

# **Section 4: Quantization**

**CSE 493G1, Spring 2025**

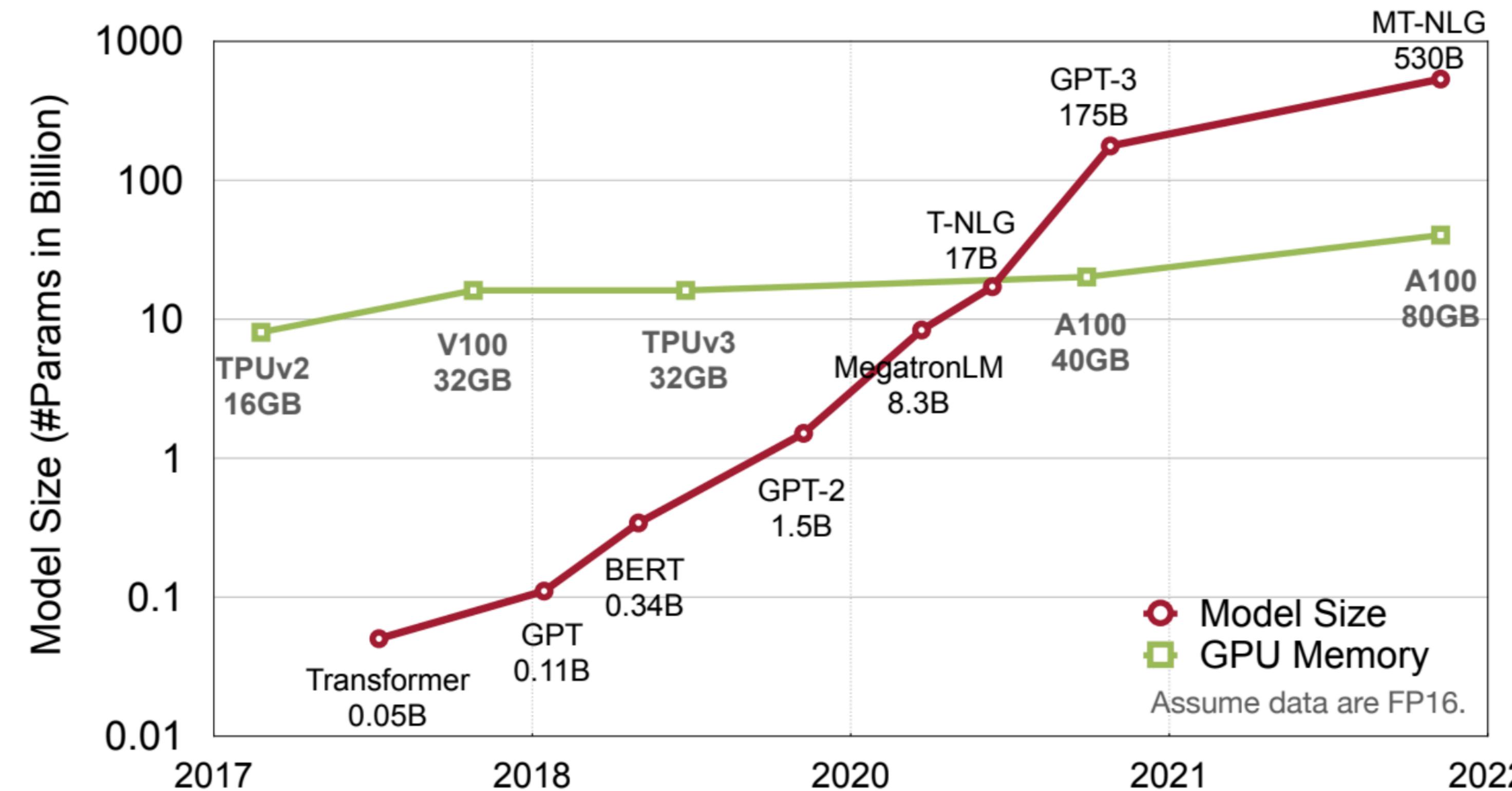
**Tanush Yadav**

# Course Logistics

- A2 due Sunday (4/27)
- Project Proposal due Tuesday (4/29)
- W2 due 5/9
- A3 due 5/11

# Motivation & Background

# DL models are outgrowing hardware

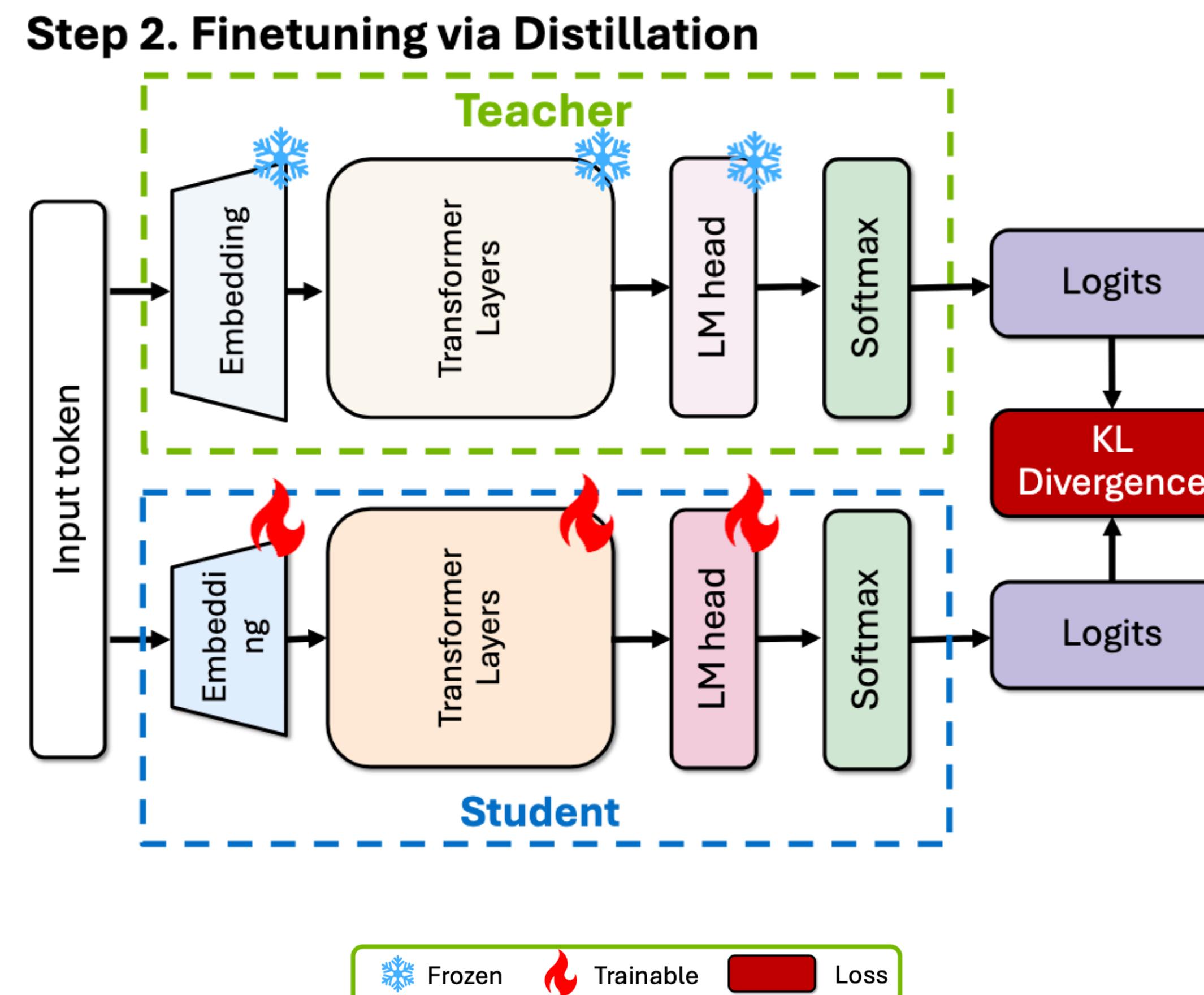


**Decreasing model size\***  
**is a worthy goal.**

\*while preserving performance

# Distillation

A core tenet of the DeepSeek paper



model size = (# of parameters) × (parameter size)

model size = (# of parameters) × (parameter size)

pruning

quantization

# Parameter Size

**How many bits does each parameter occupy? (“bit width”)**

- FP32  $\Rightarrow$  4 bytes per parameter
- FP16, BF16  $\Rightarrow$  2 bytes per parameter
- FP8, int8  $\Rightarrow$  1 byte per parameter
- int4  $\Rightarrow$  0.5 bytes per parameter

# Parameter Size

**Qwen2.5-VL has 3B, 7B, and 72B versions**



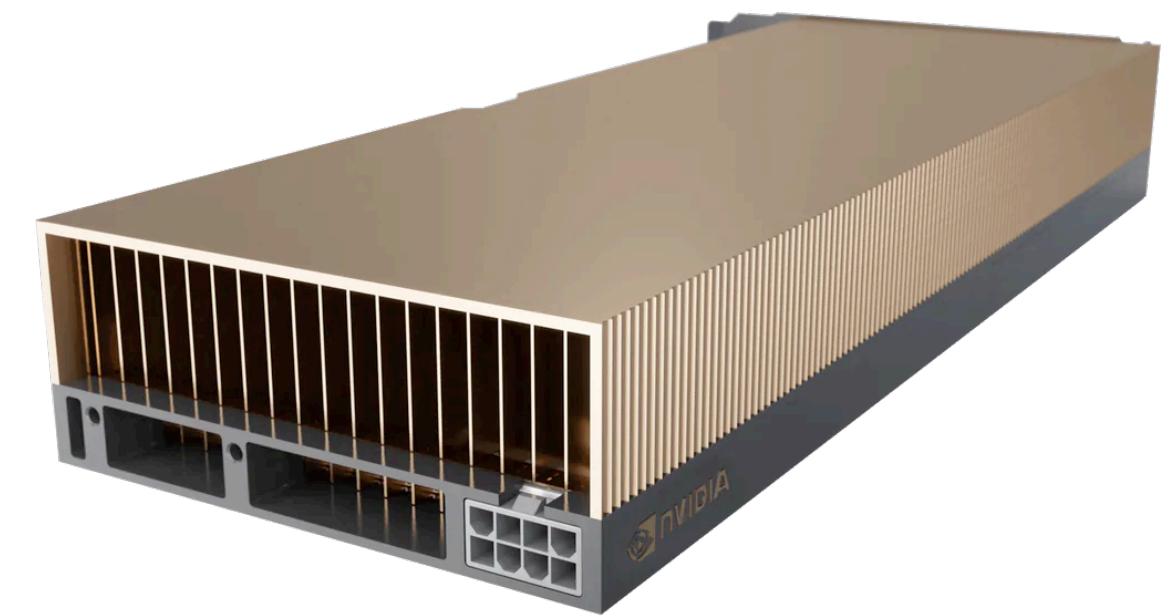
**NVIDIA 2080 Ti**

**11GB memory**



**NVIDIA RTX 6000**

**24GB memory**

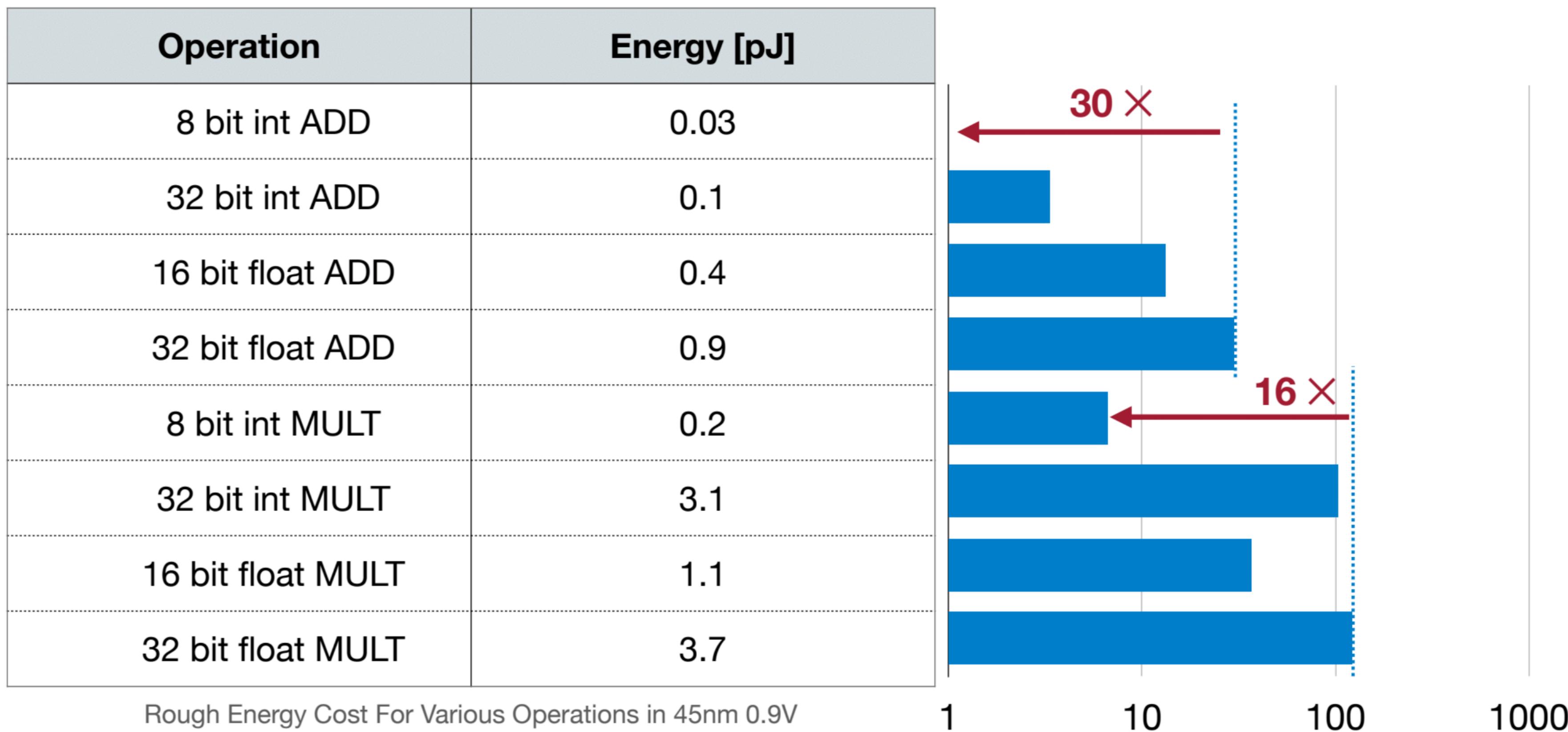


**NVIDIA A40**

**40GB memory**

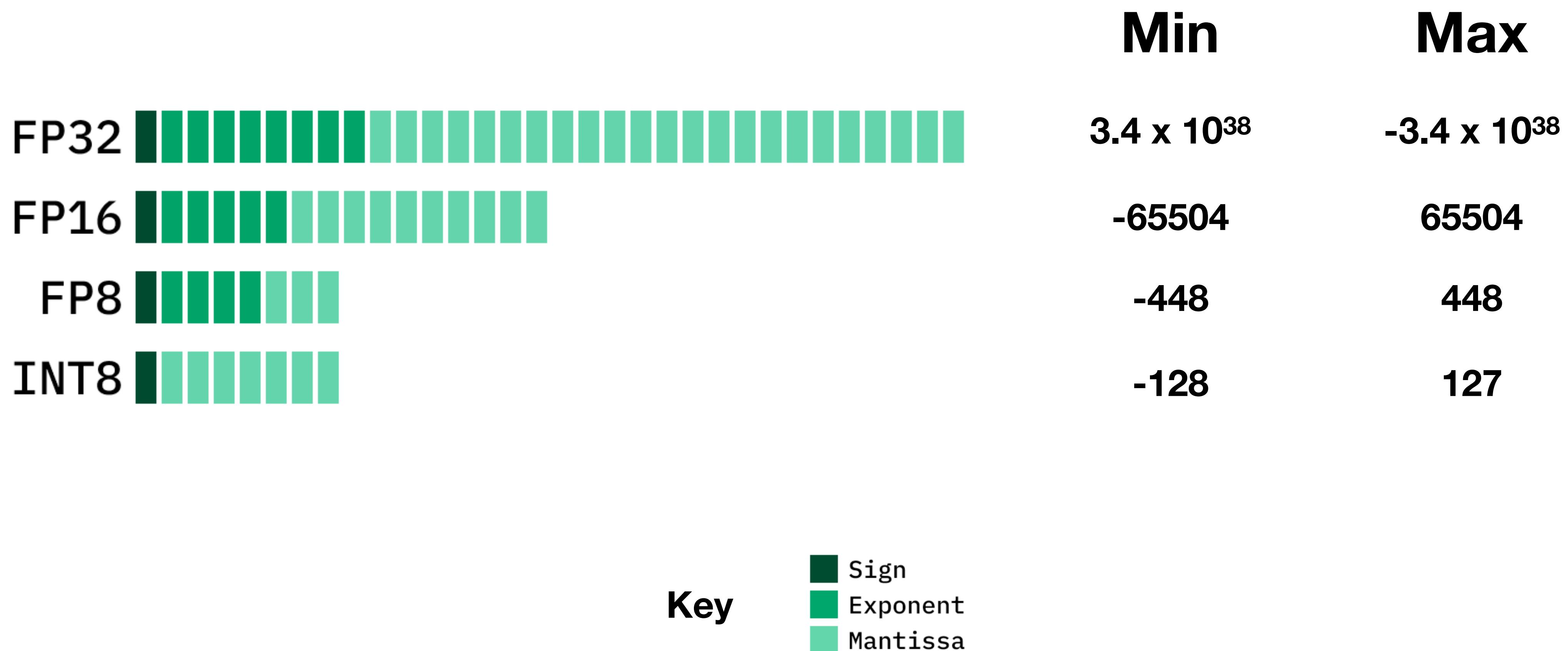
# Parameter Size

Operations are cheaper on lower bit widths



# Parameter Size

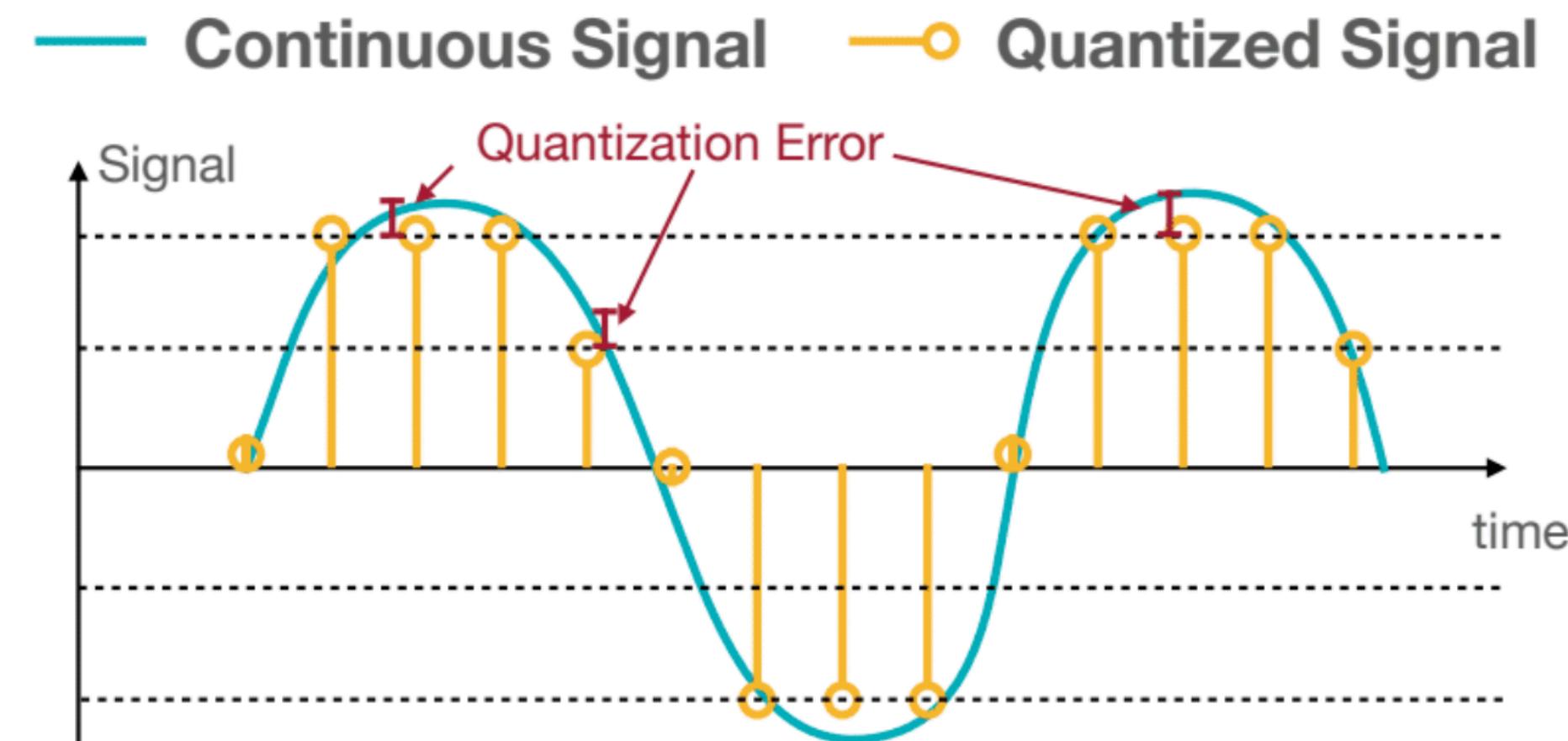
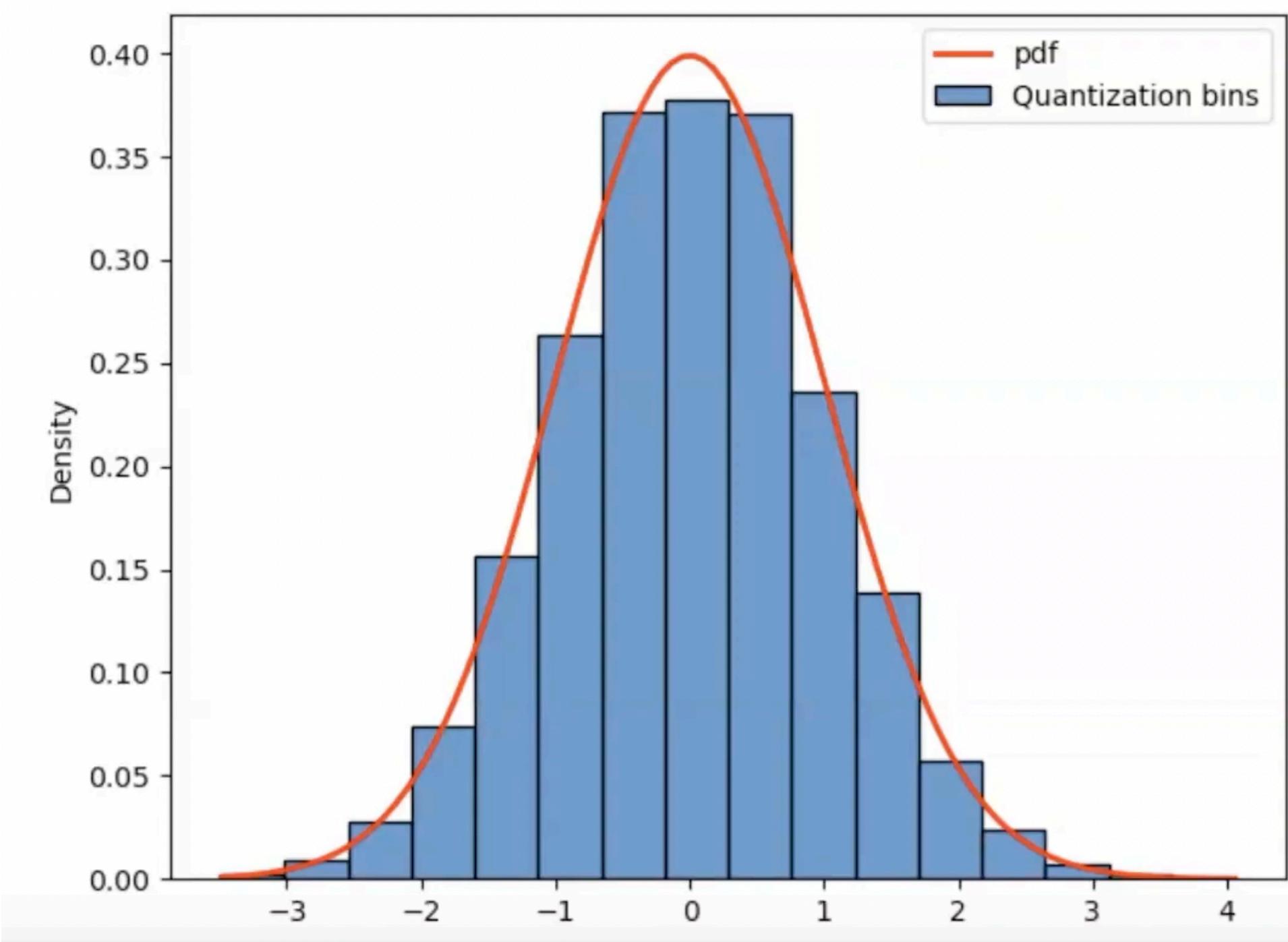
What do the underlying bits represent?



# Quantization

A well-defined tool in math & signal processing

- process of constraining an input from a continuous or otherwise large set of values to a smaller discrete set



# Quantization

## In the context of deep learning

- process of constraining an input from a continuous or otherwise large set of values to a smaller discrete set

**fp16 → int8**

while minimizing quantization error

# **Mathematical Quantization Techniques**

# Quantization

## Formalizing our problem

- Given a matrix of floating point numbers, we wish to...
  - store it as an integer matrix from which we can re-construct a float matrix
  - minimize the difference between the original matrix and the reconstructed matrix (“quantization error”)

# Quantization

## Formalizing our problem

- Let  $q$  be an integer. We wish to map it to a real number  $r$ 
  - Intuitively, we want to **shift** and **scale** our integer
  - Mathematically, we want to apply an **affine mapping**  
(that is, a transformation that preserves collinearity and ratio of distances)

$$r = S(q - Z) = Sq - Sz$$

$S$  is the scaling factor (float)

$Z$  is the zero-point (int)

# Zeropoint Quantization

Also known as “asymmetric quantization”

$$r = S(q - Z)$$

- Assume  $S = 1.5 \in \mathbb{R}$  and  $Z = 1 \in \mathbb{Z}$ , then
  - for  $q = 0$  we have  $r = 1.5(0 - 1) = 1.5(-1) = -1.5 \in \mathbb{R}$
  - for  $q = 1$  we have  $r = 1.5(1 - 1) = 1.5(0) = 0 \in \mathbb{R}$
  - for  $q = 2$  we have  $r = 1.5(2 - 1) = 1.5(1) = 1.5 \in \mathbb{R}$
  - for  $q = 3$  we have  $r = 1.5(3 - 1) = 1.5(2) = 3 \in \mathbb{R}$

# Zeropoint Quantization

We can go the other way too

$$r = S(q - Z) \quad \Rightarrow \quad q = \left\lfloor \frac{r}{S} + Z \right\rfloor$$

- Assume  $S = 1.5 \in \mathbb{R}$  and  $Z = 1 \in \mathbb{Z}$ , then

- for  $r = 1$  we have  $q = \left\lfloor \frac{1}{1.5} + 1 \right\rfloor = \lfloor 0.6 + 1 \rfloor = \lfloor 1.6 \rfloor = 2 \in \mathbb{R}$

- for  $r = 0$  we have  $q = \left\lfloor \frac{0}{1.5} + 1 \right\rfloor = \lfloor 0 + 1 \rfloor = \lfloor 1 \rfloor = 1 \in \mathbb{Z}$

# Zeropoint Quantization

We can go the other way too

$$r = S(q - Z) \quad \Rightarrow \quad q = \left\lfloor \frac{r}{S} + Z \right\rfloor$$

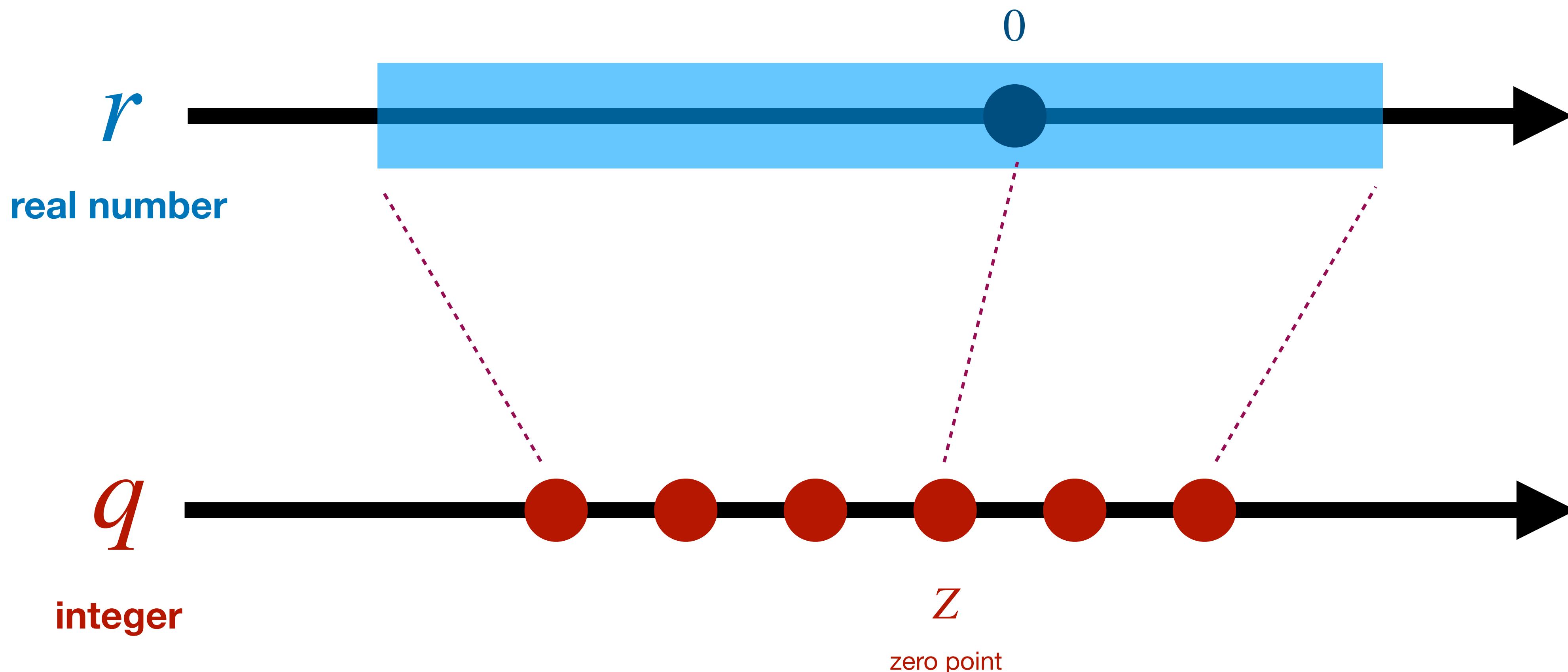
- Regardless of which  $S$  and  $Z$  we choose,  $r = 0$  maps **exactly** to  $q = Z$

$$q = \left\lfloor \frac{0}{S} + Z \right\rfloor = \lfloor Z \rfloor = Z$$

# Zeropoint Quantization

How do we determine  $S$  and  $Z$ ?

$$r = S(q - Z)$$



# Quantization

## Formalizing our problem

- Given a **specific** matrix,

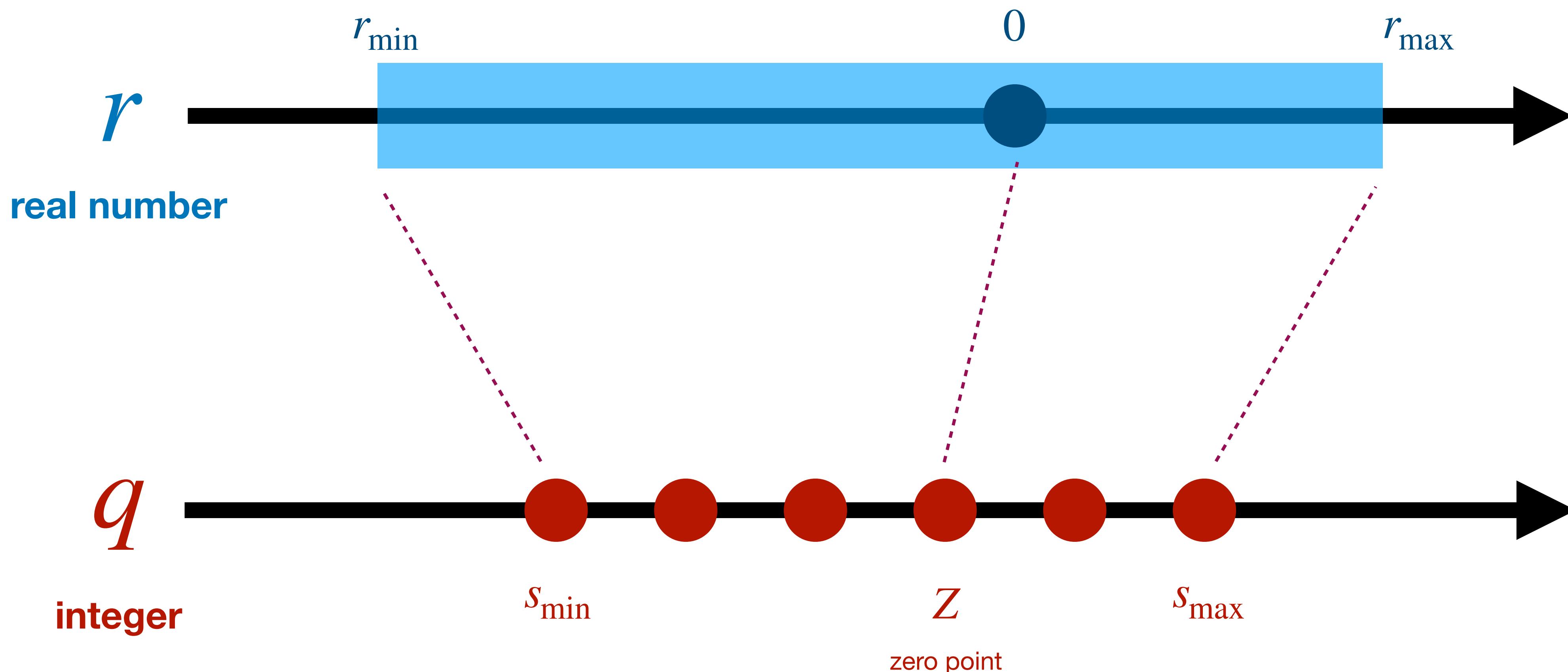
**fp16 → int8**

while minimizing quantization error

# Zeropoint Quantization

How do we determine  $S$  and  $Z$  ?

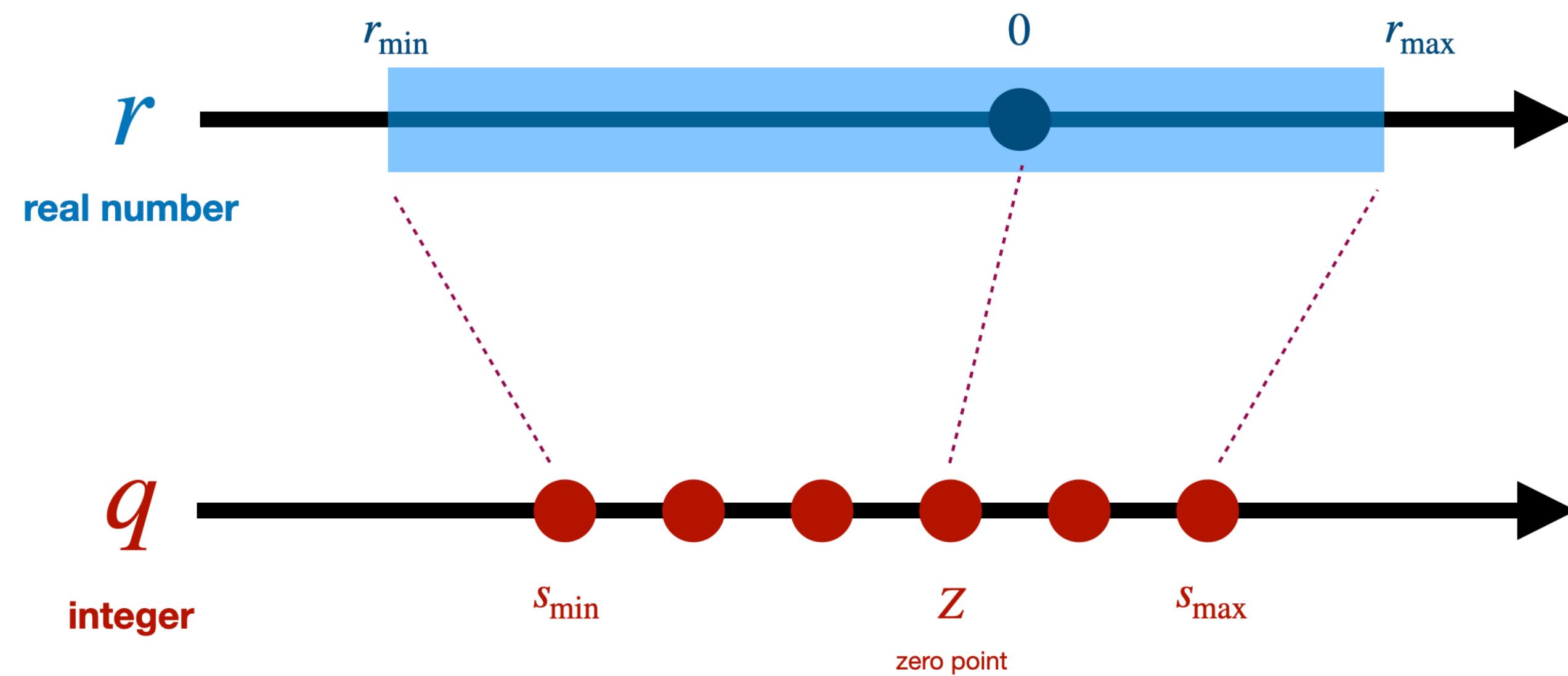
$$r = S(q - Z)$$



# Zeropoint Quantization

How do we determine  $S$  and  $Z$ ?

$$r = S(q - Z)$$



$$r_{\min} = S (q_{\min} - Z)$$

$$r_{\max} = S (q_{\max} - Z)$$

# Zeropoint Quantization

How do we determine  $S$  and  $Z$ ?

$$r_{\min} = S(q_{\min} - Z) \quad (1)$$

$$r_{\max} = S(q_{\max} - Z) \quad (2)$$

**Subtracting (1) from (2),**

$$r_{\max} - r_{\min} = S(q_{\max} - Z) - S(q_{\min} - Z)$$

**Simplifying,**

$$r_{\max} - r_{\min} = S(q_{\max} - q_{\min})$$



# Zeropoint Quantization

How do we determine  $S$  and  $Z$ ?

$$r_{\min} = S (q_{\min} - Z) \quad (1)$$

$$S = \frac{r_{\max} - r_{\min}}{q_{\max} - q_{\min}} \quad (3)$$

We now have a quantity for  $S$ , so we can solve (1) for  $Z$ ,

$$r_{\min} = S (q_{\min} - Z) \Rightarrow Z = \left\lfloor q_{\min} - \frac{r_{\min}}{S} \right\rfloor$$

# Zeropoint Quantization

Quantize the following matrix to 2-bit signed ints

Recall from CSE 351 that a  $N$ -bit signed integer has the range  $[-2^{N-1}, 2^{N-1} - 1]$

2.09	-0.98	1.48	0.09
0.05	-0.14	<b>-1.08</b>	<b>2.12</b>
0.91	1.92	0	-1.03
1.87	0	1.53	1.49

$$r_{\min} = -1.08$$

$$r_{\max} = 2.12$$

$$S = \frac{r_{\max} - r_{\min}}{q_{\max} - q_{\min}} = \frac{2.12 - (-1.08)}{1 - (-2)} = 1.07 \in \mathbb{R}$$

$$Z = \left\lfloor q_{\min} - \frac{r_{\min}}{S} \right\rfloor = \left\lfloor -2 - \frac{-1.08}{1.07} \right\rfloor = -1 \in \mathbb{Z}$$

# Zeropoint Quantization

Quantize the following matrix to 2-bit signed ints

2.09	-0.98	1.48	0.09
0.05	-0.14	-1.08	2.12
0.91	1.92	0	-1.03
1.87	0	1.53	1.49

$$r = S(q - Z) \quad \Rightarrow \quad q = \left\lfloor \frac{r}{S} + Z \right\rfloor$$

$$q = \left\lfloor \frac{r}{1.07} - 1 \right\rfloor$$

$$S = 1.07$$

$$Z = -1$$

# Zeropoint Quantization

Quantize the following matrix to 2-bit signed ints

Original Weights (float)

2.09	-0.98	1.48	0.09
0.05	-0.14	-1.08	2.12
0.91	1.92	0	-1.03
1.87	0	1.53	1.49



$$q = \left\lfloor \frac{r}{1.07} - 1 \right\rfloor$$

Quantized Weights (int2)

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

Examples

$$q = \left\lfloor \frac{r}{1.07} - 1 \right\rfloor = \left\lfloor \frac{1.48}{1.07} - 1 \right\rfloor = [1.38 - 1] = [0.38] = 0$$

$$q = \left\lfloor \frac{r}{1.07} - 1 \right\rfloor = \left\lfloor \frac{1.87}{1.07} - 1 \right\rfloor = [1.75 - 1] = [0.75] = 1$$

# Zeropoint Quantization

Determine the quantization error

Quantized Weights (int2)

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

Reconstructed Weights (float)

2.14	-1.07	1.07	0
0	0	-1.07	2.14
-1.07	2.14	0	-1.07
2.14	0	1.07	1.07



$$r = 1.07(q + 1)$$

Examples

$$r = 1.07(q + 1) = 1.07(0 + 1) = 1.07$$

$$r = 1.07(q + 1) = 1.07(1 + 1) = 2.14$$

$$r = 1.07(q + 1) = 1.07(-1 + 1) = 0$$

# Zeropoint Quantization

Determine the quantization error

**Original Weights (float)**

2.09	-0.98	1.48	0.09
0.05	-0.14	-1.08	2.12
0.91	1.92	0	-1.03
1.87	0	1.53	1.49

**Reconstructed Weights (float)**

2.14	-1.07	1.07	0
0	0	-1.07	2.14
-1.07	2.14	0	-1.07
2.14	0	1.07	1.07

**Quantization Error**

-0.05	0.09	0.41	0.09
0.05	-0.14	-0.01	-0.02
0.16	-0.22	0	0.04
-0.27	0	0.46	0.42

# Zeropoint Quantization

## Practice!

- Quantize the following matrix into 3-bit signed integers
- Reconstruct the weights from the quantized weights
- Calculate the quantization error

18	7	8
3	6	12
0	-3	-2

Recall from CSE 351 that a  $N$ -bit signed integer has the range  $[-2^{N-1}, 2^{N-1} - 1]$

# Zeropoint Quantization

## Practice!

Original Weights (float)

15	6	9
3	-3	12
0	-6	0

Quantized Weights (int3)

4	1	2
0	-2	3
-1	-3	-1



$$q = \left\lfloor \frac{r}{3} - 1 \right\rfloor$$

$$r_{\min} = -6$$

$$r_{\max} = 15$$

$$S = 3$$

$$Z = -1$$

# Zeropoint Quantization

## Practice!

Quantized Weights (int3)

4	1	2
0	-2	3
-1	-3	1

Reconstructed Weights (float)

15	6	9
3	-3	12
0	-6	0



$$r = 3(q + 1)$$

# Zeropoint Quantization

## Practice!

**Original Weights (float)**

15	6	9
3	-3	12
0	-6	0

**Reconstructed Weights (float)**

15	6	9
3	-3	12
0	-6	0

# Zeropoint Quantization

## Matmul as an integer-arithmetic-only operation

- We mostly care about the first equation
- $S$  is still stored as a float 😔

$$r = S(q - Z)$$

$$q = \left\lfloor \frac{r}{S} + Z \right\rfloor$$

# Zeropoint Quantization

## Matmul as an integer-arithmetic-only operation

$$r = S(q - Z)$$

**Let  $A, B \in \mathbb{R}^{n \times n}$ . Let  $\hat{A}, \hat{B} \in \mathbb{Z}^{n \times n}$  be their quantized versions. Let  $C = AB$ . Find  $\hat{C}$ .**

Notation: for some matrix  $M$ , we will refer to the element in its  $i$ -th row and  $j$ -th column as  $m_{ij}$ .

**By the definition of matrix multiplication,**

$$c_{ij} = \sum_{k=1}^N a_{ik} b_{kj}$$

**Applying our quantization schemes,**

$$S_C(\hat{c}_{ij} - Z_C) = \sum_{k=1}^N S_A(\hat{a}_{ij} - Z_A) S_B(\hat{b}_{ij} - Z_B)$$

# Zeropoint Quantization

Matmul as an integer-arithmetic-only operation

$$S_C (\hat{c}_{ij} - Z_C) = \sum_{k=1}^N S_A (\hat{a}_{ij} - Z_A) S_B (\hat{b}_{ij} - Z_B)$$

Factoring constants out of the sum,

$$S_C (\hat{c}_{ij} - Z_C) = S_A S_B \sum_{k=1}^N (\hat{a}_{ij} - Z_A) (\hat{b}_{ij} - Z_B)$$

Isolating  $\hat{c}_{ij}$  on the left,

$$\hat{c}_{ij} = Z_C + \frac{S_A S_B}{S_C} \sum_{k=1}^N (\hat{a}_{ij} - Z_A) (\hat{b}_{ij} - Z_B)$$

# Zeropoint Quantization

## Matmul as an integer-arithmetic-only operation

$$\hat{c}_{ij} = Z_C + \frac{S_A S_B}{S_C} \sum_{k=1}^N (\hat{a}_{kj} - Z_A) (\hat{b}_{kj} - Z_B)$$

**Define**  $M := \frac{S_A S_B}{S_C}$ . **Then,**

$$\hat{c}_{ij} = Z_C + M \sum_{k=1}^N (\hat{a}_{kj} - Z_A) (\hat{b}_{kj} - Z_B)$$

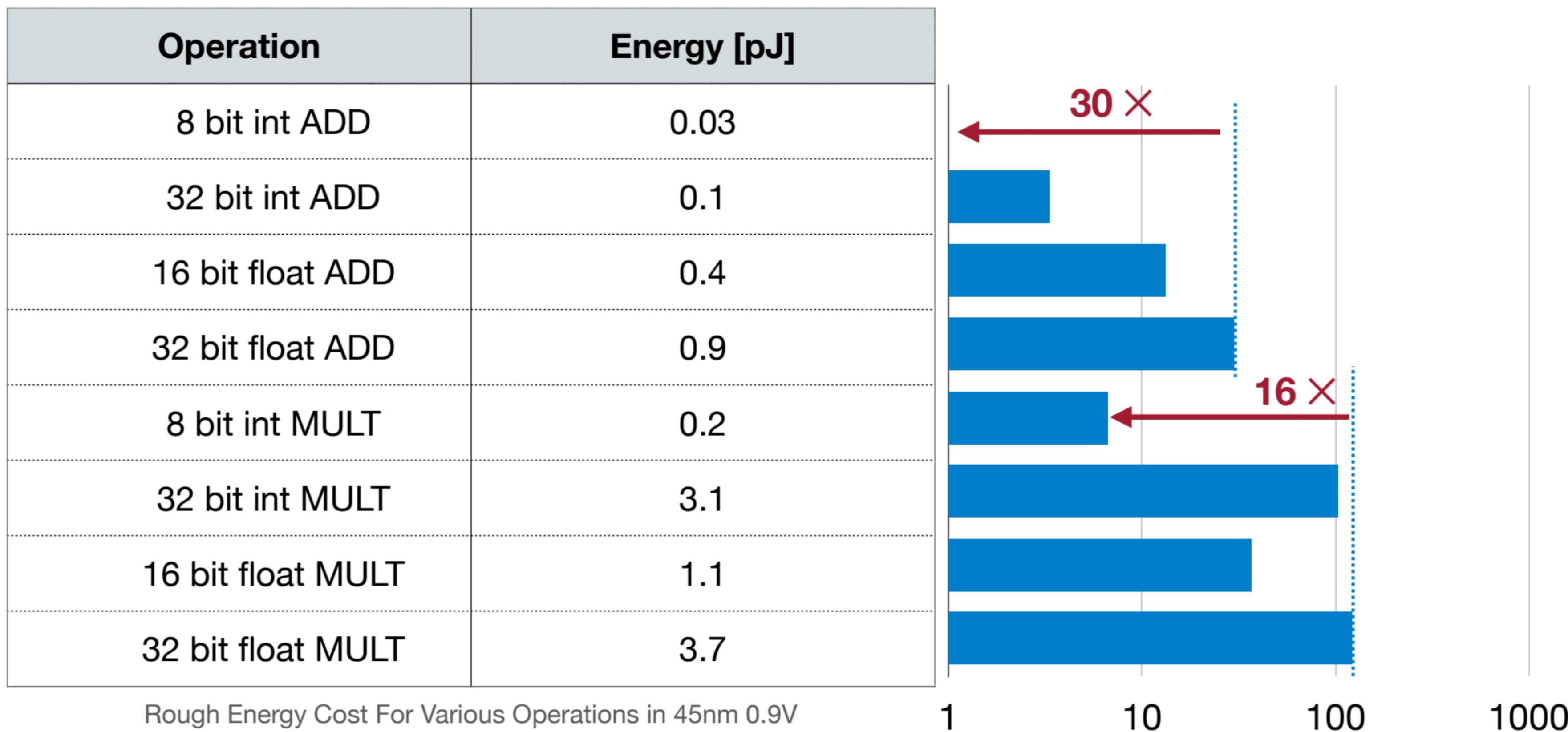
**Empirically,  $M \in (0, 1)$ . Normalizing for some  $M_0 \in [0.5, 1)$ ,**

$$M = \underbrace{2^{-n}}_{\text{bit shift}} \cdot \underbrace{M_0}_{\text{fixed-point multiplier}}$$

Assuming int32, this will always have 30 bits of relative accuracy. Why?

# Zeropoint Quantization

We can now do integer-only matmul!



# Zeropoint Quantization

## Practical Considerations

$$\hat{c}_{ij} = Z_C + M \sum_{k=1}^N (\hat{a}_{ij} - Z_A) (\hat{b}_{ij} - Z_B)$$

The bulk of the computation occurs in the sum:

$$(\hat{a}_{ij} - Z_A) (\hat{b}_{ij} - Z_B) = \sum_{k=1}^N \hat{a}_{ij} \hat{b}_{ij} - \hat{a}_{ij} Z_B - \hat{b}_{ij} Z_A + Z_A Z_B$$

int8

int16

CPUs contain PMADDUBSW, but most GPUs/TPUs don't.

The zeropoint slows us down tremendously!

# Absolute Maximum Quantization

Also known as “symmetric quantization”

- What if get rid of the zeropoint?

$$r = S(q - Z) \quad \Rightarrow \quad r = Sq$$

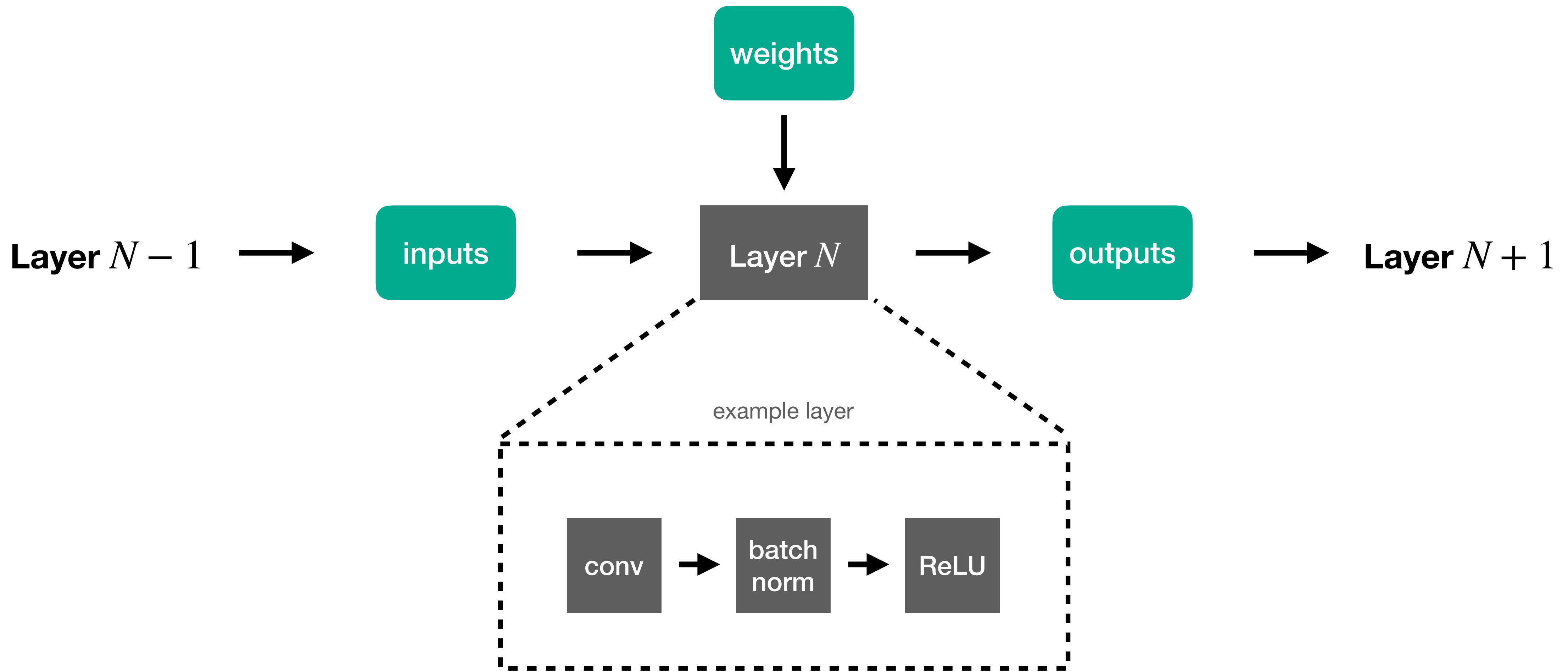
- This forces a symmetric quantization scheme. **Why?**

$$S = \frac{r_{\max} - r_{\min}}{q_{\max} - q_{\min}} \approx \frac{2 \cdot r_{absmax}}{2 \cdot q_{\max}} = \frac{r_{absmax}}{q_{\max}}$$

# Deep Learning Quantization Paradigms

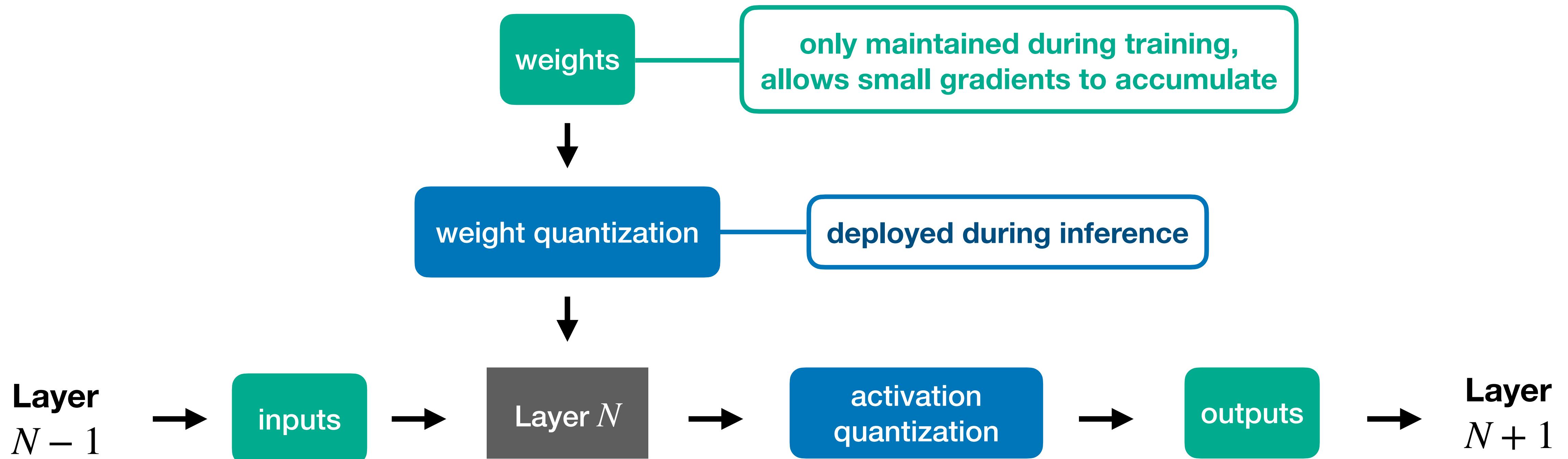
# Quantization-Aware Training

## Simulated / Fake Quantization



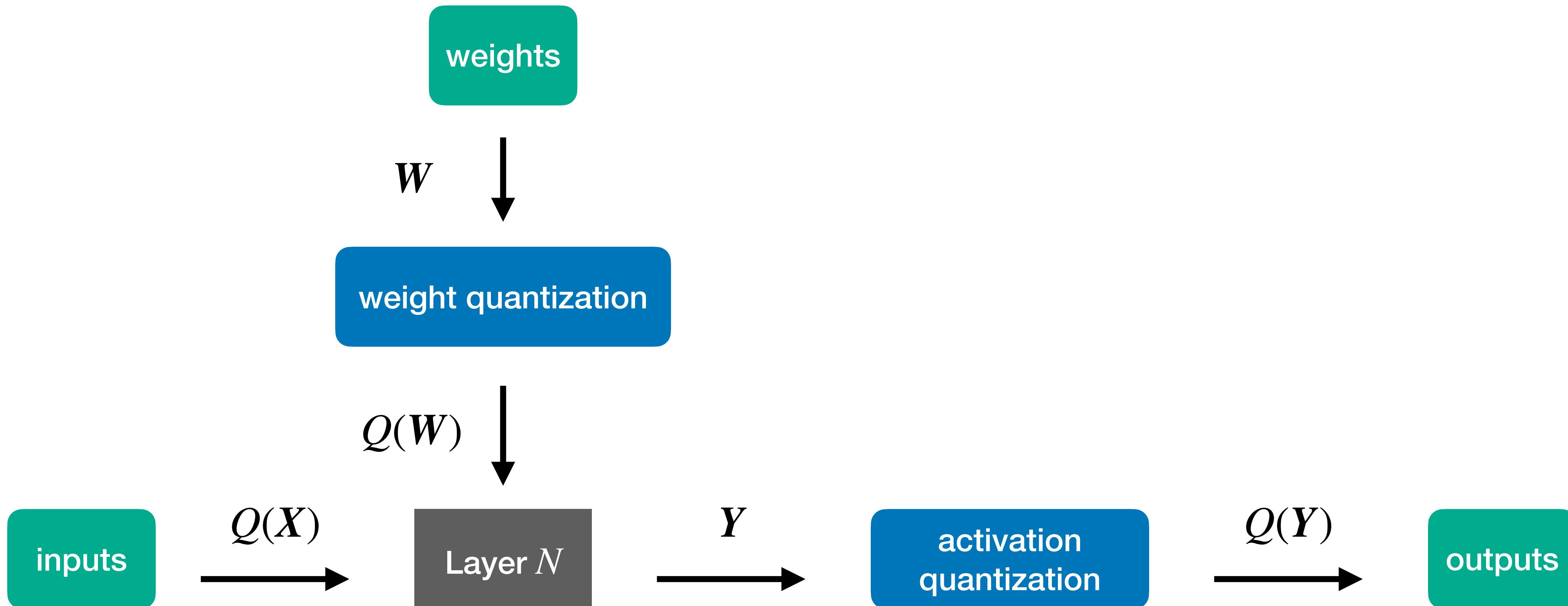
# Quantization-Aware Training

## Simulated / Fake Quantization



# Quantization-Aware Training

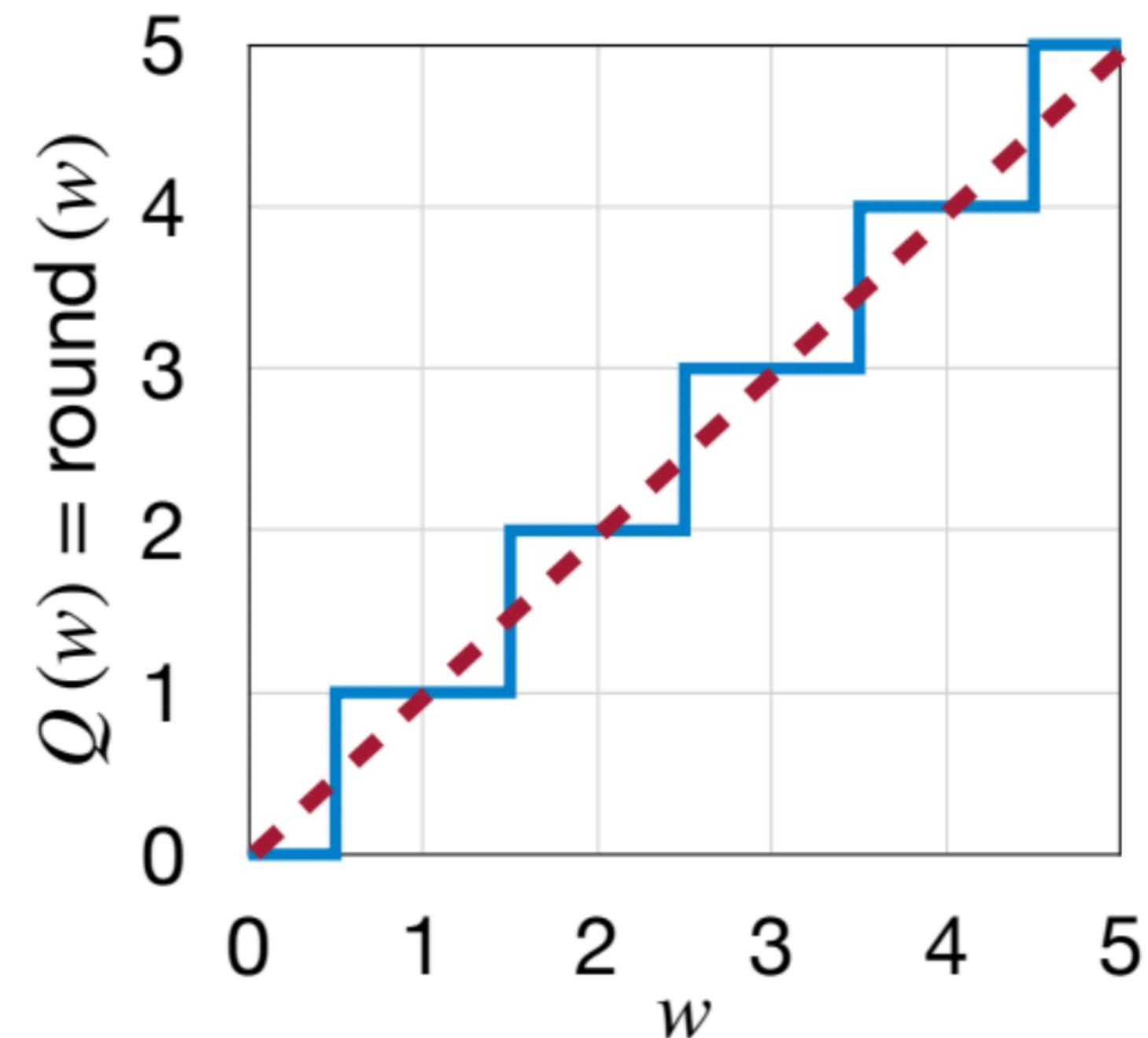
## Simulated / Fake Quantization



# Quantization-Aware Training

## Straight-Through Estimator (STE)

- Quantization is a step function  $\Rightarrow$  gradient is almost always 0



# Quantization-Aware Training

## Straight-Through Estimator (STE)

- Quantization is a step function  $\Rightarrow$  gradient is almost always 0
- STE passes gradients along “as if it had been the identity function”

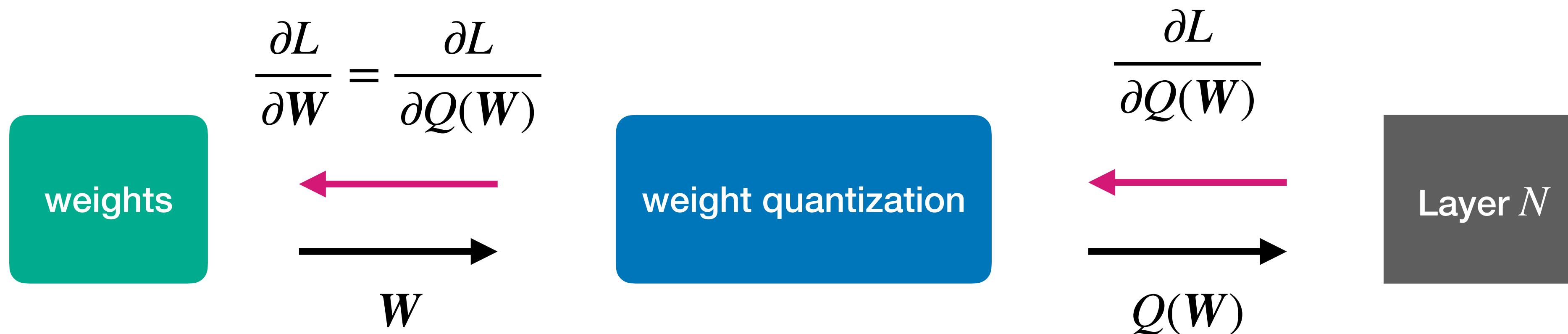
$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial Q(W)} \cdot \frac{\partial Q(W)}{\partial W} = \frac{\partial L}{\partial Q(W)} \cdot 1 = \frac{\partial L}{\partial Q(W)}$$

chain rule      straight-through estimator

# Quantization-Aware Training

## Straight-Through Estimator (STE)

- Quantization is a step function  $\Rightarrow$  gradient is almost always 0
- STE passes gradients along “as if it had been the identity function”



# Quantization-Aware Training

## Practical Considerations

- Quite effective, especially useful for really large models
  - Requires training (or re-training) the model 😔
  - Also called “Mixed Precision Training”

# Post-Training Quantization

## Practical Considerations

<b>Model Size</b>	<b>Model Name</b>	<b>Asymmetric PTQ</b>	<b>Symmetric PTQ</b>	<b>Asymmetric QAT</b>	<b>Symmetric QAT</b>	<b>Floating Point</b>
3.5M	MobileNet-v2	0.1%	69.8%	70.9%	71.1%	71.9%
25M	ResNet-50	75%	75%	75%	75%	75.6%
60M	ResNet-152	76.1%	76%	76%	76%	77.8%

**Smaller models suffer from PQT, likely because they have a smaller representational capacity**

# **Quantization Libraries**

**& their underlying PTQ methods**

# LLM.int8()

## Seminal work in quantization for LLMs

- Implemented in bitsandbytes
- Integrated with 😊 Transformers

Parameters	125M	1.3B	2.7B	6.7B	13B
32-bit Float	25.65	15.91	14.43	13.30	12.45
Int8 absmax	87.76	16.55	15.11	14.59	19.08
Int8 zeropoint	56.66	16.24	14.76	13.49	13.94
Int8 absmax row-wise	30.93	17.08	15.24	14.13	16.49
Int8 absmax vector-wise	35.84	16.82	14.98	14.13	16.48
Int8 zeropoint vector-wise	25.72	15.94	14.36	13.38	13.47
Int8 absmax row-wise + decomposition	30.76	16.19	14.65	13.25	12.46
Absmax LLM.int8() (vector-wise + decomp)	25.83	15.93	14.44	13.24	12.45
Zeropoint LLM.int8() (vector-wise + decomp)	<b>25.69</b>	<b>15.92</b>	<b>14.43</b>	<b>13.24</b>	<b>12.45</b>

perplexity scores

previously discussed methods

what we're about to discuss

# LLM.int8()

## Absmax Quantization

$$r = S(q - Z) \Rightarrow r = Sq \Rightarrow q = \left\lfloor \frac{r}{S} \right\rfloor$$

$$S = \frac{r_{\max} - r_{\min}}{q_{\max} - q_{\min}} \approx \frac{r_{absmax}}{q_{\max}} = \frac{r_{absmax}}{127}$$

$\text{int8} \Rightarrow q_{\max} = 127$

# LLM.int8()

## Absmax Quantization

$$r = S(q - Z) \Rightarrow r = Sq \Rightarrow q = \left\lfloor \frac{r}{S} \right\rfloor$$

$$S = \frac{r_{\max} - r_{\min}}{q_{\max} - q_{\min}} \approx \frac{r_{absmax}}{q_{\max}} = \frac{r_{absmax}}{127}$$

we will refer to this value as  $s$

quantized  $\rightarrow$  real

$$\frac{r_{absmax}}{127}$$

real  $\rightarrow$  quantized

$$\frac{127}{r_{absmax}}$$

# LLM.int8()

## Absmax Quantization

- $\frac{1}{s}$  : quantized  $\rightarrow$  real
- $s$  : real  $\rightarrow$  quantized

quantized  $\rightarrow$  real

$$\frac{r_{absmax}}{127}$$

real  $\rightarrow$  quantized

$$\frac{127}{r_{absmax}}$$

# LLM.int8()

## int8 matmul with fp16 inputs and outputs

- Given inputs  $X_{f16} \in \mathbb{R}^{s \times h}$  and  $W_{f16} \in \mathbb{R}^{h \times o}$ , we compute  $C_{f16} \in \mathbb{R}^{s \times o}$

$$X_{f16} W_{f16} = C_{f16}$$



from previous slide

# LLM.int8()

## int8 matmul with fp16 inputs and outputs

- Given inputs  $X_{f16} \in \mathbb{R}^{s \times h}$  and  $W_{f16} \in \mathbb{R}^{h \times o}$ , we compute  $C_{f16} \in \mathbb{R}^{s \times o}$

$$\begin{aligned} X_{f16} W_{f16} = C_{f16} &\approx \frac{1}{s_{X_{f16}} s_{W_{f16}}} \cdot C_{i32} \\ &\approx \frac{1}{s_{X_{f16}} s_{W_{f16}}} \cdot X_{i8} W_{i8} \end{aligned}$$

**Why int32?**

# LLM.int8()

## int8 matmul with fp16 inputs and outputs

- Given inputs  $X_{f16} \in \mathbb{R}^{s \times h}$  and  $W_{f16} \in \mathbb{R}^{h \times o}$ , we compute  $C_{f16} \in \mathbb{R}^{s \times o}$

$$\begin{aligned} X_{f16} W_{f16} = C_{f16} &\approx \frac{1}{S_{X_{f16}} S_{W_{f16}}} \cdot C_{i32} \\ &\approx \frac{1}{S_{X_{f16}} S_{W_{f16}}} \cdot X_{i8} W_{i8} \\ &= \frac{1}{S_{X_{f16}} S_{W_{f16}}} \cdot Q(X_{f16}) Q(W_{f16}) \end{aligned}$$

# **LLM.int8()**

## **Vector-wise Quantization**

- More scaling constants  $\Rightarrow$  improved quantization error
- Matrix multiplication is a series of dot products
- Idea: one scaling factor for each row of  $X$  and each column of  $W$

# LLM.int8()

## Vector-wise Quantization

- Given inputs  $X_{f16} \in \mathbb{R}^{s \times h}$  and  $W_{f16} \in \mathbb{R}^{h \times o}$ , we compute  $C_{f16} \in \mathbb{R}^{s \times o}$ 
  - Assign a different scaling constant  $s_{x_{f16}}$  to each row of  $X_{f16}$
  - Assign a different scaling constant  $s_{w_{f16}}$  to each col of  $W_{f16}$

# LLM.int8()

## Vector-wise Quantization

**first entry of  $C = (\text{first row of } X) \cdot (\text{first col of } W)$**

$$\approx \frac{1}{s \text{ for first row of } X \cdot s \text{ for first col of } W} \cdot Q(\text{first row of } X) \cdot Q(\text{first col of } W)$$

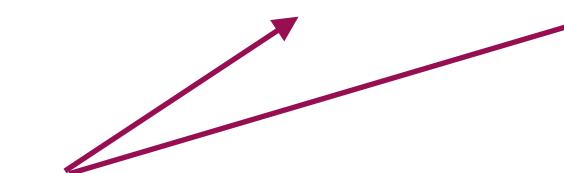
**first entry of  $C =$**

2	6	-15
---	---	-----

7
2
-4

$$\approx \frac{1}{8.466 \cdot 18.143} \cdot Q(\boxed{2 \quad 6 \quad -15}) \cdot Q(\boxed{\begin{array}{c} 7 \\ 2 \\ -4 \end{array}})$$

derived via  $\frac{127}{absmax}$  formula



# LLM.int8()

## Vector-wise Quantization

$$= \frac{1}{8.46\bar{6} \cdot 18.143} \cdot Q(\begin{array}{|c|c|c|}\hline 2 & 6 & -15 \\ \hline\end{array}) \cdot Q(\begin{array}{|c|}\hline 7 \\ \hline 2 \\ \hline -4 \\ \hline\end{array})$$

$$= \frac{1}{8.46\bar{6} \cdot 18.143} \cdot \begin{array}{|c|c|c|}\hline [2 \cdot 8.46\bar{6}] & [6 \cdot 8.46\bar{6}] & [-15 \cdot 8.46\bar{6}] \\ \hline\end{array} \cdot \begin{array}{|c|}\hline [18.143 \cdot 7] \\ \hline [18.143 \cdot 2] \\ \hline [18.143 \cdot -4] \\ \hline\end{array}$$

$$= \frac{1}{153.61} \cdot \begin{array}{|c|c|c|}\hline 17 & 51 & -127 \\ \hline\end{array} \cdot \begin{array}{|c|}\hline 127 \\ \hline 36 \\ \hline -73 \\ \hline\end{array}$$

$$= \frac{1}{153.61} \cdot (17(127) + 51(36) - 127(-73)) = \frac{1}{153.61} \cdot (13266) = 86.36$$

for reference,  
the correct result is 86

# LLM.int8()

## Vector-wise Quantization

$X_{f16}$

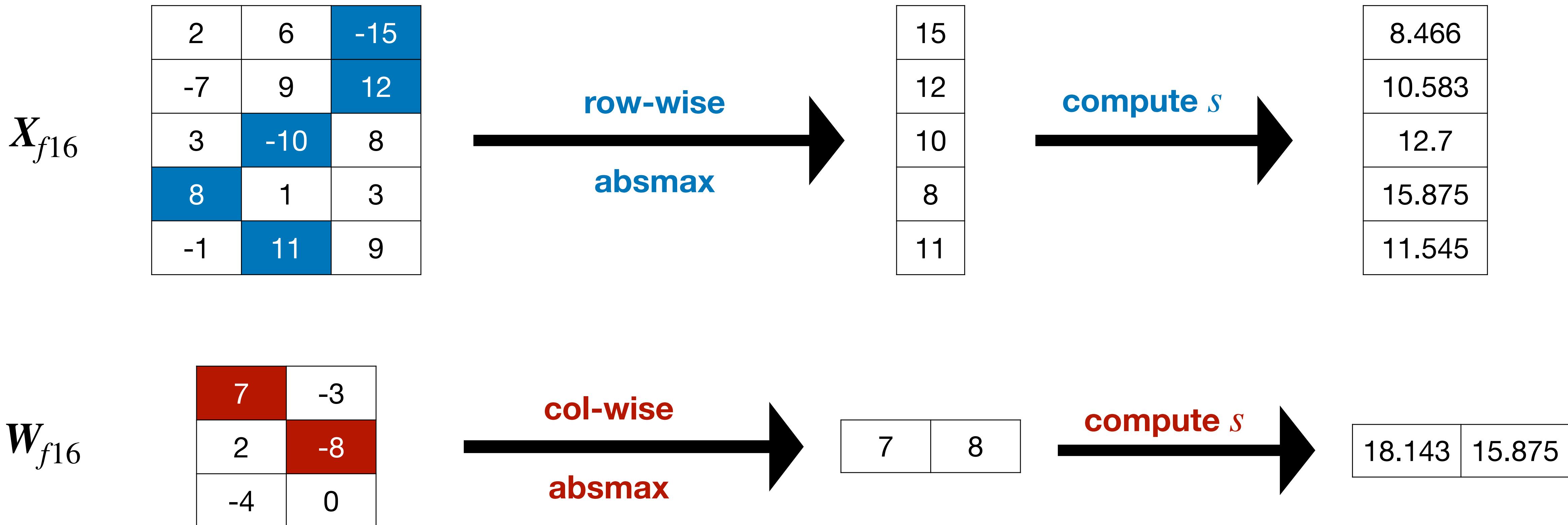
2	6	-15
-7	9	12
3	-10	8
8	1	3
-1	11	9

$W_{f16}$

7	-3
2	8
-4	0

# LLM.int8()

## Vector-wise Quantization



real  $\rightarrow$  quantized

$$s_{f16} = \frac{127}{r_{absmax}}$$

# LLM.int8()

## Vector-wise Quantization

$X_{f16}$

2	6	-15
-7	9	12
3	-10	8
8	1	3
-1	11	9

$s_{f16}$

8.466
10.583
12.7
15.875
11.545

$$Q(X_{f16}) = \left\lfloor s_{f16} X_{f16} \right\rfloor$$

17	51	-127
-74	95	127
38	-127	102
127	16	3
-12	127	104

$W_{f16}$

7	-3
2	-8
-4	0

$s_{f16}$

18.143	15.875
--------	--------

$$Q(W_{f16}) = \left\lfloor s_{f16} W_{f16} \right\rfloor$$

127	-48
36	-127
-73	0

# LLM.int8()

## Vector-wise Quantization

- Recall our dequantization equation

$$C_{f16} \approx \frac{1}{s_{x_{f16}} s_{w_{f16}}} \cdot Q(X_{f16}) Q(W_{f16})$$

- We now have multiple  $s$
- How can we efficiently dequantize?

# LLM.int8()

## Vector-wise Quantization

- How can we efficiently dequantize?
  - $ij$ -th entry of  $C$  is the dot product of the  $i$ -th row of  $X$  & the  $j$ -th col of  $W$
  - $i$ -th row of  $X$  was quantized with  $i$ -th entry of  $s_{f16}$
  - $j$ -th col of  $W$  was quantized with  $j$ -th entry of  $s_{f16}$

What should the  $ij$ -th entry of  $C$  be dequantized by?

# LLM.int8()

## Vector-wise Quantization

- How can we efficiently dequantize?
  - $ij$ -th entry of  $C$  is the dot product of the  $i$ -th row of  $X$  & the  $j$ -th col of  $W$
  - $i$ -th row of  $X$  was quantized with  $i$ -th entry of  $s_{f16}$
  - $j$ -th col of  $W$  was quantized with  $j$ -th entry of  $s_{f16}$

What should  $C$  be dequantized by?

$$s_{f16} \otimes s_{f16}$$

# LLM.int8()

## Vector-wise Quantization

$$S_{f16} \otimes S_{f16} = \begin{array}{|c|} \hline 8.466 \\ \hline 10.583 \\ \hline 12.7 \\ \hline 15.875 \\ \hline 11.545 \\ \hline \end{array} \otimes \begin{array}{|c|c|} \hline 18.143 & 15.875 \\ \hline \end{array} = \begin{array}{|c|c|} \hline 153.610 & 134.408 \\ \hline 192.012 & 168.010 \\ \hline 230.414 & 201.613 \\ \hline 288.018 & 252.016 \\ \hline 209.468 & 183.284 \\ \hline \end{array}$$

# LLM.int8()

## Vector-wise Quantization

- Given inputs  $X_{f16} \in \mathbb{R}^{b \times h}$  and  $W_{f16} \in \mathbb{R}^{h \times o}$ , we compute  $C_{f16} \in \mathbb{R}^{s \times o}$

$$C_{f16} = \frac{1}{S_{f16} \otimes S_{f16}}$$

$\vdots$

$$Q(X_{f16}) Q(W_{f16})$$

**b × o matrix**      **b × o matrix**      **b × o matrix**

**elementwise  
multiplication**

# LLM.int8()

## Vector-wise Quantization

- Given inputs  $X_{f16} \in \mathbb{R}^{b \times h}$  and  $W_{f16} \in \mathbb{R}^{h \times o}$ , we compute  $C_{f16} \in \mathbb{R}^{s \times o}$

$$C_{f16} = \frac{1}{s_{f16} \otimes s_{f16}} \cdot Q(X_{f16}) Q(W_{f16})$$

$\left[ \begin{array}{c} s_{f16} X_{f16} \\ s_{f16} W_{f16} \end{array} \right] \quad \parallel \quad \parallel$

# LLM.int8()

## Vector-wise Quantization

- Given inputs  $X_{f16} \in \mathbb{R}^{b \times h}$  and  $W_{f16} \in \mathbb{R}^{h \times o}$ , we compute  $C_{f16} \in \mathbb{R}^{s \times o}$

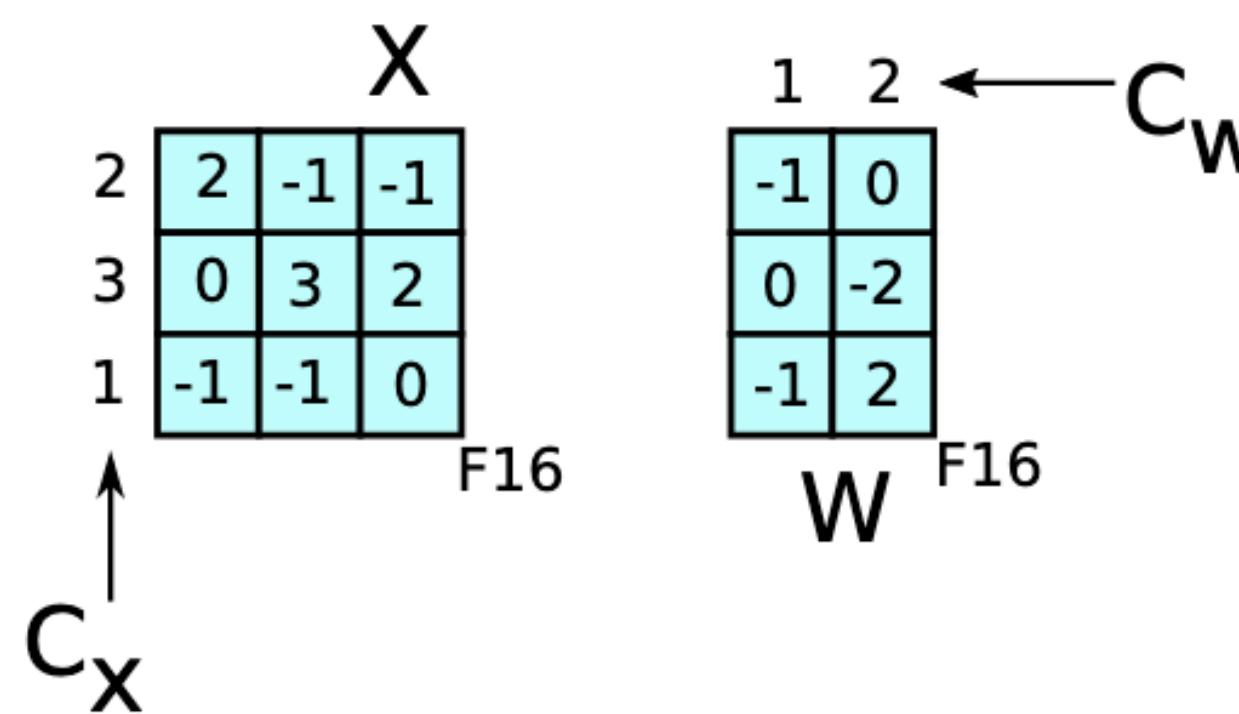
$$C_{f16} = s \cdot Q(X_{f16}) \cdot Q(W_{f16})$$

$\left[ \begin{array}{c} s_{f16} X_{f16} \\ \parallel \\ s_{f16} W_{f16} \end{array} \right]$

# LLM.int8()

## Vector-wise Quantization

(1) Find vector-wise constants:  $C_w$  &  $C_x$



(2) Quantize

$$X_{\text{F16}} * (127/C_x) = X_{\text{I8}}$$
$$W_{\text{F16}} * (127/C_w) = W_{\text{I8}}$$

(3) Int8 Matmul

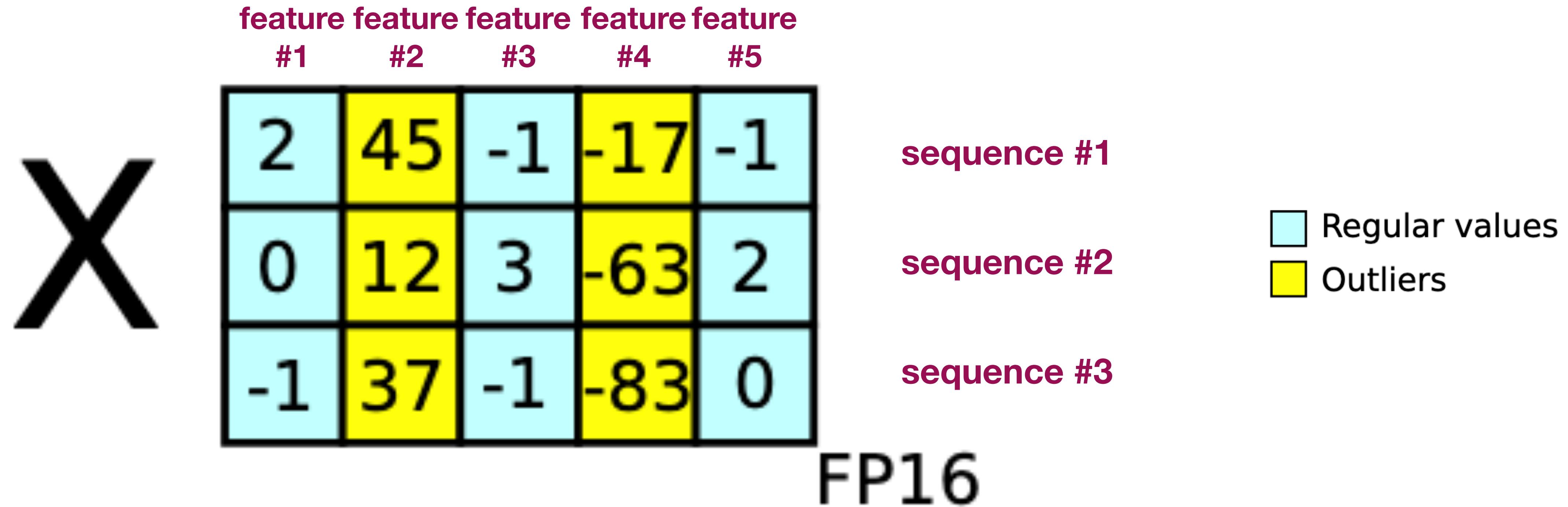
$$X_{\text{I8}} \cdot W_{\text{I8}} = \text{Out}_{\text{I32}}$$

(4) Dequantize

$$\frac{\text{Out}_{\text{I32}} * (C_x \otimes C_w)}{127*127} = \text{Out}_{\text{F16}}$$

# LLM.int8()

## Outliers



# LLM.int8()

## Outliers

- Empirically, large (>6B) transformer models have **outliers**:
  - Large magnitude features (columns)
  - Extremely important for performance
  - Require high quantization precision

$$\mathbf{X} \begin{array}{|c|c|c|c|c|} \hline 2 & 45 & -1 & -17 & -1 \\ \hline 0 & 12 & 3 & -63 & 2 \\ \hline -1 & 37 & -1 & -83 & 0 \\ \hline \end{array}$$

FP16

# **LLM.int8()**

## **Outliers**

- Empirically, large transformer models have...
  - 99.9% regular features – medium quantization precision OK
  - 0.01% outlier features – require high quantization precision

**Great, we already employ vector-wise quantization, right?**

# LLM.int8()

## Outliers

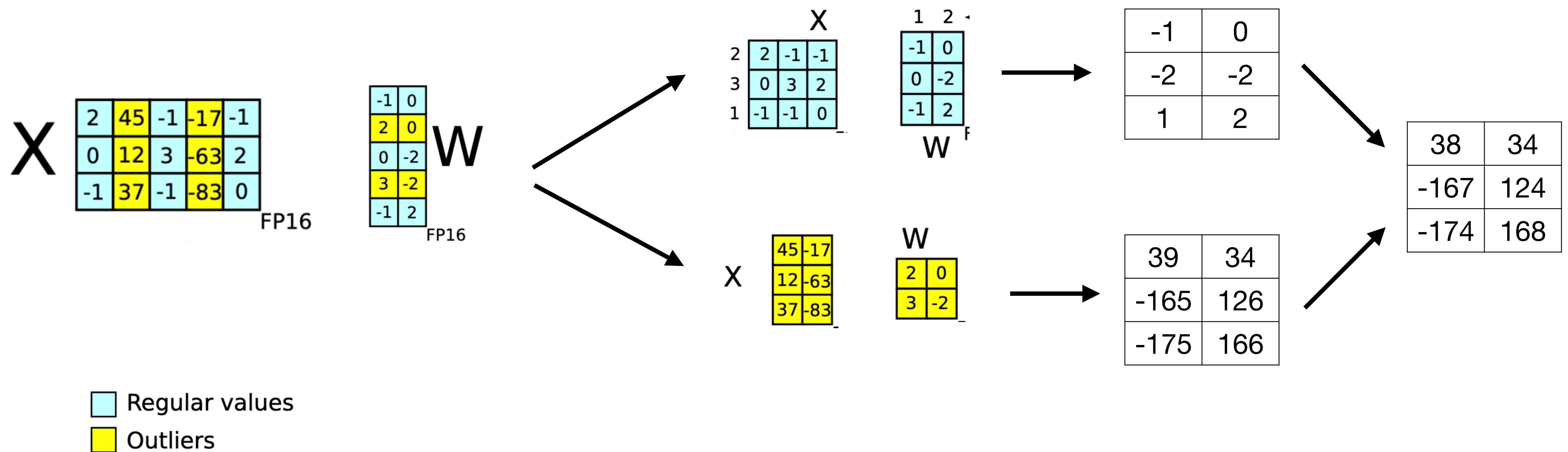
- Vector-wise quantization assigns different scaling factors for rows in  $X$
- Outliers occur in columns in  $X$

**Vector-wise quantization doesn't help**

# LLM.int8()

## Mixed-precision Decomposition

- What if we handled the normal features and outlier features separately?



# Aside: Einstein Notation

This is non-standard



$$X^2 W^2 + X^4 W^4$$

# LLM.int8()

## Mixed-precision Decomposition

- Notation for handling outliers  $O$  “separately” (still with normal fp16 matmul)

$$C_{f16} = \sum_{h \in O} X_{f16}^h W_{f16}^h + \sum_{h \notin O} X_{f16}^h W_{f16}^h$$

# LLM.int8()

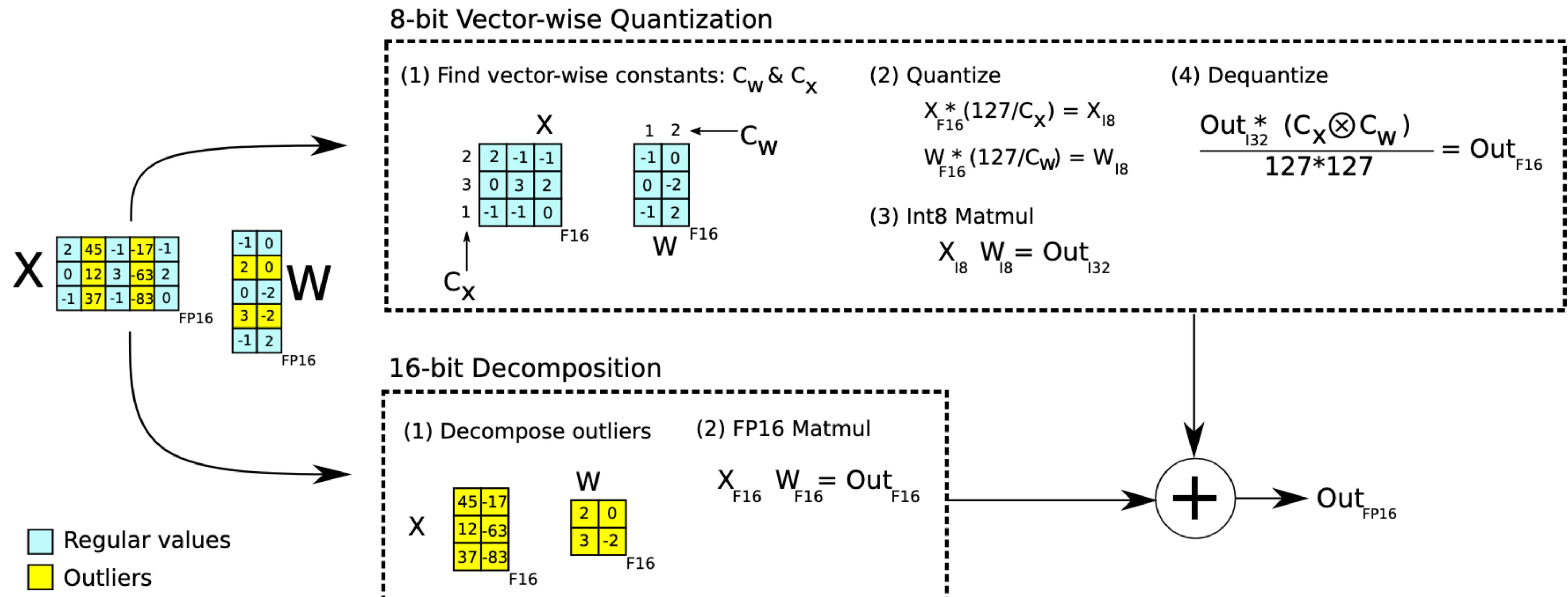
## Mixed-precision Decomposition

$$C_{f16} = S \cdot Q(X_{f16}) Q(W_{f16})$$

- What if we handled the normal features and outlier features separately?

$$C_{f16} = \sum_{h \in O} X_{f16}^h W_{f16}^h + S \cdot \sum_{h \notin O} X_{i8}^h W_{i8}^h$$

# LLM.int8()



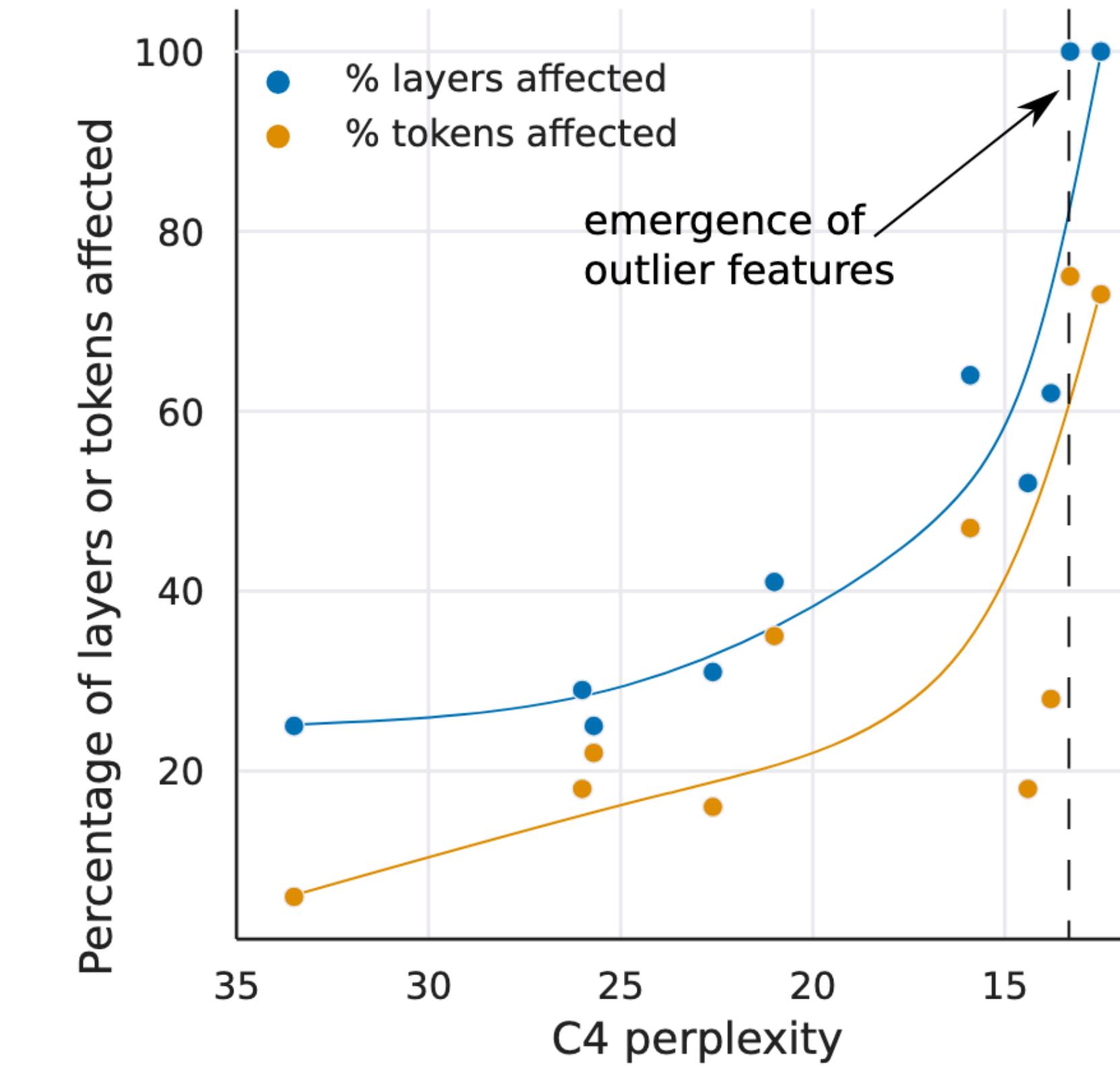
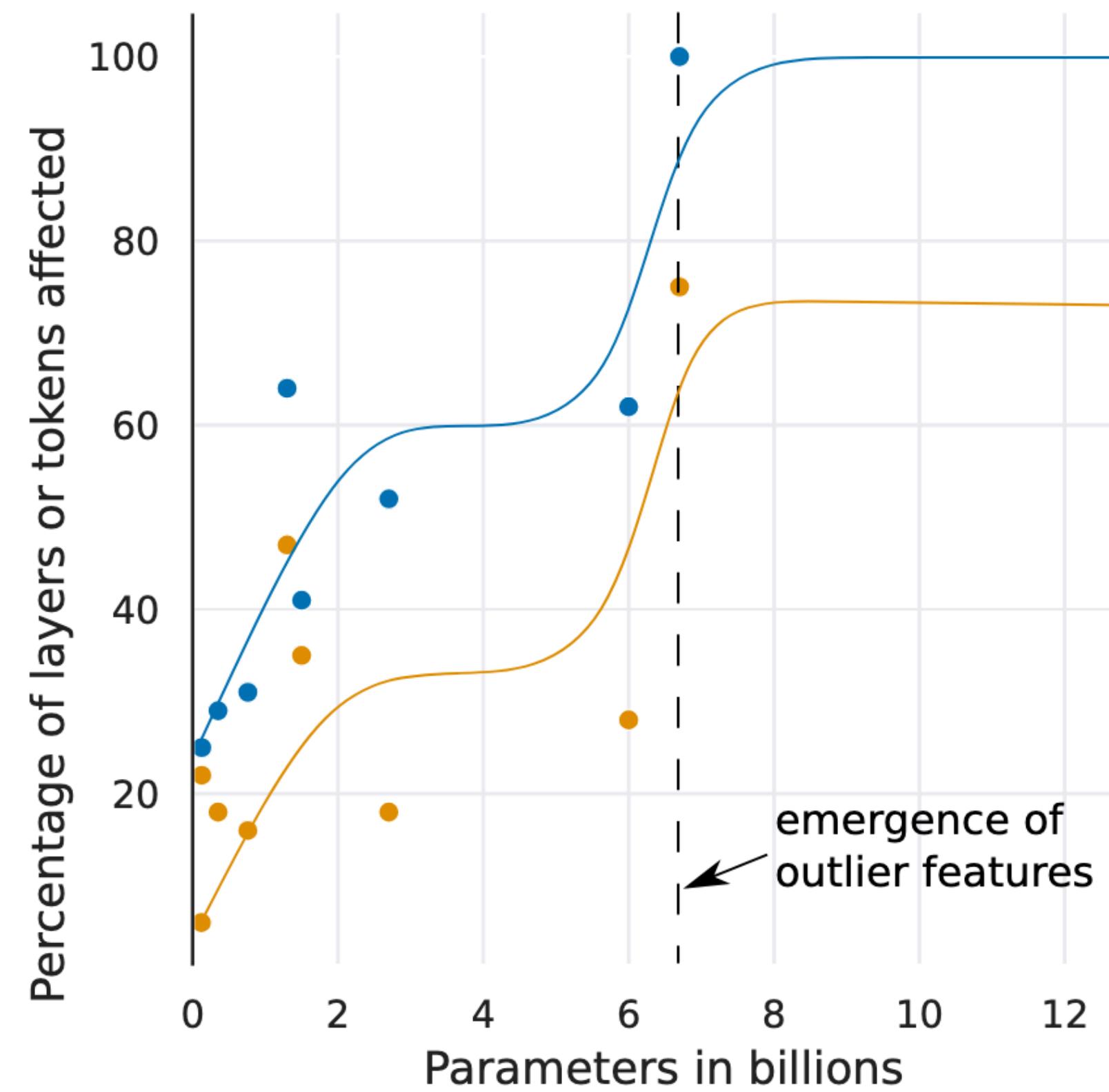
# LLM.int8()

## Outliers are critical to transformer performance

- Removing **<7 outlier features** causes...
  - top-1 softmax probability cut in half
  - validation perplexity increases 6-10x
- Removing **7 random features** causes...
  - top-1 softmax probability decreases by 0.3%
  - validation perplexity increases by 0.1%

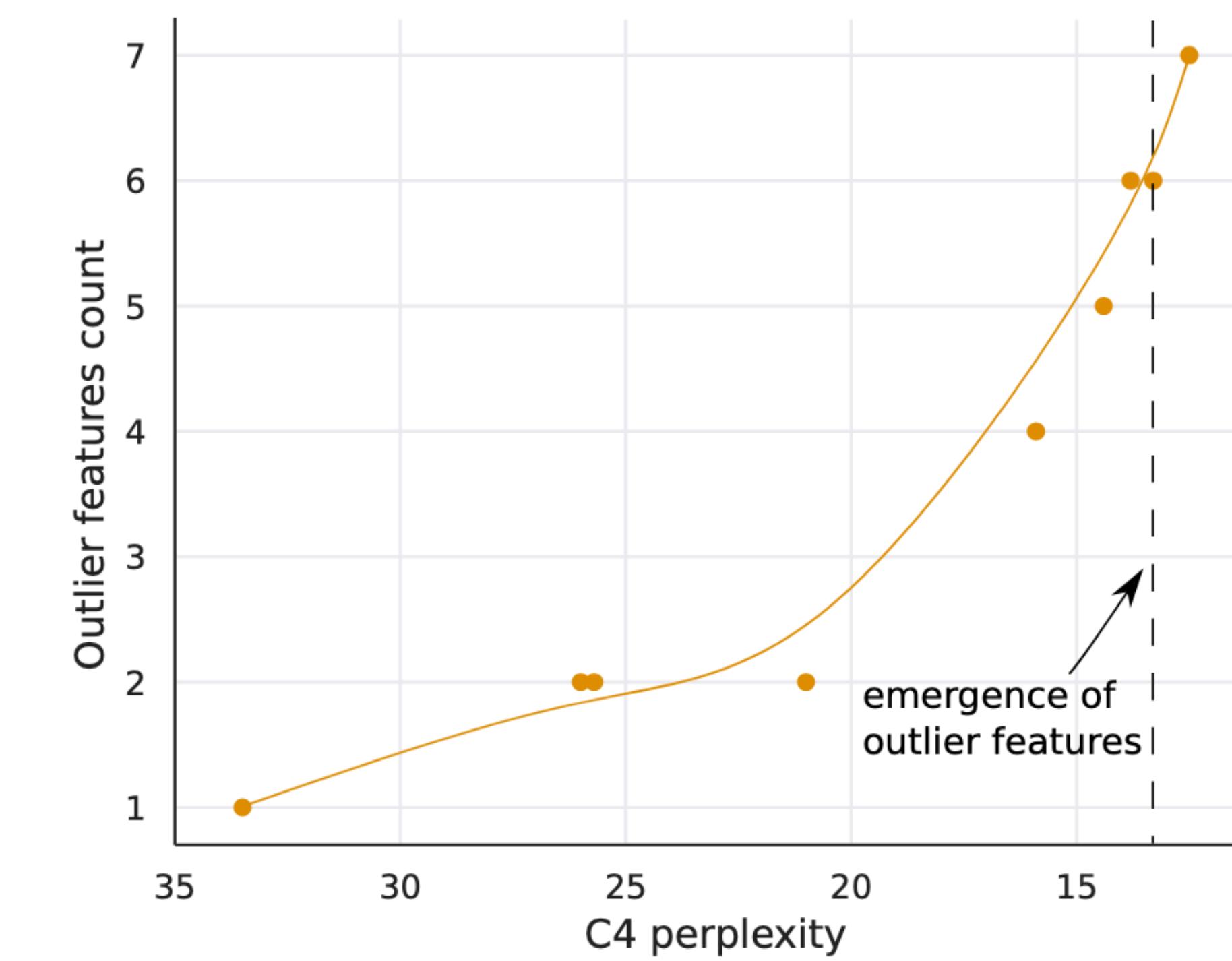
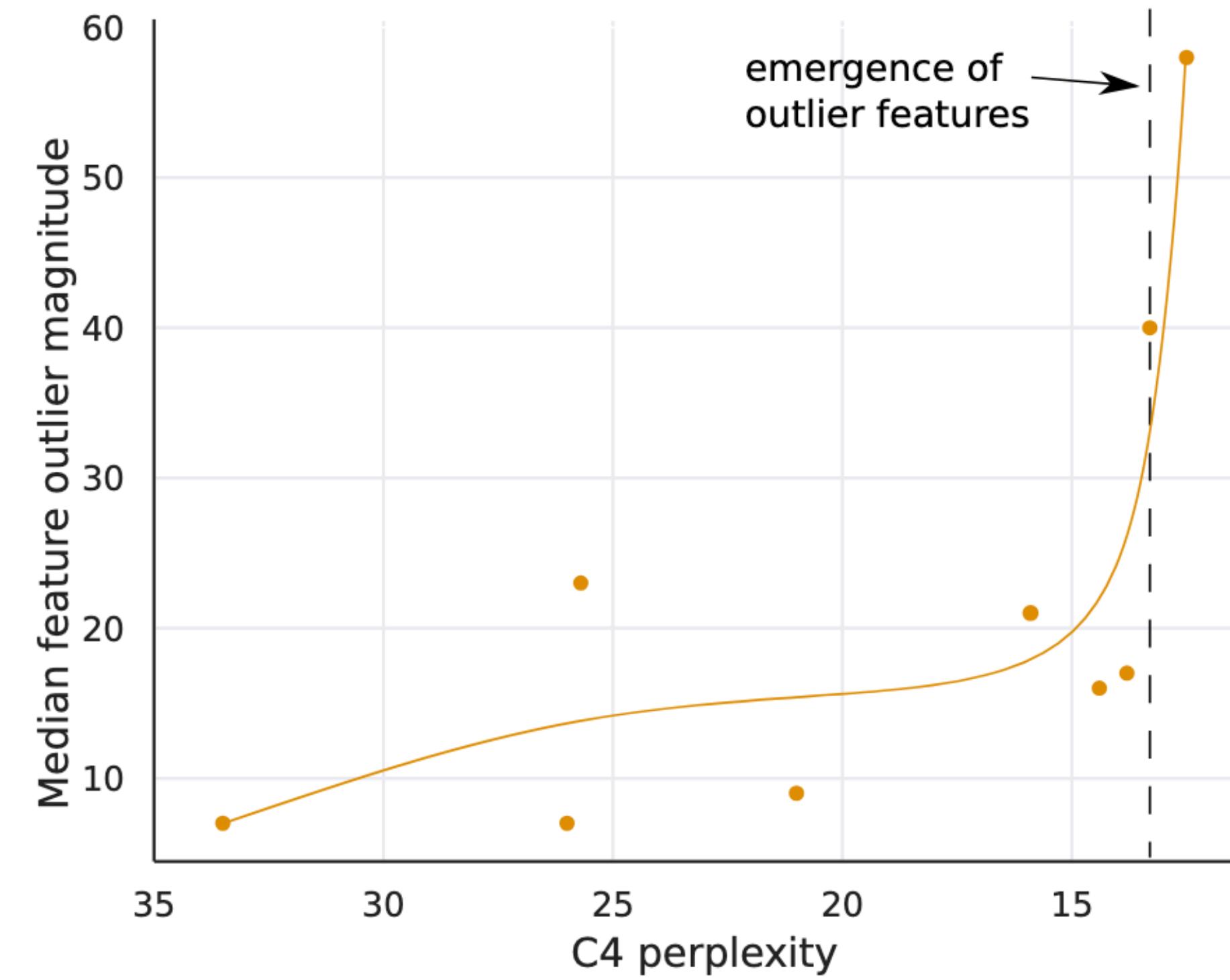
# LLM.int8()

## Phase Shift



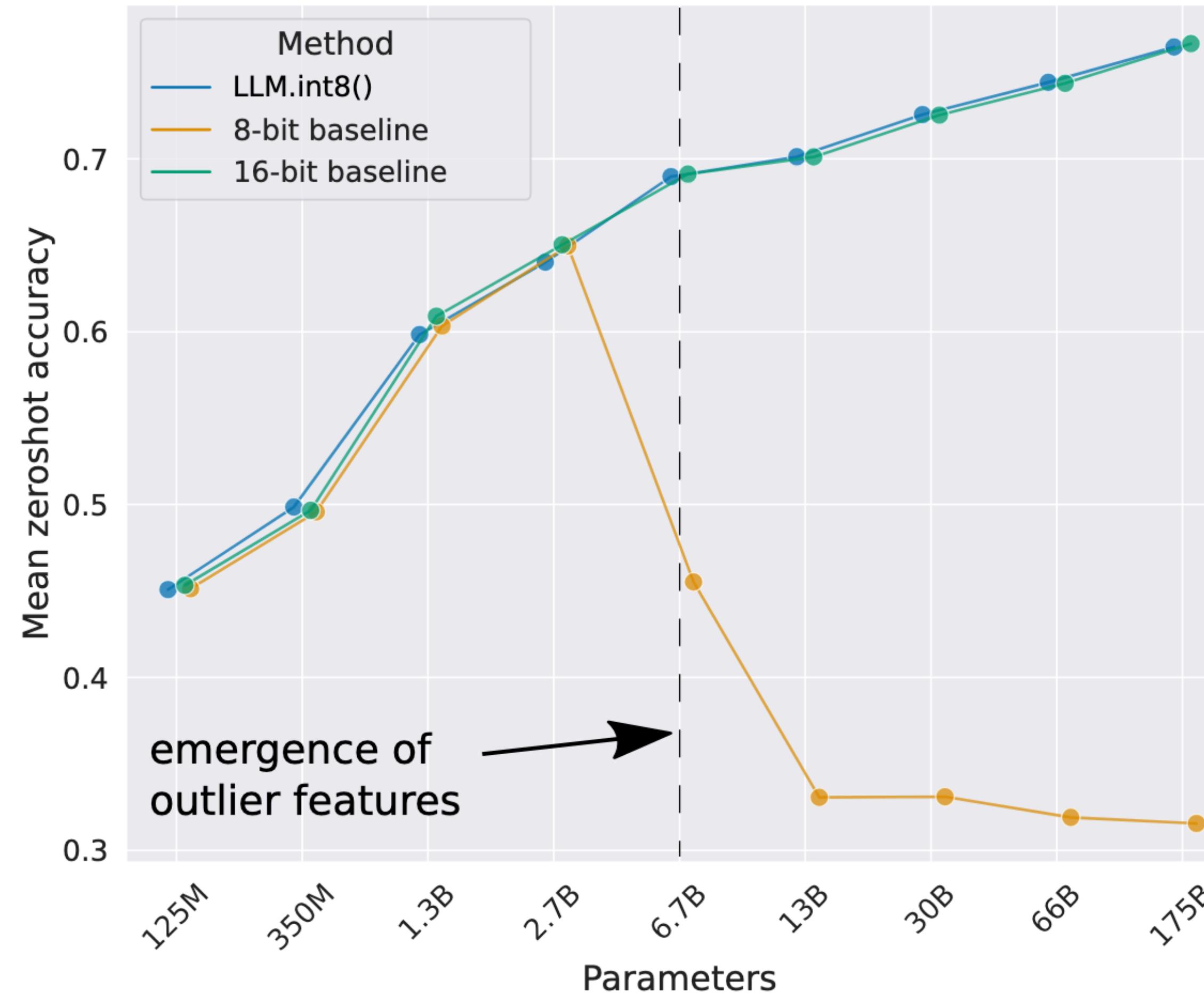
# LLM.int8()

## Phase Shift



# LLM.int8()

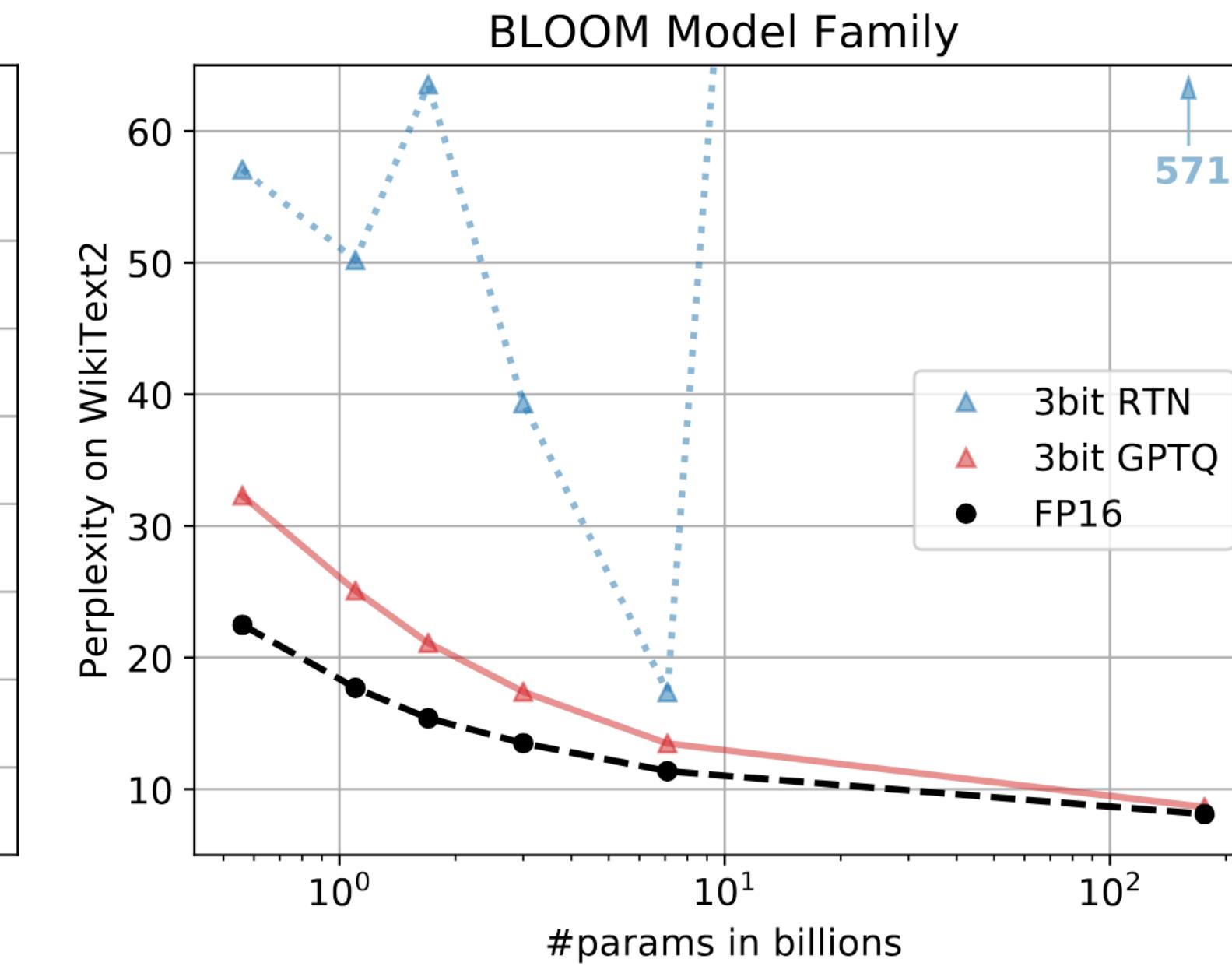
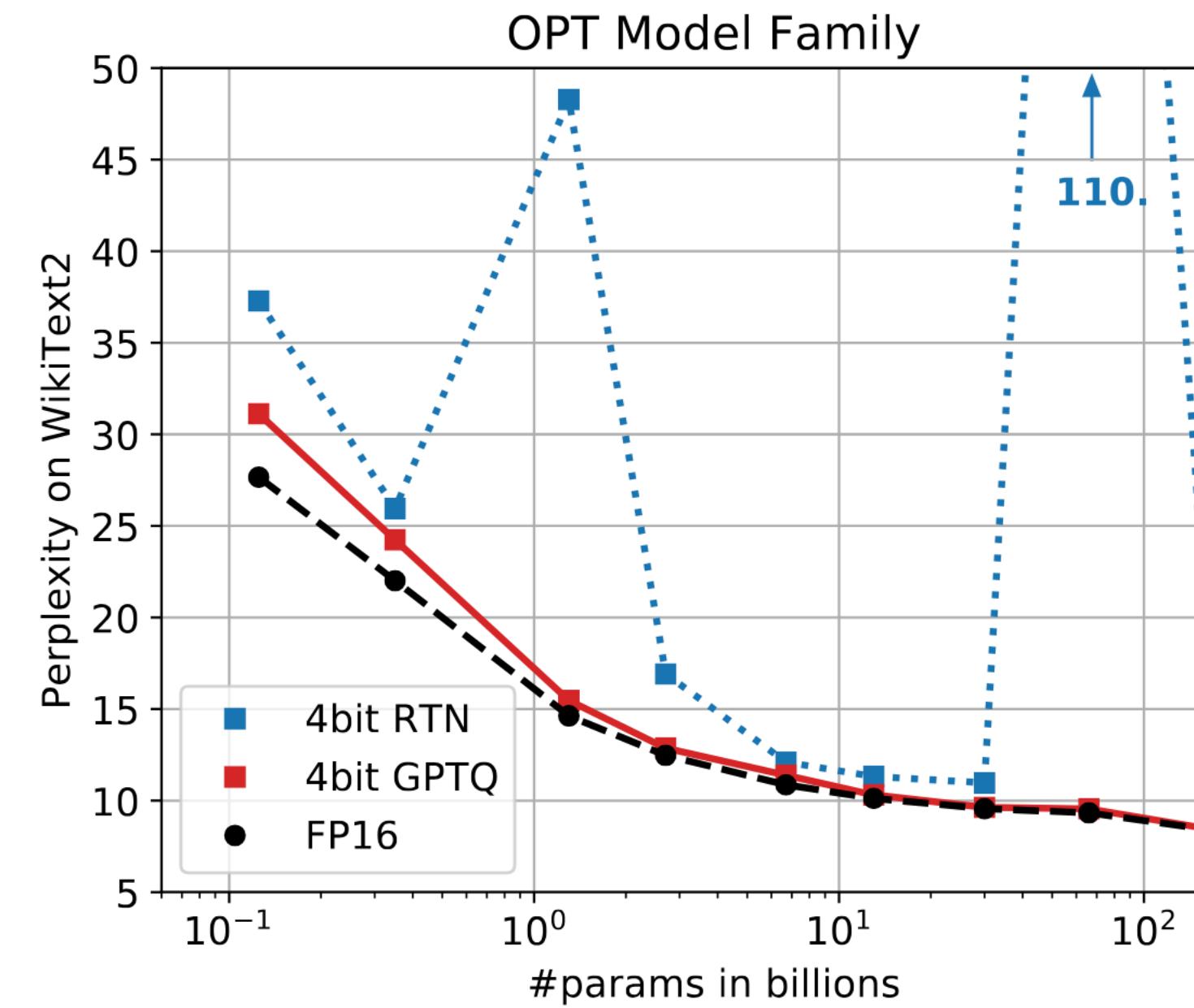
## Phase Shift



# GPTQ

## 4 Bit Quantization

- Integrated with 😊 Transformers
- First paper to break the 8 bit quantization barrier



# AWQ

## 4 Bit Quantization... but simpler :)

- 1% of weights are “salient weights”
  - Similar insight to LLM.int8()
- Can be found via activation distribution instead of weights
  - Intuition: input features with large magnitude are important
- Neat trick: multiply weight by  $s$ , input by  $s^{-1}$

# **Conclusion**

# **Summary**

**Quantization is like printers.**

- Powerful tool to reduce memory footprint of models with many use-cases
  - Frontier labs use quantization to train very large models
  - Academics & home users use quantization to fit large models on their GPUs
- Awe-inspiring results

# References

## Recommended further reading is bolded

- Song Han, MIT 6.5940 TinyML and Efficient Deep Learning Computing, 2024, <https://hanlab.mit.edu/courses/2024-fall-65940>
- Sreenivas et al., “LLM Pruning and Distillation in Practice: The Minitron Approach,” December 10, 2024, <https://arxiv.org/pdf/2408.11796>
- Pankaj Gupta and Philip Kiely, “FP8: Efficient model inference with 8-bit floating point numbers,” <https://www.baseten.co/blog/fp8-efficient-model-inference-with-8-bit-floating-point-numbers/>
- Tim Dettmers, “8-bit Methods for Efficient Deep Learning,” March 2023, <https://www.youtube.com/watch?v=jyOqtw4ry2w>
- **Jacob et al., “Quantization and Training of Neural Networks for Efficient Integer-Arithmetic-Only Inference,” CPVR 2018, <https://arxiv.org/pdf/1712.05877>**
- Raghuraman Krishnamoorthi, “Quantizing deep convolutional networks for efficient inference: A whitepaper,” June 2018, <https://arxiv.org/pdf/1806.08342>
- **Dettmers et al., “LLM.int8(): 8-bit Matrix Multiplication for Transformers at Scale,” NeurIPS 2022, <https://arxiv.org/pdf/2208.07339>**
- Yoshua Bengio, Nicholas Léonard and Aaron Courville, “Estimating or Propagating Gradients Through Stochastic Neurons for Conditional Computation,” August 15, 2013, <https://arxiv.org/pdf/1308.3432>
- Frantar et al., “GPTQ: Accurate Post-Training Quantization for Generative Pre-trained Transformers,” ICLR 2023, <https://arxiv.org/pdf/2210.17323>
- **Lin et al., “AWQ: Activation-aware Weight Quantization for LLM Compression and Acceleration,” MLSys 2024, <https://arxiv.org/pdf/2306.00978>**