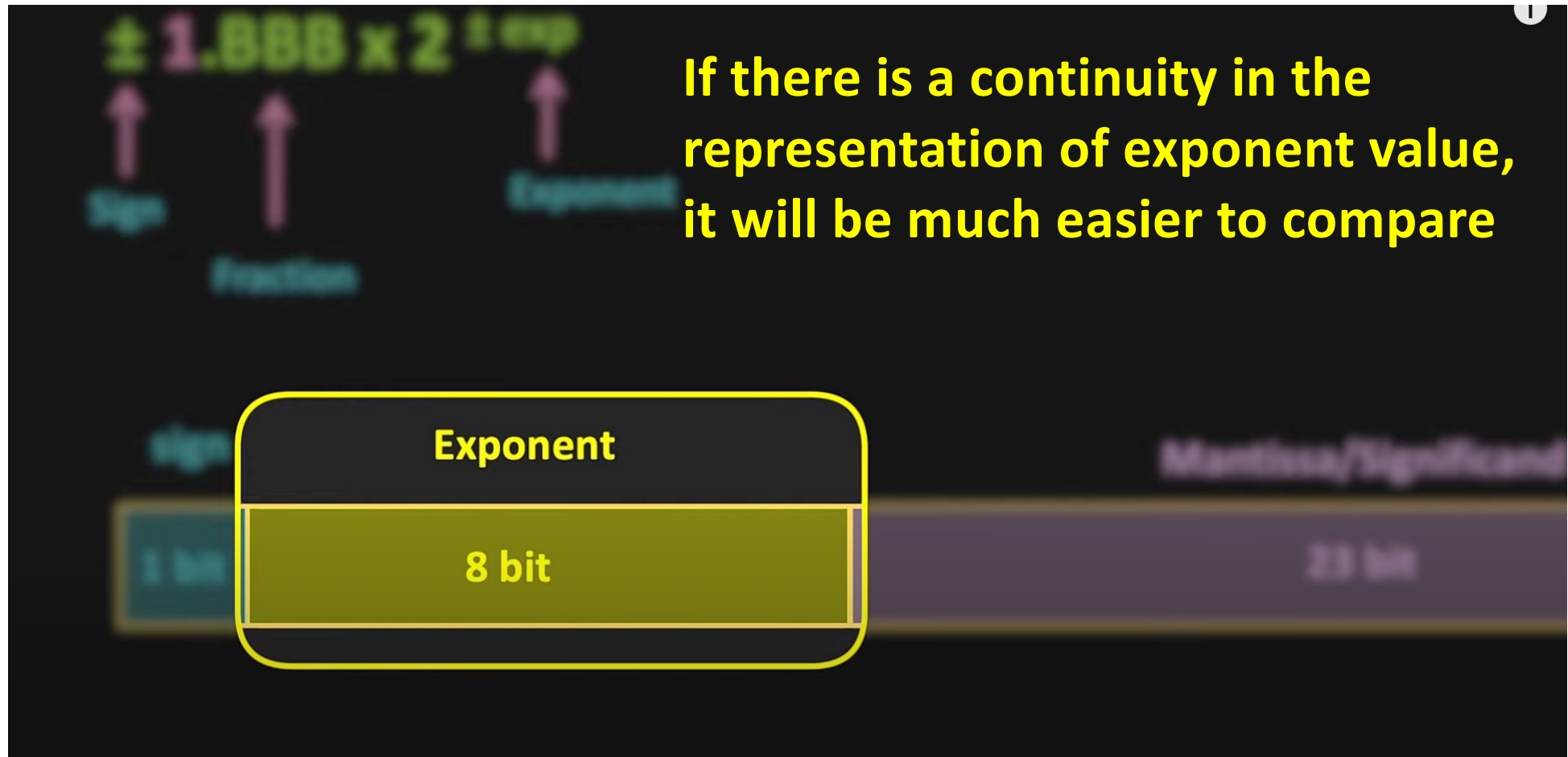
↑ ↑ ↑
Sign Fraction Exponent

How to compare two numbers?

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



Biased Representation



Biased Representation



Let's compare these two floating point numbers

Number 1

sign	Exponent	Mantissa/Significand
0	0001 0110	010101000000000000000000

Number 2

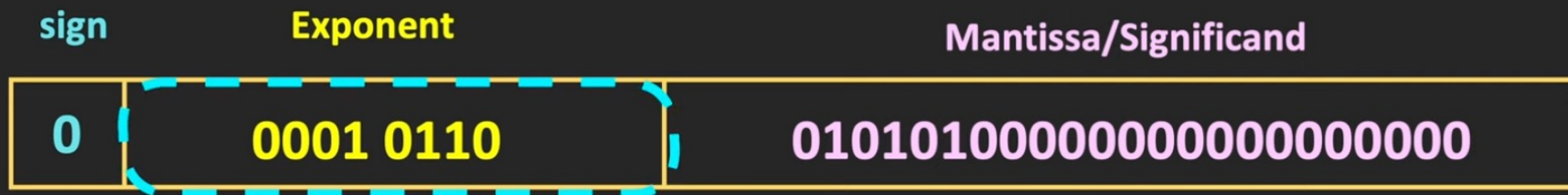
sign	Exponent	Mantissa/Significand
0	1011 0110	111101000000000000000000

Biased Representation

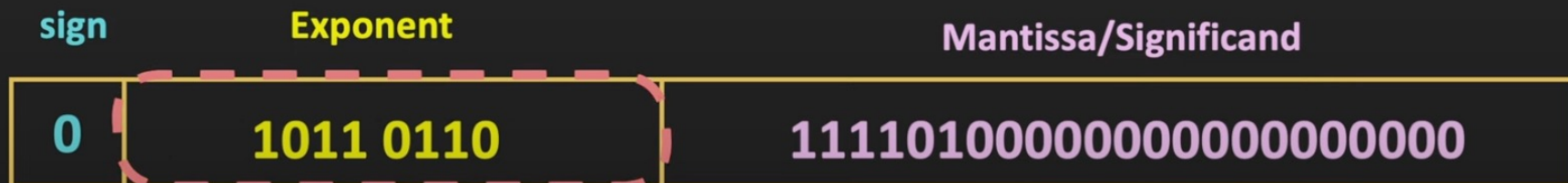


By comparing the exponent parts, we can easily say that number 2 is greater than number 1

Number 1



Number 2



Biased Representation



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**

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

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

The MSB is 0. So this is a positive number

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

Actual value of the exponent

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

Actual value of the exponent

1000 0101 \longrightarrow 133

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

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

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

Exponent: 2^6

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

Exponent: 2^6

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

1. 001111000000000000000000

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

Exponent: 2^6

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

1. 001111

Significand

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

Exponent: 2^6

Actual normalized
binary number:

$$1.001111 \times 2^6$$

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

1.001111
Significand

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

Exponent: 2^6

Significand: 1.001111

Normalized binary number: 1.001111×2^6

Actual binary number: 1001111

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

Exponent: 2^6

Significand: 1.001111

Normalized binary number: 1.001111×2^6

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

IEEE 754 – Single Precision Format



(Example 2)

sign	Exponent	Mantissa/Significand
1	1000 0011	110011000000000000000000

?

IEEE 754 – Single Precision Format



Let's try to represent a number

IEEE 754 – Single Precision Format



(Example 3)

$(12.625)_{10}$

IEEE 754 – Single Precision Format



(Example 3)

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

IEEE 754 – Single Precision Format



(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

IEEE 754 – Single Precision Format



(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

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

IEEE 754 – Single Precision Format



(Example 3)

1.100101 x 2³

sign

Exponent

Mantissa/Significand

--	--	--

IEEE 754 – Single Precision Format



(Example 3)

+1.100101 x 2³

↑
Sign

sign

Exponent

Mantissa/Significand

0

IEEE 754 – Single Precision Format



(Example 3)

~~+~~**X.100101** x 2³

↑
Sign

sign

Exponent

Mantissa/Significand

0		
---	--	--

IEEE 754 – Single Precision Format



(Example 3)

~~+~~**X.100101** x 2³

↑
Sign

↑
Fraction

sign

Exponent

Mantissa/Significand

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

IEEE 754 – Single Precision Format



(Example 3)

~~+~~**X.100101** x 2³

↑
Sign

↑
Fraction

sign

Exponent

Mantissa/Significand

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

IEEE 754 – Single Precision Format



(Example 3)

~~+~~**X.100101** x 2³

↑
Sign

↑
Fraction

sign

Exponent

Mantissa/Significand

0

100101000000000000000000

IEEE 754 – Single Precision Format



(Example 3)

$+ \cancel{X}.100101 \times 2^3$

↑ ↑ ↑
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



Thank you

Any Question?

