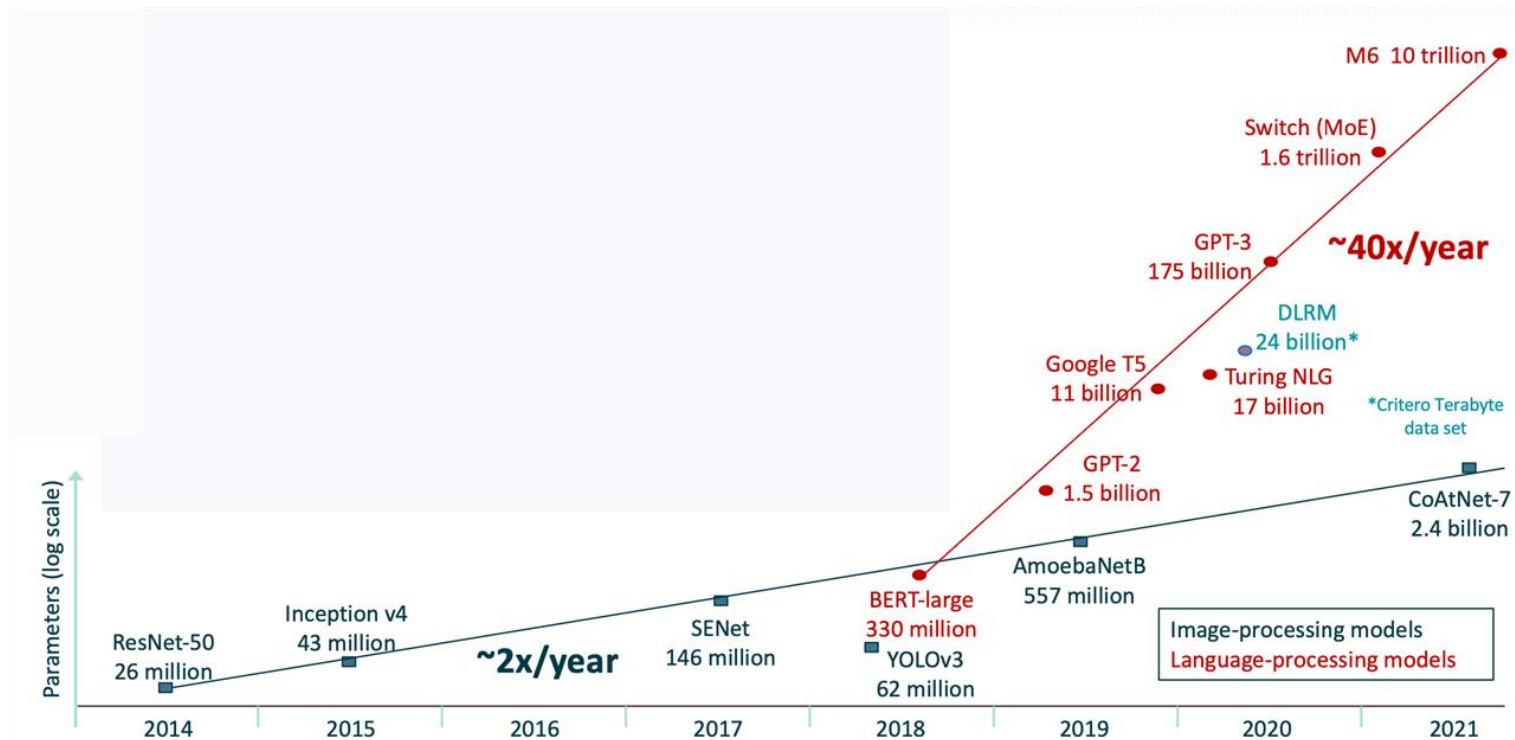


# Model Quantization

**Bach-Hoang Ngo**

# Motivation



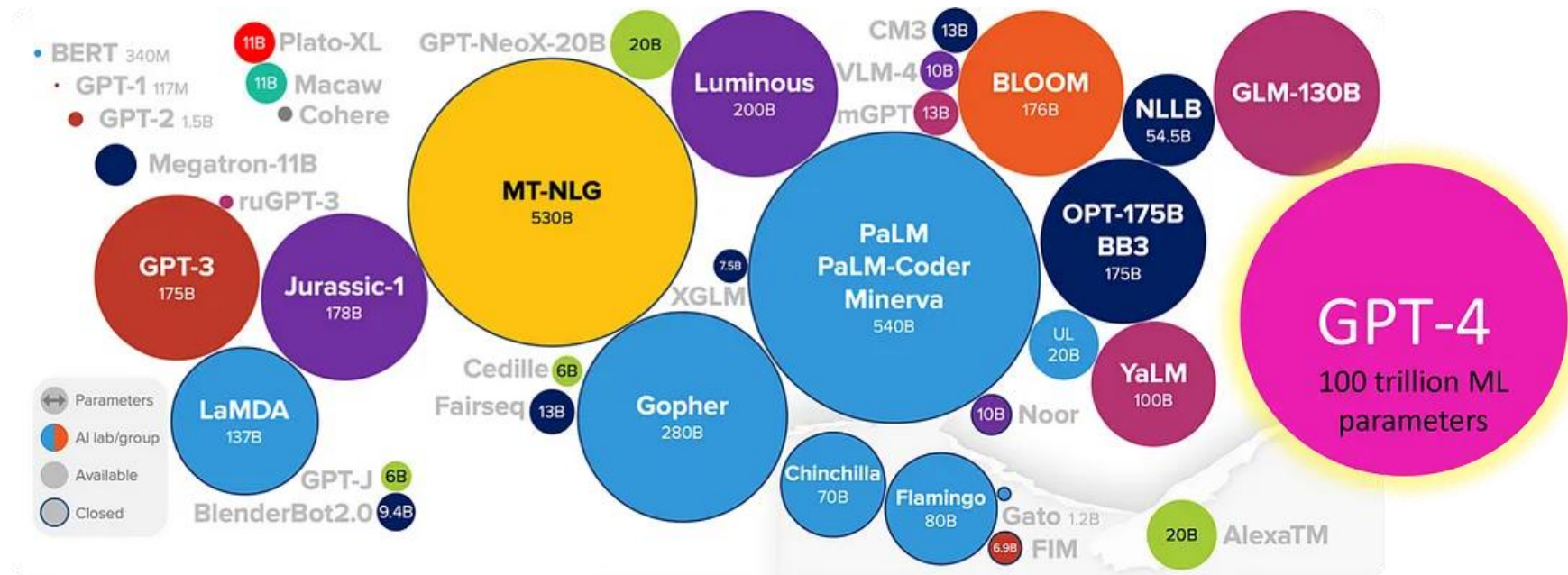
# Outline

- **Motivation**
- **Floating Point**
- **Uniform Quantization**
- **Non-Uniform Quantization**

# Outline

- **Motivation**
- **Floating Point**
- **Uniform Quantization**
- **Non-Uniform Quantization**

# Motivation

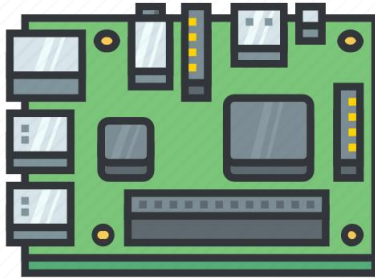
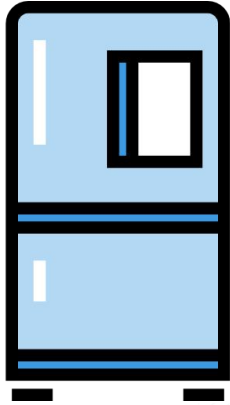


VRAM Usage: Model Size \* 2 (Loading with FP16)

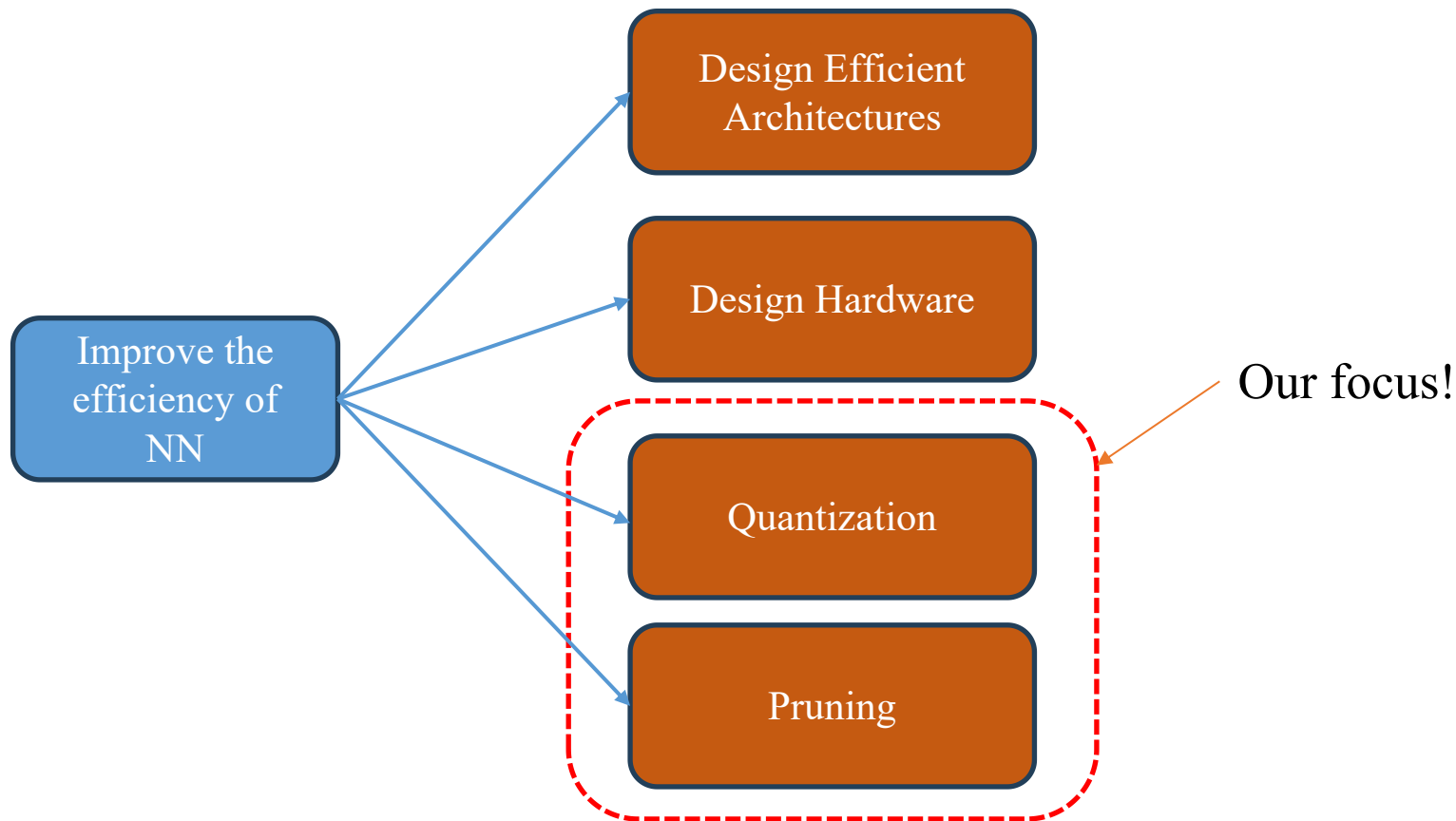
Model	Size	VRAM Usage
GPT3	175B	350 GB
Bloom	176B	352GB
Llama-2-70B	70B	140GB
Mistral	7B	14GB

# Motivation

---



# Motivation





# Quantization

Floating Point

-0.2	1	0.3
0.1	-0.6	-0.7
1.2	0.4	0

32 bit



Integer

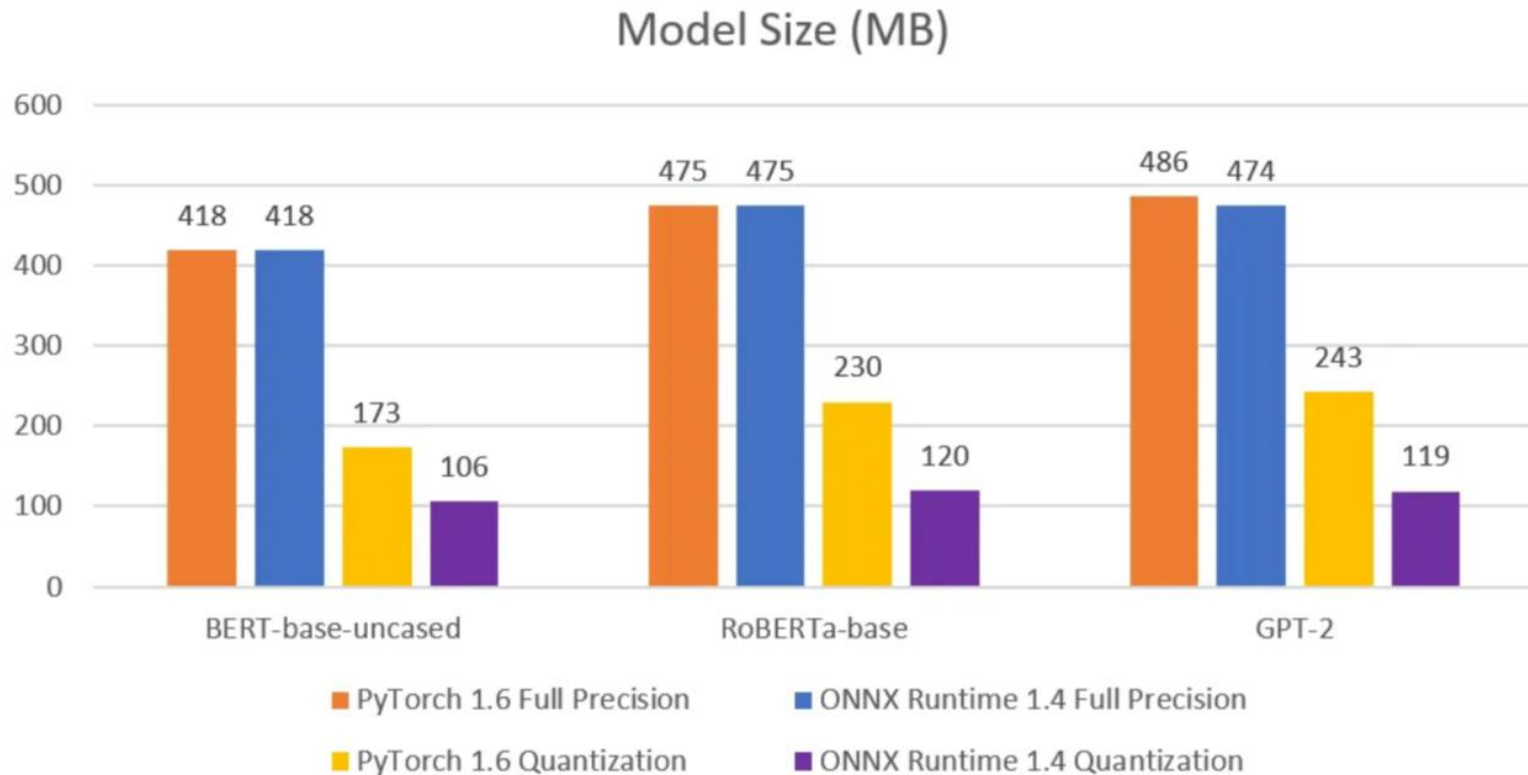
1	3	2
1	0	0
3	2	1

8 bit

Quantization

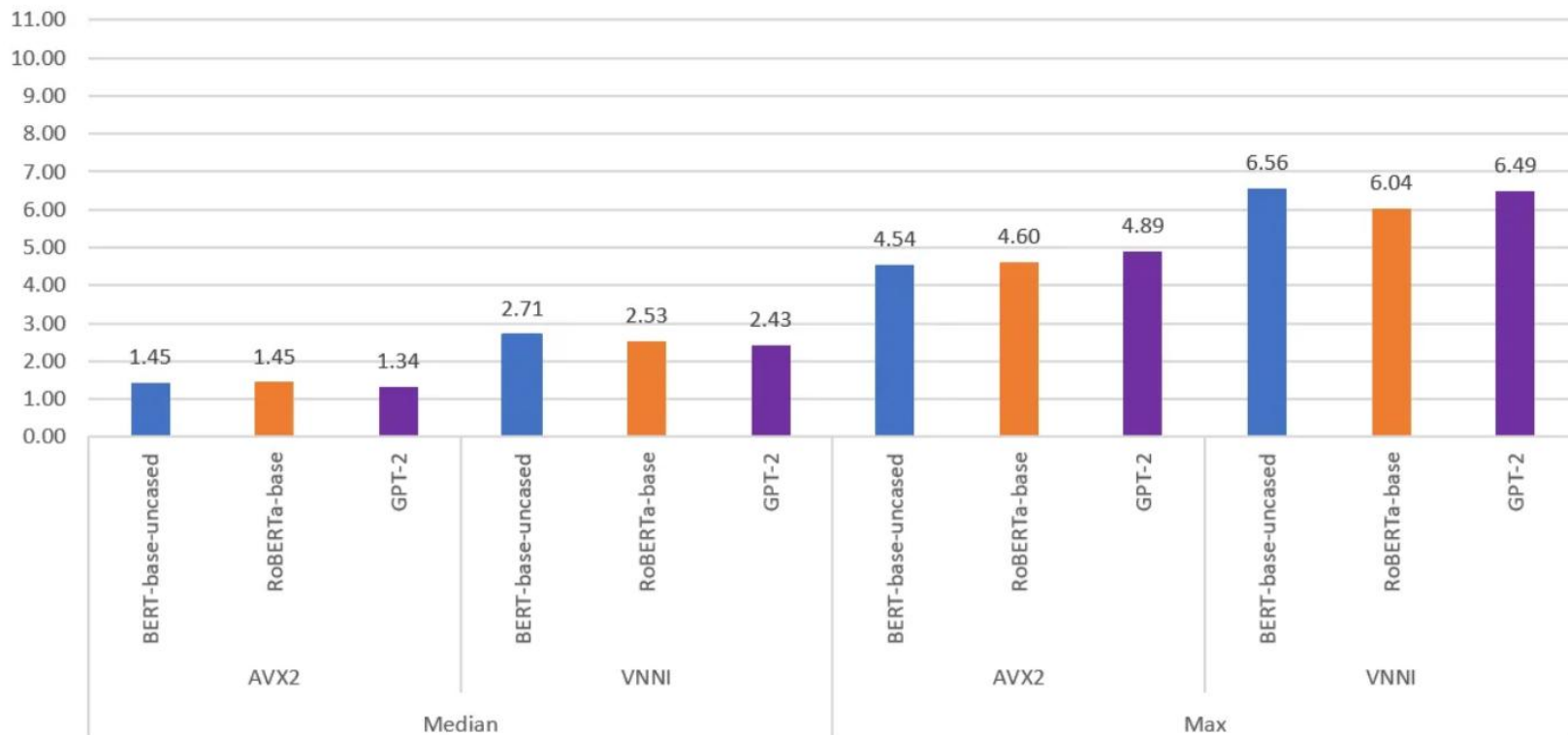


# Quantization



# Quantization

Inference Speedup with ONNX Runtime INT8 Quantization  
(Compared to ONNX Runtime FP32)

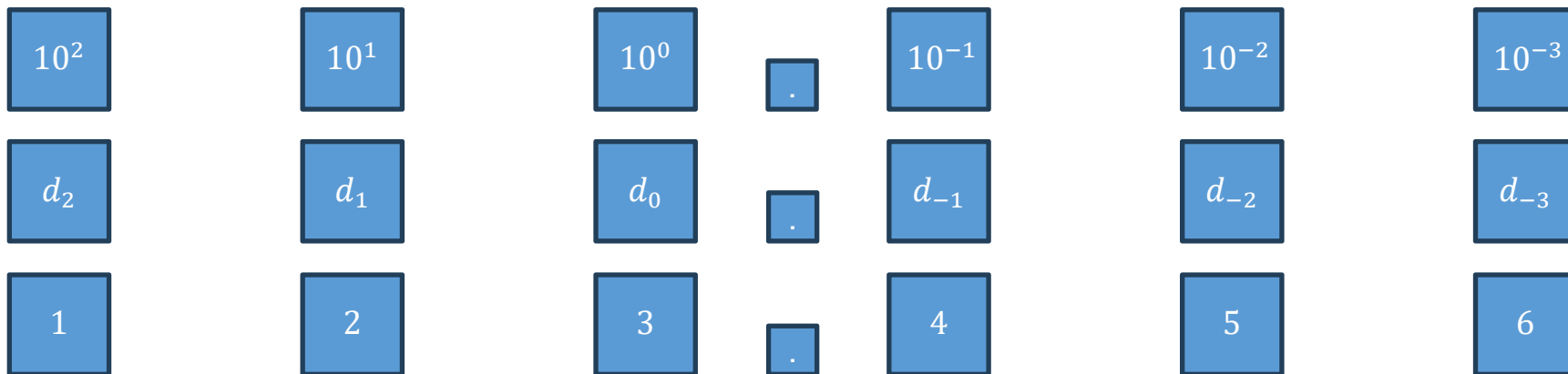


# Quantization

- Motivation
- **Floating Point**
- Uniform Quantization
- Non-Uniform Quantization

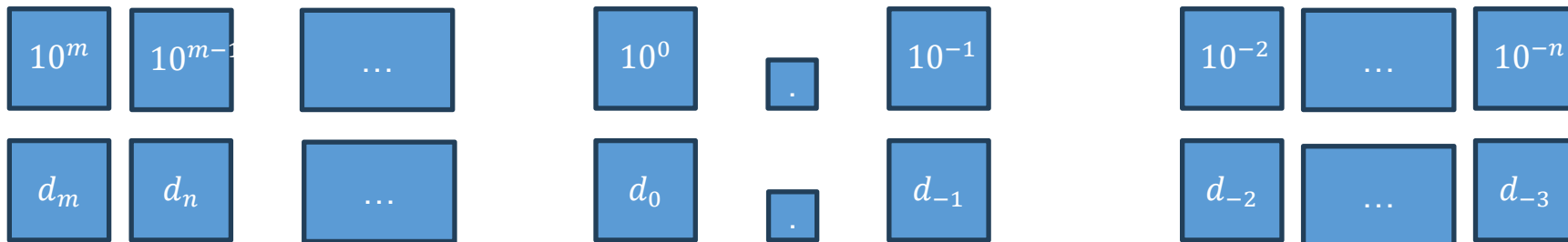
# Representing numbers

## Representing Real Numbers: Decimal



# Representing numbers

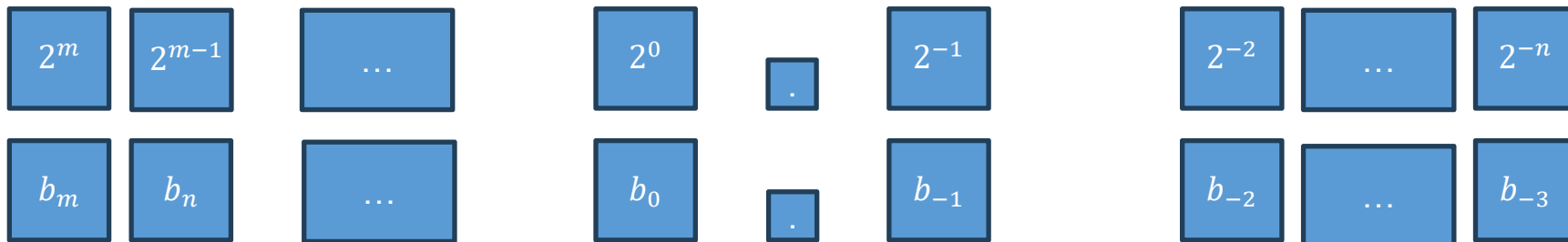
## Representing Real Numbers: Decimal



$$d = \sum_{i=-n}^m 10^i * d_i$$

# Representing numbers

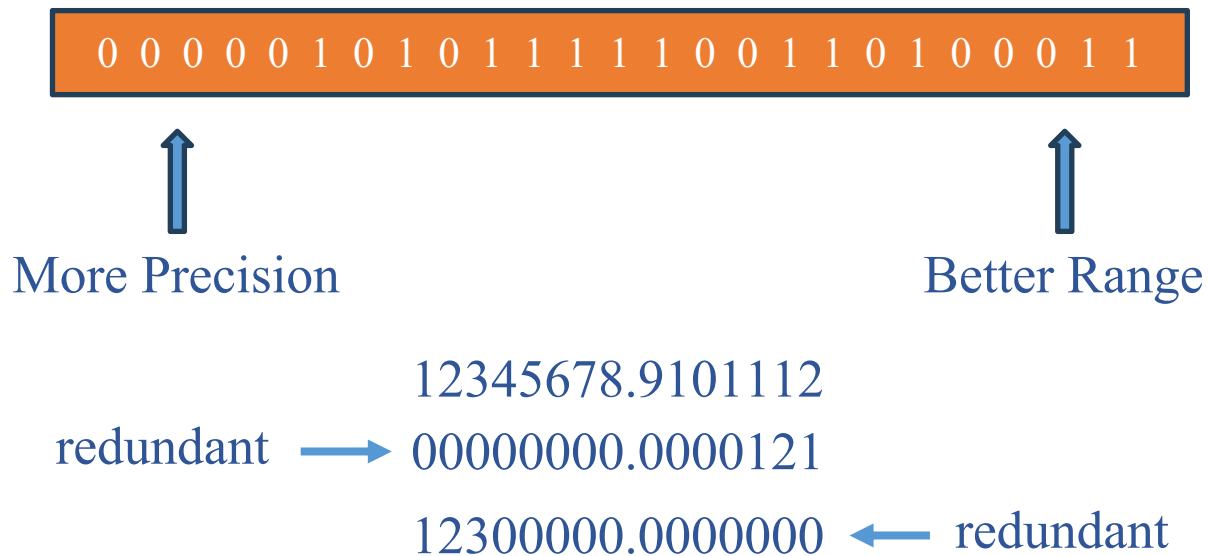
## Representing Real Numbers: Binary



$$b = \sum_{i=-n}^m 2^i * b_i$$

# Fixed-point

Where to put the point?



Any way to make the point moving?



# Fixed-point

Where to put the point?

12300000.0000000	→	$1.23 * 10^7$
00000000.0000121	→	$1.21 * 10^{-5}$

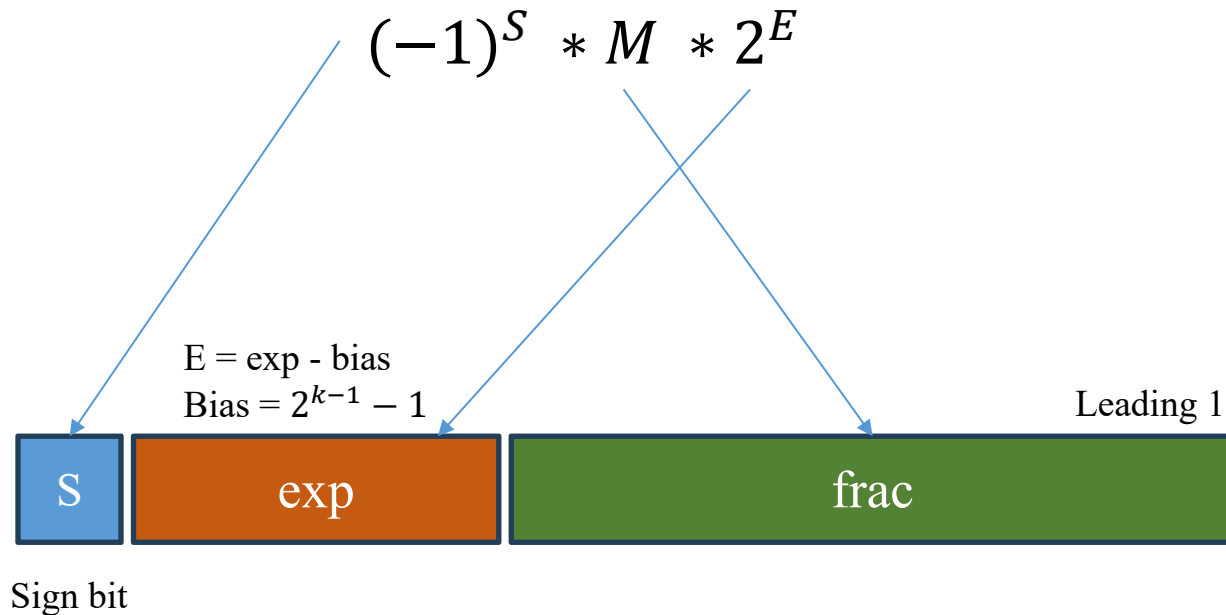
For Binary

11000000.0000000	→	$1.1 * 2^7$
00000000.0000011	→	$1.1 * 2^{-6}$

Indicates where to put the point

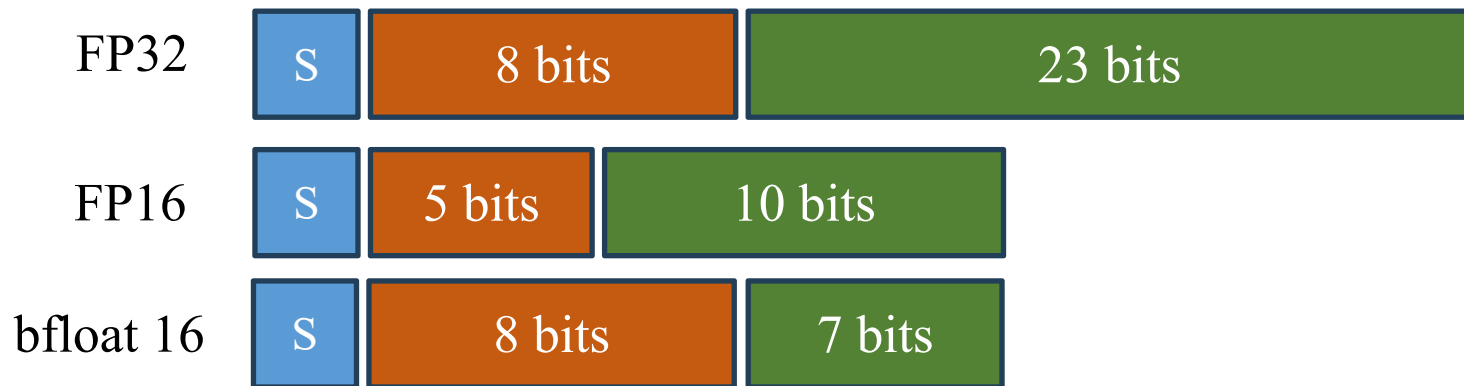
$$(-1)^S * M * 2^E$$

# Floating Point



# Floating Point

$$(-1)^S * M * 2^E$$



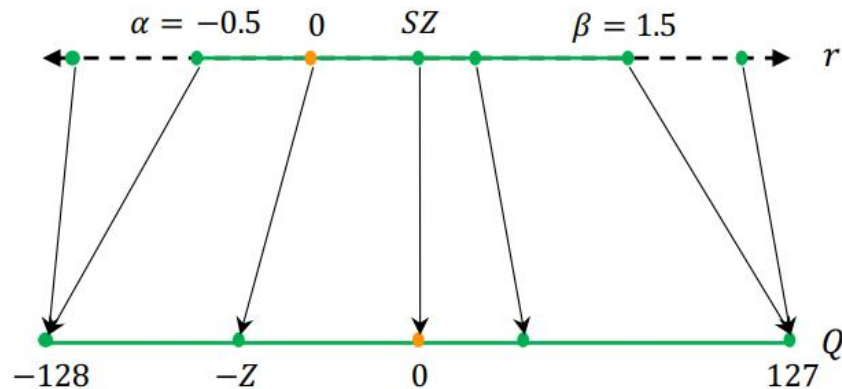
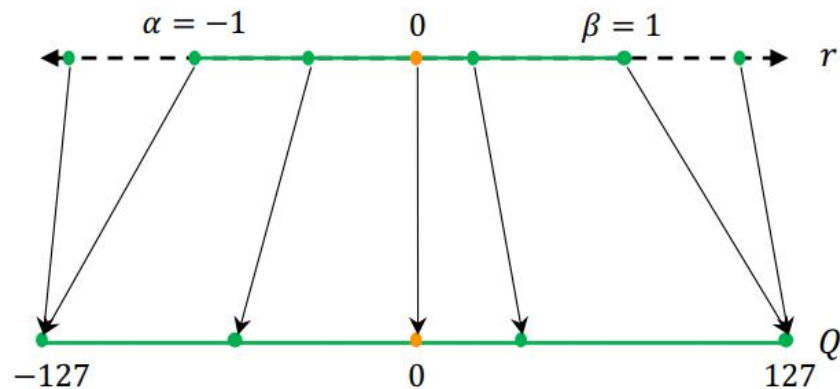
# Floating Point

Data type	Range	Precision
FP32	$1.18 \times 10^{-38}$ to $3.40 \times 10^{38}$	6 – 9 decimal digits
FP16	$4.88 \times 10^{-4}$ to $6.55 \times 10^4$	5 – 6 decimal digits
Bfloat16	$1.18 \times 10^{-38}$ to $3.40 \times 10^{38}$	3 decimal digits

# Quantization

- Motivation
- Floating Point
- **Uniform Quantization**
- Non-Uniform Quantization

# Uniform Quantization



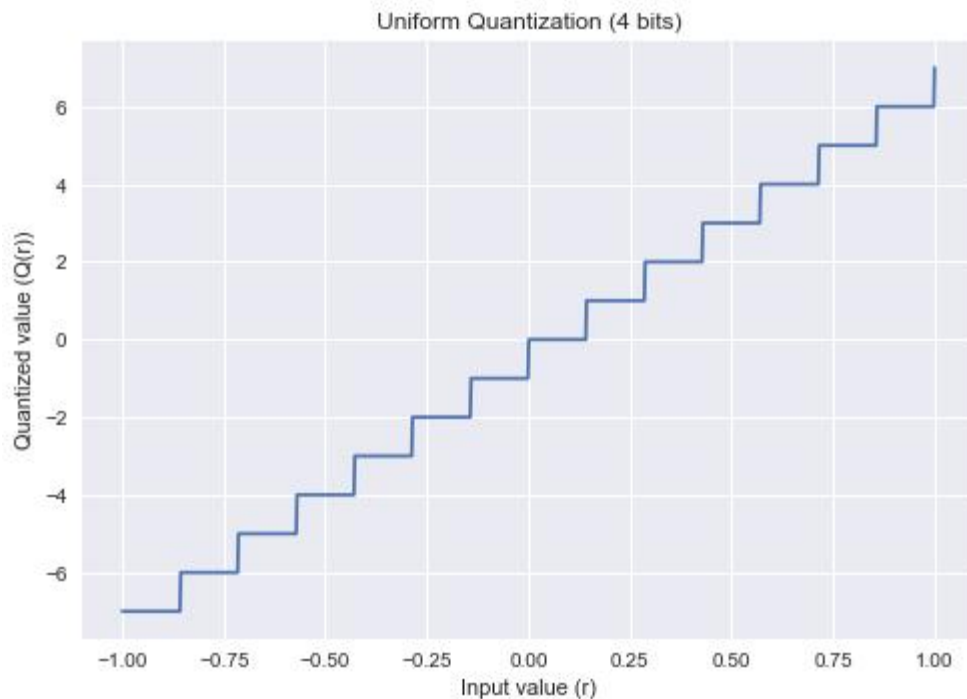
$$Q(r) = \text{Int}\left(\frac{r}{S}\right) - Z$$

Scaling Factor

Zero Value

# Uniform Quantization

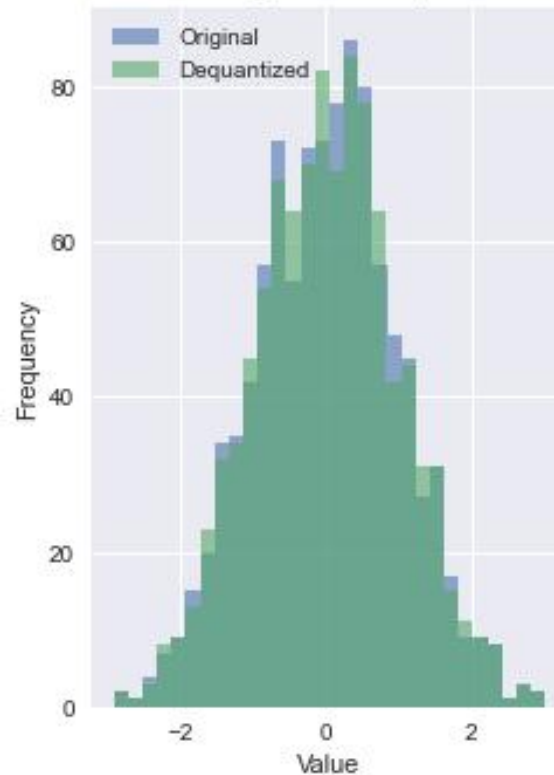
Example  $S = 1/127$ ,  $Z = 0$



# Dequantization

$$Q(r) = \text{Int}\left(\frac{r}{S}\right) - Z$$
$$\rightarrow Q(r) + Z = \text{Int}\left(\frac{r}{S}\right)$$
$$\rightarrow r \approx S * (Q(r) + Z)$$

Distribution of Original and Dequantized Values





# Scaling Factor

$$Q(r) = \text{Int} \left( \frac{r}{S} \right) - Z \quad S = \frac{\beta - \alpha}{2^b - 1}$$

$S$ : scaling factor  
 $[\alpha, \beta]$ : clipping range  
 $b$ : bit width

For 8-bit quantization

$$S = \frac{\beta - \alpha}{255}$$

How to find  $\beta$  and  $\alpha$  ?

Calibration: the process of choosing the clipping range

Simple approach: Choosing the min/max of the signal

Assymetric Quantization

$$\alpha = r_{min}$$

$$\beta = r_{max}$$

$$Z \neq 0$$

Symetric Quantization

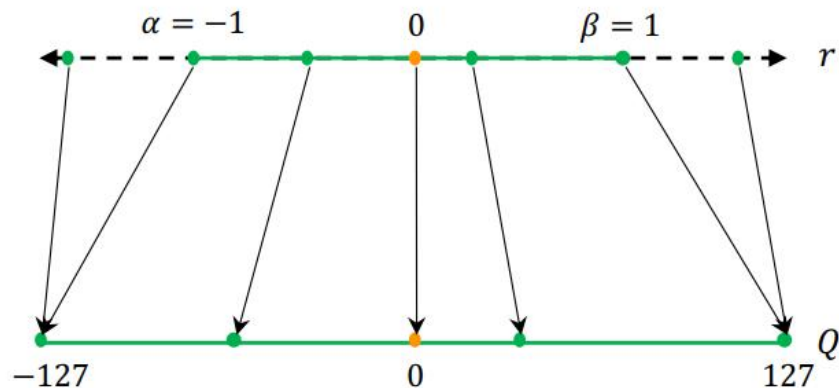
$$-\alpha = \beta$$

$$= \max(|r_{max}|, |r_{min}|)$$

$$Z = 0$$

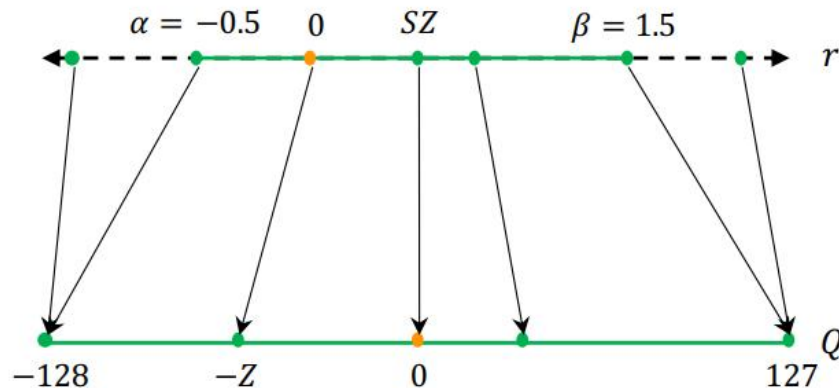
# Uniform Quantization

Symmetric Quantization



$$Q(r) = \text{Int}\left(\frac{r}{S}\right)$$

Asymmetric Quantization



$$Q(r) = \text{Int}\left(\frac{r}{S}\right) - Z$$

# Symmetric Quantization

```
import torch

def absmax_quantize(x):
    # Calculate scale
    scale = 127 / torch.max(torch.abs(x))

    # Quantize
    x_quant = (scale * x).round()

    # Dequantize
    x_dequant = x_quant / scale

    return x_quant.to(torch.int8), x_dequant
```

$$Q(r) = \text{Int}\left(\frac{r}{S}\right)$$

For symmetry  $S = \text{absmax}(r)/127$

# Symmetric Quantization

Quantize

0.3	0.6	-0.3	-0.5	0.0	0.7	-0.8	0.2
-----	-----	------	------	-----	-----	------	-----

Absolute max: 0.8

$$\text{Scale: } \frac{1}{s} = \frac{127}{0.8} = 158.75$$

$$Q(r) = \text{Int}\left(\frac{r}{s}\right)$$

48	95	-48	-79	0	111	-127	32
----	----	-----	-----	---	-----	------	----

# Symmetric Quantization

Dequantize

48	95	-48	-79	0	111	-127	32
----	----	-----	-----	---	-----	------	----

Absolute max: 0.8

$$\text{Scale: } \frac{1}{s} = \frac{127}{0.8} = 158.75$$

$$\hat{r} = Q(r) * S$$

0.30	0.49	-0.30	-0.49	0.0	0.69	-0.8	0.20
------	------	-------	-------	-----	------	------	------

$$Q(r) = \text{Int}\left(\frac{r}{S}\right) - Z$$

$$S = \frac{\beta - \alpha}{255}$$

Calculate zero point

$$\begin{aligned} r_{min} &\rightarrow -128 \\ -128 &= \text{Int}\left(\frac{r_{min}}{S}\right) - Z \\ Z &= \text{Int}\left(\frac{r_{min}}{S}\right) + 128 \end{aligned}$$

```
def zeropoint_quantize(X):
    # Calculate value range (denominator)
    x_range = torch.max(X) - torch.min(X)
    x_range = 1 if x_range == 0 else x_range

    # Calculate scale
    scale = x_range / 255

    # Shift by zero-point
    zeropoint = (torch.min(X) / scale).round() + 128

    # Scale and round the inputs
    X_quant = torch.clip((X / scale - zeropoint).round(), -128, 127)

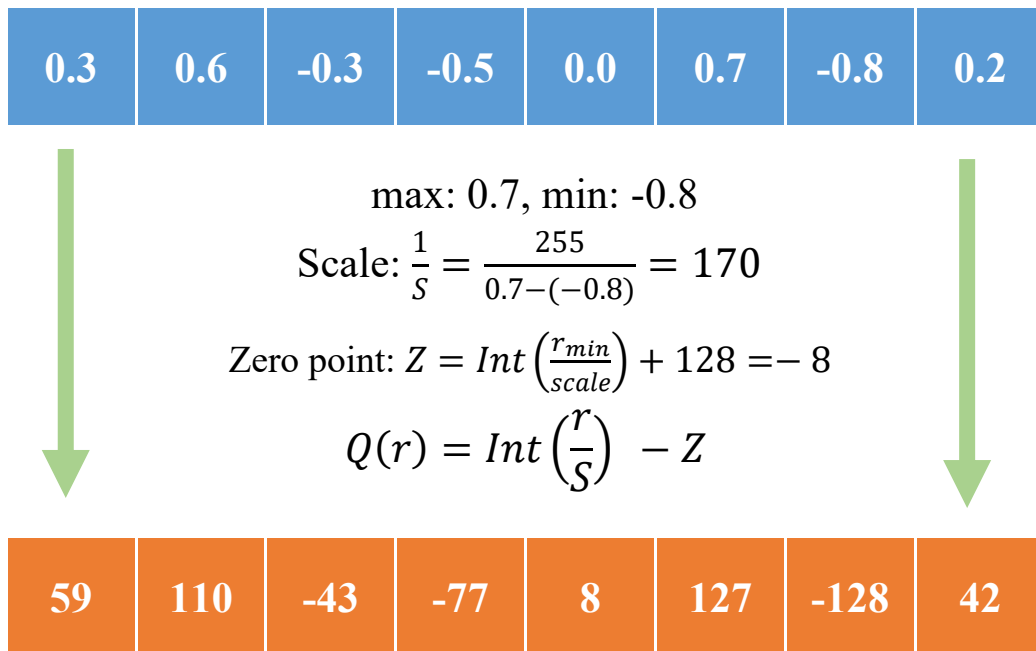
    # Dequantize
    X_dequant = (X_quant + zeropoint) * scale

    return X_quant.to(torch.int8), X_dequant
```



# Symmetric Quantization

Quantize



# Symmetric Quantization

Dequantize

59	110	-43	-77	8	127	-128	42
----	-----	-----	-----	---	-----	------	----

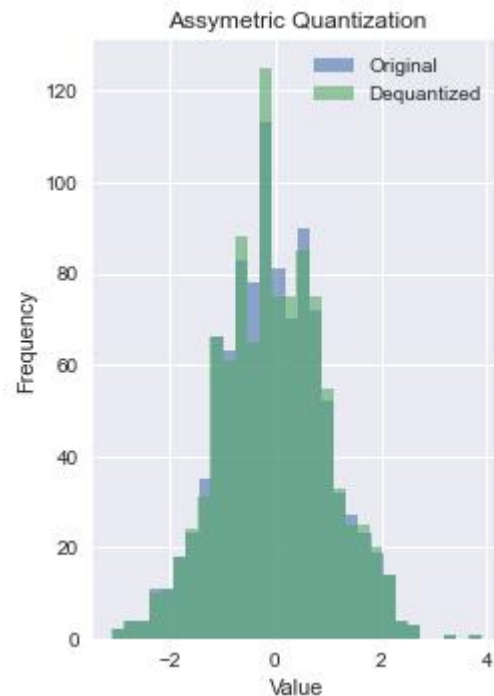
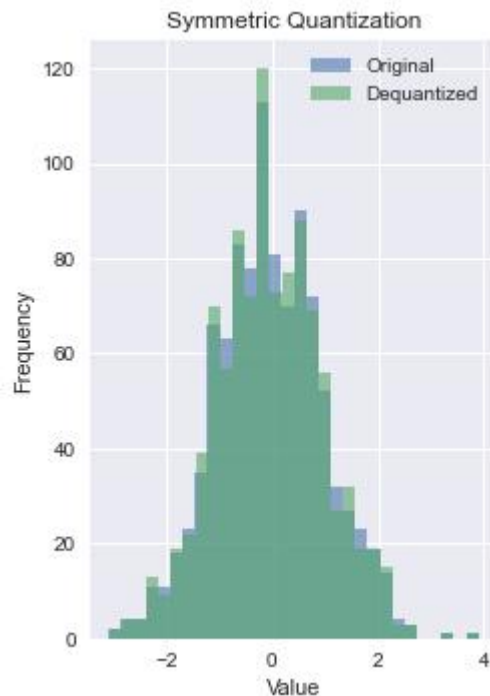


$$\begin{aligned}\text{max: } 0.7, \text{ min: } -0.8 \\ \text{Scale: } \frac{1}{S} &= \frac{255}{0.7 - (-0.8)} = 170 \\ \text{Zero point: } Z &= \text{Int}\left(\frac{r_{\min}}{\text{scale}}\right) + 128 = -8 \\ \hat{r} &= S * (Q(r) + Z)\end{aligned}$$



0.30	0.60	-0.30	-0.50	0.0	0.70	-0.80	0.20
------	------	-------	-------	-----	------	-------	------

# Asymmetric vs. Symmetric

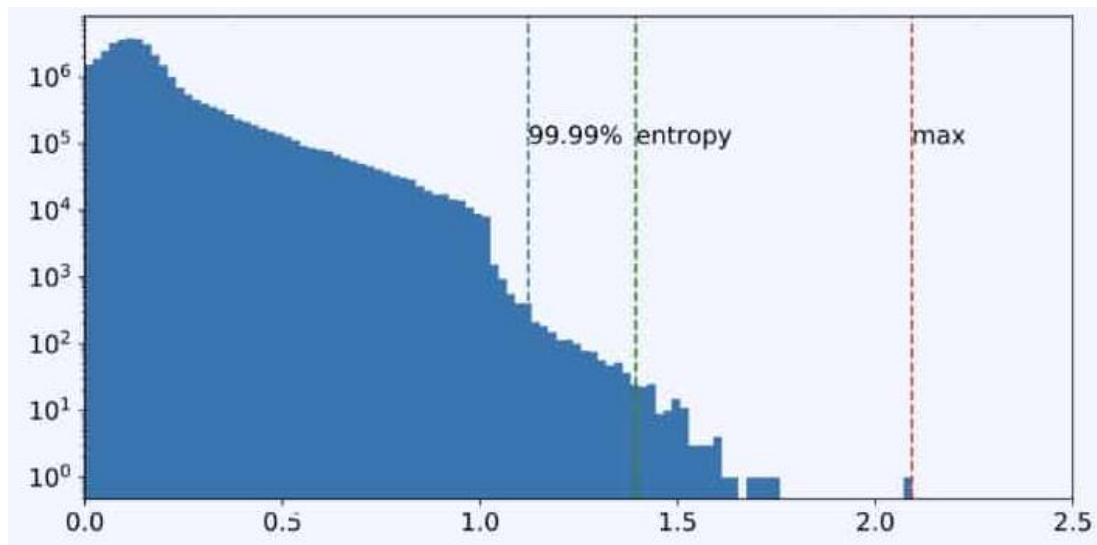


# Asymmetric vs. Symmetric

Feature	Symmetric Quantization	Asymmetric Quantization
Implementation	Simple Implementation	More complex
Computational cost	Simpler calculation	Slightly higher
Efficiency for non-uniform distributions	Less efficient	More efficient

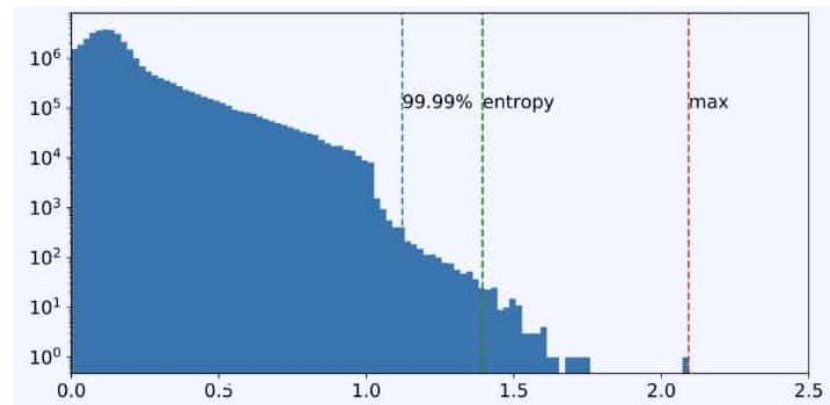
# Calibration

Is choosing max/min values a good option?

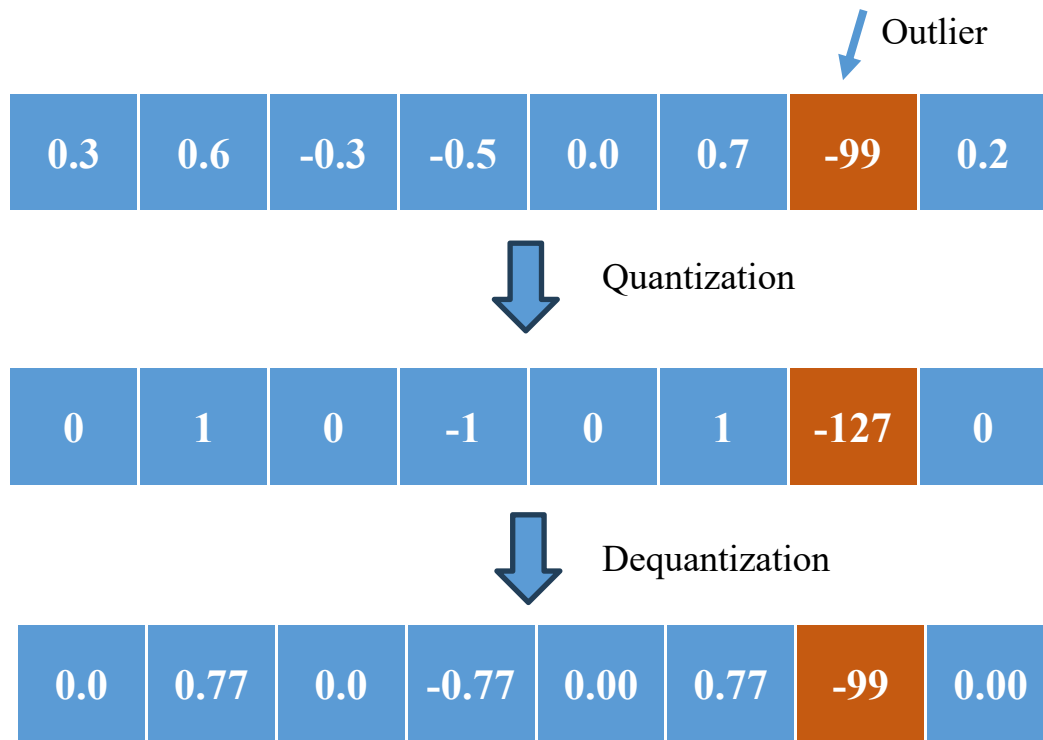


Choose max absolute value  $\rightarrow$  vulnerable to outliers.

- Entropy: KL divergence to minimize information loss.
- Percentile: Set a range to a percentile of the distribution.



# Calibration

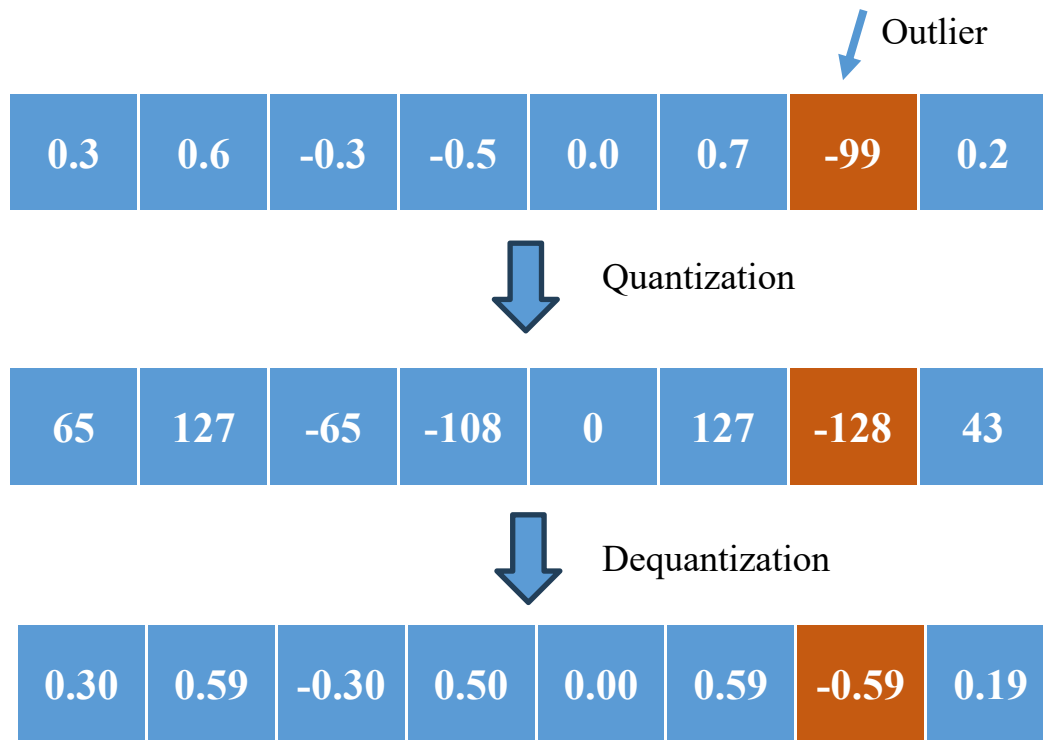


## Percentile Quantization

```
def percentile_quantize(x, percentile=70):  
    # Calculate the 99th percentile of the absolute values in x  
    scale = torch.quantile(torch.abs(x), percentile/100.) / 127  
  
    # Quantize  
    x_quant = (x / scale).round()  
  
    # Clamp the values to ensure they are within the int8 range after rounding  
    x_quant = torch.clamp(x_quant, min=-128, max=127)  
  
    # Dequantize  
    x_dequant = x_quant * scale  
  
    return x_quant.to(torch.int8), x_dequant
```



# Calibration



# QUIZ TIME!!!



**Câu 1: Quá trình quantization là gì?**

- A) Chuyển đổi giá trị từ dạng phân số sang số nguyên
- B) Chuyển đổi giá trị liên tục thành giá trị rời rạc / giảm số bit để biểu diễn 1 giá trị.
- C) Tăng độ chính xác của dữ liệu số
- D) Giảm tốc độ xử lý của hệ thống số

**Câu 2 Trong machine learning, quantization thường được sử dụng để:**

- A) Tăng độ phức tạp của mô hình
- B) Tăng độ chính xác của mô hình
- C) Giảm kích thước và yêu cầu tính toán của mô hình
- D) Tăng số lượng lớp trong mạng neural

# QUIZ TIME!!!

**4) Phương pháp quantization nào cung cấp các đặc tính sau:**

- **Dễ dàng triển khai (easy to implement) hơn.**
  - **Khoảng cách bằng nhau (equal spacing) giữa các mức lượng hóa.**
- A) Lượng hóa đồng nhất (Uniform quantization)  
B) Lượng hóa không đồng nhất (Non-uniform quantization)  
C) Cả hai lượng hóa đồng nhất và không đồng nhất  
D) Không phải lượng hóa đồng nhất cũng không phải không đồng nhất

**5) Loại lượng hóa nào phù hợp hơn cho các tín hiệu có hầu hết các giá trị tập trung ở phía dương, xét về cả tính đơn giản và hiệu quả?**

- a) Lượng hóa đối xứng (Symmetric quantization)  
b) Lượng hóa không đối xứng (Non-symmetric quantization)  
c) Cả hai đều hoạt động tốt như nhau.  
d) Không phương pháp nào phù hợp.

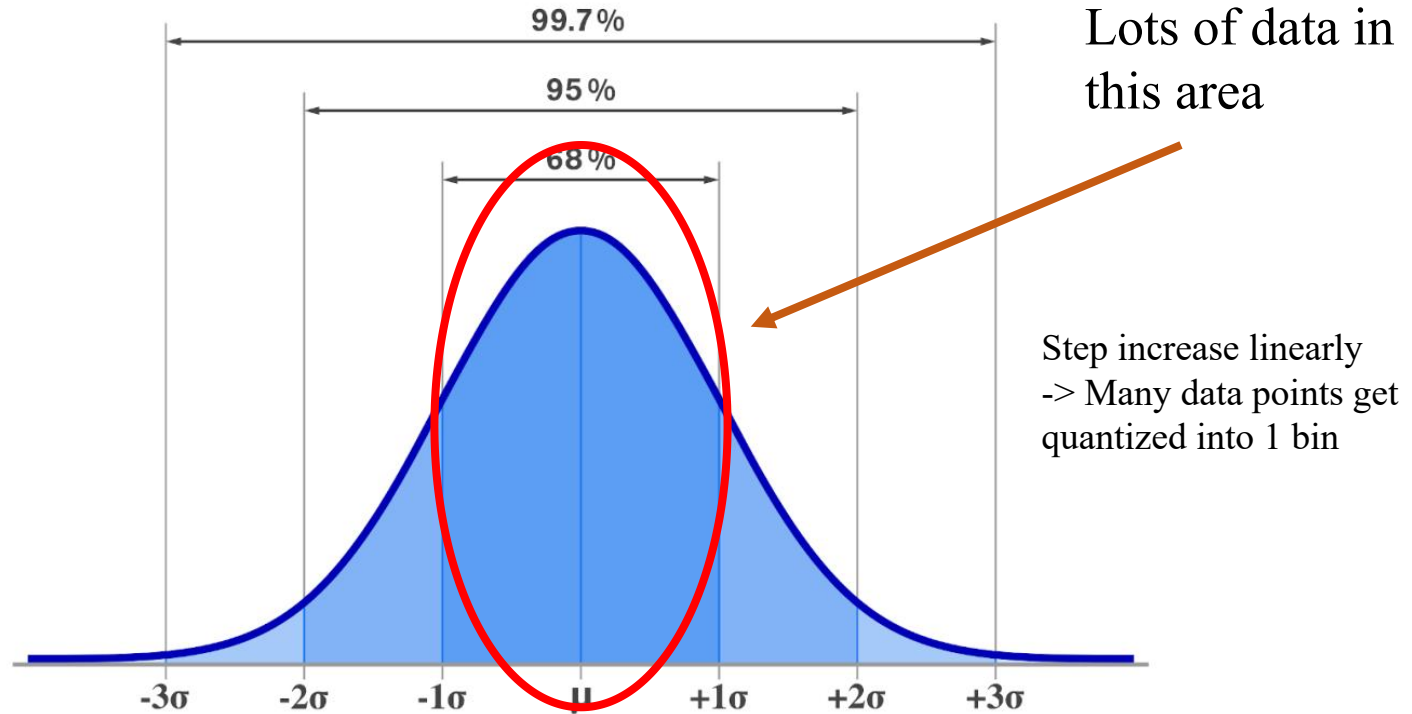
**Câu 5: Quá trình hiệu chỉnh (calibration) trong quantization được sử dụng để làm gì?**

- a) Tăng độ chính xác của mô hình sau quantization
- b) Giảm kích thước của mô hình mà không cần quantization
- c) Xác định phạm vi giá trị tối ưu cho các tham số để giảm thiểu lỗi quantization
- d) Tăng tốc độ tính toán của mô hình mà không ảnh hưởng đến kích thước mô hình

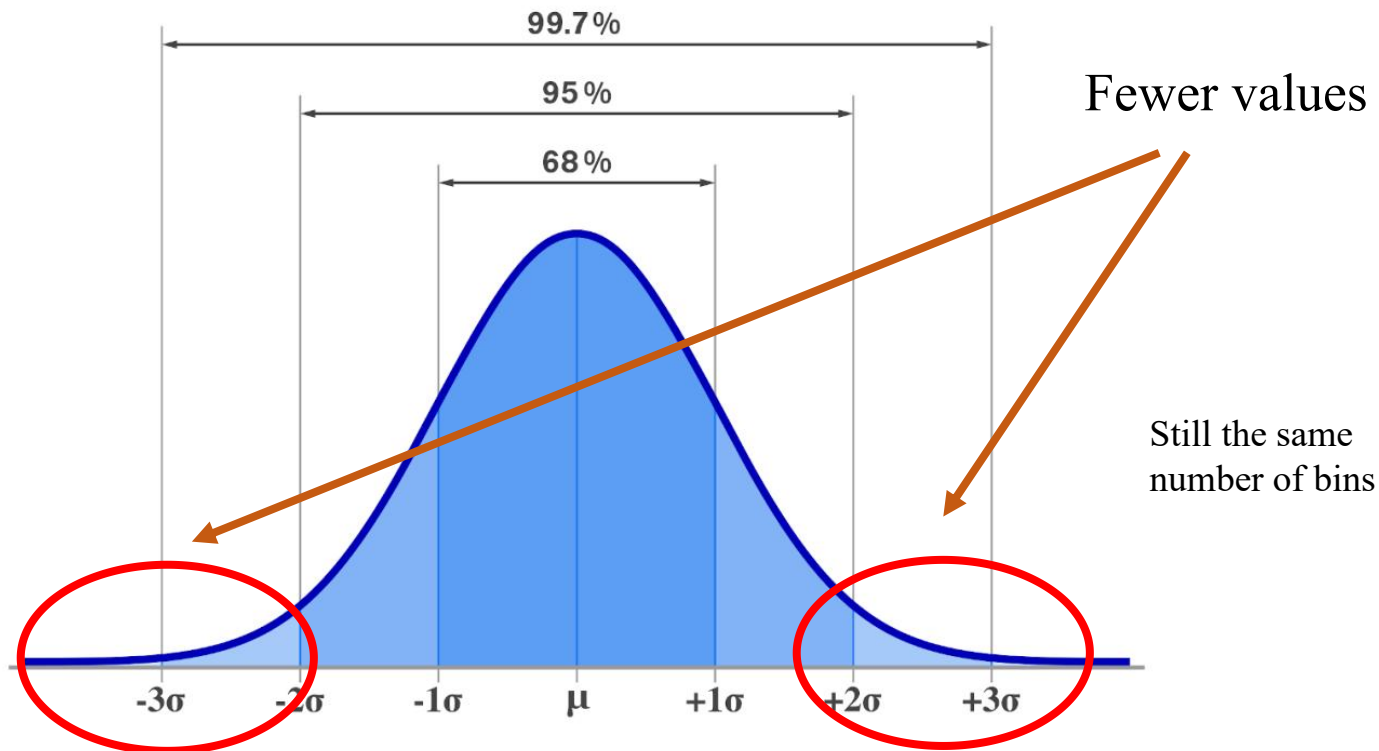
# Quantization

- Motivation
- Floating Point
- Uniform Quantization
- **Non-Uniform Quantization**

# Non-uniform Quantization



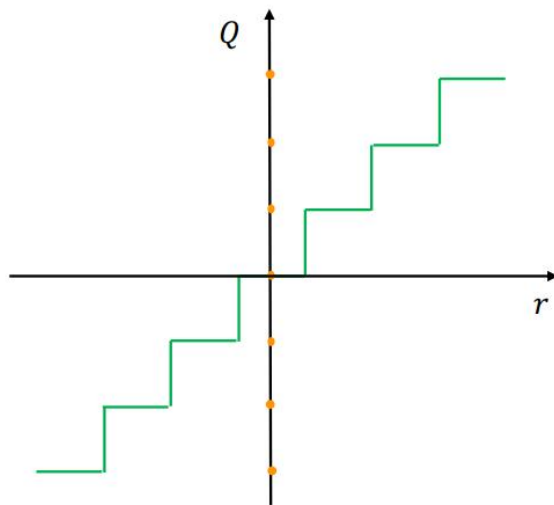
# Non-uniform Quantization



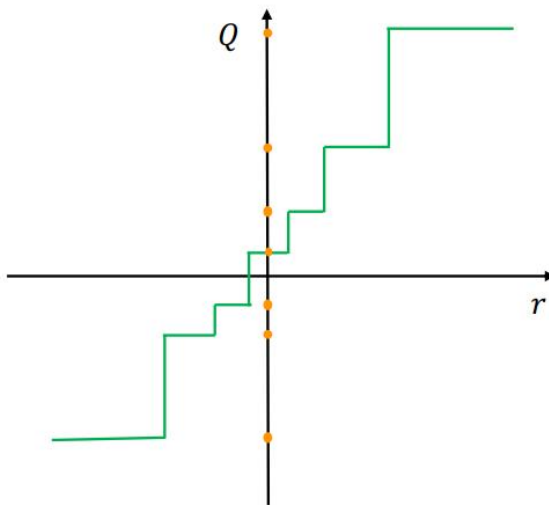


# Non-uniform Quantization

Uniform quantization



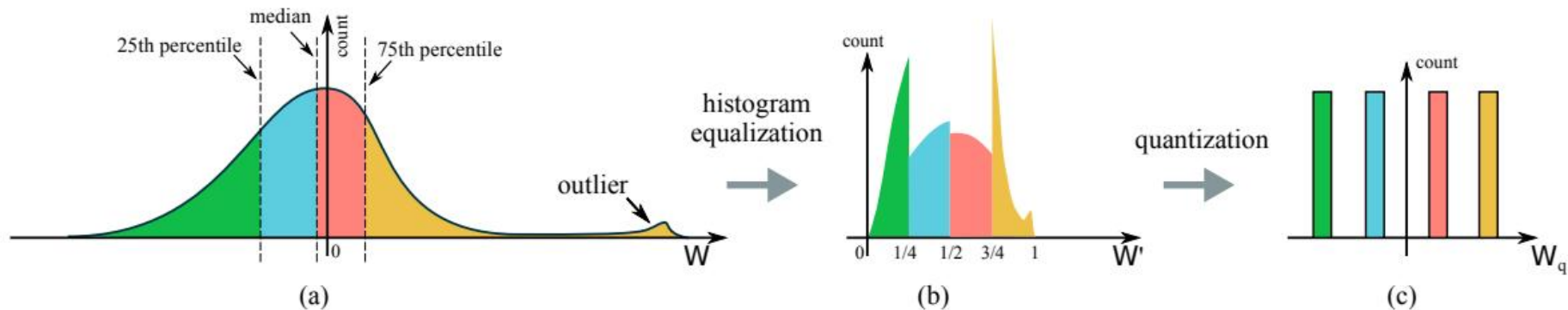
Non-uniform quantization



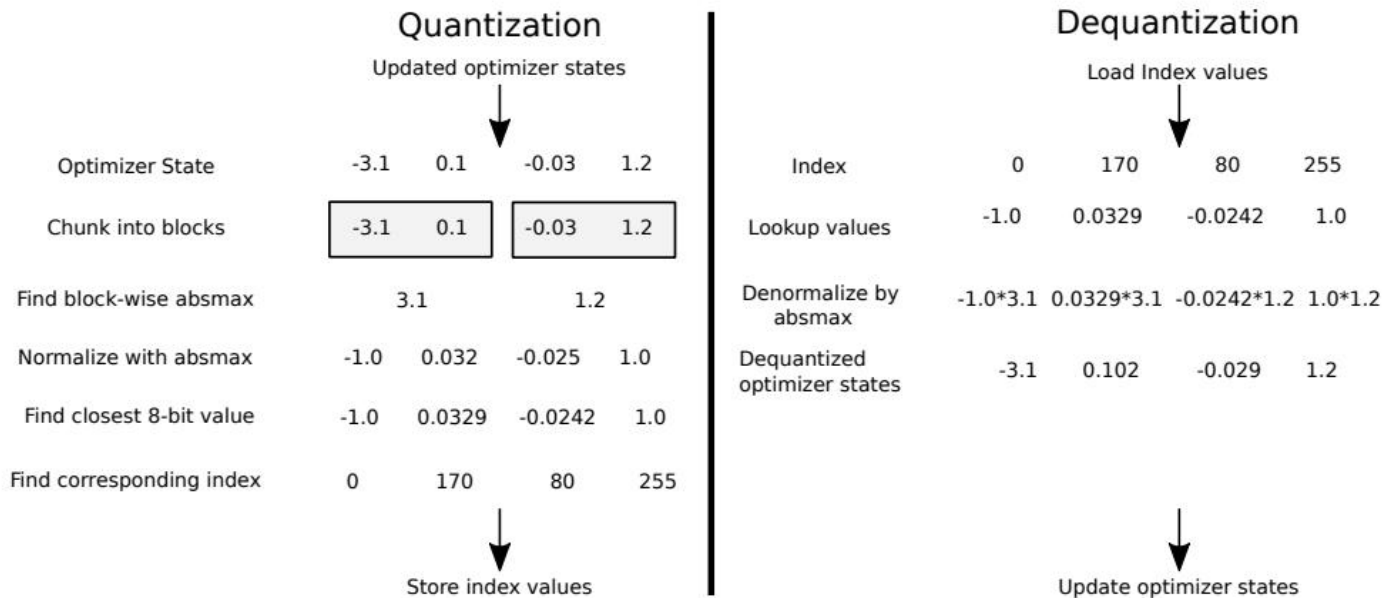
$$Q(r) = X_i, \text{ if } r \in [\Delta_i, \Delta_{i+1})$$

$$Q(r) = X_i, \text{ if } r \in [\Delta_i, \Delta_{i+1})$$

Balance Quantization

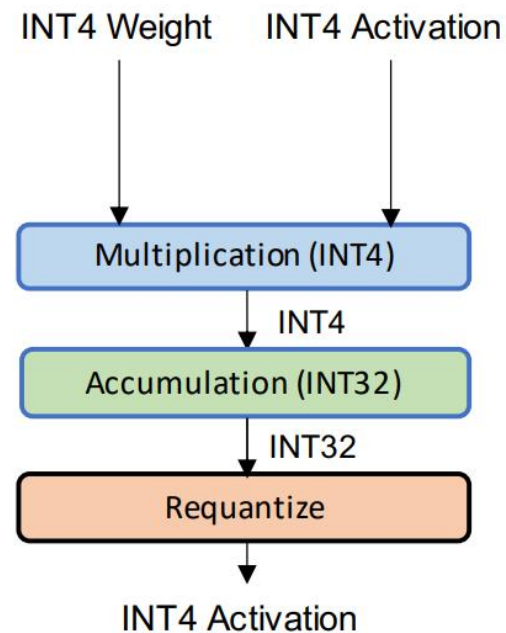
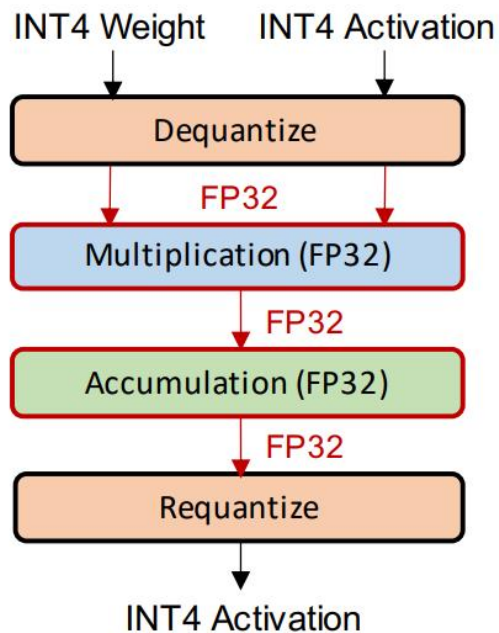
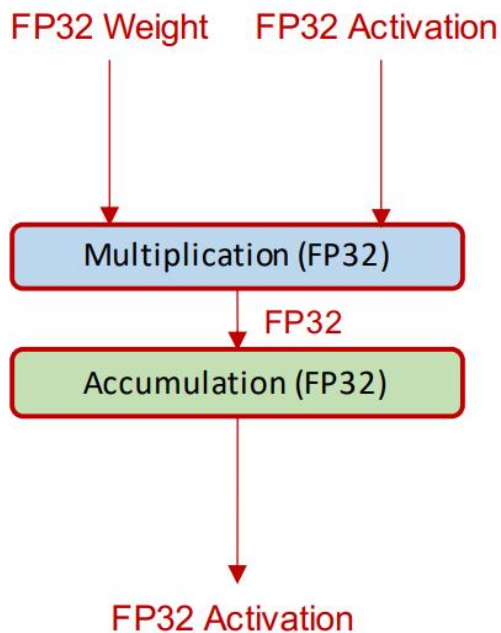


## Block-wise Quantization



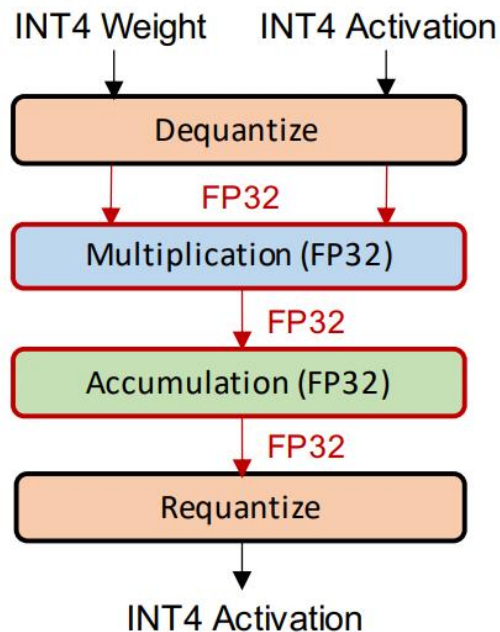
# Non-uniform Quantization

Feature	Uniform Quantization	Non-uniform Quantization
Spacing	Equally spaced levels	Unequally spaced levels (based on signal characteristics)
Simplicity	Simple and straightforward	More complex
Efficiency (uniform distribution)	Efficient	Less efficient for uniform distributions
Efficiency (non-uniform distribution)	May waste bits, poor accuracy for non-uniform distributions	More efficient and accurate for non-uniform distributions
Complexity	Low	High

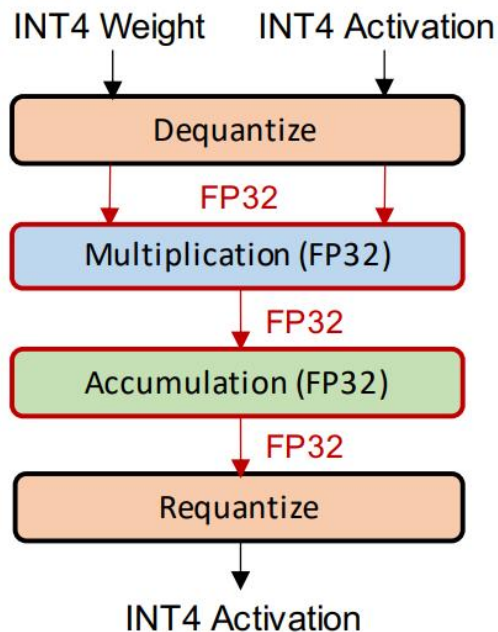


# Simulated Quantization

## Simulated Quantization



## Simulated Quantization



-64	32	-16
96	-48	112
-12	60	-56



-0.5	0.25	0.12
0.75	-0.37	0.88
-0.09	0.47	-0.44

32	50	16
16	32	44



0.25	0.39	0.12
0.12	0.25	0.34

Range  $[-1, 1]$ ,  $S = 1/127$

# Simulated Quantization

-0.5	0.25	0.12
0.75	-0.37	0.88
-0.09	0.47	-0.44

MatMul



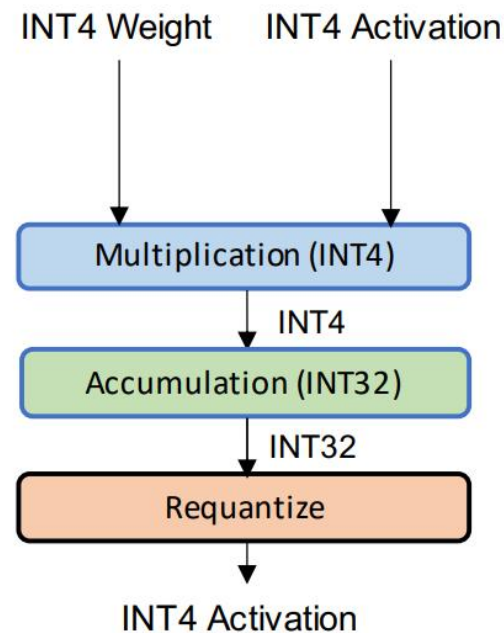
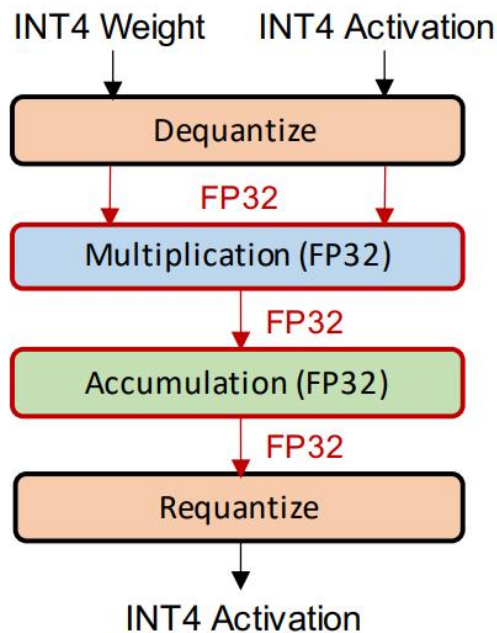
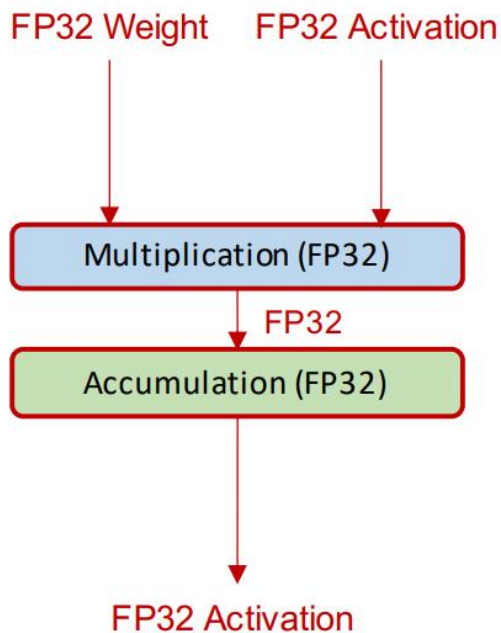
-0.01	0.04
0.15	0.30
0.10	-0.04

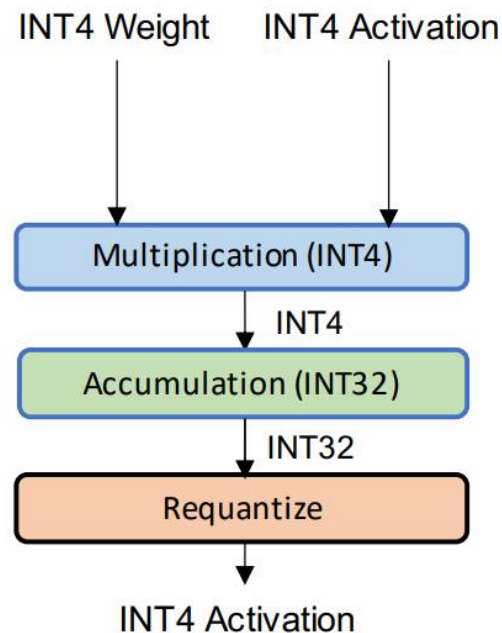
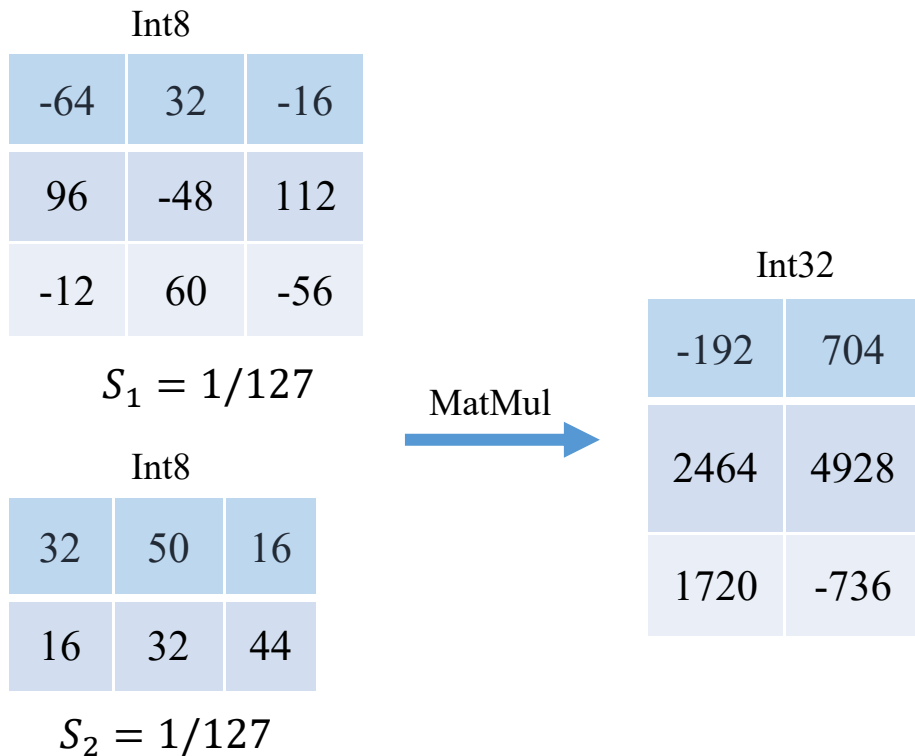
0.25	0.39	0.12
0.12	0.25	0.34

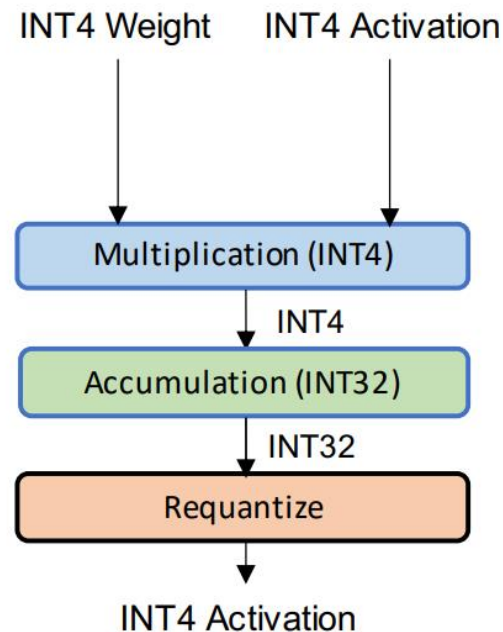
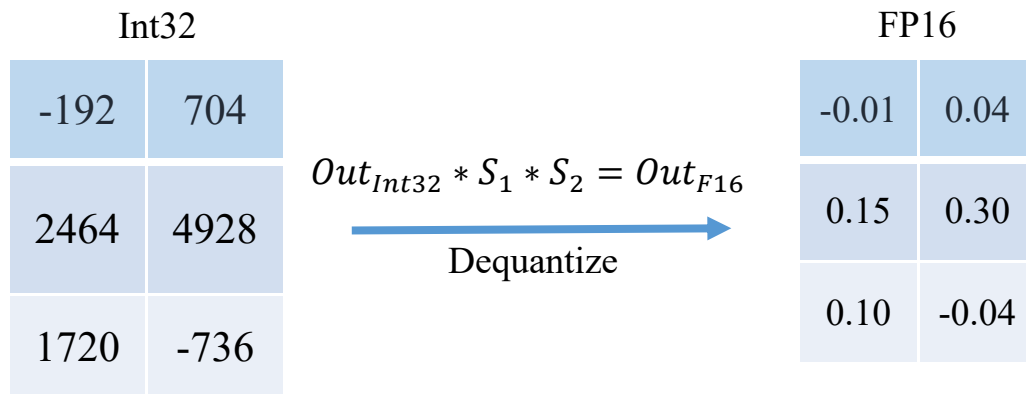
- Reduce Model Size
- Slow inference time

Range [-1, 1], S = 1/127



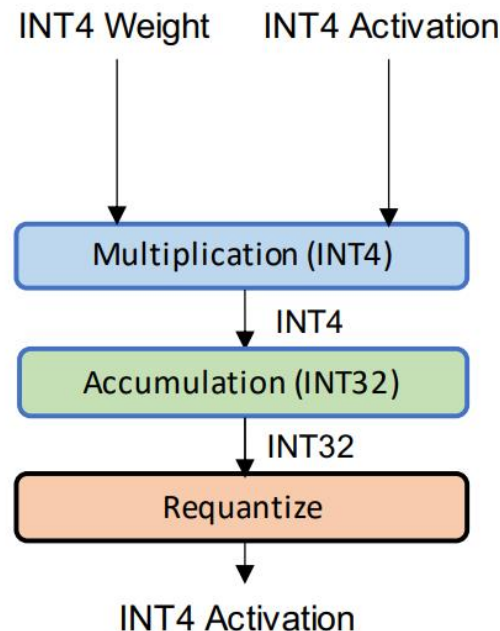




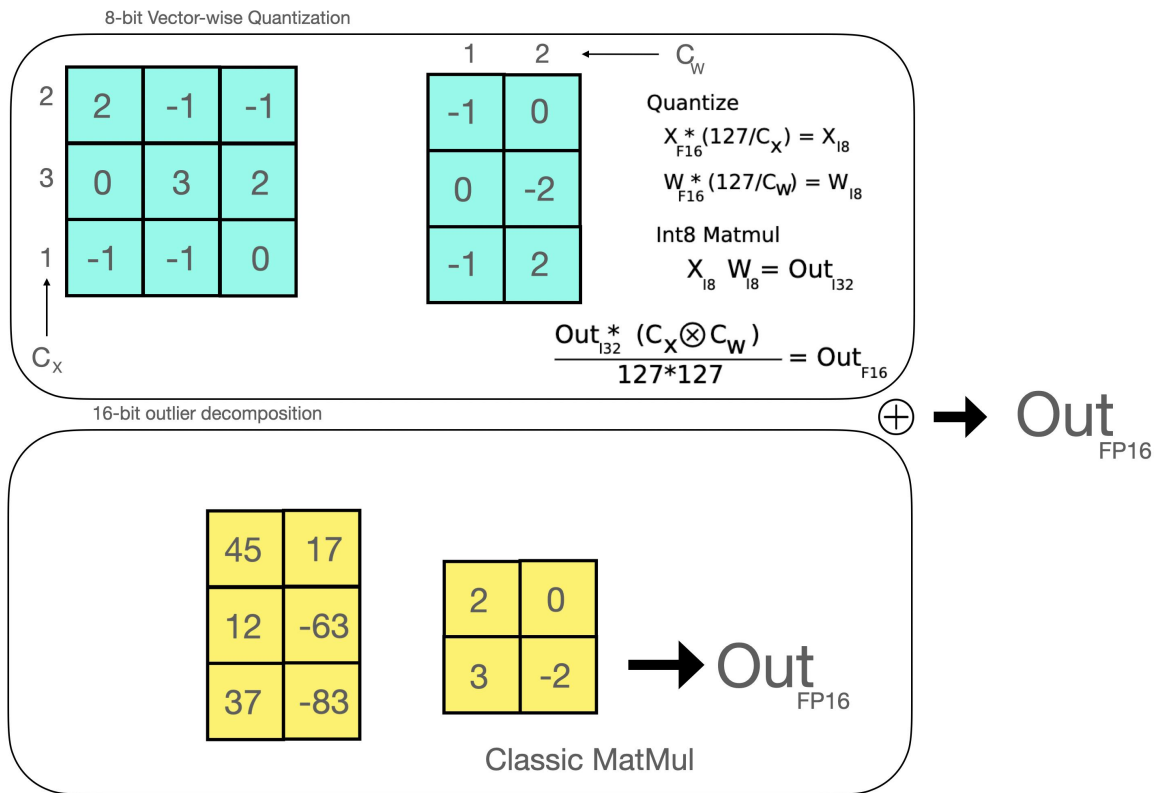


All operations -> Using low-precision integer arithmetic

- Latency
- Power consumption
- Area efficiency



# Inferencing Quantized Models



# Summary

