# LLM Quantization

Xin Yu   2023.11.01

# GPT-4

## Write python codes

## Write poems

Prompt: Can you write a proof that there are many primes, with every line that rhymes?

GPT-4:

Yes, I think I can,
Though it might take a clever plan.
I'll start by noting Euclid's proof,
Which shows that primes aren't just aloof.
Assume we have a finite list,
Of primes, and that none have been missed.
Multiply them all together,
And add one, just to be clever.
The result will be a number,
That has no prime factors, oh wonder!
But every number has a prime,
That divides it, it's just a matter of time.
So we've found a contradiction,
And our finite list needs eviction.
There must be infinitely many primes,
And that's the end of my rhyming lines.s

## Write LaTeX codes

Prompt: Draw a unicorn in TiKZ.

GPT-4: [Produces LaTeX compiling to following picture.]
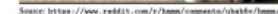


## Explain jokes



## Compose music



## Do math problems

Prompt: Suppose $g(x) = f^{-1}(x)$, $g(0) = 5$, $g(4) = 7$, $g(3) = 2$, $g(7) = 9$, $g(9) = 6$ what is $f(f(f(6)))$?

GPT-4: There are different ways to approach this problem, but one possible method is to use the fact that $g(x) = f^{-1}(x)$ implies that $f(g(x)) = x$ for any $x$ in the domain of $g$. That is, applying $f$ to both sides of the equation $g(x) = f^{-1}(x)$ gives $f(g(x)) = f(f^{-1}(x)) = x$. This means that we can undo the effect of $f$ by applying $g$, and vice versa. Using this idea, we can try to find $f(6)$ by looking for the value of $x$ such that $g(x) = 6$. From the given information, we see that $g(9) = 6$, so $f(6) = f(g(9)) = 9$. Similarly, we can find $f(9)$ by looking for the value of $x$ such that $g(x) = 9$. We see that $g(7) = 9$, so $f(9) = f(g(7)) = 7$. Finally, we can find $f(7)$ by looking for the value of $x$ such that $g(x) = 7$. We see that $g(4) = 7$, so $f(7) = f(g(4)) = 4$. Therefore, $f(f(f(6))) = f(f(9)) = f(7) = 4$.

NLP's Moore's Law: Every year model size increases by 10x

NLP model size and computation are increasing **exponentially**

Ref: https://indiaai.gov.in/article/the-future-of-large-language-models-llms-strategy-opportunities-and-challenges

# 2023-2024 OPTIMAL LANGUAGE MODELS

SEP/2023

**GPT-4**
Estimate 1.76T MoE

**Chinchilla** 70B 20:1

**MPT** 7B 143:1 — 7B

**StableLM** 65B 24:1

**LLaMA-1** 65B 22:1

**MOSS** 20B 27:1

**Cerebras-GPT** 13B 20:1 — 13B

**Inflection-1**

**InternLM** 104B 16:1

**PaLM 2** 340B 11:1

**Llama 2** 70B 29:1

**Claude 2**

**Falcon 180B** 180B 20:1

**Olmo** 70B 20:1

**Next...**

**Google DeepMind Gemini**
**Anthropic Claude-Next**
**OpenAI GPT-5**
**Other stealth projects**
>2T 20:1

Parameters
AI lab/group
Available
Closed

Bubble plot, sizes linear to scale. Selected highlights only. *Chinchilla scale means tokens:parameters ratio ≥11:1. https://lifearchitect.ai/chinchilla/ Alan D. Thompson. 2023. https://lifearchitect.ai/

🔗 LifeArchitect.ai/models

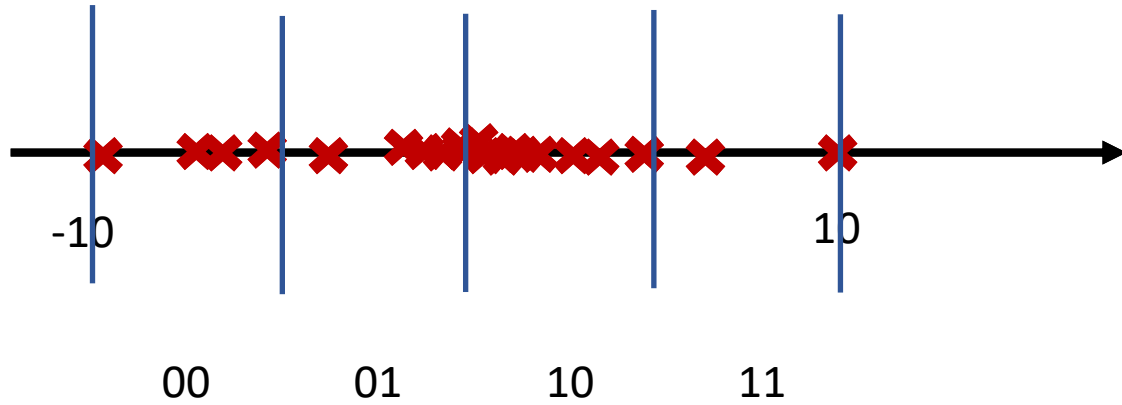Ref: https://lifearchitect.ai/models/

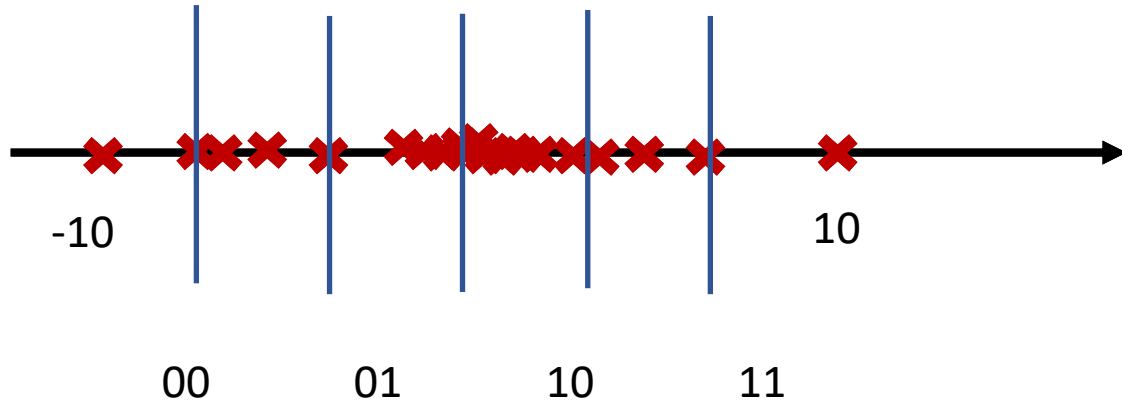# Quantization: Data with less bits

# Quantization: Data with less bits

RTN: round to nearest
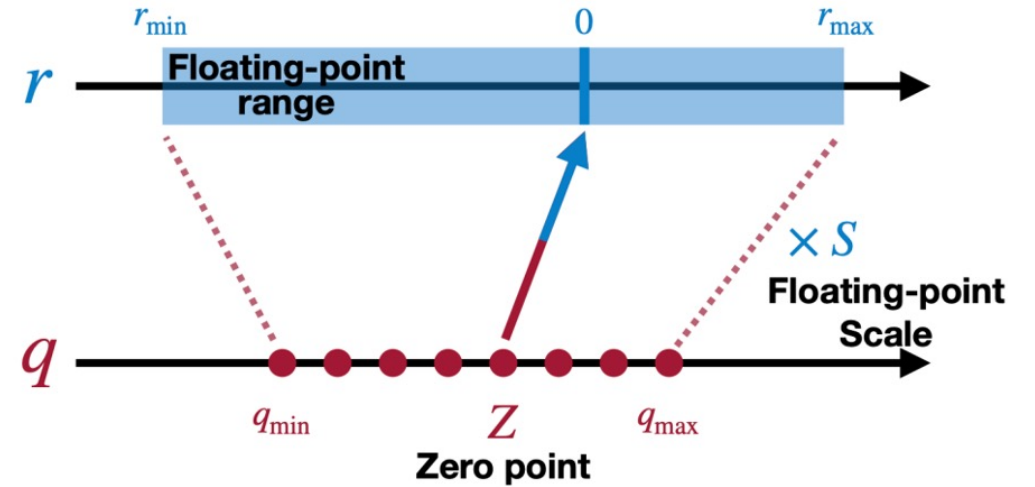
# Quantization: Data with less bits

RTN: round to nearest
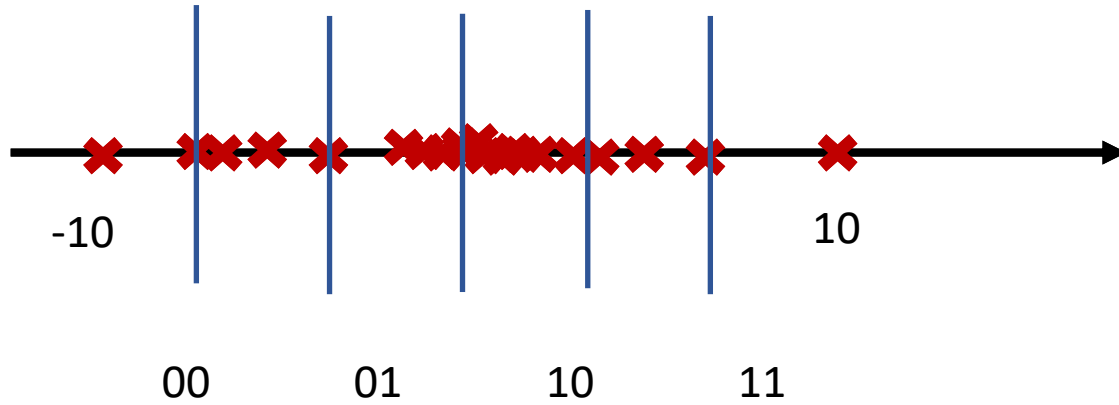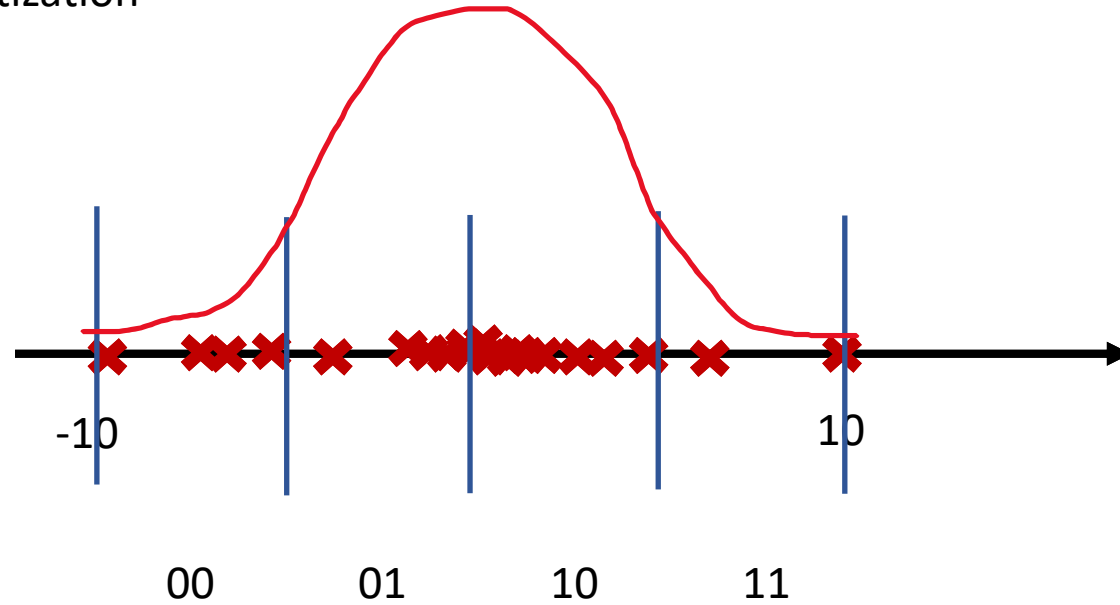


-10

10

00    01    10    11

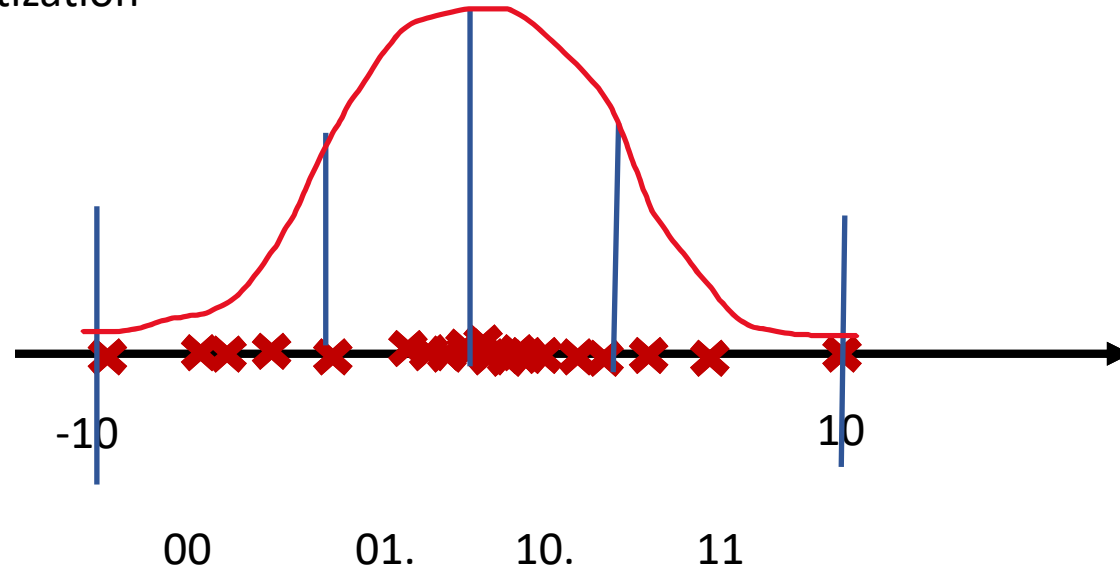# Quantization: Data with less bits

RTN: round to nearest

# Quantization: Data with less bits

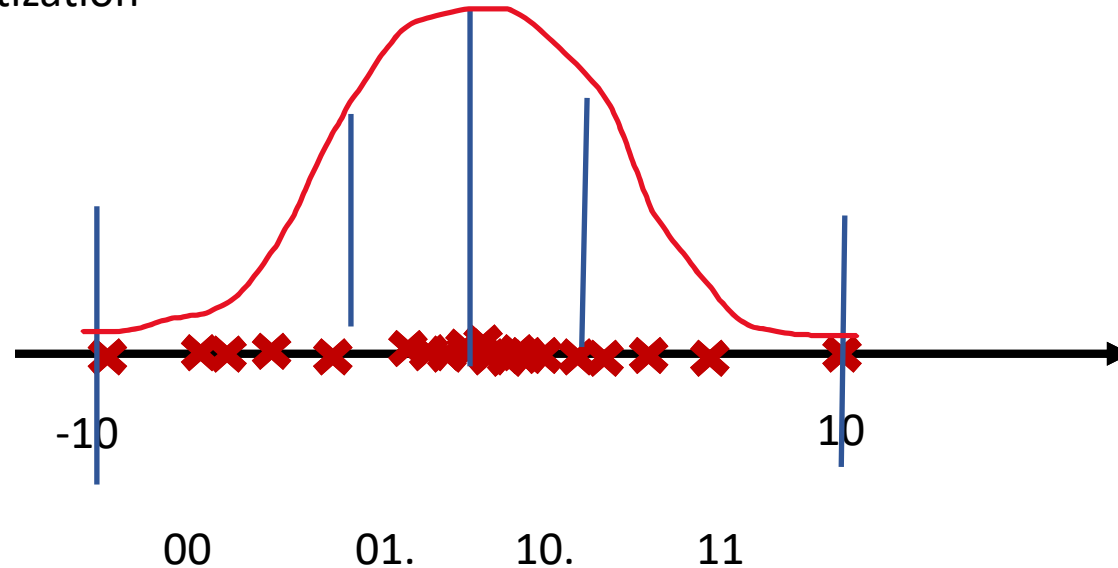Normal Float Quantization

# Quantization: Data with less bits

Normal Float Quantization

Quantile Quantization



-10

10

00          01.          10.          11

# Quantization: Data with less bits



Normal Float Quantization

Quantile Quantization

-10

10

00    01.    10.    11
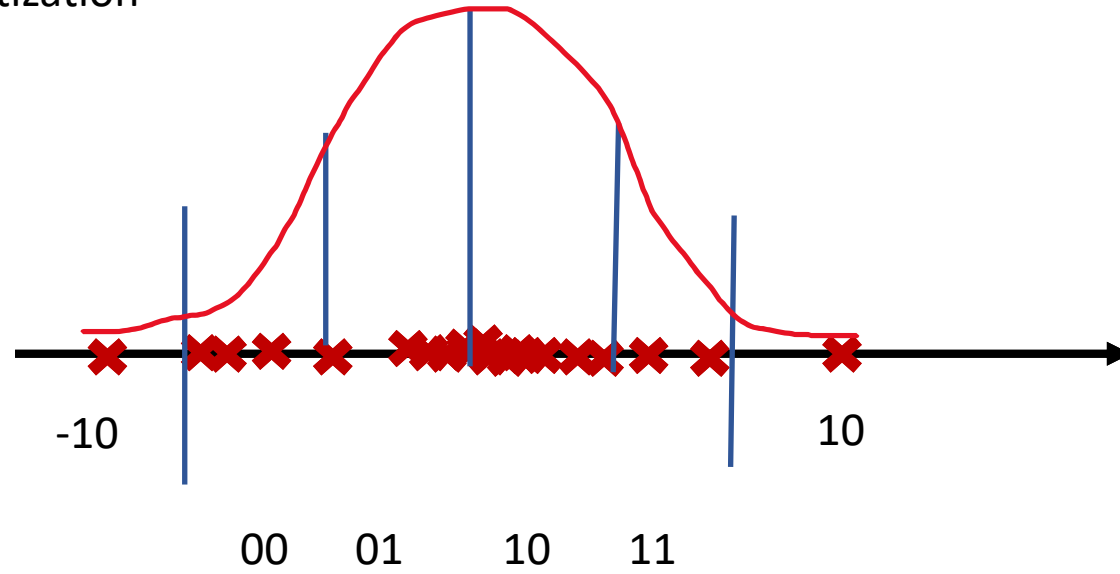
The exact values of the NF4 data type are as follows:

[-1.0, -0.6961928009986877, -0.5250730514526367, -0.39491748809814453, -0.28444138169288635, -0.18477343022823334, -0.09105003625154495, 0.0, 0.07958029955625534, 0.16093020141124725, 0.24611230194568634, 0.33791524171829224, 0.44070982933044434, 0.5626170039176941, 0.7229568362236023, 1.0]

# Quantization: Data with less bits
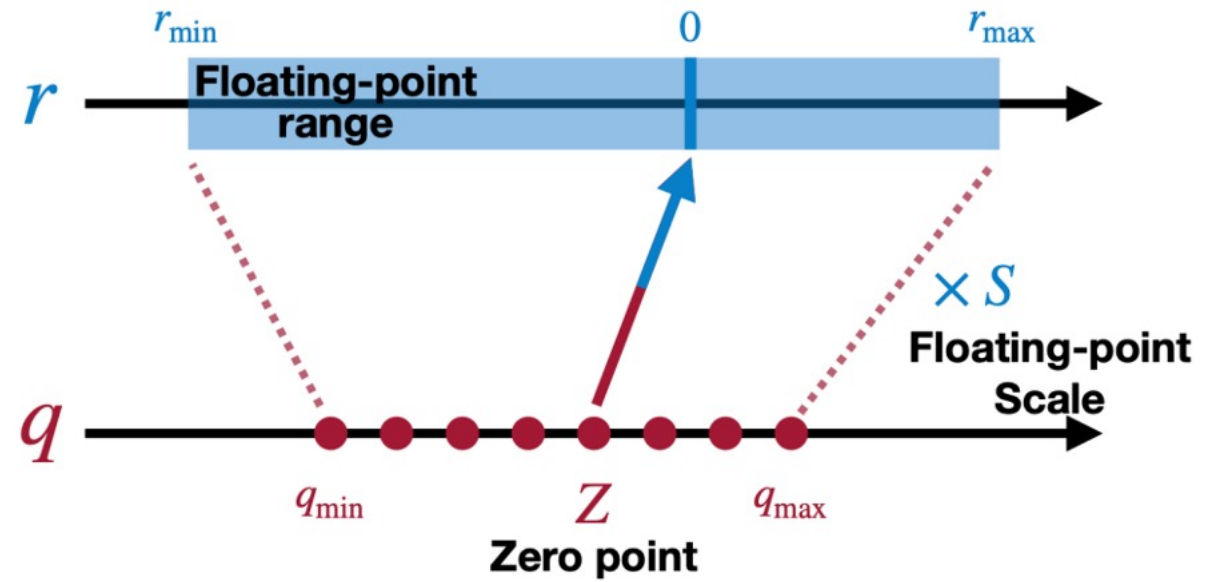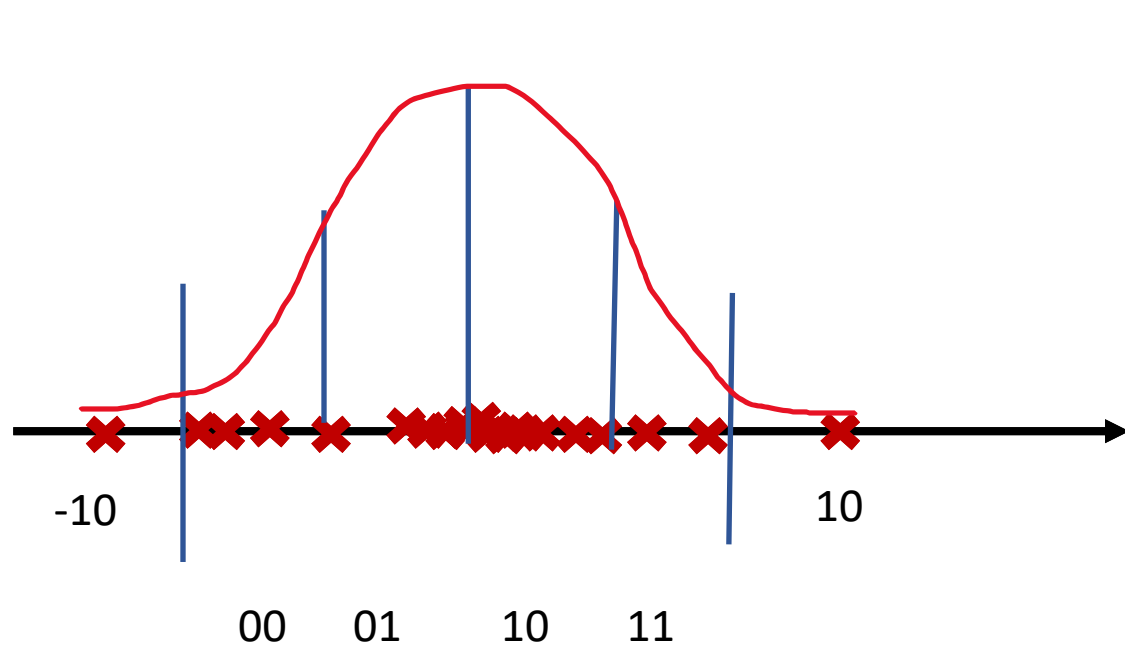
Normal Float Quantization



-10    00    01    10    11    10

Truncate the outliers!

# Quantization: Data with less bits

# Quantization: Data with less bits
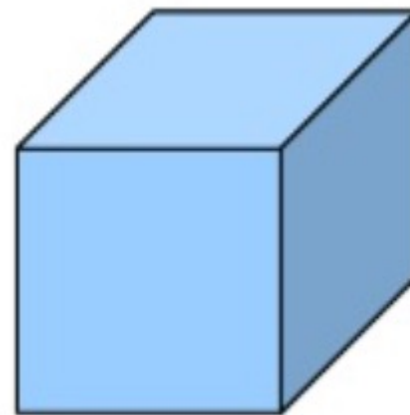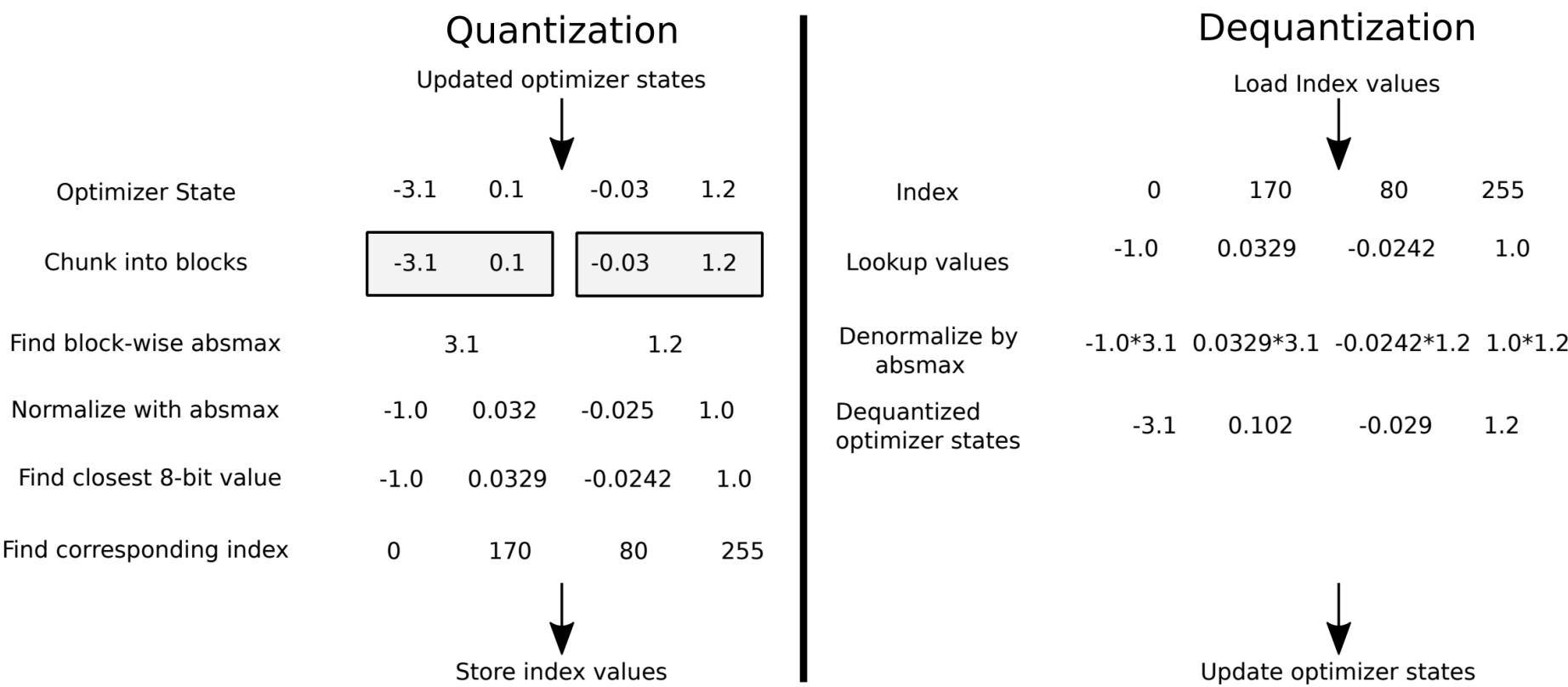
Block-wise dynamic quantization



1d-tensor          2d-tensor          3d-tensor

# Quantization: Data with less bits
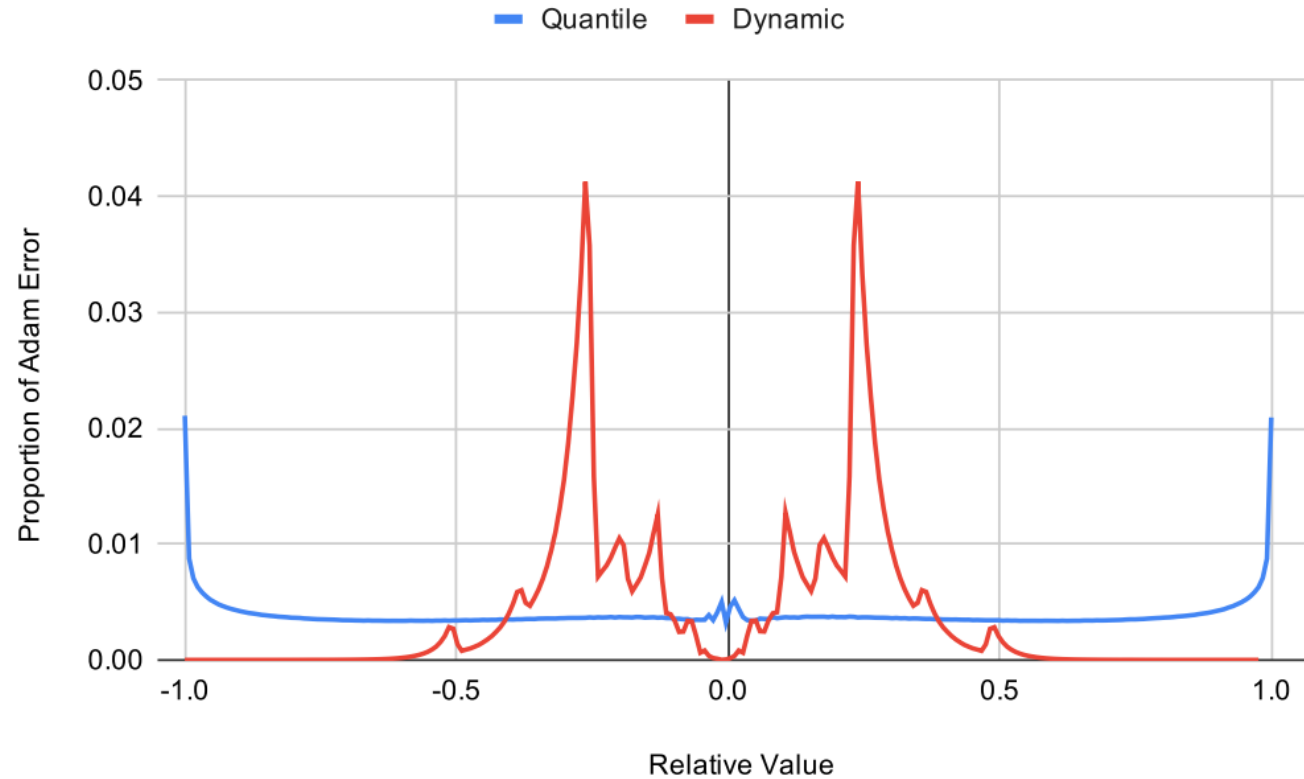
Block-wise dynamic quantization



| | Quantization | | | |
|---|---|---|---|---|
| | **Updated optimizer states** ↓ | | | |
| Optimizer State | -3.1 | 0.1 | -0.03 | 1.2 |
| Chunk into blocks | -3.1 | 0.1 | -0.03 | 1.2 |
| Find block-wise absmax | 3.1 | | 1.2 | |
| Normalize with absmax | -1.0 | 0.032 | -0.025 | 1.0 |
| Find closest 8-bit value | -1.0 | 0.0329 | -0.0242 | 1.0 |
| Find corresponding index | 0 | 170 | 80 | 255 |
| | ↓ Store index values | | | |

| | Dequantization | | | |
|---|---|---|---|---|
| | **Load Index values** ↓ | | | |
| Index | 0 | 170 | 80 | 255 |
| Lookup values | -1.0 | 0.0329 | -0.0242 | 1.0 |
| Denormalize by absmax | -1.0*3.1 | 0.0329*3.1 | -0.0242*1.2 | 1.0*1.2 |
| Dequantized optimizer states | -3.1 | 0.102 | -0.029 | 1.2 |
| | ↓ Update optimizer states | | | |

Ref: 8-BIT OPTIMIZERS VIA BLOCK-WISE QUANTIZATION

# Quantization: Data with less bits

Block-wise dynamic quantization



Ref: 8-BIT OPTIMIZERS VIA BLOCK-WISE QUANTIZATION

# Minimize quantization error of each layer!
Reduce the accumulated error cross the network

$$\text{argmin}_{\widehat{\mathbf{W}}} \ ||\mathbf{W}\mathbf{X} - \widehat{\mathbf{W}}\mathbf{X}||_2^2.$$

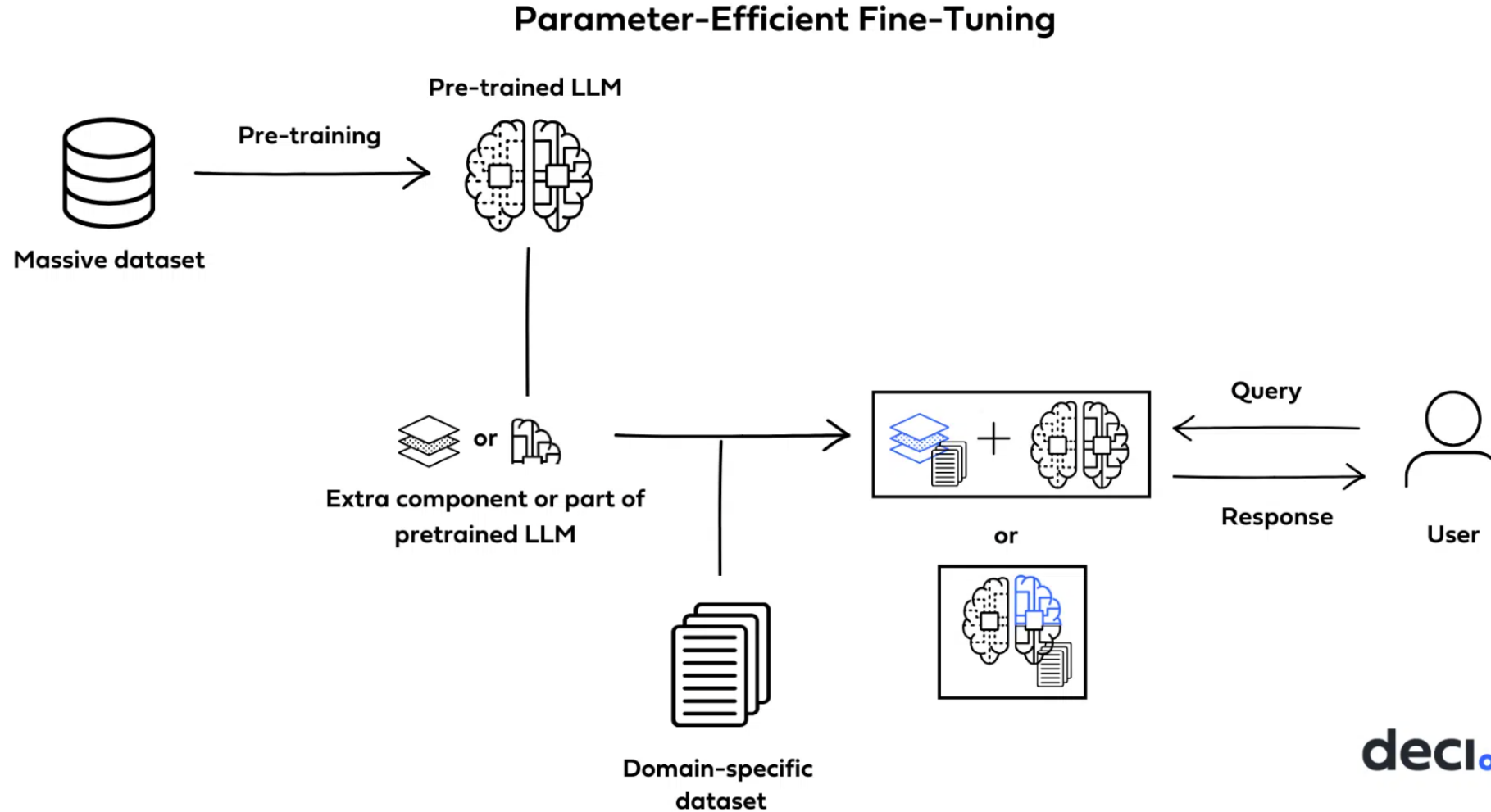# *Method-1* **QLORA: Efficient Finetuning of Quantized LLMs**

Full fine-tuning vs. Parameter-Efficient finetuning
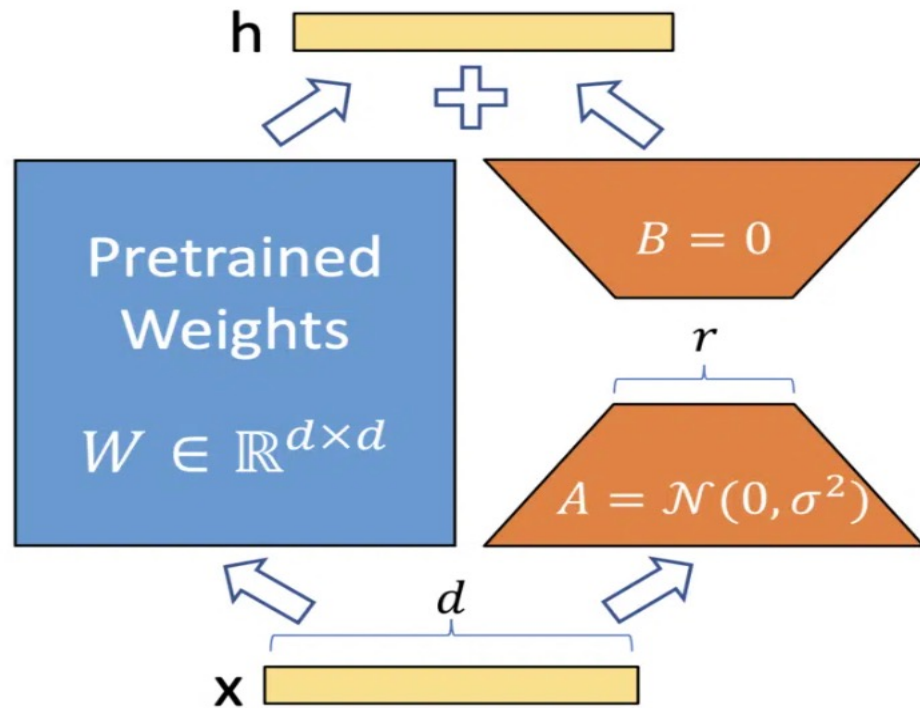


## Full Fine Tuning

# *Method-1* QLORA: Efficient Finetuning of Quantized LLMs

Full fine-tuning vs. Parameter-Efficient finetuning

# *Method-1* QLORA: Efficient Finetuning of Quantized LLMs

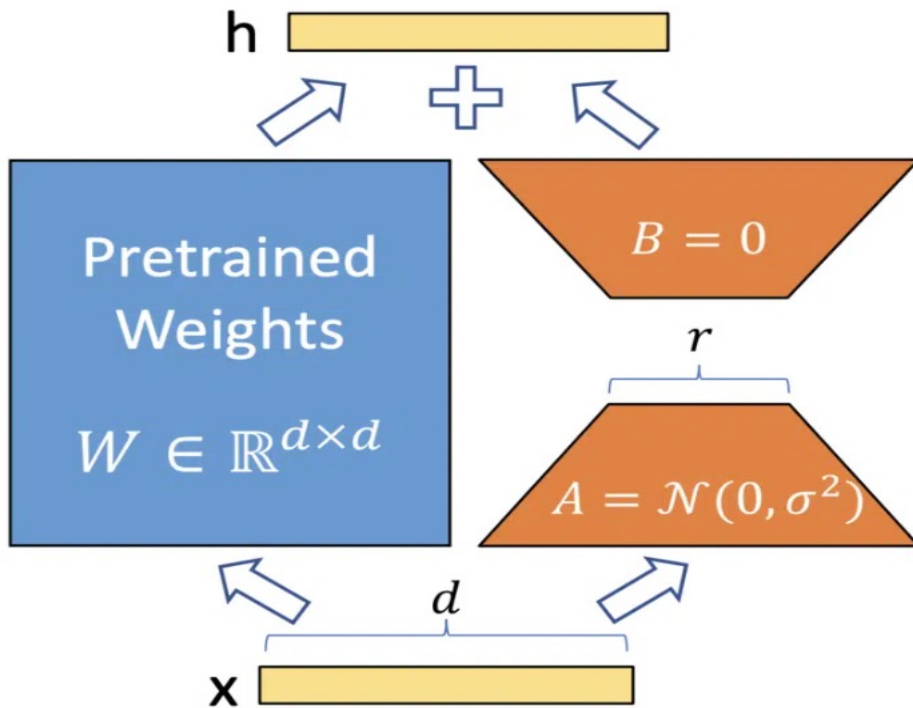LoRA: was introduced in early 2023, the most commonly used PEFT method



Figure 1: LoRA reparameterization trains only A and B. (Credit: https://arxiv.org/abs/2106.09685)

$$W' = W + A \text{ @ } B$$

# *Method-1* QLORA: Efficient Finetuning of Quantized LLMs

QLoRA: the most commonly used PEFT Quantization



Figure 1: LoRA reparameterization trains only A and B. (Credit: https://arxiv.org/abs/2106.09685)

$$W' = quant(W) + A @ B$$

$$W' = double\_quantize(W) + \mathbf{A \text{ @ } B}$$

$$W' = double\_normal\_float(W) + \mathbf{A @ B}$$

**QLoRA.** Using the components described above, we define QLoRA for a single linear layer in the quantized base model with a single LoRA adapter as follows:
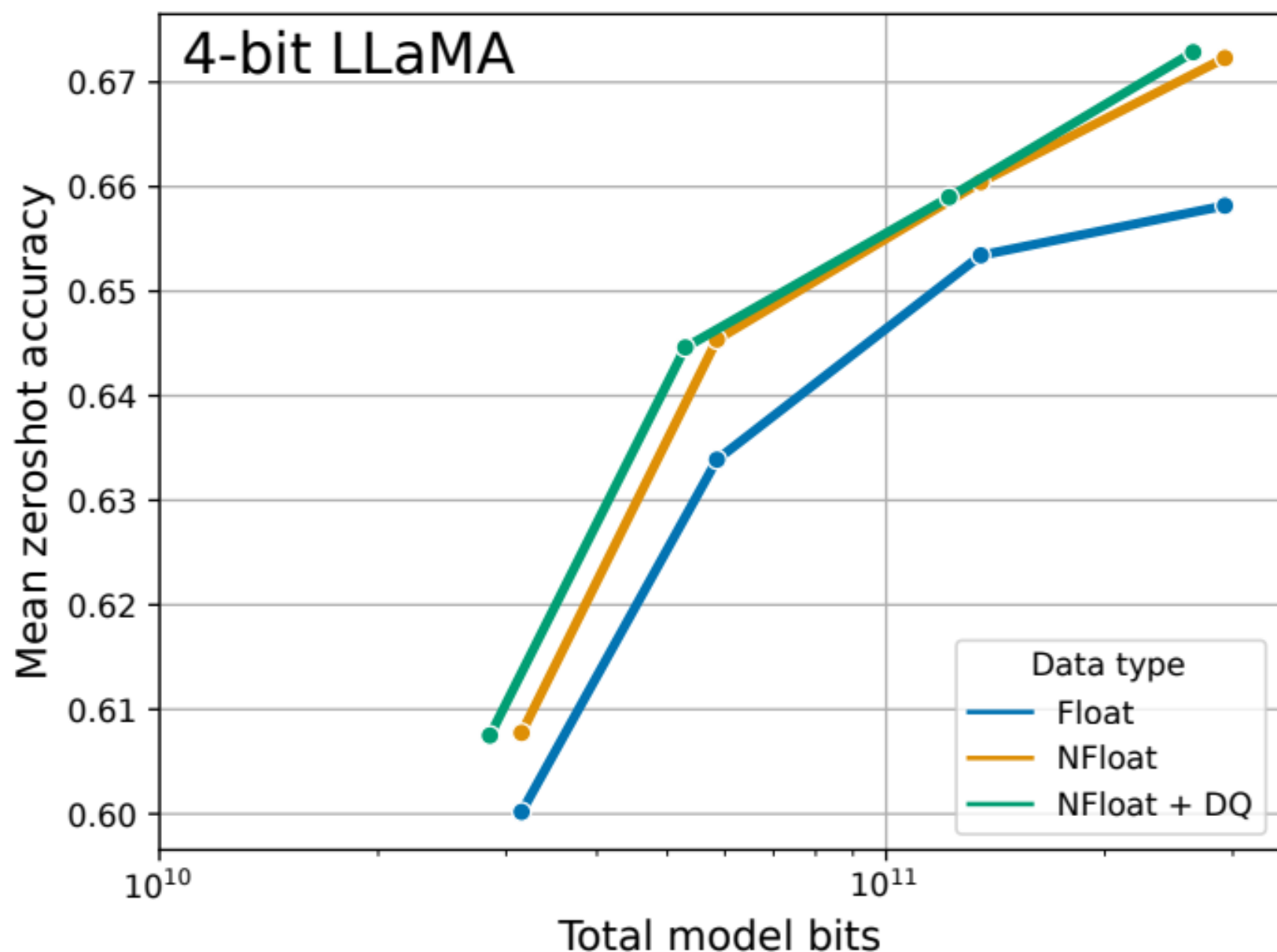
$$\mathbf{Y}^{\text{BF16}} = \mathbf{X}^{\text{BF16}}\text{doubleDequant}(c_1^{\text{FP32}}, c_2^{\text{k-bit}}, \mathbf{W}^{\text{NF4}}) + \mathbf{X}^{\text{BF16}}\mathbf{L}_1^{\text{BF16}}\mathbf{L}_2^{\text{BF16}}, \quad (5)$$

where doubleDequant($\cdot$) is defined as:

$$\text{doubleDequant}(c_1^{\text{FP32}}, c_2^{\text{k-bit}}, \mathbf{W}^{\text{k-bit}}) = \text{dequant}(\text{dequant}(c_1^{\text{FP32}}, c_2^{\text{k-bit}}), \mathbf{W}^{\text{4bit}}) = \mathbf{W}^{\text{BF16}}, \quad (6)$$

We use NF4 for $\mathbf{W}$ and FP8 for $c_2$. We use a blocksize of 64 for $\mathbf{W}$ for higher quantization precision and a blocksize of 256 for $c_2$ to conserve memory.
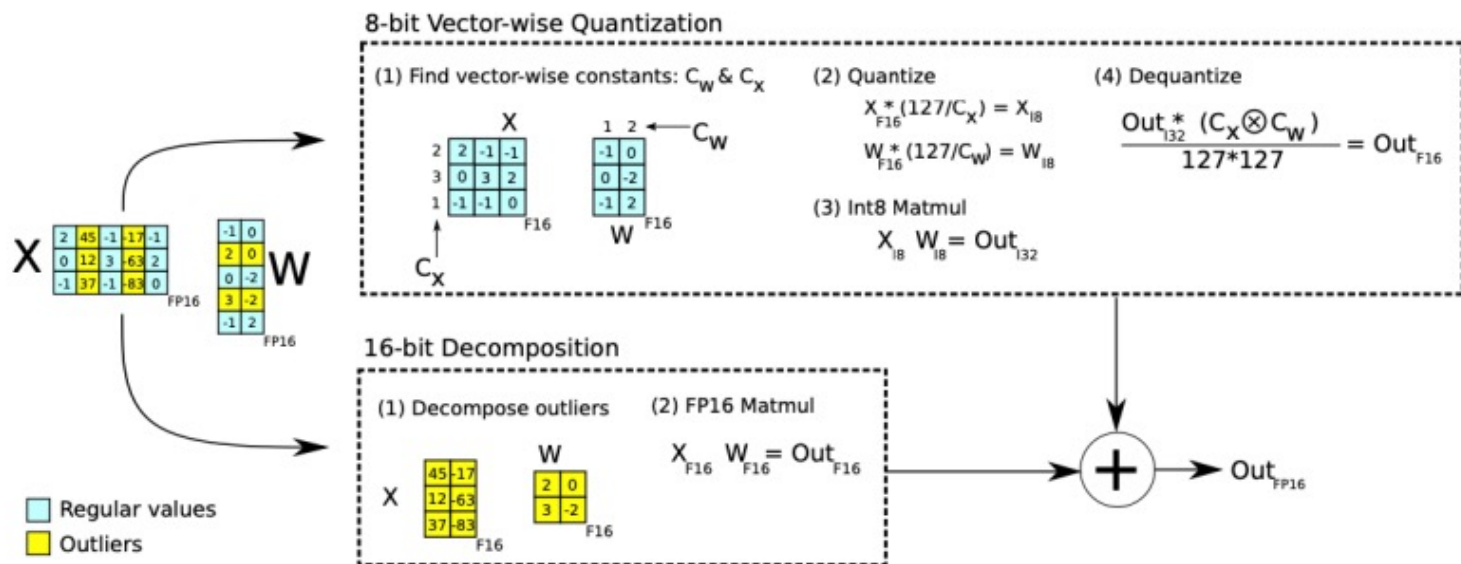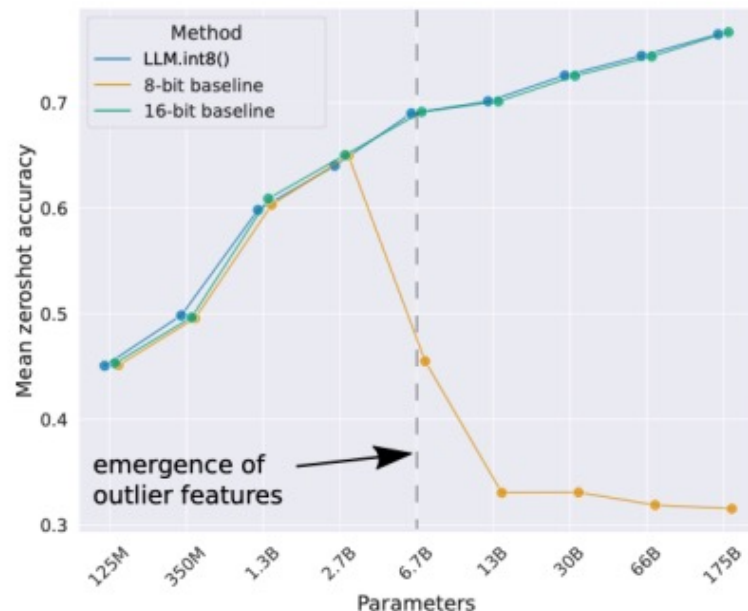
**Table 3:** Experiments comparing 16-bit BrainFloat (BF16), 8-bit Integer (Int8), 4-bit Float (FP4), and 4-bit NormalFloat (NF4) on GLUE and Super-NaturalInstructions. QLoRA replicates 16-bit LoRA and full-finetuning.

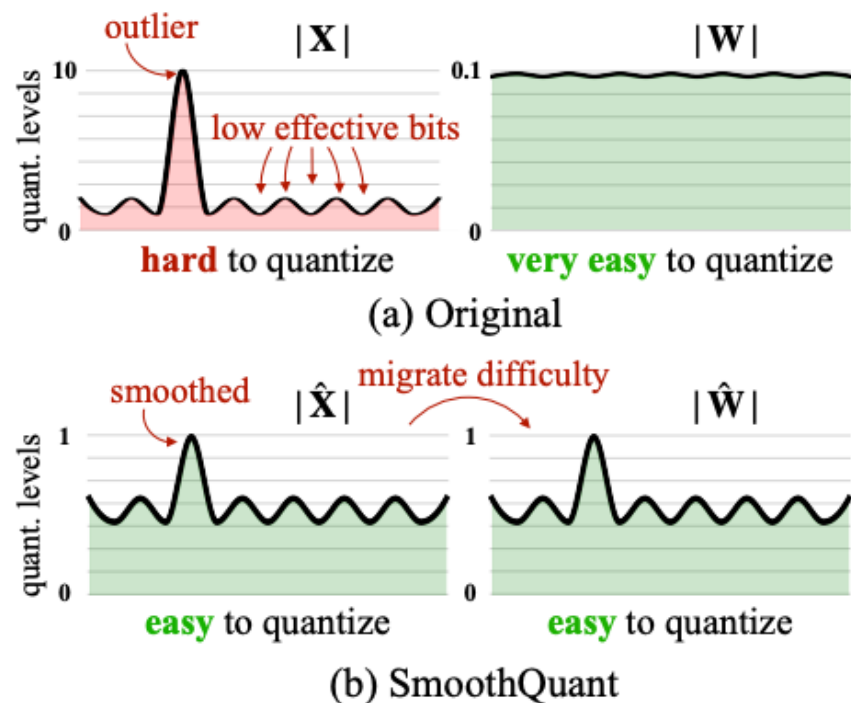| Dataset | GLUE (Acc.) | Super-NaturalInstructions (RougeL) | | | | |
|---|---|---|---|---|---|---|
| Model | RoBERTa-large | T5-80M | T5-250M | T5-780M | T5-3B | T5-11B |
| BF16 | 88.6 | 40.1 | 42.1 | 48.0 | 54.3 | 62.0 |
| BF16 replication | 88.6 | 40.0 | 42.2 | 47.3 | 54.9 | - |
| LoRA BF16 | 88.8 | 40.5 | 42.6 | 47.1 | 55.4 | 60.7 |
| QLoRA Int8 | 88.8 | 40.4 | 42.9 | 45.4 | 56.5 | 60.7 |
| QLoRA FP4 | 88.6 | 40.3 | 42.4 | 47.5 | 55.6 | 60.9 |
| QLoRA NF4 + DQ | - | 40.4 | 42.7 | 47.7 | 55.3 | 60.9 |

# Quantization — LLM.int8()

## Mixed-Precision Decomposition



- **Motivation**: Transformers have outlier features that have **large values** (especially large models).
  - They occur in particular hidden dimensions, leading to large quantization error.
- **Key idea**: Separate outlier features into a **separate FP16 MM**, quantize the other values to Int8.
  - Outlier: At least one feature dimension with a magnitude larger than the threshold (6).
  - Token-wise scale factor (for X) and (output) channel-wise scale factor (for W).

LLM.int8(): 8-bit Matrix Multiplication for Transformers at Scale [Dettmers *et al.*, 2022]

# SmoothQuant: Accurate and Efficient Post-Training Quantization for LLMs

| | LLM (100B+) Accuracy | Hardware Efficiency |
|---|---|---|
| ZeroQuant | ✗ | ✔ |
| Outlier Suppression | ✗ | ✔ |
| LLM.int8() | ✔ | ✗ |
| **SmoothQuant** | ✔ | ✔ |



(a) Original

(b) SmoothQuant

- We propose SmoothQuant, an **accurate** and **efficient** post-training-quantization (PTQ) method to enable 8-bit weight, 8-bit activation (**W8A8**) quantization for LLMs.

- Since **weights are easy** to quantize while **activations are not**, SmoothQuant smooths the activation outliers by **migrating the quantization difficulty from activations to weights** with a mathematically equivalent transformation.