# Quantization in llms:

Quantization is the process of converting the data from higher memory format to lower memory format.



*Figure 1: Quantization of Numbers from FP32 into INT8*

# Why do we need Quantization ?

We need quantization to store high-parameter models in our RAM. Without the concept of quantization, it's impossible to store higher GB models in our computer RAM or GPU RAM.

By default, the weights and biases (parameters) of the LLM are in memory format of 32bit or FP32 datatype. Suppose we want to load a language model with 7 billion parameters into our RAM, let's see how much memory in RAM we require:

**Lets Visualize It**:

- 1 GB = 8 × 10^9 bits
- 1 bit = 1/8 × 10^9 GB

**Each parameter requires 32 bits of memory, so:**

- 1 parameter = 32 bits

**For 7 billion parameters:**

- 7 billion parameters = 7 × 10^9

**Memory required for 7 billion parameters in bits:**

$$= (7 \times 10^9) \times 32 \text{ bits}$$
$$= 2.24 \times 10^{11} \text{ bits}$$

**To convert 2.24 × 10^11 bits into gigabytes (GB):**

$$= \frac{2.24 \times 10^{11} \text{ bits}}{8 \times 10^9 \text{ bits/GB}}$$
$$= 28 \text{ GB}$$

We need 28 GB of RAM to store the model for inference which is impossible for normal use case computer devices.

# Some Concepts and Terminologies in Quantization and Memory Structures :
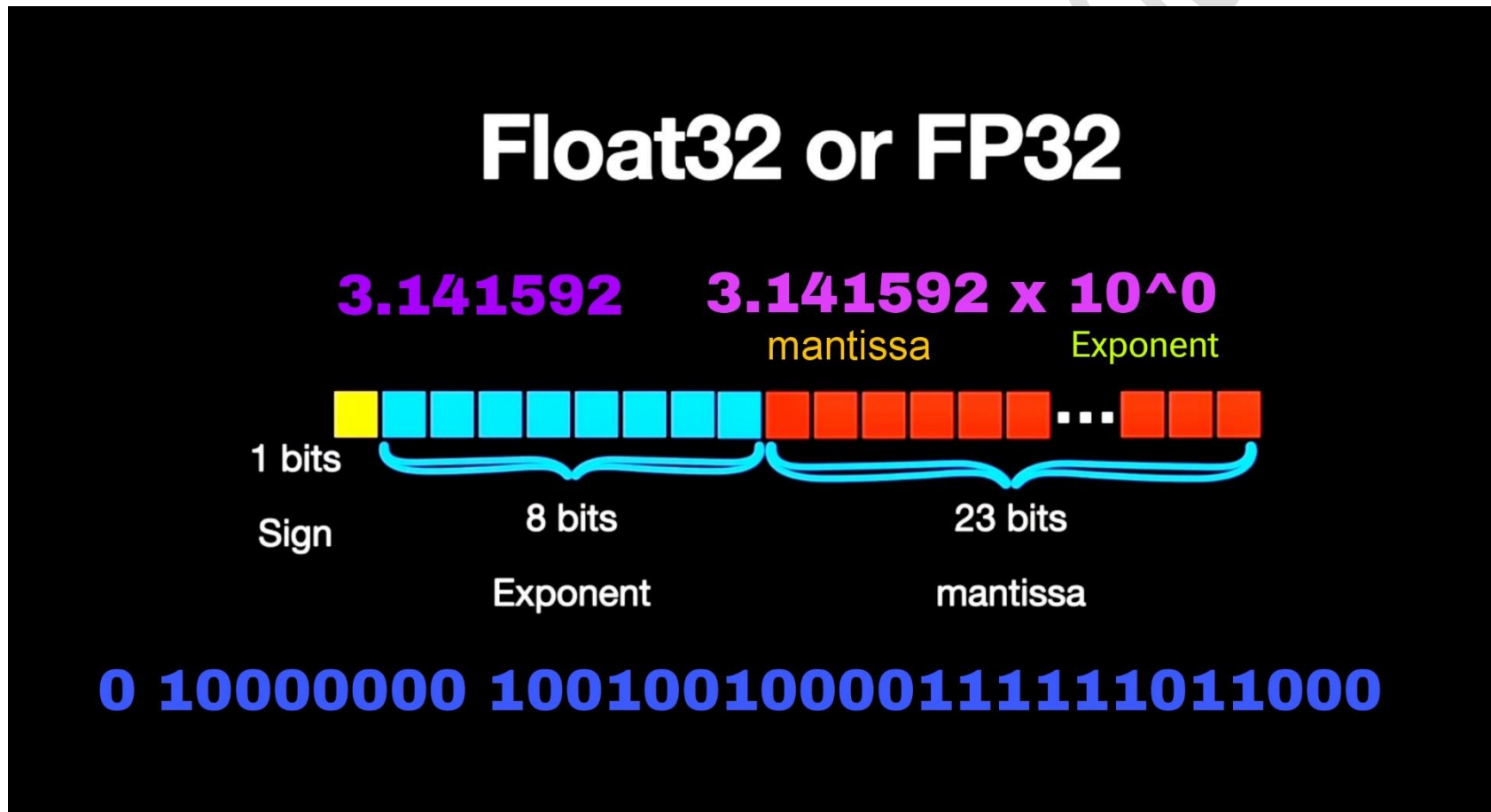
## Types of Memory Formats:

- **FP32 :** (Floating Point 32), also known as Single/Full Precision, uses 32 bits of memory.

- **FP16 :** (Floating Point 16), also known as Half Precision, uses 16 bits of memory.

- **BF16 :** (Brain Floating Point 16), also known as Half Precision, similar to FP16, also utilizes 16 bits of memory but is optimized for specific computational tasks, particularly in artificial intelligence and machine learning applications.

- **Int8 :** (Integer 8) uses 8 bits of memory for numbers. It is useful for applications where memory efficiency and a relatively small range of values are important, such as in neural networks quantization and certain signal processing tasks.

- **FP4 :** (Floating Point 4) uses 4 bits of memory to store floating numbers.

- **NF4 :** (Normalized Float 4) normalizes floating-point numbers and uses 4 bits of memory to store them. NF4 performs better than FP4 experimentally.
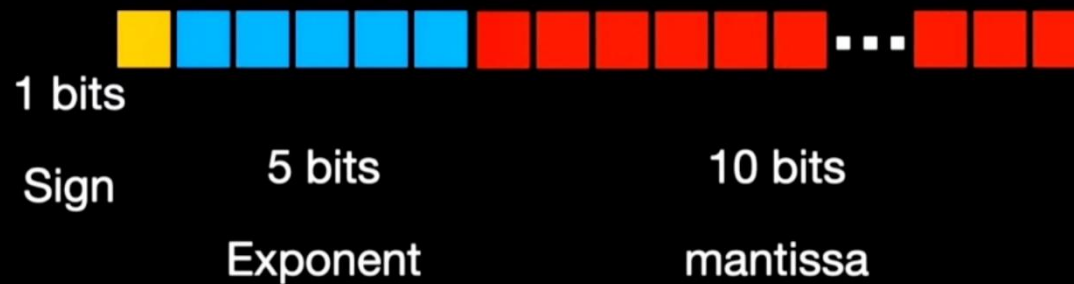
# Lets see how does this memory Format Store data :

*FP32 (Floating Point 32) :*

## FP16 (Floating Point 16) :

**BF16  (Brain Floating Point 16) :**
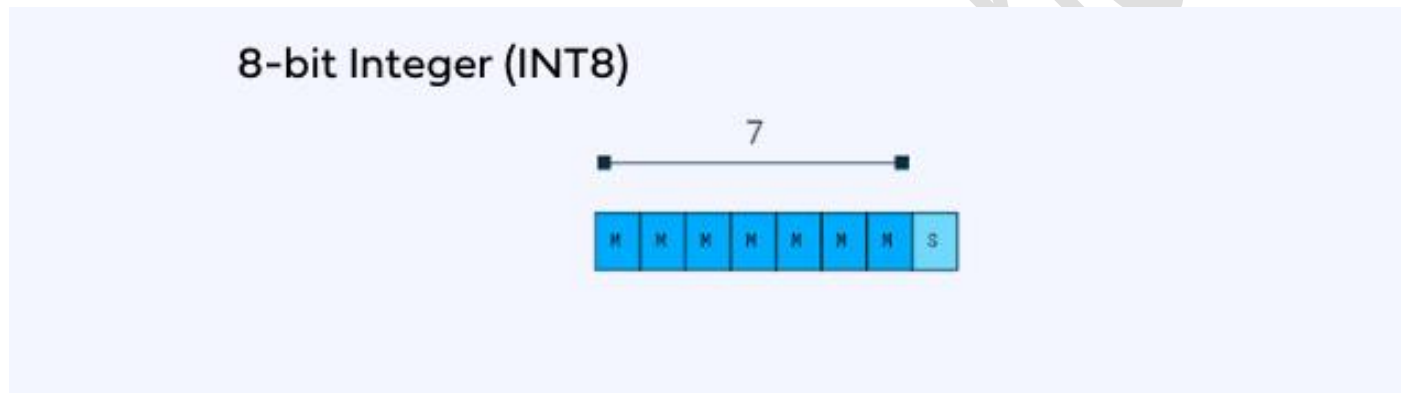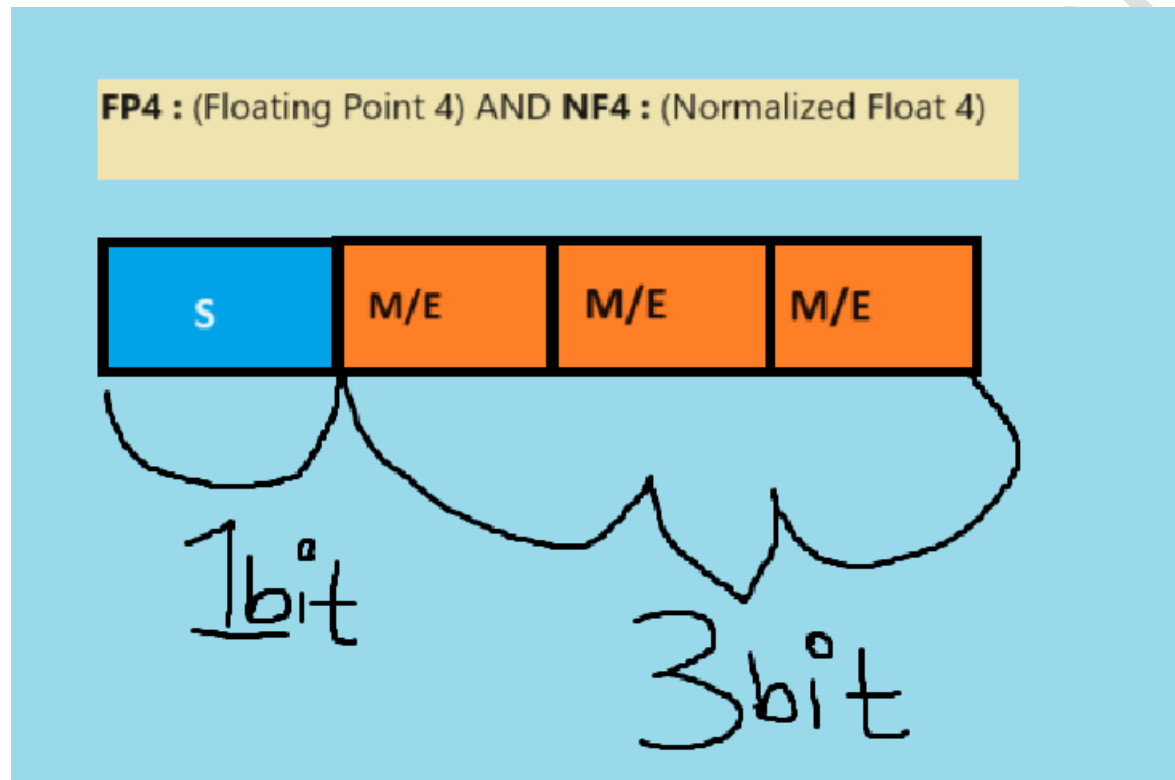
## Int8 (Integer 8) :



**Figure 2: Here m is for mantissa and s is for sign**

*FP4 (Floating Point 4) and NF4 (Normalized Float 4):*



In the FP4 and NF4 memory formats, one bit is reserved for the sign, and 3 bits are used by the mantissa or exponent. It's not fixed; sometimes, 2 bits are taken by the mantissa, and the remaining one bit is used by the exponent. There is no fixed format, memory try itself best combinations of different mantissa/exponent combinations.

# Now Lets see how Quantization Occurs Mathematically:

How to perform Quantization:

1) Symmetric Quantization:

In this quantization the min value of our set is perfectly align with the min value of quantization set.
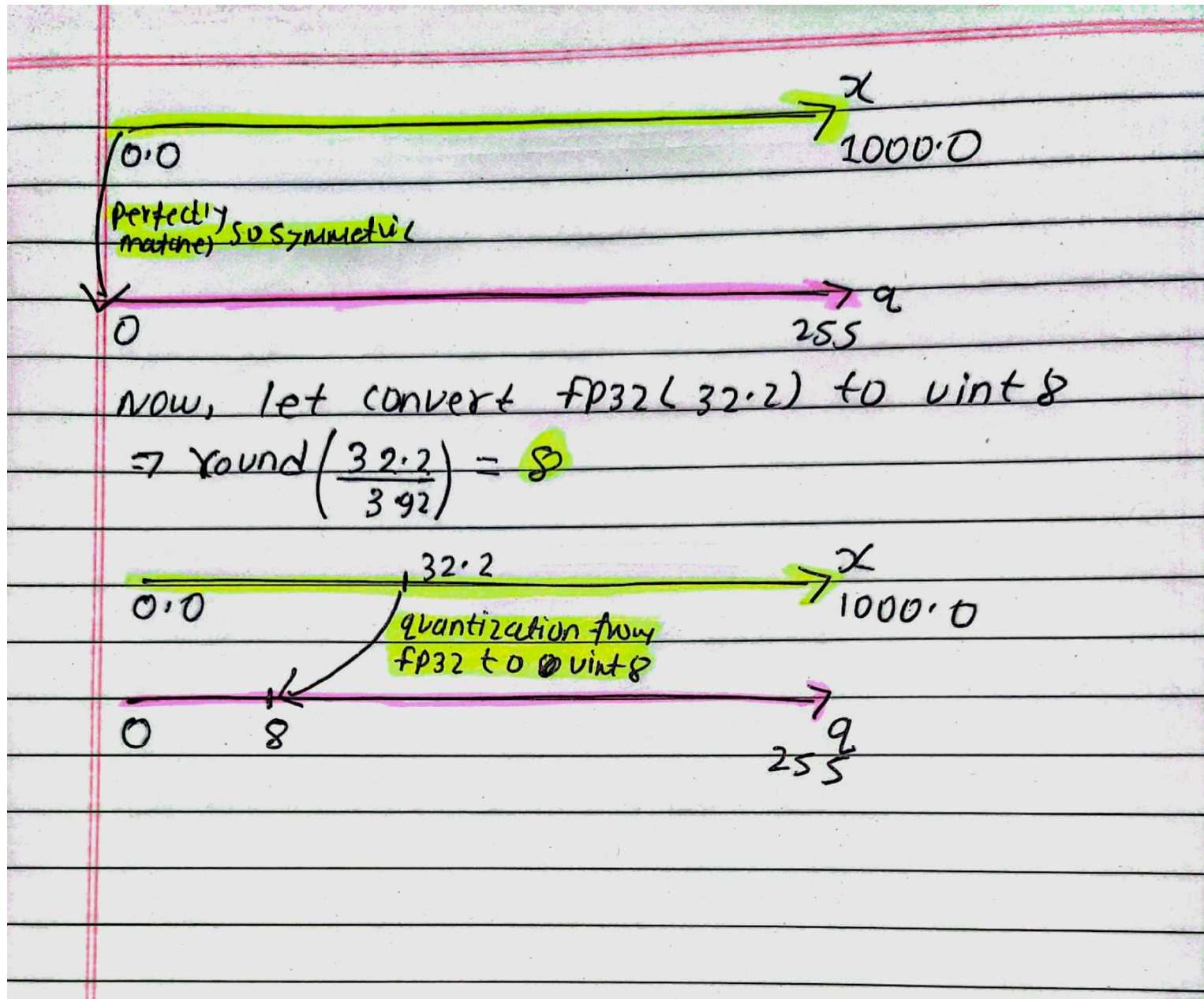
let's perform quantization:

$[0.0 — 1000.0] \rightarrow uint8$

The range of number in uint8 is $[0-255]$

So,

$[0.0 — 1000.0] \rightarrow [0-255]$

   FP 32             uint8

                        $x$

$0.0$                     $1000.0$

                        $q$

$O$                     $255$

Now,

let's calculate scale factor first,

$$Scale = \frac{Xmax - Xmin}{a\,max - a\,min} = \frac{1000.0 - 0.0}{255 - 0} = 3.92$$

let's check It's symmetric or not

$$\Rightarrow round\left(\frac{min\,no\,of\,ourset}{scale\,factor}\right) + (no\,uf\,req\,no\,to\,get\,min\,value\,of\,quantization\,set)$$

$$\Rightarrow round\left(\frac{0.0}{3.92}\right) + (num)$$

$$= O + \boxed{O} \leftarrow O\,is\,added\,since\,we\,already\,got\,(q)min\,value.$$

$$= It\,is\,symmetric$$

$x$

0.0 ———————————————→ 1000.0

Perfectly
match(ed) so symmetric

0 ———————————————→ $q$

255

Now, let convert fp32 ( 32.2) to uint 8

$\Rightarrow$ round $\left(\dfrac{32.2}{3.92}\right) = 8$

32.2

0.0 ———————————————→ $x$

1000.0

quantization from
fp32 to uint8

0 ———8———————————→

255 $q$

Date ___
Page ___

## ② Asymmetric Uint 8 quantization:

In this quantization the $\wedge$ min value — of our

$$\frac{}{\text{min value}}$$ of

Set is not perfectly align with the $\frac{}{\text{min value}}$ of

the quantization set.

eg:

**x**

—————————— 1000.0

-20.0

—————————— 255 **q**

0

let's perform quantization:

$$\frac{[-20.0 \quad 1000.0]}{fp32} \longrightarrow \frac{[0 - 255]}{\text{uInt 8}}$$

x

————————→ 1000.0

-20.0

————————→ q²55

0

Now,

$$Scale = \frac{x_{max} - x_{min}}{q_{max} - q_{min}} = \frac{1000.0 - (-20.0)}{255 + 20.0} = 4.0$$
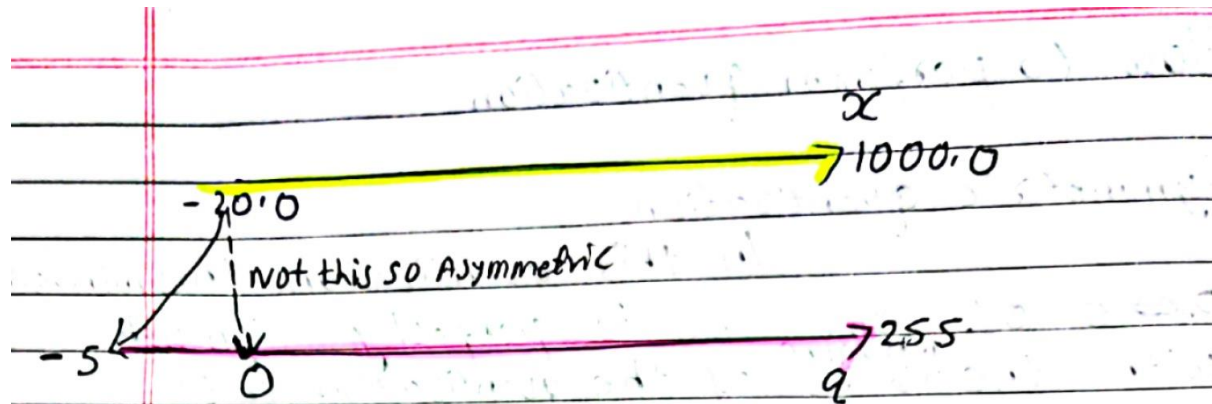
## lets check its is symmetric or not

$$\Rightarrow round\left(\frac{\text{min no of our set}}{\text{scale factor}}\right) \pm (\text{no that is need to get min value of quantization set})$$

$$\Rightarrow round\left(\frac{-20.0}{4.0}\right) \pm (\text{num})$$

$$\Rightarrow \qquad \boxed{-5} + \boxed{5} \quad (\text{adding +5 to get min value of quantization set} \therefore \text{zero point = 5})$$

Asymmetric

$x$

$\rightarrow 1000.0$

$-20.0$

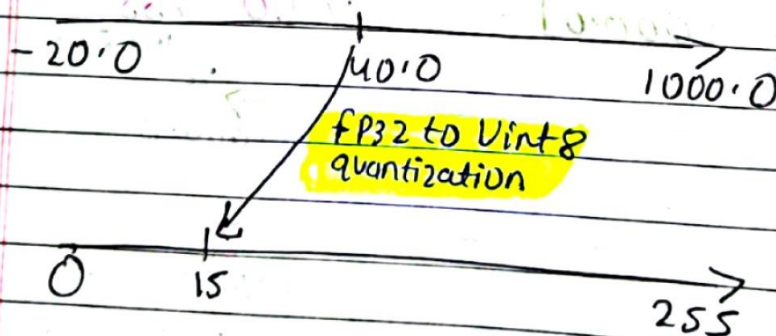| not this so Asymmetric

$\rightarrow 255$

$-5$     $0$                                    $q$

Now, let convert 40.0 (fP32) to uint8

$$\Rightarrow round\left(\frac{num}{scale factor}\right) + zeropoint$$

$$\Rightarrow round\left(\frac{40.0}{4.0}\right) + 5$$

$$\Rightarrow 15$$

## In figure

$-20.0$       $40.0$           $1000.0$

fP32 to Uint8
quantization

$0$     $15$

$255$

# Different ranges numbers that comes under different memory format

| Memory Format | Range of Values |
| --- | --- |
| uint8 | 0 to 255 |
| int8 | -128 to 127 |
| uint16 | 0 to 65,535 |
| int16 | -32,768 to 32,767 |
| uint32 | 0 to 4,294,967,295 |
| int32 | -2,147,483,648 to 2,147,483,647 |
| uint64 | 0 to 18,446,744,073,709,551,615 |
| int64 | -9,223,372,036,854,775,808 to 9,223,372,036,854,775,807 |
| float16 | Approximately $\pm6.1E-05$ to $\pm6.55E+04$ (3 significant digits) |
| float32 | Approximately $\pm1.4E-45$ to $\pm3.4E+38$ (7 significant digits) |
| float64 | Approximately $\pm5E-324$ to $\pm1.8E+308$ (15 significant digits) |