

# CSCI 5541: Natural Language Processing

## Lecture 15: LLM Compute efficiency and engineering

James Mooney

With slides borrowed from Song Han (MIT)

# What Is Efficiency and Why Does It Matter?

- ❑ Efficiency for NLP is concerned with delivering faster, cheaper, smaller, less energy intensive solutions to problems involving natural language
- ❑ Faster models means LLM model services (GPT3.5, Claude 2.0, etc.) can meet the demands of many clients more quickly
- ❑ Cheaper models reduce costs for LLM model service providers
- ❑ Smaller model sizes allow for service providers to use fewer resources and can allow for individuals to deploy LLMs to their own (smaller) devices
- ❑ Less energy intensive means lower cost and easier to deploy at the edge, where energy is harder to come by

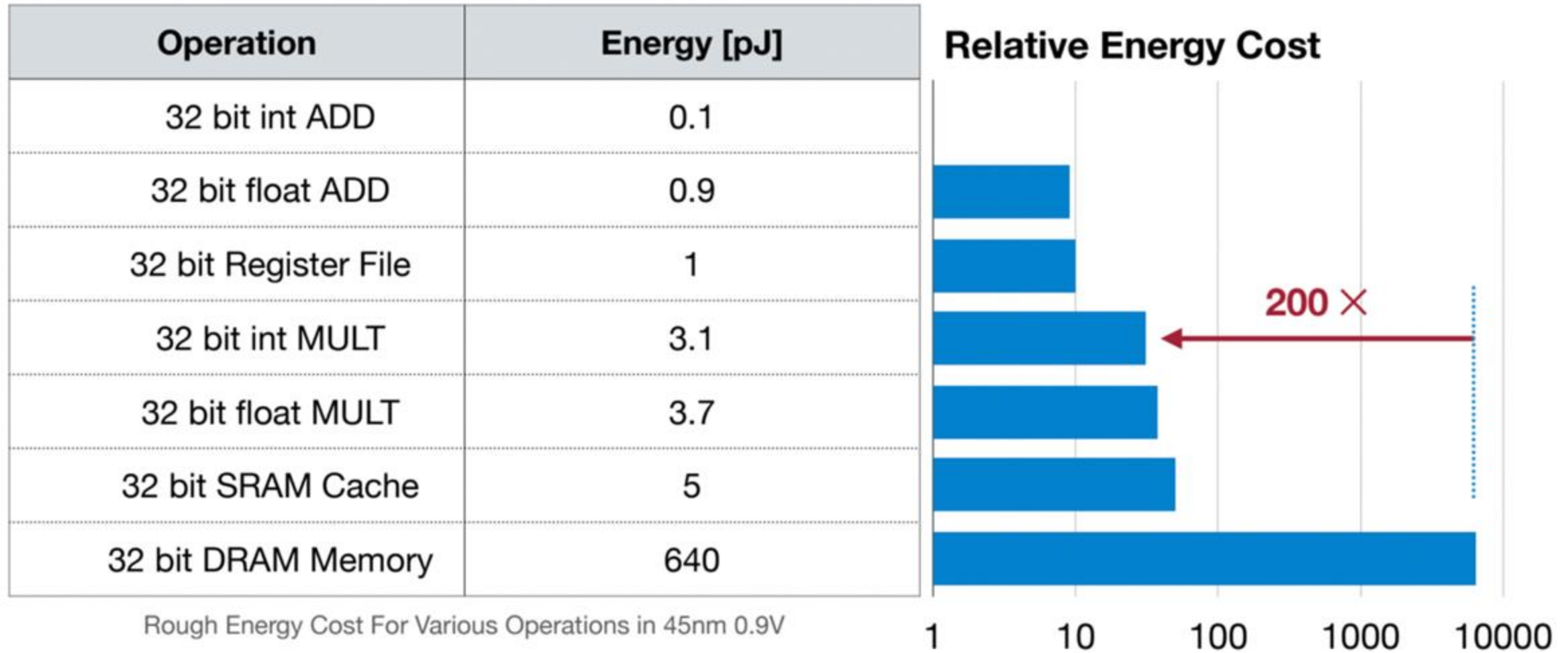


# What Is Efficiency and Why Does It Matter?

- ❑ Efficiency for NLP is concerned with delivering **faster, cheaper, smaller, less energy** intensive solutions to problems involving natural language
- ❑ **Faster** models means LLM model services (GPT3.5, Claude 2.0, etc.) can meet the demands of many clients more quickly
- ❑ **Cheaper** models reduce costs for LLM model service providers
- ❑ **Smaller** model sizes allow for service providers to use fewer resources and can allow for individuals to deploy LLMs to their own (smaller) devices
- ❑ **Less energy** intensive means lower cost and easier to deploy at the edge, where energy is harder to come by



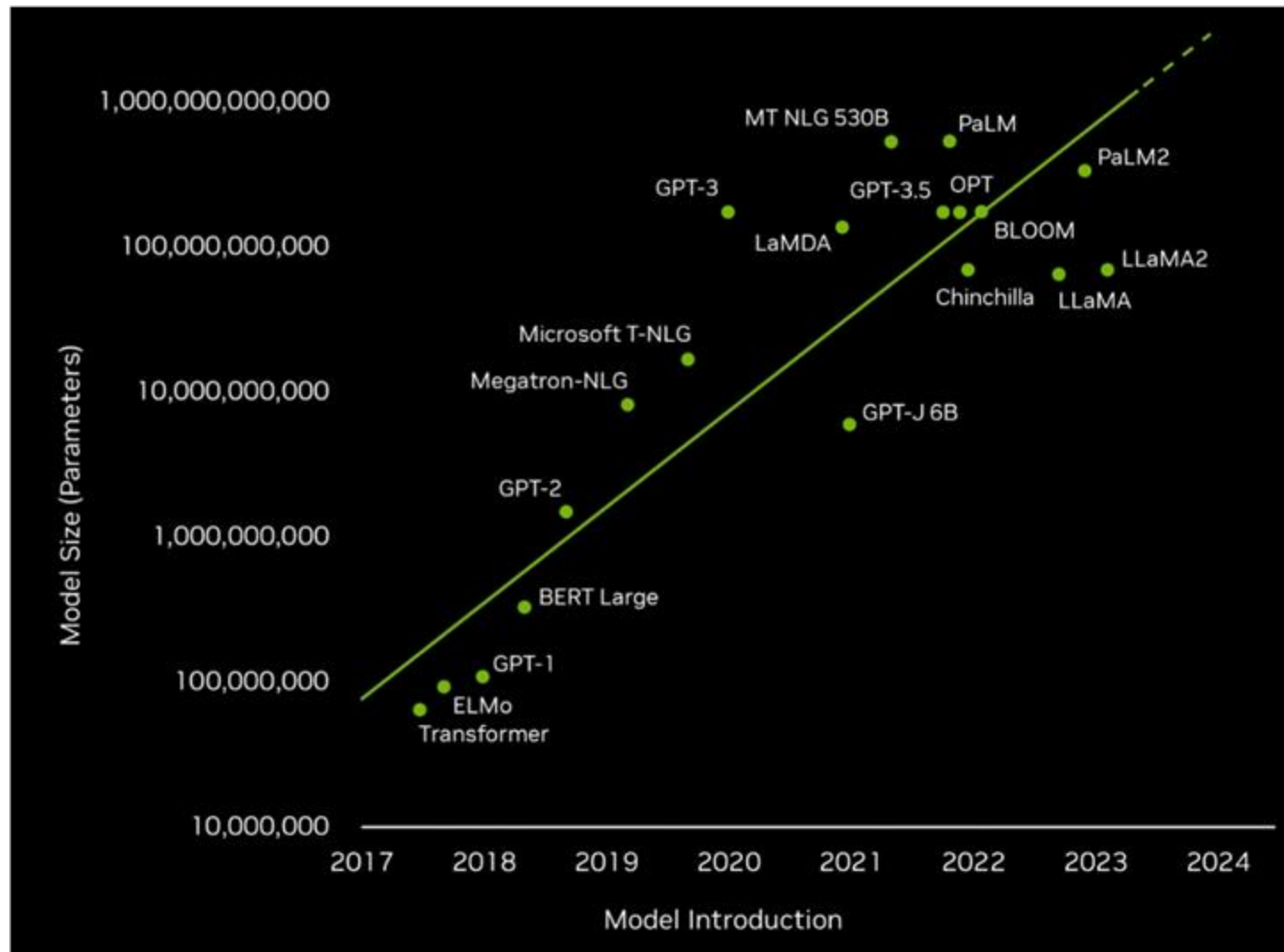
# Model Energy Use



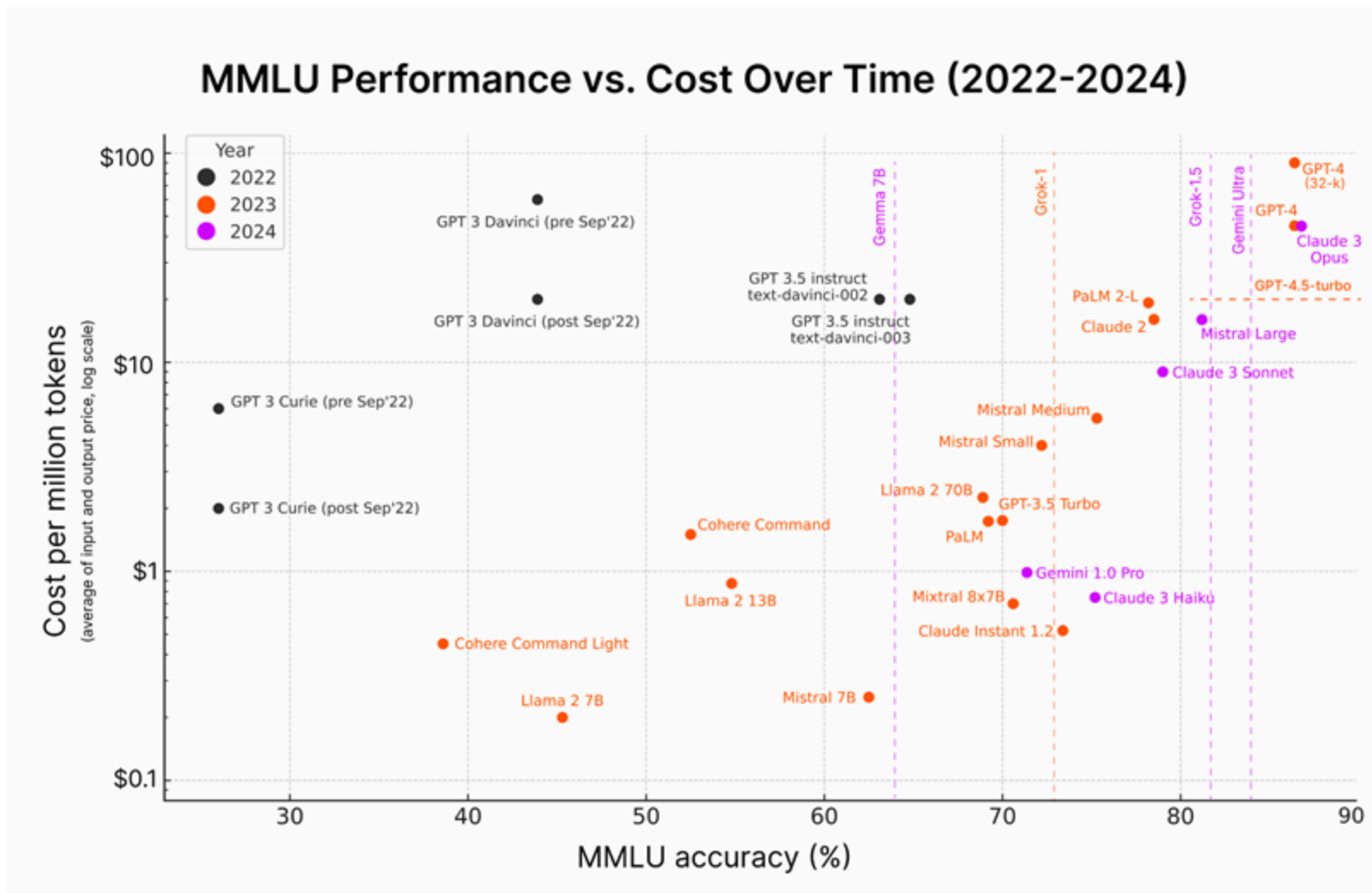
Computing's Energy Problem (and What We Can Do About it) [Horowitz, M., IEEE ISSCC 2014]



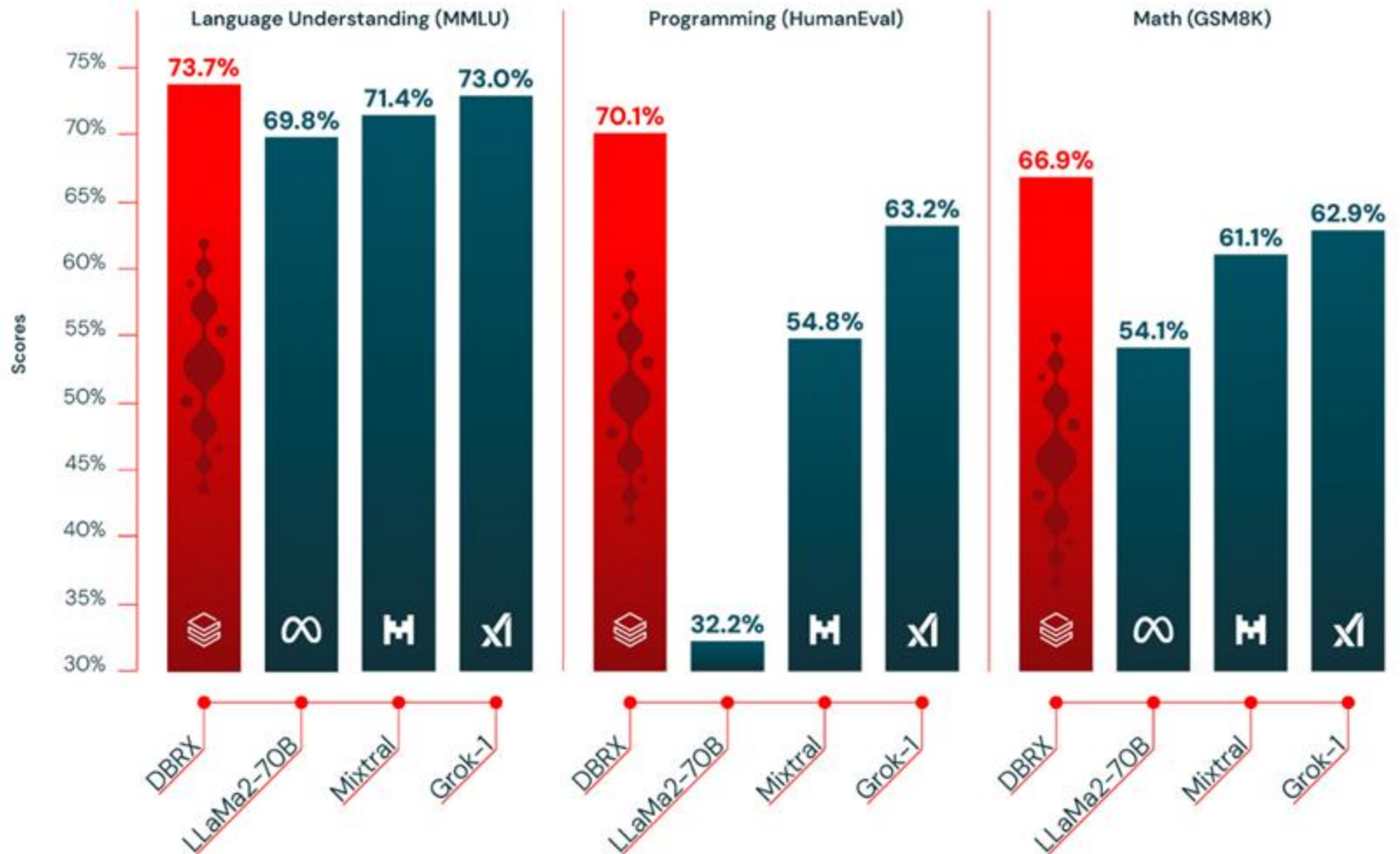
# Model Size



# Model Cost



# Development Speed

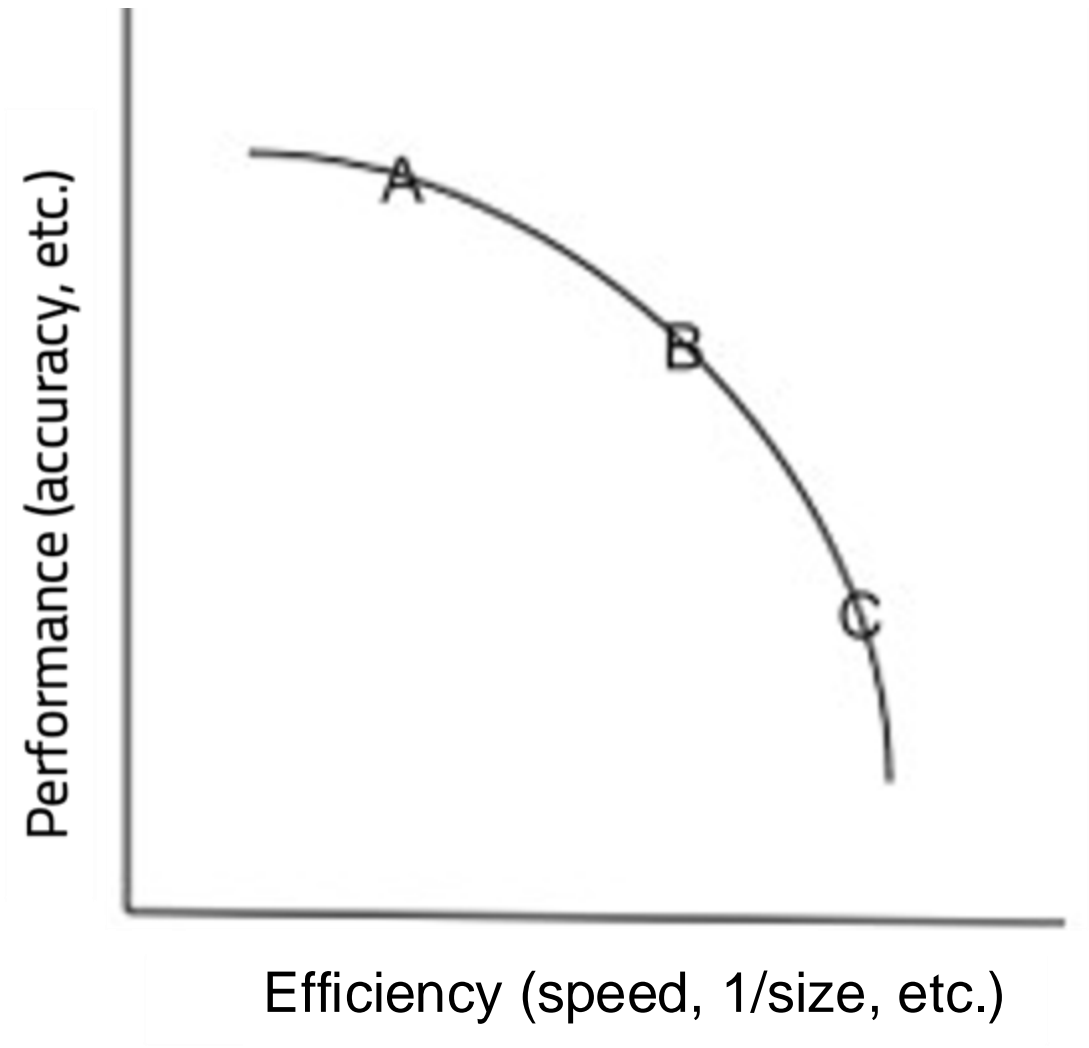


<https://www.databricks.com/blog/introducing-dbrx-new-state-art-open-llm>



# Efficiency Tradeoff

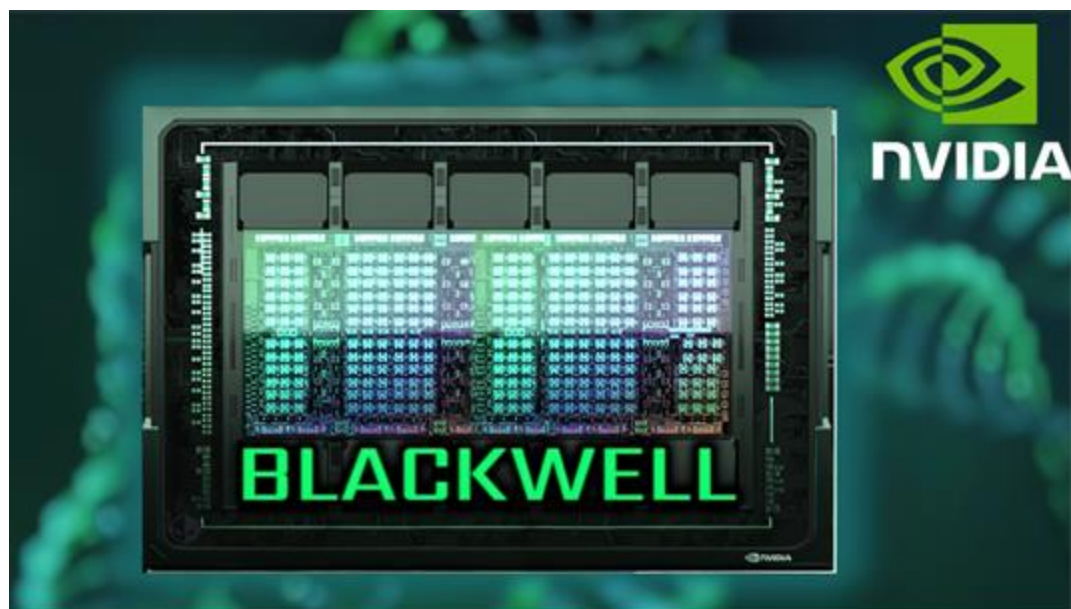
- ❑ More efficient models (smaller, faster) typically come at a cost of some performance of the model itself
- ❑ In the other direction, getting more performance from a model architecture likely means it will be larger, and require more computation



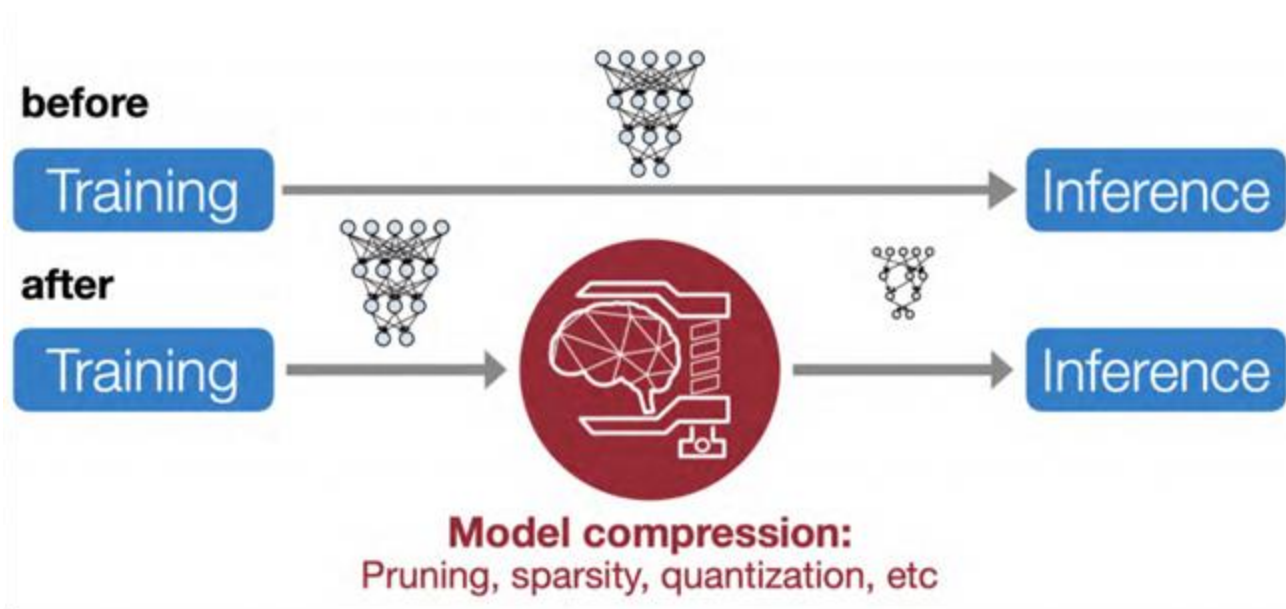


# How to Improve Model Efficiency?

Hardware



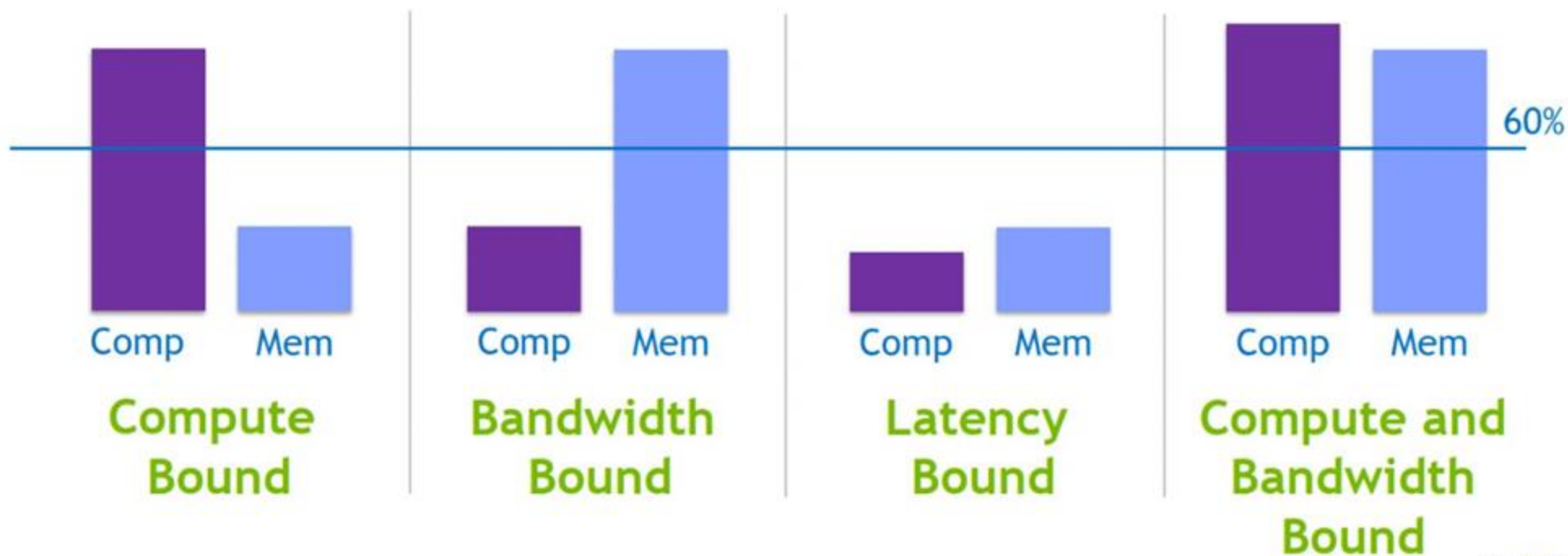
Software



# What Makes a Language Model Slow

## Memory Utilization vs Compute Utilization

Four possible combinations:



17 NVIDIA



# Efficient LLMs

## □ Quantization

- Background
- K-Means vs. Linear Quantization
- Quantization Granularity
- Quantization Aware Training (QAT) vs Post-Training Quantization (PTQ)
- LLM Quantization (LLM.int8(), SmoothQuant, AWQ, 1-bit LLMs)

## □ Sparsity (Mixture of Experts, Deja Vu: Contextual Sparsity)

## □ Long Context (PagedAttention, StreamingLLM, MHA/GQA/MQA, H<sub>2</sub>O)

## □ Speculative Decoding

## □ Parameter Efficient Fine-Tuning (BitFit, Adapter, Prompt Tuning, LoRA)



# Efficient LLMs

## □ Quantization

- **Background**
- K-Means vs. Linear Quantization
- Quantization Granularity
- Quantization Aware Training (QAT) vs Post-Training Quantization (PTQ)
- LLM Quantization (LLM.int8(), SmoothQuant, AWQ, 1-bit LLMs)

## □ Sparsity (Mixture of Experts, Deja Vu: Contextual Sparsity)

## □ Long Context (PagedAttention, StreamingLLM, MHA/GQA/MQA, H<sub>2</sub>O)

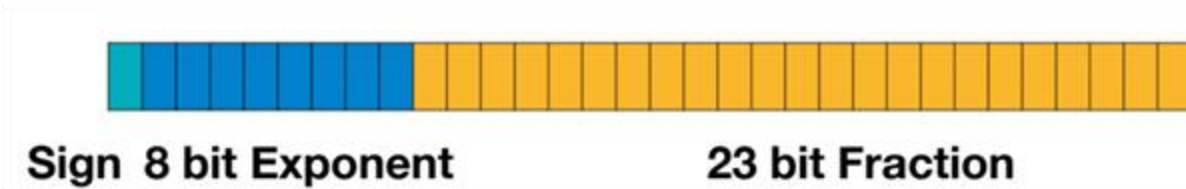
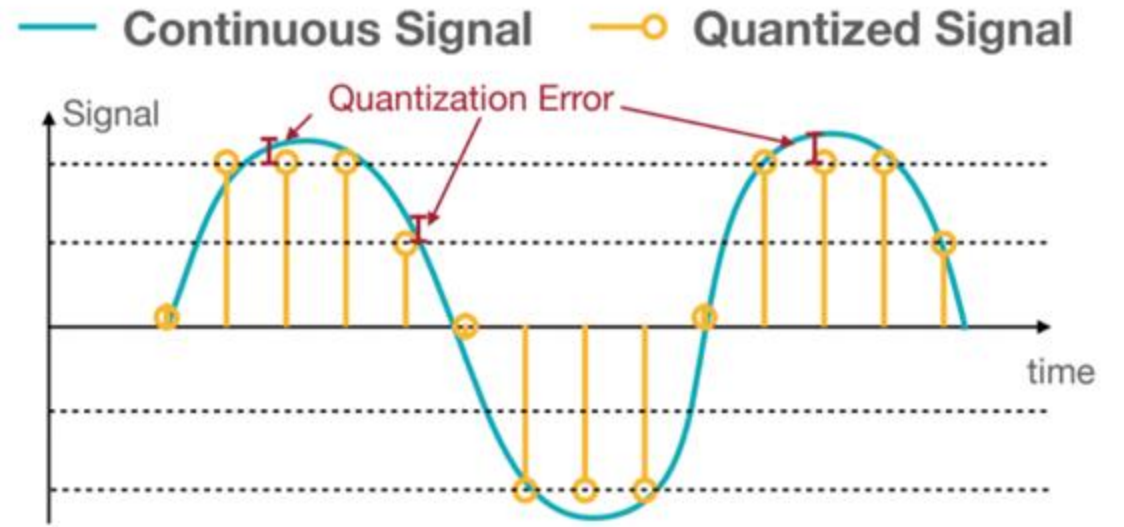
## □ Speculative Decoding

## □ Parameter Efficient Fine-Tuning (BitFit, Adapter, Prompt Tuning, LoRA)

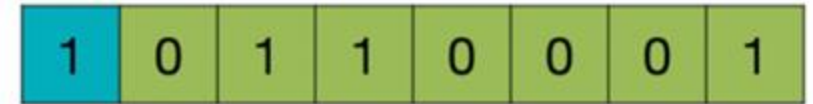


# Quantization

Reduce model size by replacing high bit-width representations with low bit-width representations



## Sign Bit



[IEEE 754](#) Half Precision 16-bit Float (IEEE FP16)



[Google Brain Float](#) (BF16)



# Efficient LLMs

## □ Quantization

- Background
- **K-Means vs. Linear Quantization**
- Quantization Granularity
- Quantization Aware Training (QAT) vs Post-Training Quantization (PTQ)
- LLM Quantization (LLM.int8(), SmoothQuant, AWQ, 1-bit LLMs)

## □ Sparsity (Mixture of Experts, Deja Vu: Contextual Sparsity)

## □ Long Context (PagedAttention, StreamingLLM, MHA/GQA/MQA, H<sub>2</sub>O)

## □ Speculative Decoding

## □ Parameter Efficient Fine-Tuning (BitFit, Adapter, Prompt Tuning, LoRA)



# K-Means Quantization vs Linear Quantization

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

3	0	2	1	3:	2.00
1	1	0	3	2:	1.50
0	3	1	0	1:	0.00
3	1	2	2	0:	-1.00

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

$(-1) \times 1.07$

**K-Means-based  
Quantization**

**Linear  
Quantization**

<b>Storage</b>	Floating-Point Weights	Integer Weights; Floating-Point Codebook	Integer Weights
<b>Computation</b>	Floating-Point Arithmetic	Floating-Point Arithmetic	Integer Arithmetic



# K-Means Quantization vs Linear Quantization

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

3	0	2	1	3:	2.00
1	1	0	3	2:	1.50
0	3	1	0	1:	0.00
3	1	2	2	0:	-1.00

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

$(-1) \times 1.07$

**K-Means-based  
Quantization**

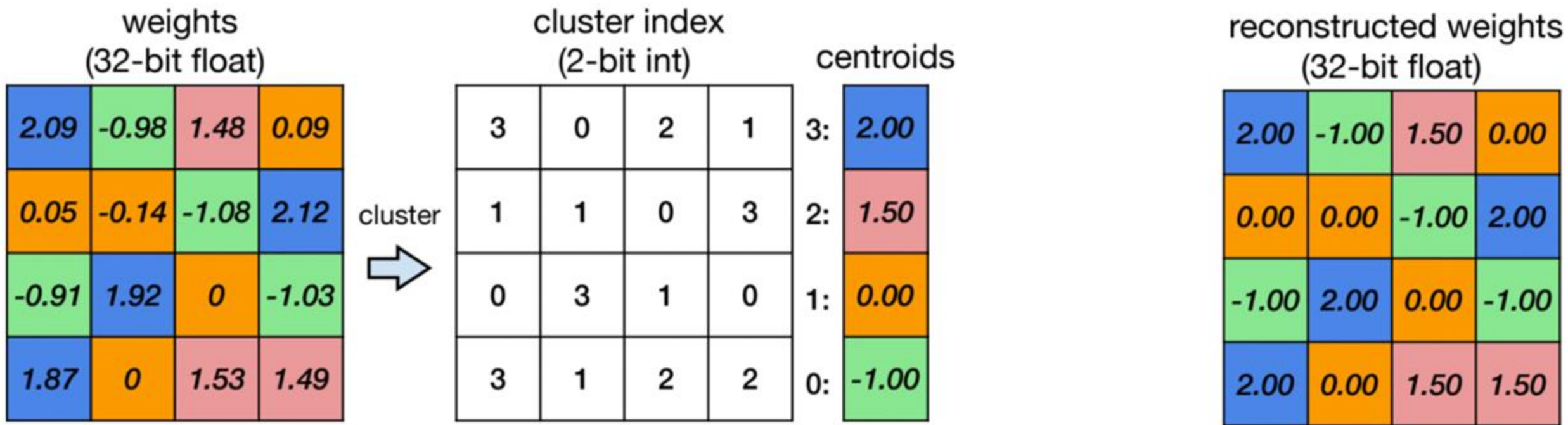
**Linear  
Quantization**

<b>Storage</b>	Floating-Point Weights	<b>Integer Weights; Floating-Point Codebook</b>	Integer Weights
<b>Computation</b>	Floating-Point Arithmetic	<b>Floating-Point Arithmetic</b>	Integer Arithmetic





# K-Means Quantization



Deep Compression [Han et al., ICLR 2016]



# K-Means Quantization

Original weights

weights (32-bit 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

cluster  
➔

cluster index (2-bit int)				centroids
3	0	2	1	3: 2.00
1	1	0	3	2: 1.50
0	3	1	0	1: 0.00
3	1	2	2	0: -1.00

reconstructed weights (32-bit float)			
2.00	-1.00	1.50	0.00
0.00	0.00	-1.00	2.00
-1.00	2.00	0.00	-1.00
2.00	0.00	1.50	1.50

Deep Compression [Han et al., ICLR 2016]



# K-Means Quantization

Stored weights after clustering

weights  
(32-bit 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

cluster  
➔

cluster index (2-bit int)				centroids
3	0	2	1	3: 2.00
1	1	0	3	2: 1.50
0	3	1	0	1: 0.00
3	1	2	2	0: -1.00

reconstructed weights  
(32-bit float)

2.00	-1.00	1.50	0.00
0.00	0.00	-1.00	2.00
-1.00	2.00	0.00	-1.00
2.00	0.00	1.50	1.50

Deep Compression [Han et al., ICLR 2016]



# K-Means Quantization

Retrieved weights to  
be used at inference  
time

weights  
(32-bit 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

cluster  
➔

cluster index  
(2-bit int)

3	0	2	1
1	1	0	3
0	3	1	0
3	1	2	2

centroids

3:	2.00
2:	1.50
1:	0.00
0:	-1.00

reconstructed weights  
(32-bit float)

2.00	-1.00	1.50	0.00
0.00	0.00	-1.00	2.00
-1.00	2.00	0.00	-1.00
2.00	0.00	1.50	1.50

Deep Compression [Han et al., ICLR 2016]



# K-Means Quantization vs Linear Quantization

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

3	0	2	1	3:	2.00
1	1	0	3	2:	1.50
0	3	1	0	1:	0.00
3	1	2	2	0:	-1.00

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

$(-1) \times 1.07$

**K-Means-based  
Quantization**

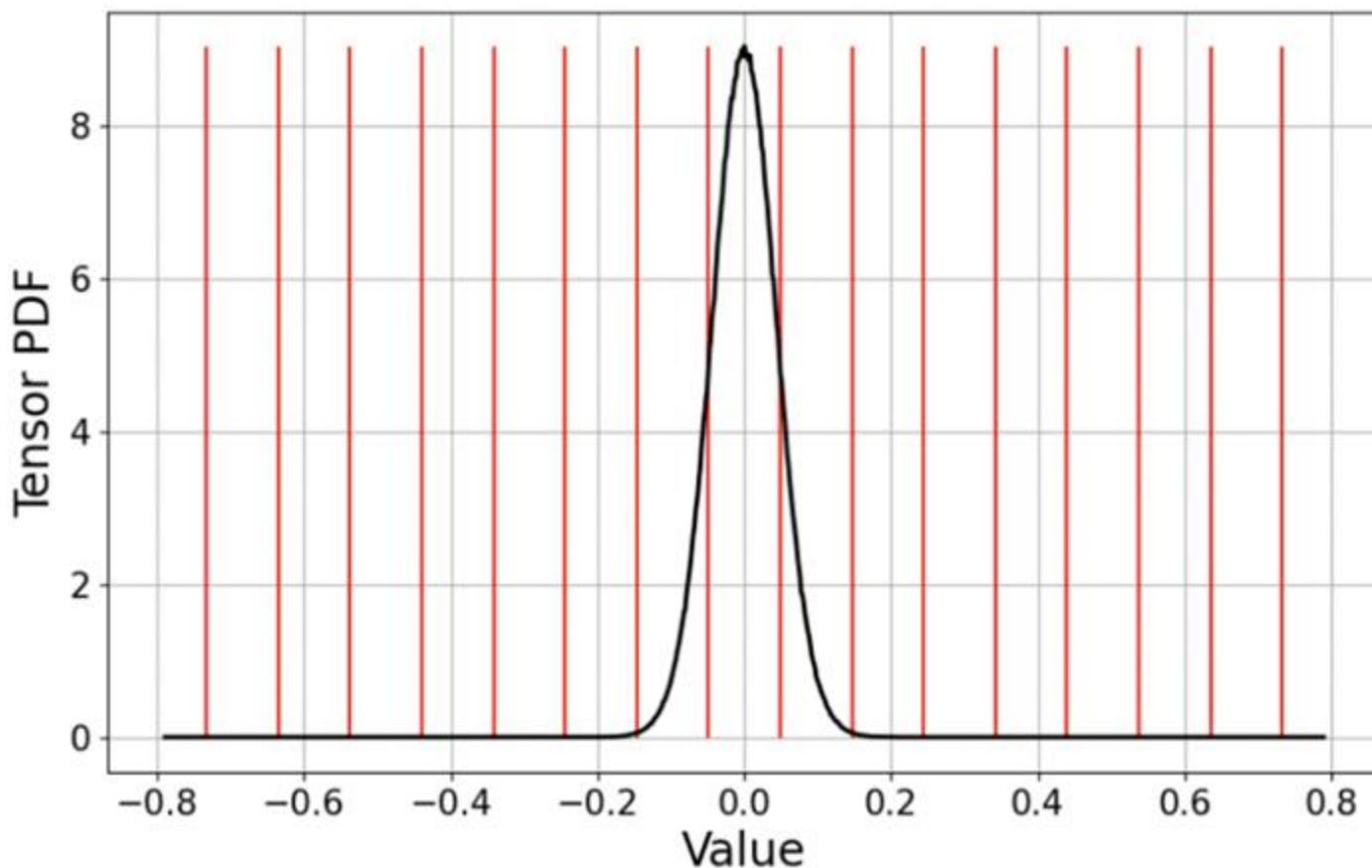
**Linear  
Quantization**

<b>Storage</b>	Floating-Point Weights	Integer Weights; Floating-Point Codebook	Integer Weights
<b>Computation</b>	Floating-Point Arithmetic	Floating-Point Arithmetic	Integer Arithmetic



# Linear Quantization

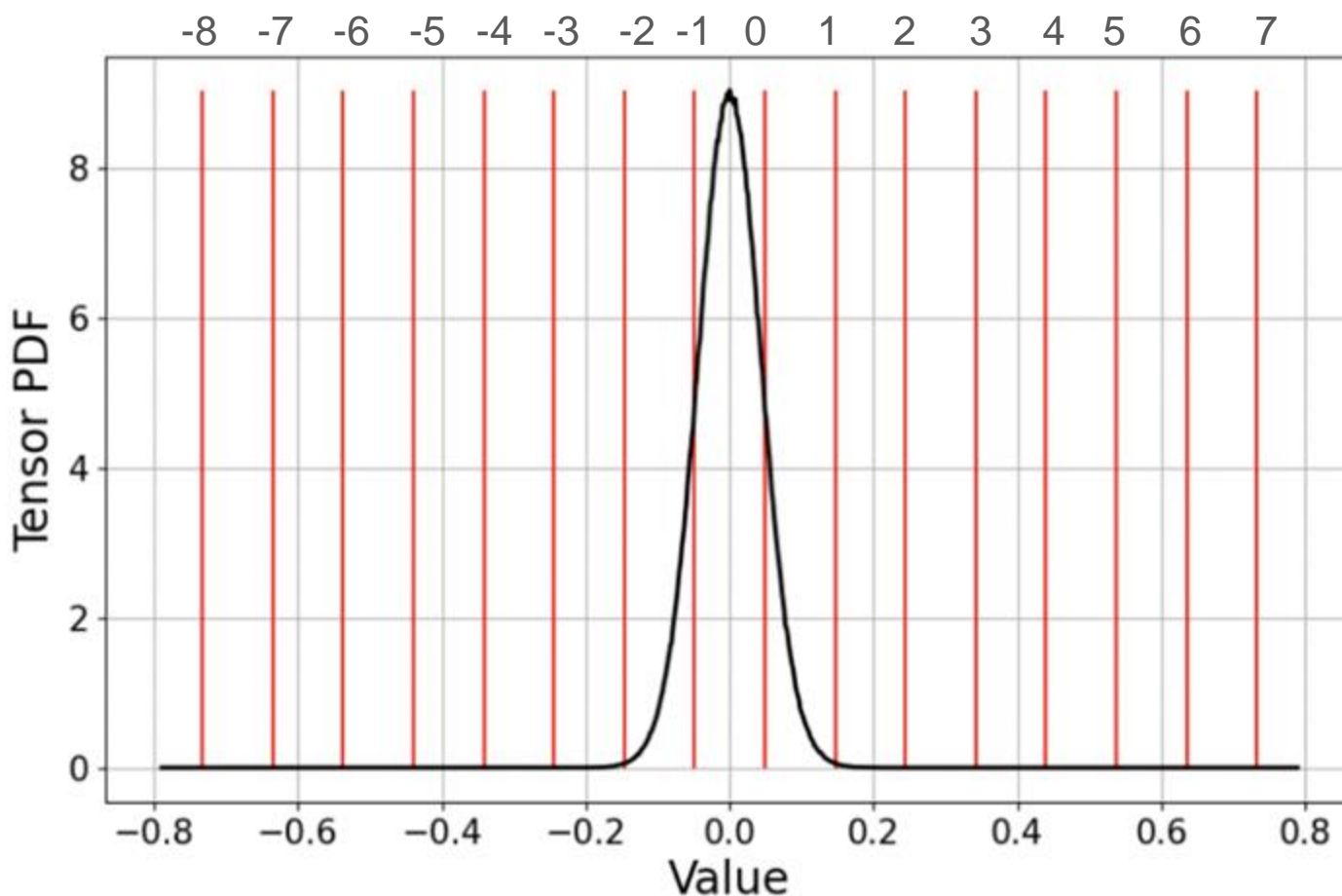
- Apply linear function on weights and hidden state activations from floating point values ( $r$ ) to integer values ( $q$ )
- Original weights (black), Quantized bins (red)
- Black weights are mapped to one of the vertical red lines



# Linear Quantization

- Apply linear function on weights and hidden state activations from floating point values ( $r$ ) to integer values ( $q$ )
- Original weights (black), Quantized bins (red)
- Black weights are mapped to one of the vertical red lines

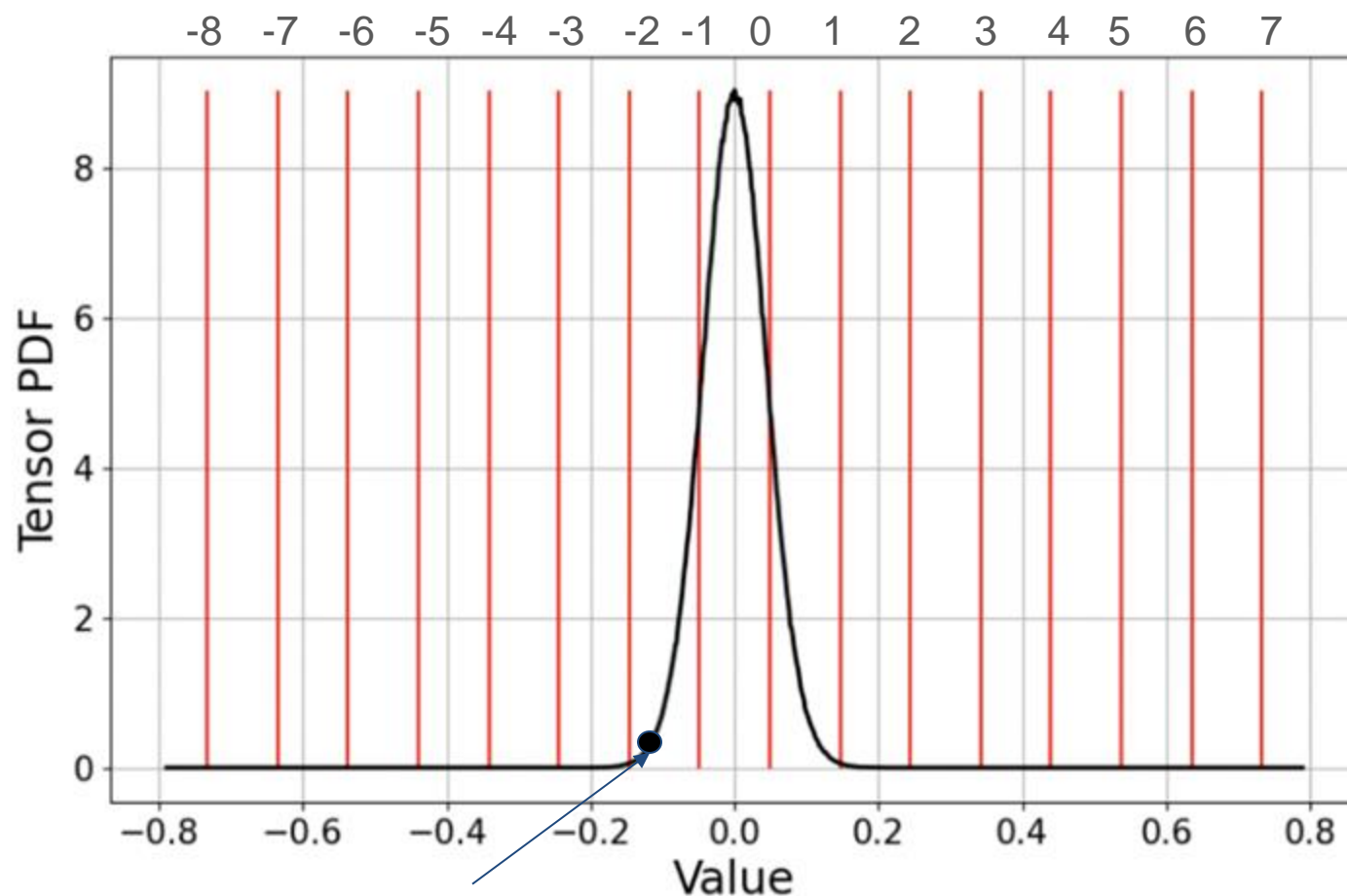
32-bit float to 4-bit int



# Linear Quantization

- Apply linear function on weights and hidden state activations from floating point values ( $r$ ) to integer values ( $q$ )
- Original weights (black), Quantized bins (red)
- Black weights are mapped to one of the vertical red lines

32-bit float to 4-bit int

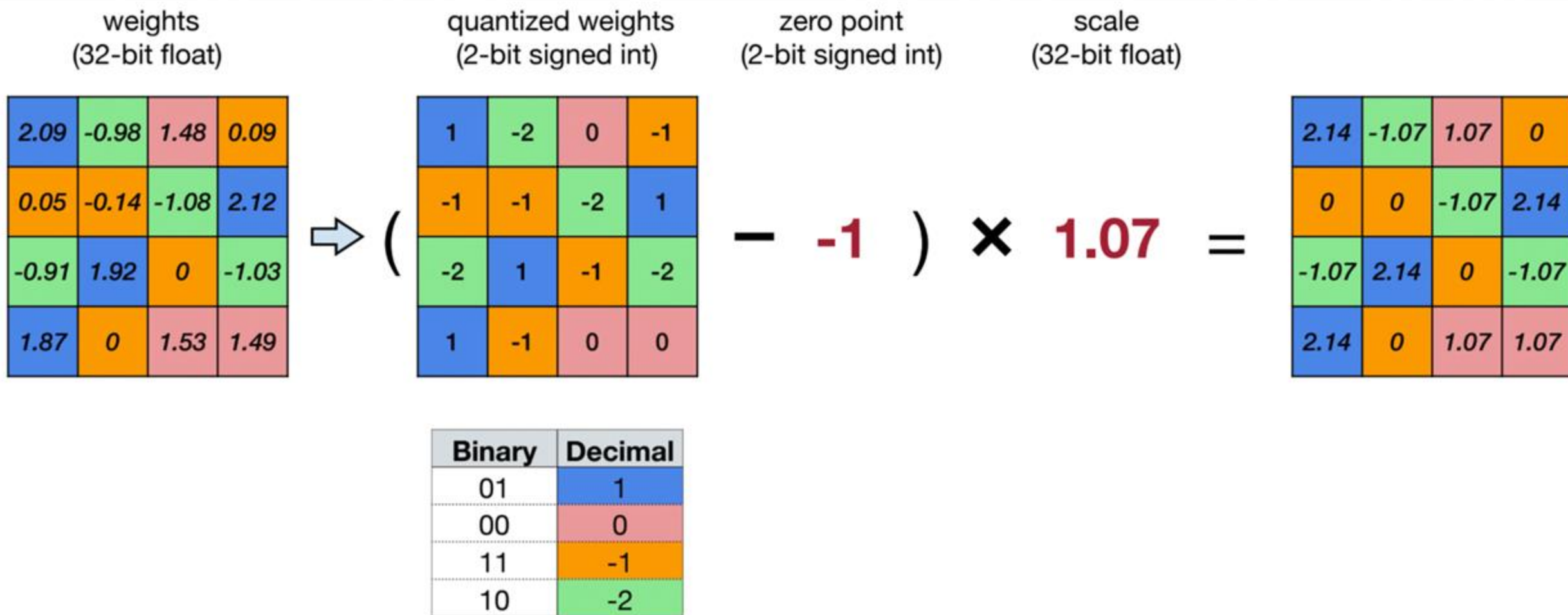


Before Quantization:  $-.14$

After Quantization:  $-2$



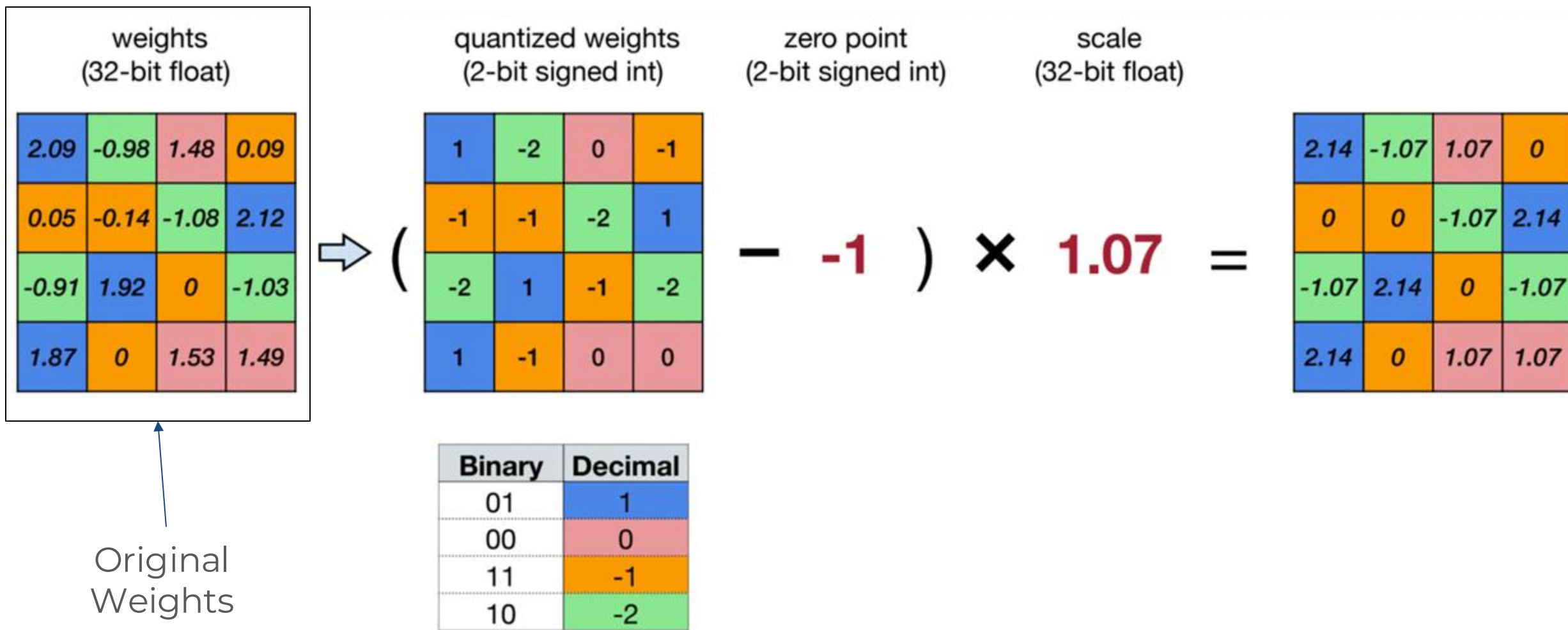
# Linear Quantization



Quantization and Training of Neural Networks for Efficient Integer-Arithmetic-Only Inference [Jacob et al., CVPR 2018]



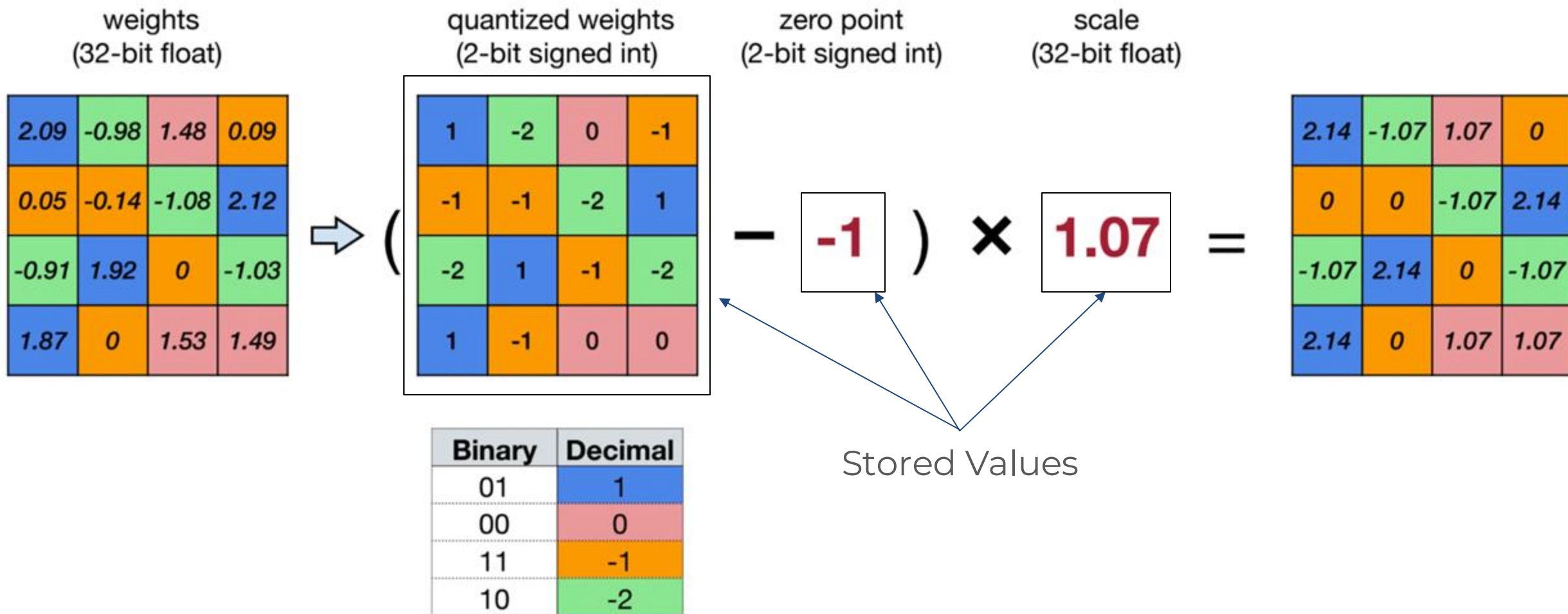
# Linear Quantization



Quantization and Training of Neural Networks for Efficient Integer-Arithmetic-Only Inference [Jacobet et al., CVPR 2018]



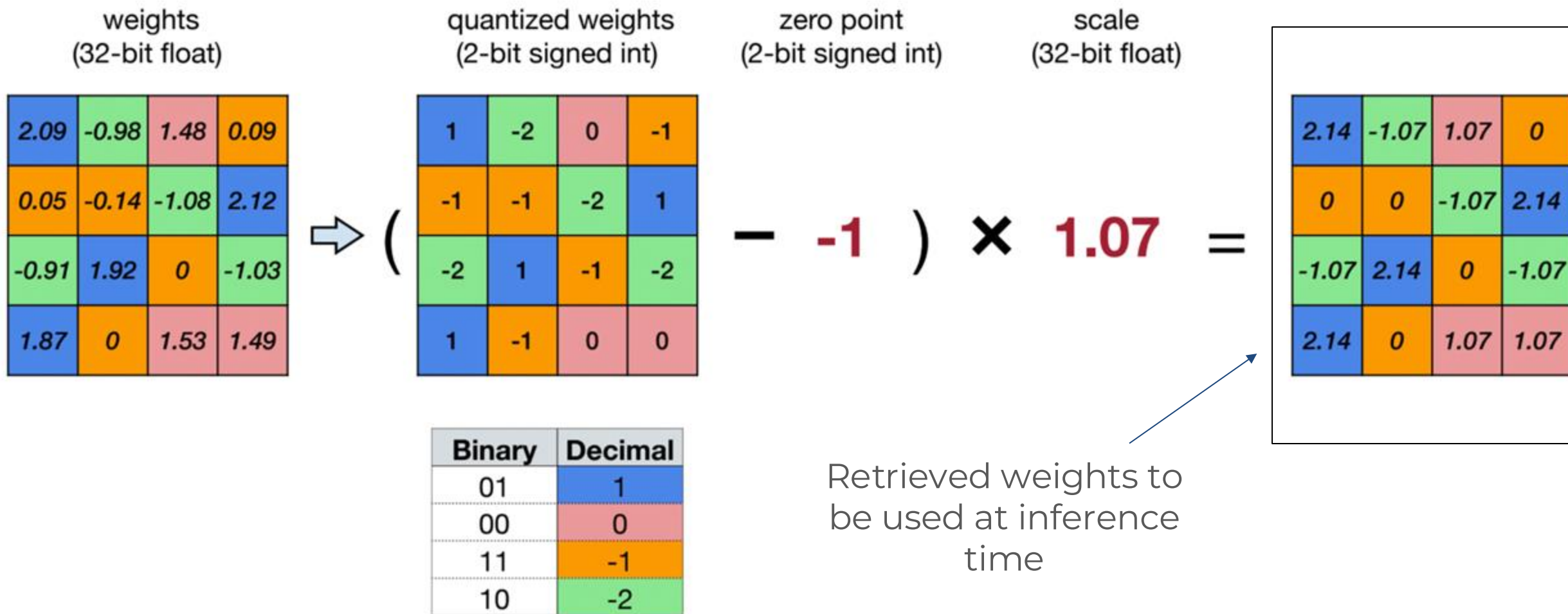
# Linear Quantization



Quantization and Training of Neural Networks for Efficient Integer-Arithmetic-Only Inference [Jacobet et al., CVPR 2018]



# Linear Quantization



Quantization and Training of Neural Networks for Efficient Integer-Arithmetic-Only Inference [Jacob et al., CVPR 2018]



# Efficient LLMs

## □ Quantization

- Background
- K-Means vs. Linear Quantization
- **Quantization Granularity**
- Quantization Aware Training (QAT) vs Post-Training Quantization (PTQ)
- LLM Quantization (LLM.int8(), SmoothQuant, AWQ, 1-bit LLMs)

## □ Sparsity (Mixture of Experts, Deja Vu: Contextual Sparsity)

## □ Long Context (PagedAttention, StreamingLLM, MHA/GQA/MQA, H<sub>2</sub>O)

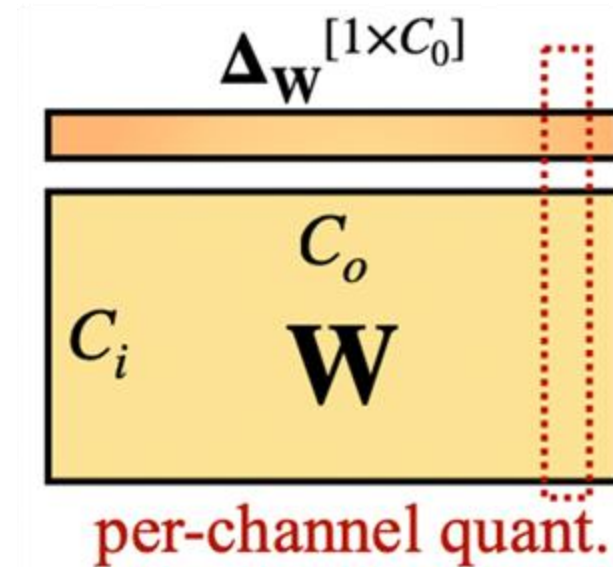
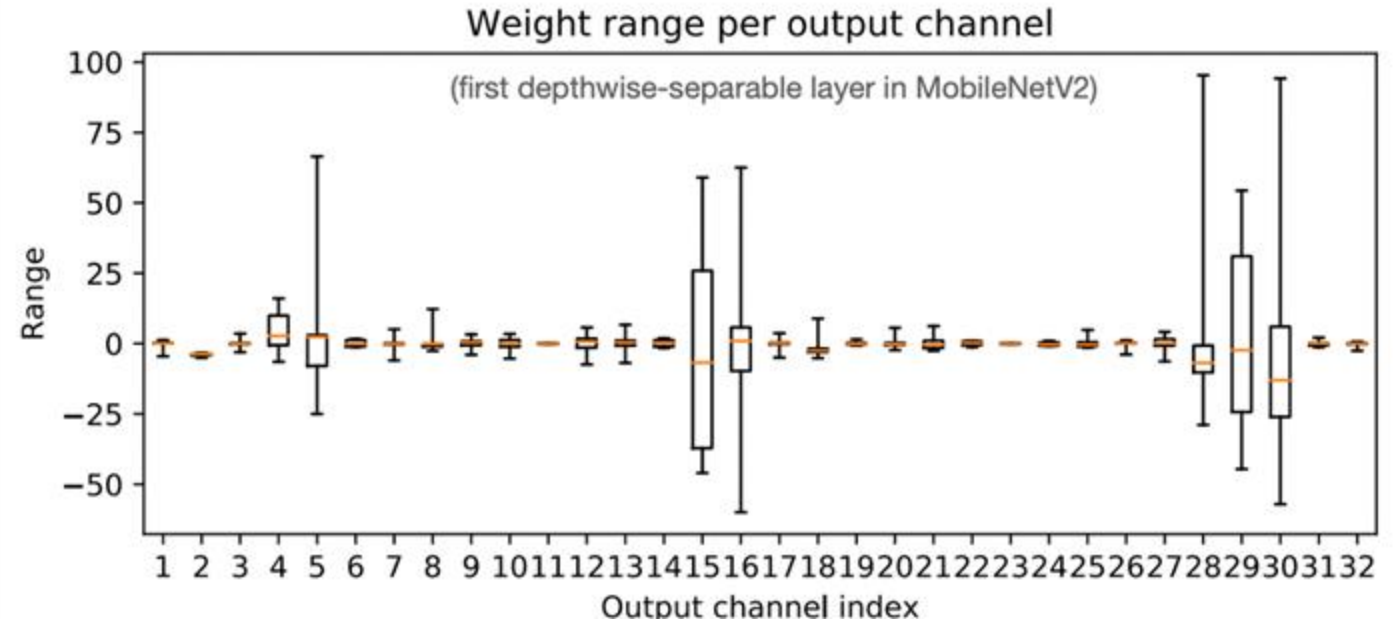
## □ Speculative Decoding

## □ Parameter Efficient Fine-Tuning (BitFit, Adapter, Prompt Tuning, LoRA)



# Weight Granularity

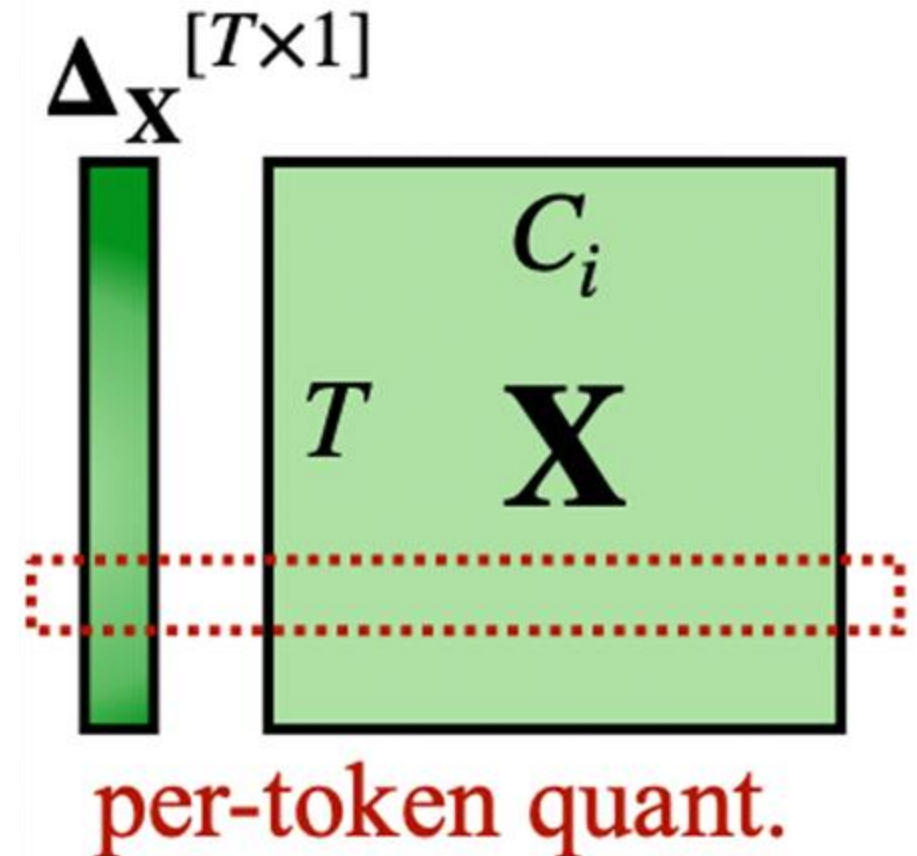
- ❑ Weight matrices will often have different variances along each output channel
- ❑ High variance in weights means that applying linear quantization will result in large performance degradation
- ❑ To fix this, we can perform linear quantization along each channel of the weight tensor separately



SmoothQuant: Accurate and Efficient Post-Training Quantization for Large Language Models[Xiao et. al., ICML 2023]

# Activation Granularity

- ❑ Activations can have a similar problem whereby the variance by channel can be quite different
- ❑ The variance by token can also differ dramatically
- ❑ When applying quantization, we should split up channels, tokens to take this into account



# Efficient LLMs

## □ Quantization

- Background
- K-Means vs. Linear Quantization
- Quantization Granularity
- **Quantization Aware Training (QAT) vs Post-Training Quantization (PTQ)**
- LLM Quantization (LLM.int8(), SmoothQuant, AWQ, 1-bit LLMs)

## □ Sparsity (Mixture of Experts, Deja Vu: Contextual Sparsity)

## □ Long Context (PagedAttention, StreamingLLM, MHA/GQA/MQA, H<sub>2</sub>O)

## □ Speculative Decoding

## □ Parameter Efficient Fine-Tuning (BitFit, Adapter, Prompt Tuning, LoRA)





# Quantization Aware Training (QAT)

Quantize while training

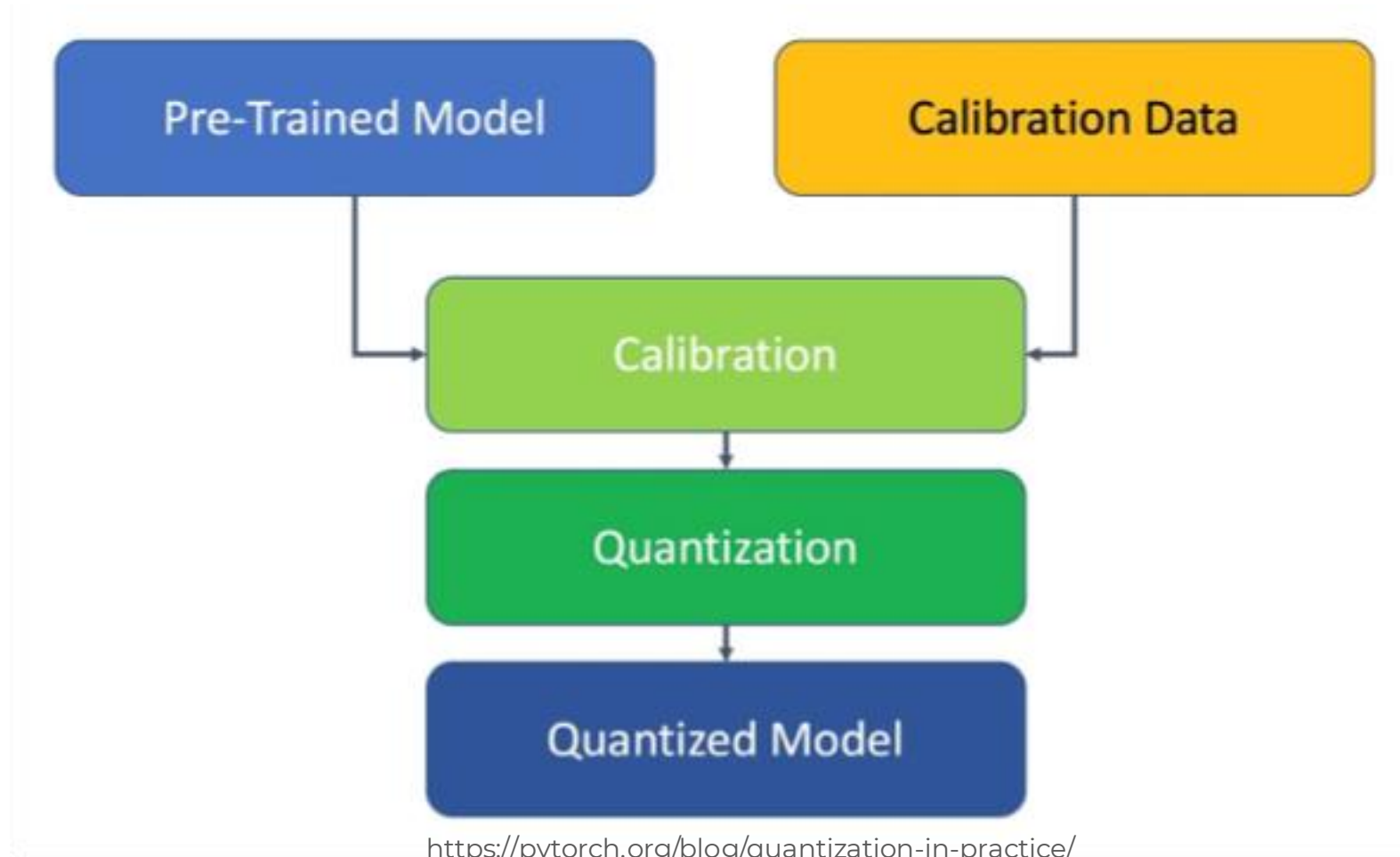


<https://pytorch.org/blog/quantization-in-practice/>



# Post Training Quantization (PTQ)

Quantize after training



# Efficient LLMs

## □ Quantization

- Background
- K-Means vs. Linear Quantization
- Quantization Granularity
- Quantization Aware Training (QAT) vs Post-Training Quantization (PTQ)
- **LLM Quantization (LLM.int8(), SmoothQuant, AWQ, 1-bit LLMs)**

## □ Sparsity (Mixture of Experts, Deja Vu: Contextual Sparsity)

## □ Long Context (PagedAttention, StreamingLLM, MHA/GQA/MQA, H<sub>2</sub>O)

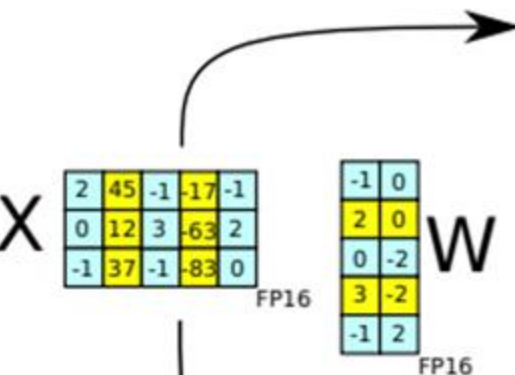
## □ Speculative Decoding

## □ Parameter Efficient Fine-Tuning (BitFit, Adapter, Prompt Tuning, LoRA)

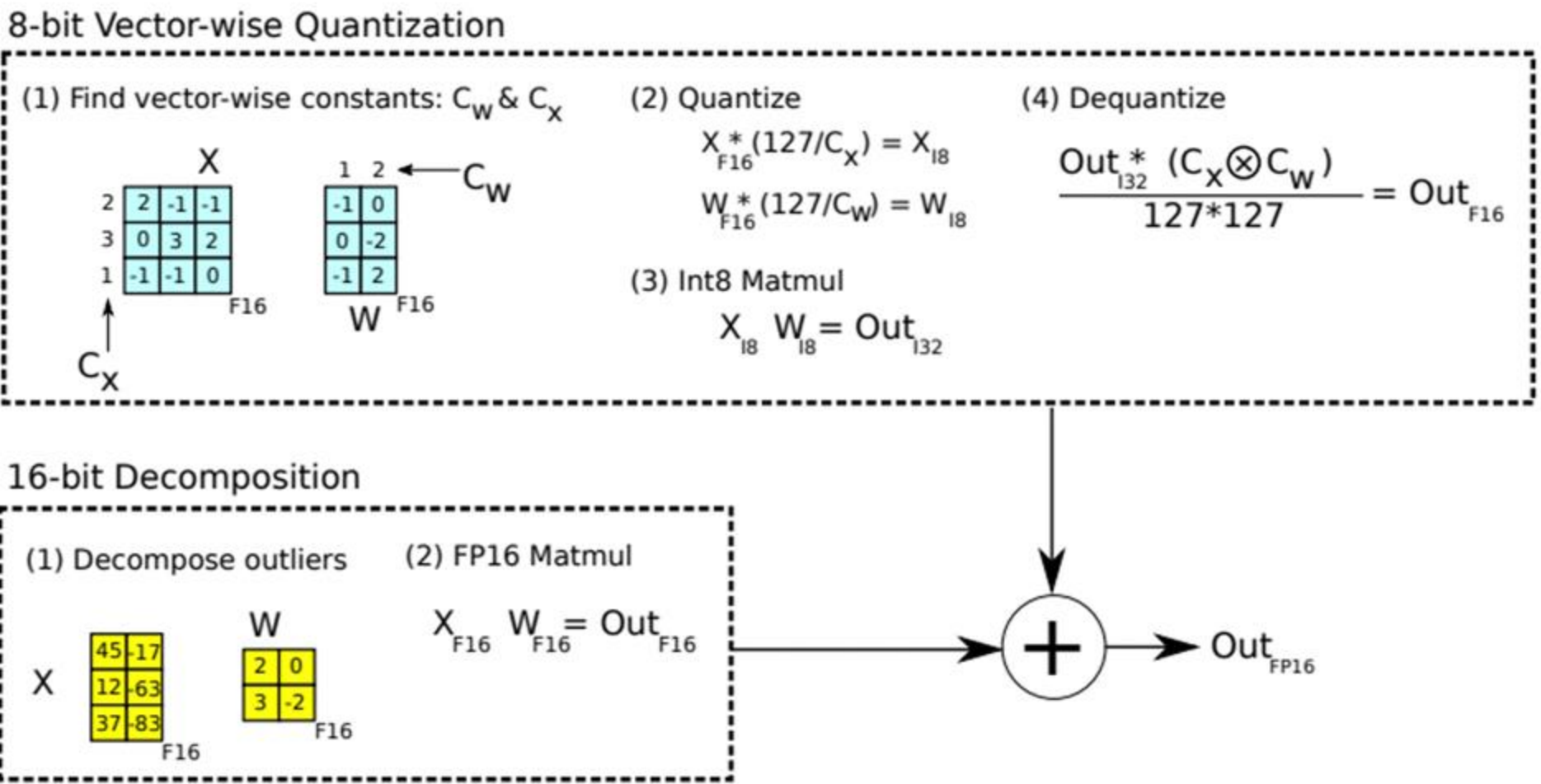


# LLM.int8() (W8A8)

## LLM.int8()



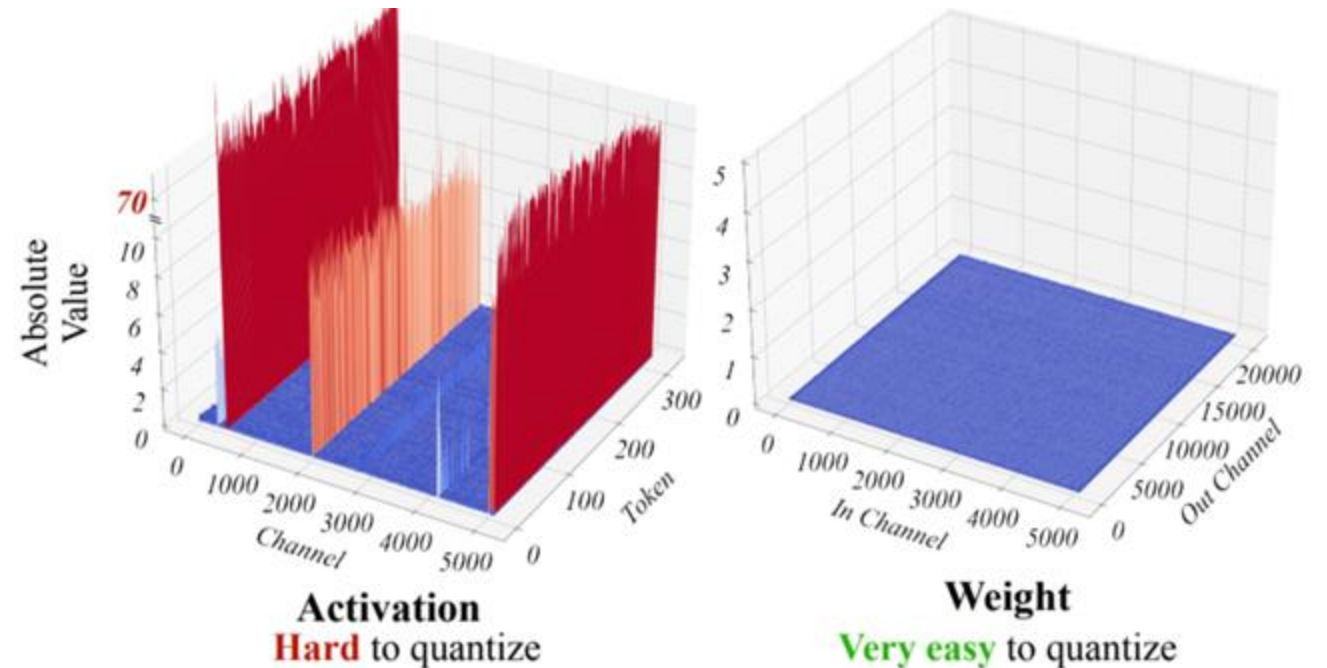
Regular values  
 Outliers



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

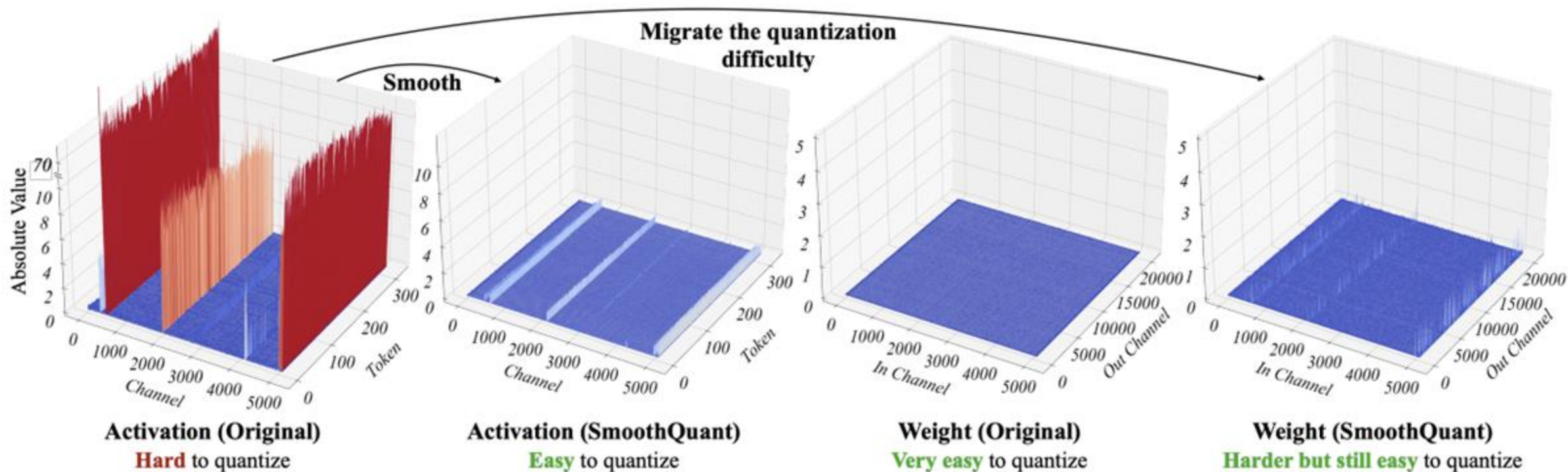
# SmoothQuant (W8A8)

Observation: High variance channels are fixed in activations in LLM FFN layers-weights have relatively little difference in variance



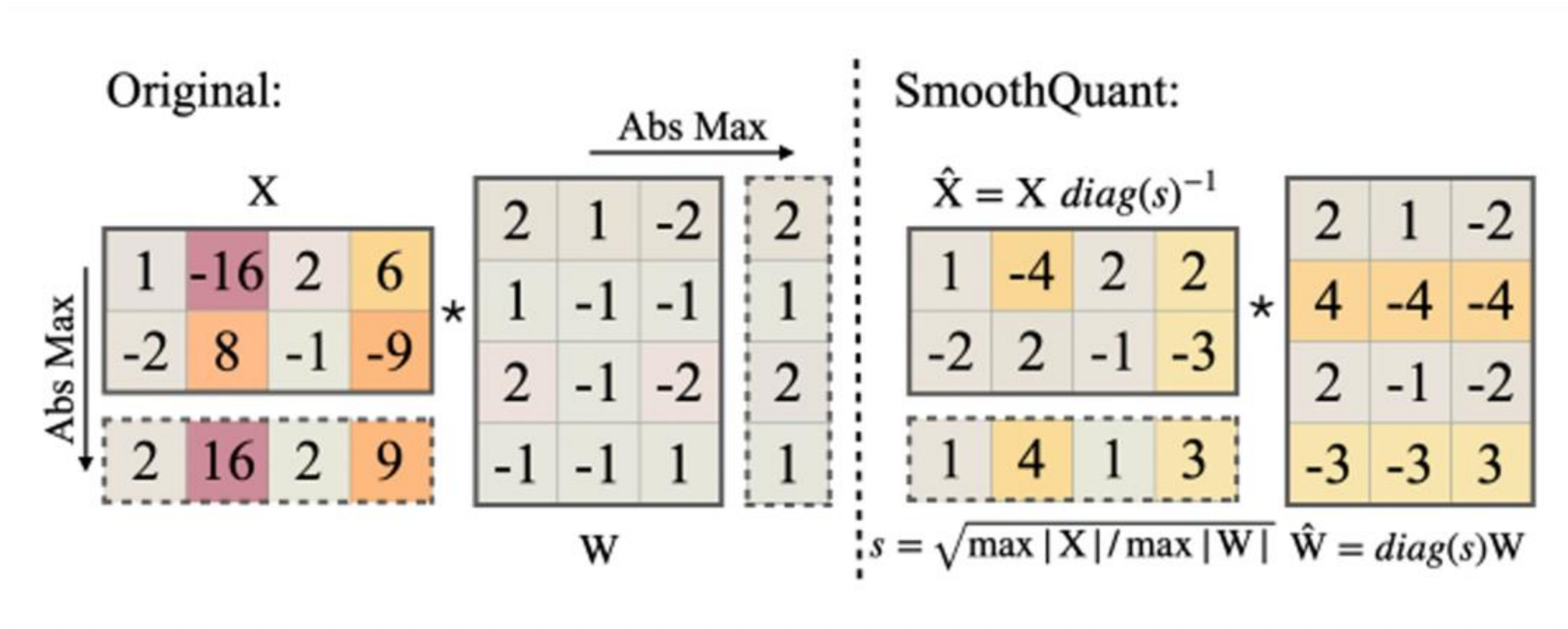
SmoothQuant: Accurate and Efficient Post-Training Quantization for Large Language Models[Xiao et. al., ICML 2023]

# SmoothQuant (W8A8)



SmoothQuant: Accurate and Efficient Post-Training Quantization for Large Language Models[Xiao et. al., ICML 2023]

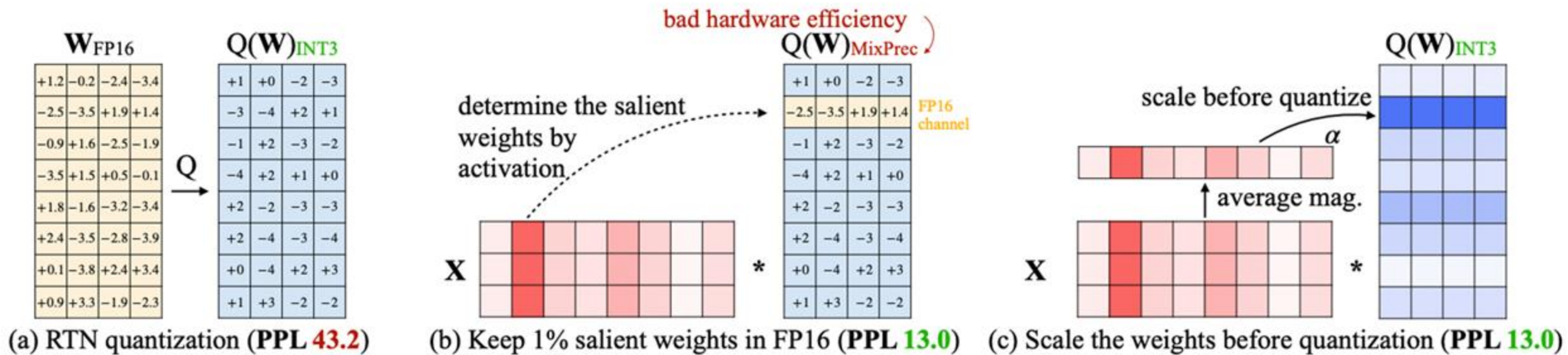
# SmoothQuant (W8A8)



SmoothQuant: Accurate and Efficient Post-Training Quantization for Large Language Models[Xiao et. al., ICML 2023]



# AWQ (W4A16)



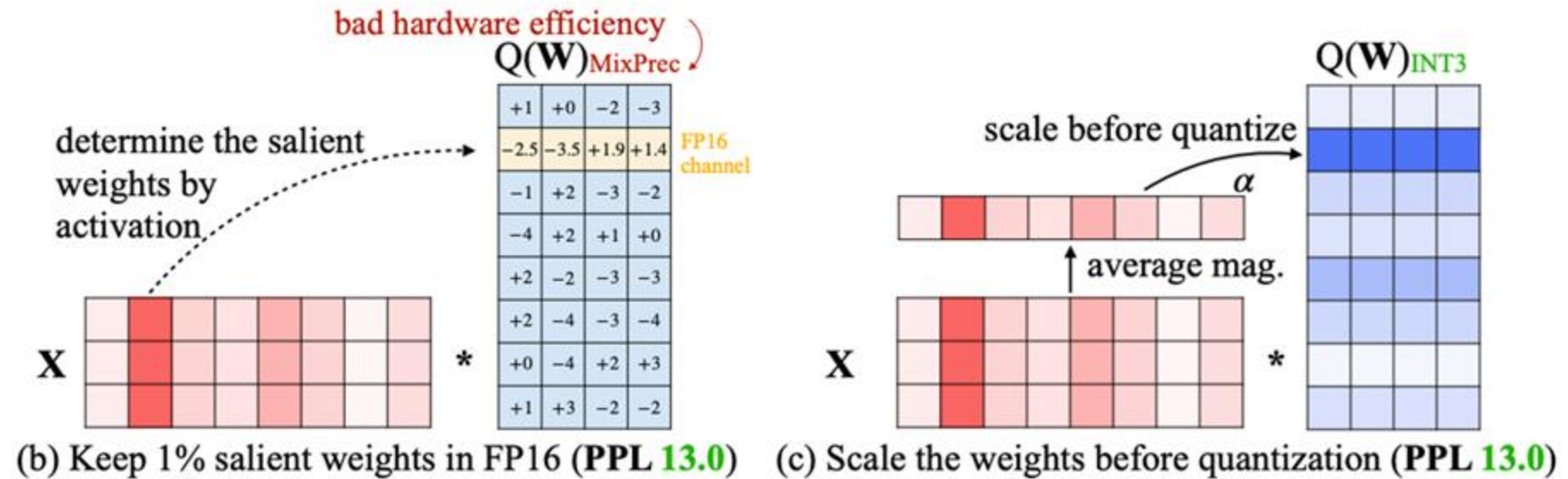
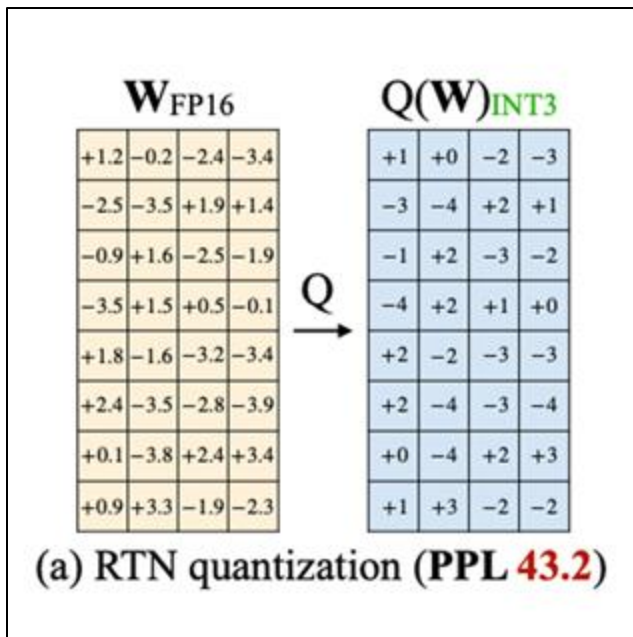
AWQ: Activation-aware Weight Quantization for LLM Compression and Acceleration [Lin et. al., arxiv 2023]





# AWQ (W4A16)

Normal quantization on LLMs performs poorly due to outliers in the model's hidden state

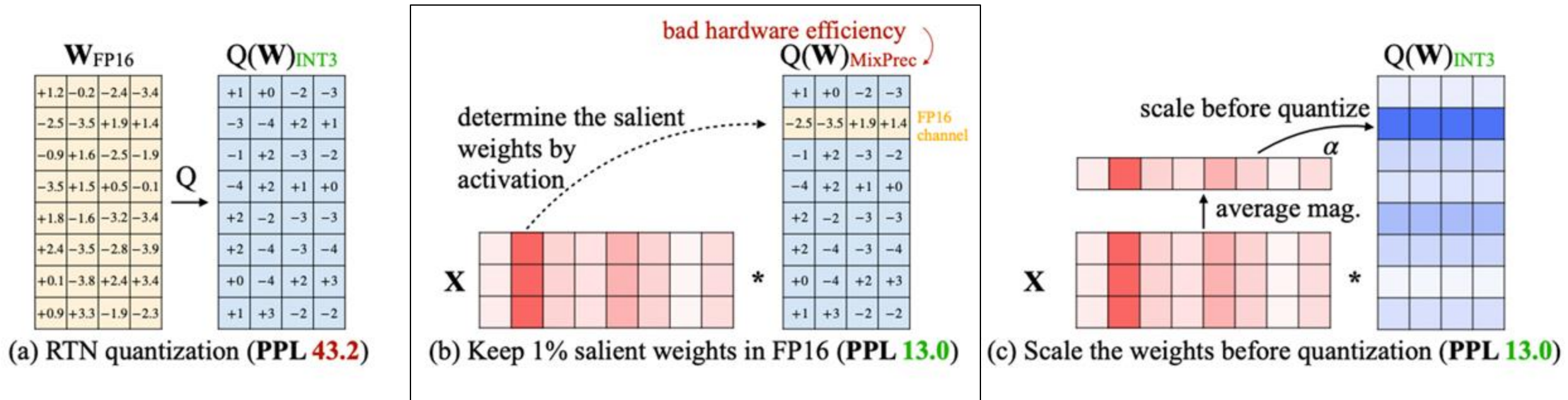


AWQ: Activation-aware Weight Quantization for LLM Compression and Acceleration [Lin et. al., arxiv 2023]



# AWQ (W4A16)

LLM.int8() can resolve these issues, but mixed precision matrix multiplication is hardware inefficient

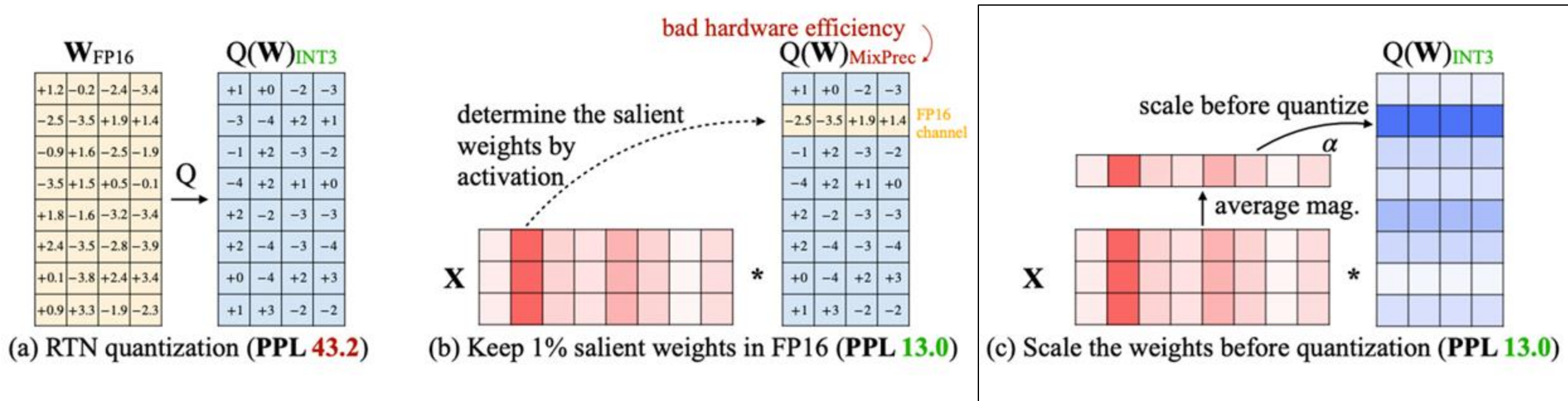


AWQ: Activation-aware Weight Quantization for LLM Compression and Acceleration [Lin et. al., arxiv 2023]



# AWQ (W4A16)

As in SmoothQuant, we can resolve this issue by shifting the difficulty to the weights using a scaling factor.

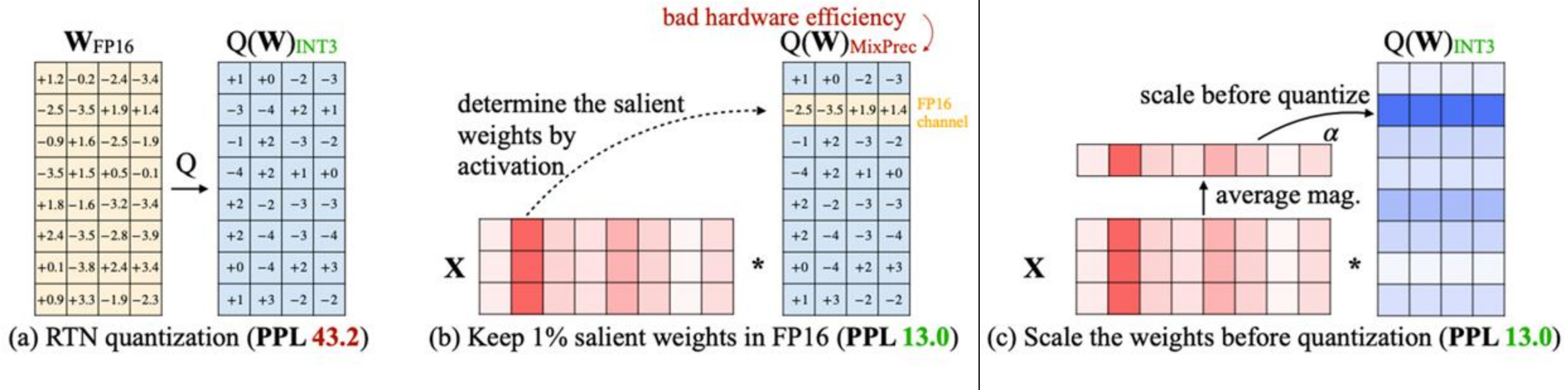


AWQ: Activation-aware Weight Quantization for LLM Compression and Acceleration [Lin et. al., arxiv 2023]



# AWQ (W4A16)

Where Smoothquant quantizes both activations and weights, AWQ only quantizes the weights

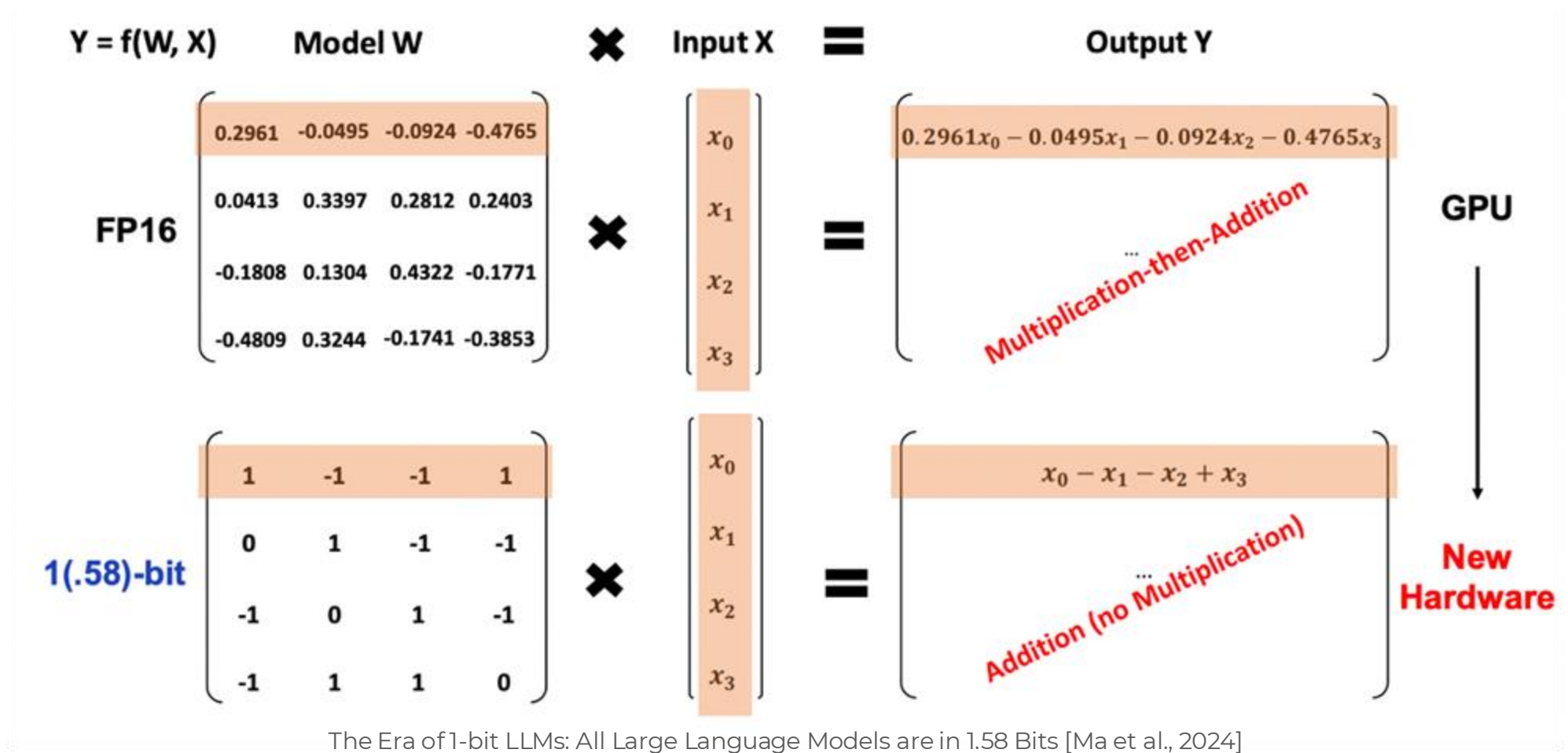


AWQ: Activation-aware Weight Quantization for LLM Compression and Acceleration [Lin et. al., arxiv 2023]



# Era of 1-bit LLMs (W1.58A8)

Weight-only QAT algorithm that uses only weights in  $\{-1, 0, 1\}$



# Era of 1-bit LLMs (W1.58A8)

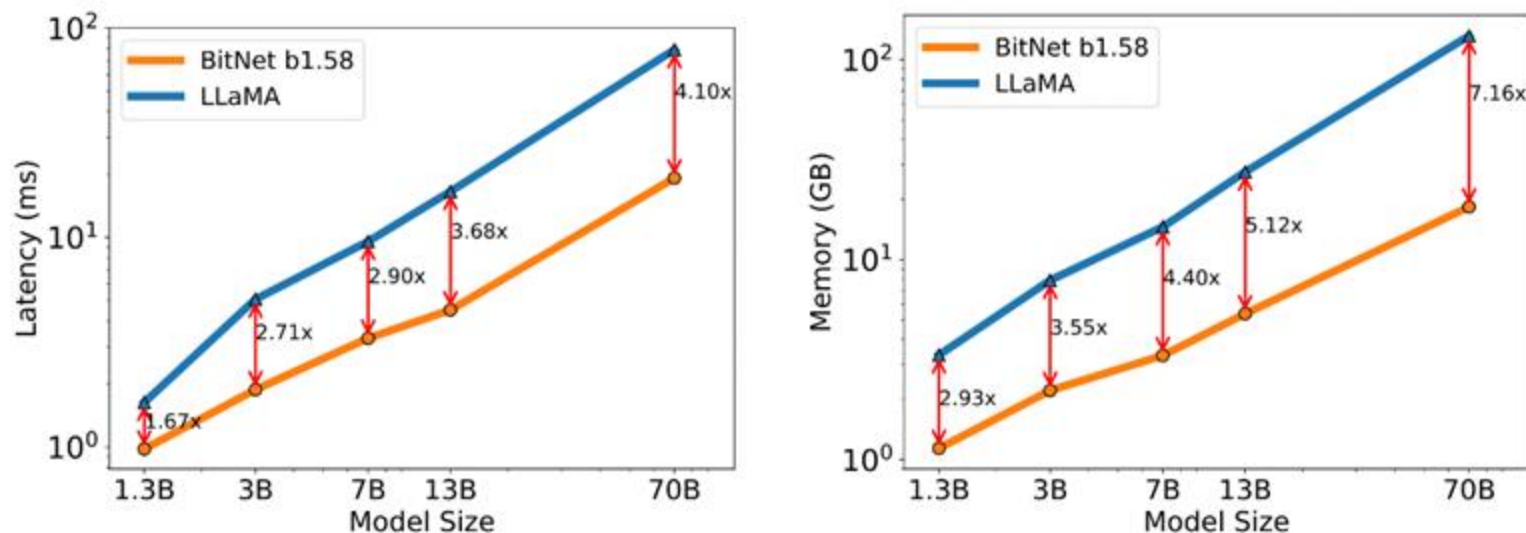


Figure 2: Decoding latency (Left) and memory consumption (Right) of BitNet b1.58 varying the model size.

Models	Size	Max Batch Size	Throughput (tokens/s)
LLaMA LLM	70B	16 (1.0x)	333 (1.0x)
<b>BitNet b1.58</b>	70B	<b>176 (11.0x)</b>	<b>2977 (8.9x)</b>

Table 3: Comparison of the throughput between BitNet b1.58 70B and LLaMA LLM 70B.

The Era of 1-bit LLMs: All Large Language Models are in 1.58 Bits [Ma et al., 2024]



# Efficient LLMs

## □ Quantization

- Background
- K-Means vs. Linear Quantization
- Quantization Granularity
- Quantization Aware Training (QAT) vs Post-Training Quantization (PTQ)
- LLM Quantization (LLM.int8(), SmoothQuant, AWQ, 1-bit LLMs)

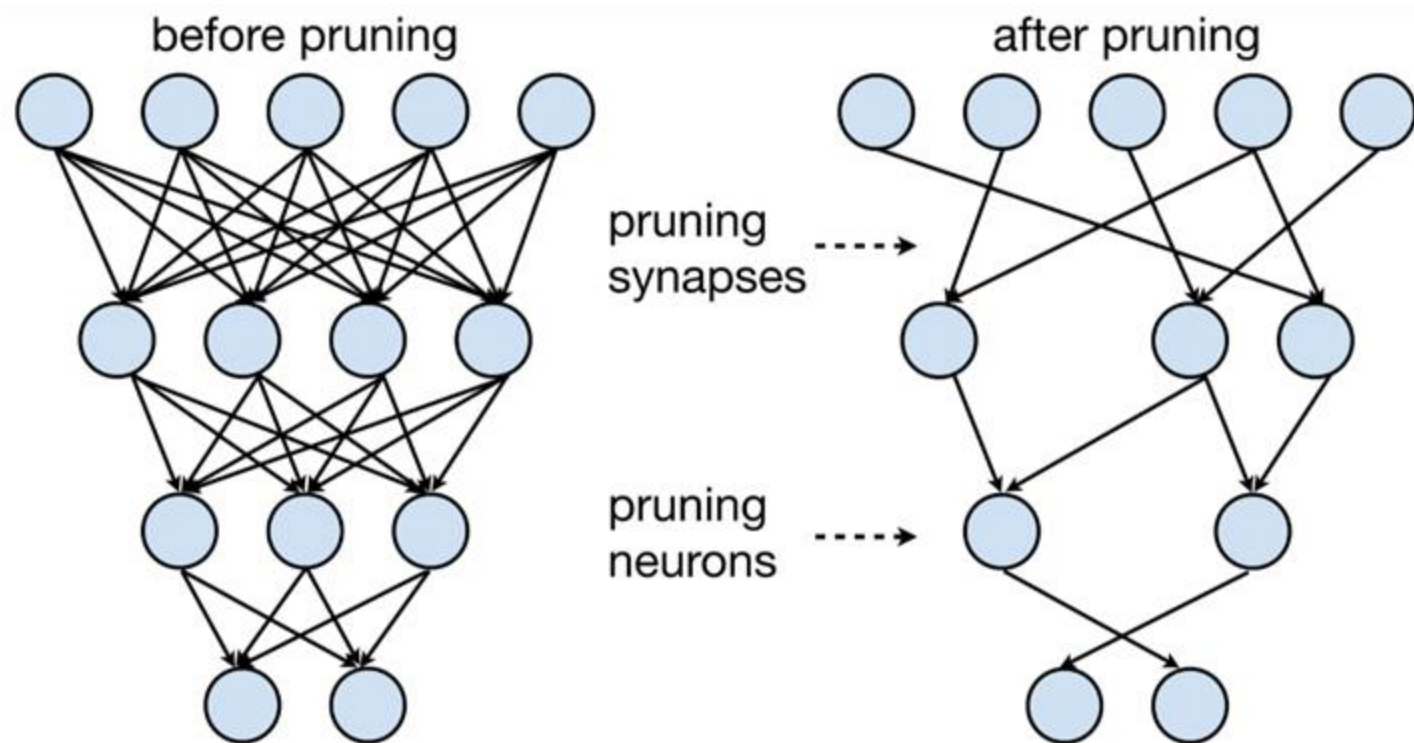
## □ **Sparsity (Mixture of Experts, Deja Vu: Contextual Sparsity)**

- Long Context (PagedAttention, StreamingLLM, MHA/GQA/MQA, H<sub>2</sub>O)
- Speculative Decoding
- Parameter Efficient Fine-Tuning (BitFit, Adapter, Prompt Tuning, LoRA)



# Sparsity

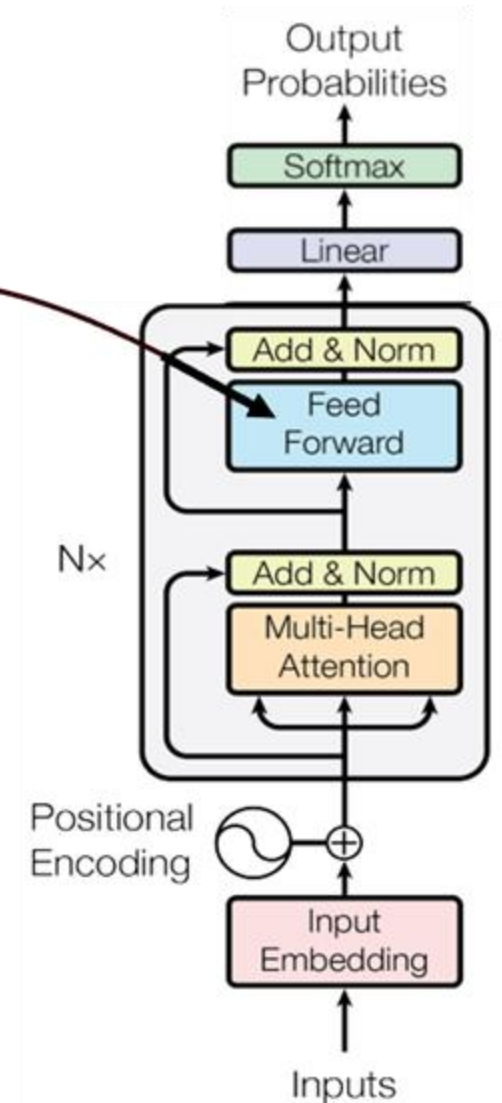
Even though our model may have many parameters, we can get speedups by only using a much smaller number of those parameters for a given instance



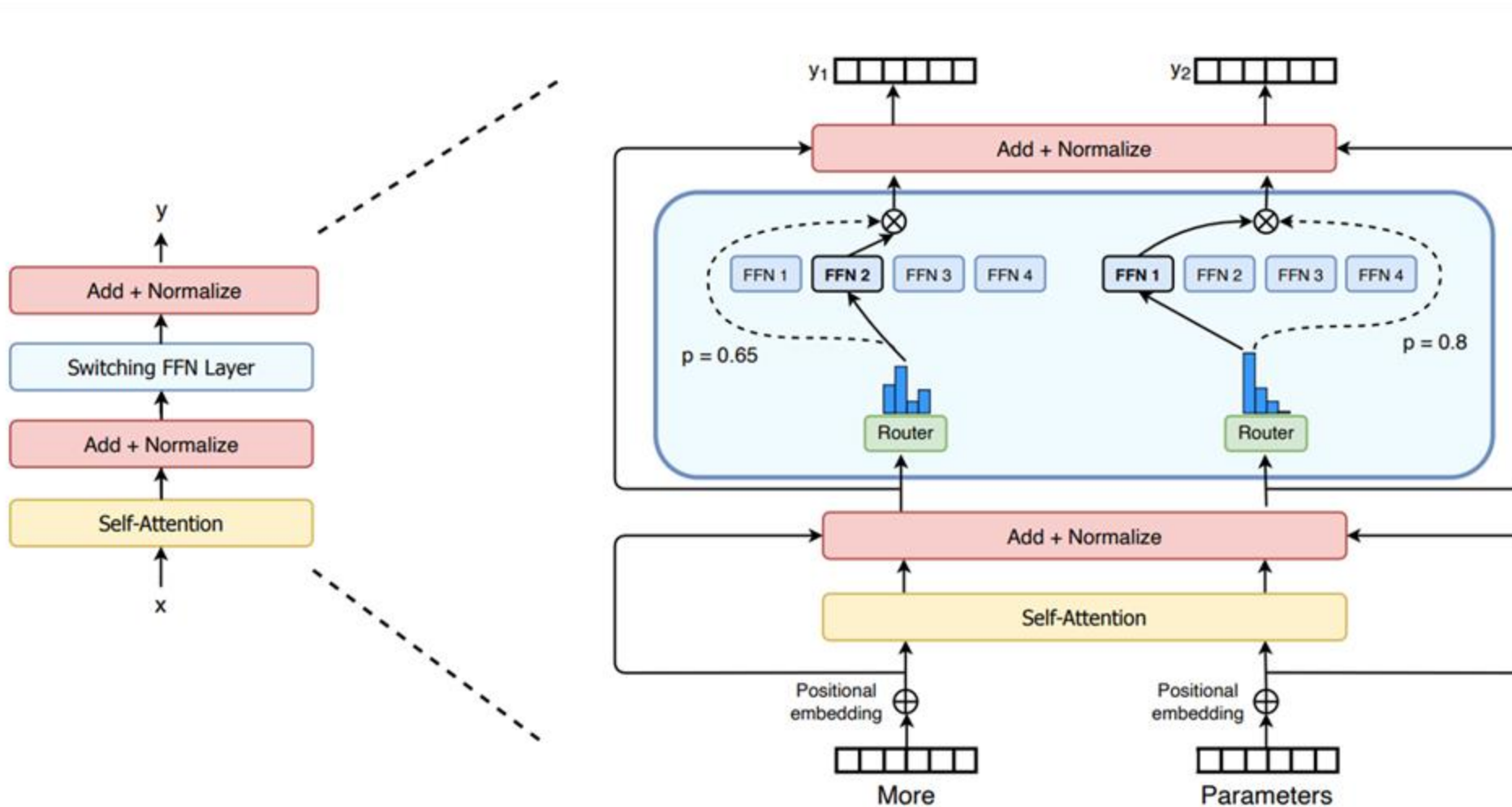


# Mixture of Experts (MoE)

Replace FFN layers in traditional transformers with a switching FFN layer (more generally called an MoE layer)



# Mixture of Experts (MoE)

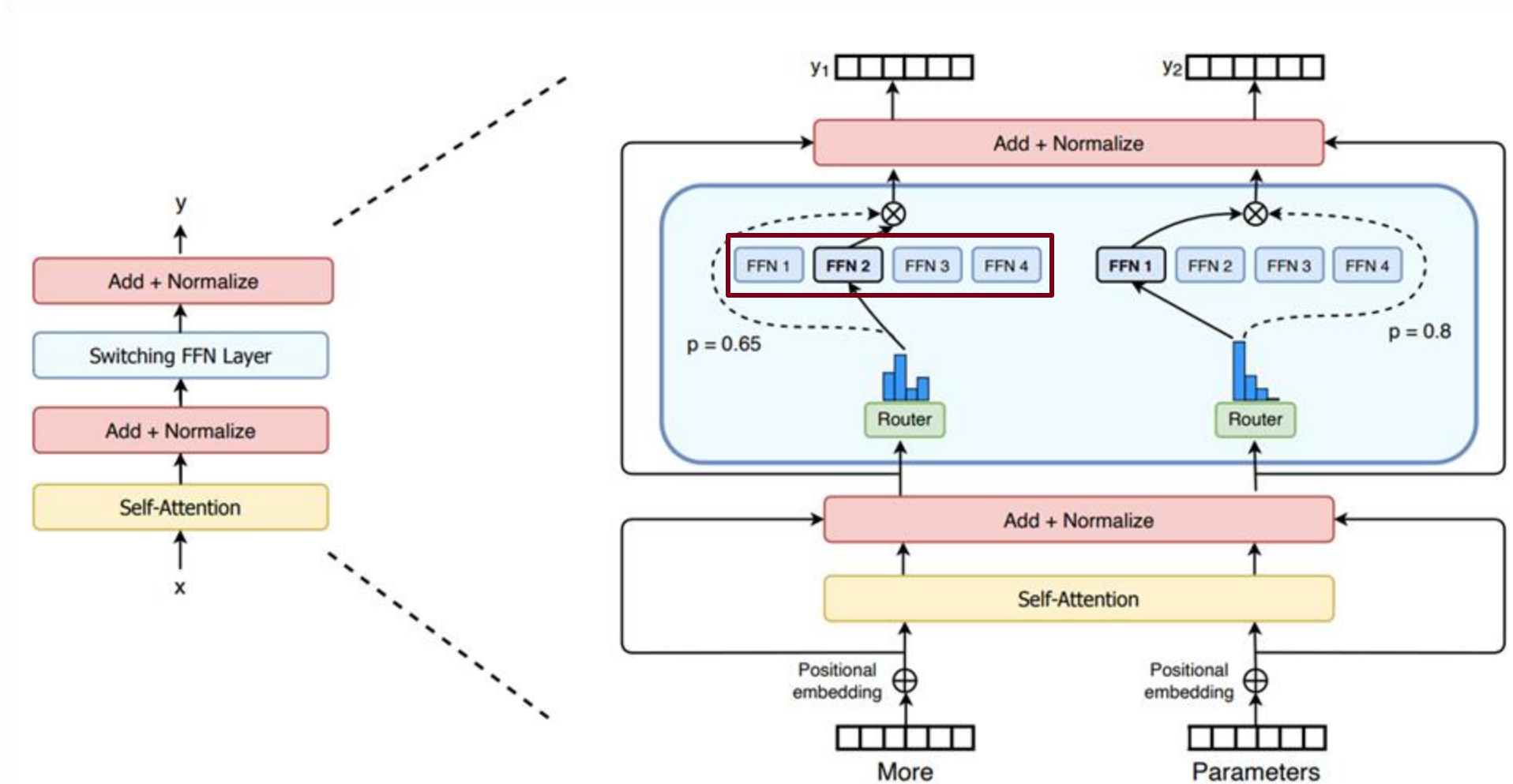


Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity [Fedus et al., CoRR 2021]



# Mixture of Experts (MoE)

Four FFN layers

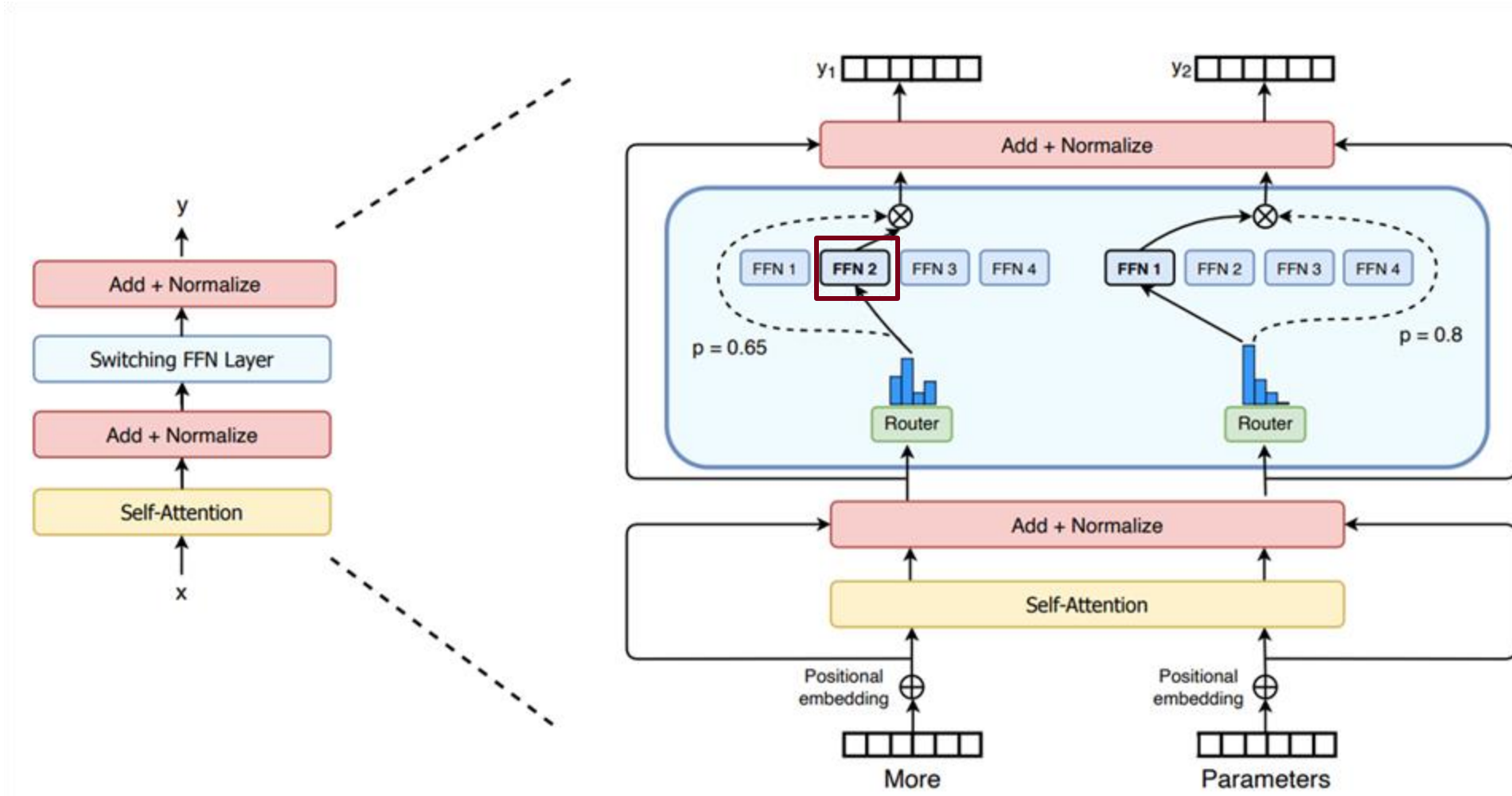


Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity [Fedus et al., CoRR 2021]



# Mixture of Experts (MoE)

Only one is used per token

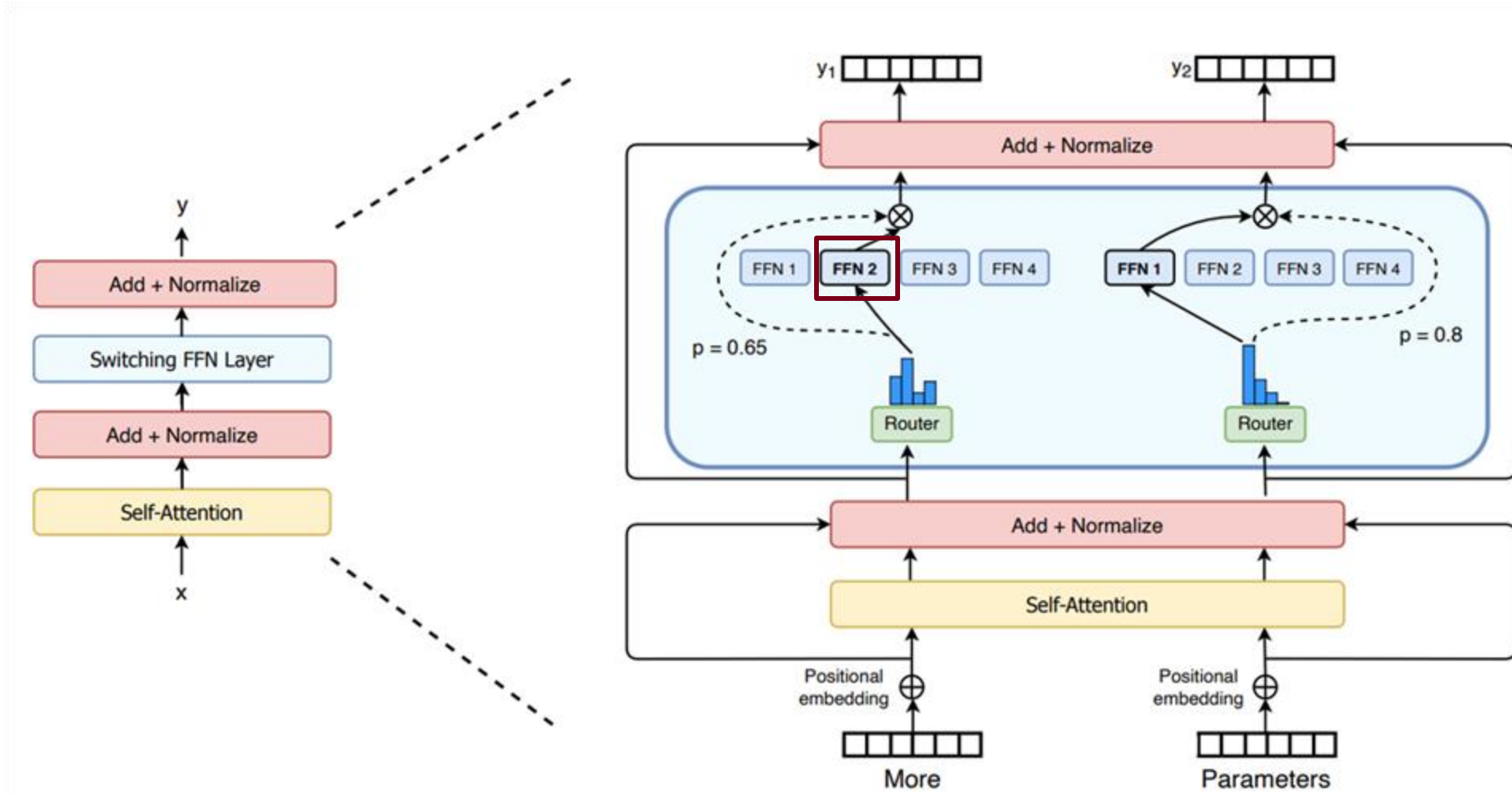


Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity [Fedus et al., CoRR 2021]



# Mixture of Experts (MoE)

Only 25% of the FFN parameters are used for a single token

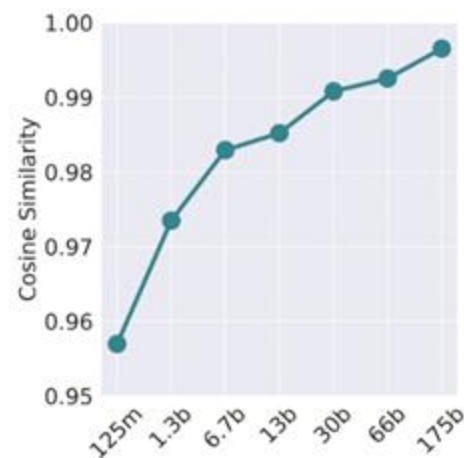


Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity [Fedus et al., CoRR 2021]

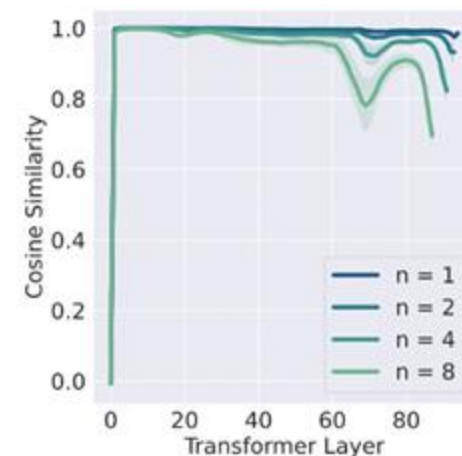


# Deja Vu: Contextual Sparsity

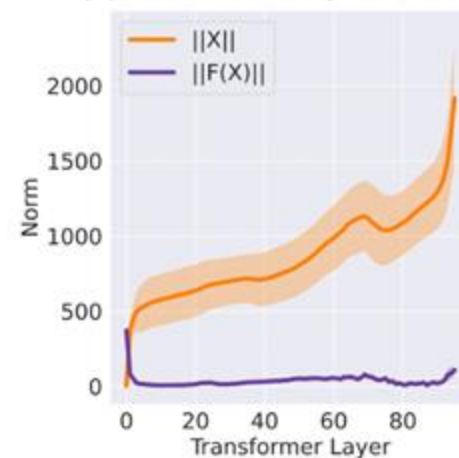
Observation 1: Model activations change very little between consecutive layers of a network



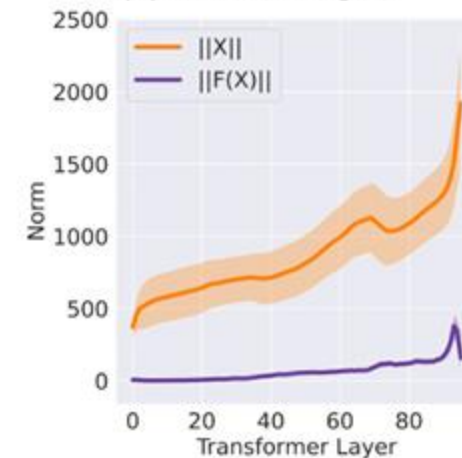
(a) Model Comparison



(b) Across Layer



(c) Residual Around Attention



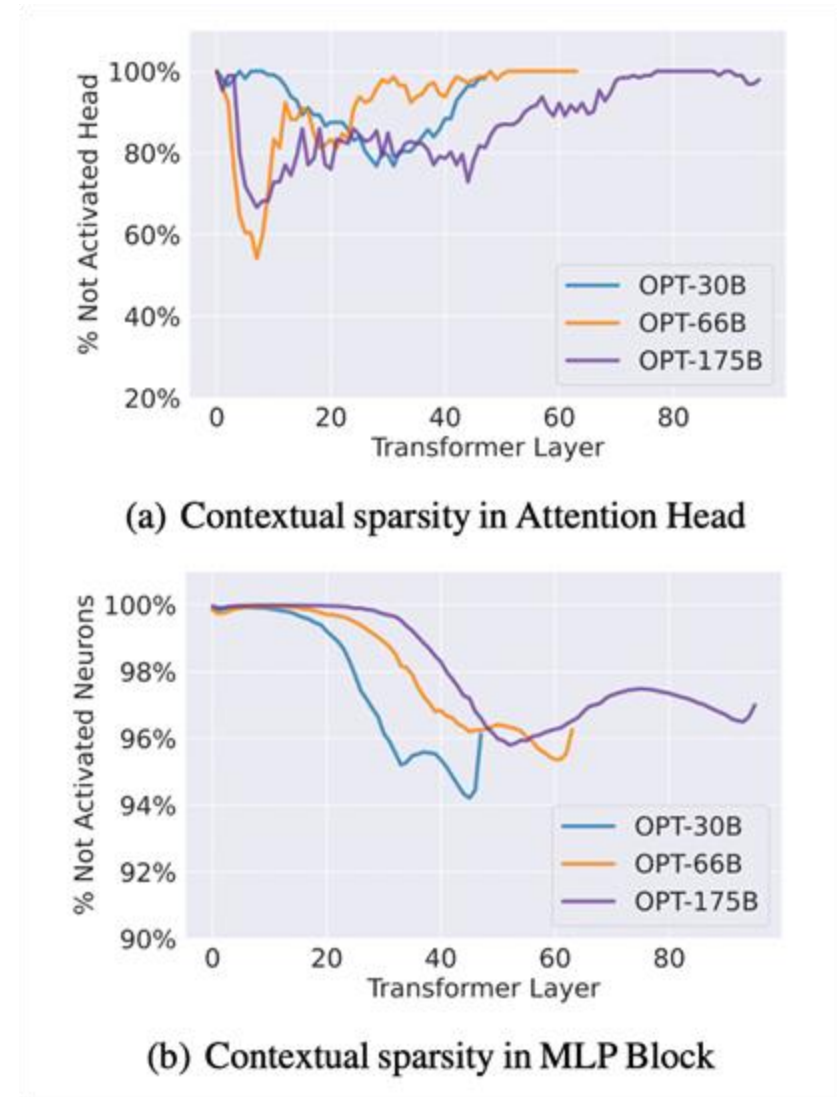
(d) Residual Around MLP

Deja Vu: Contextual Sparsity for Efficient LLMs at Inference Time [Liu et al., 2023]



# Deja Vu: Contextual Sparsity

Observation 2: Most attention heads and most neurons are not used

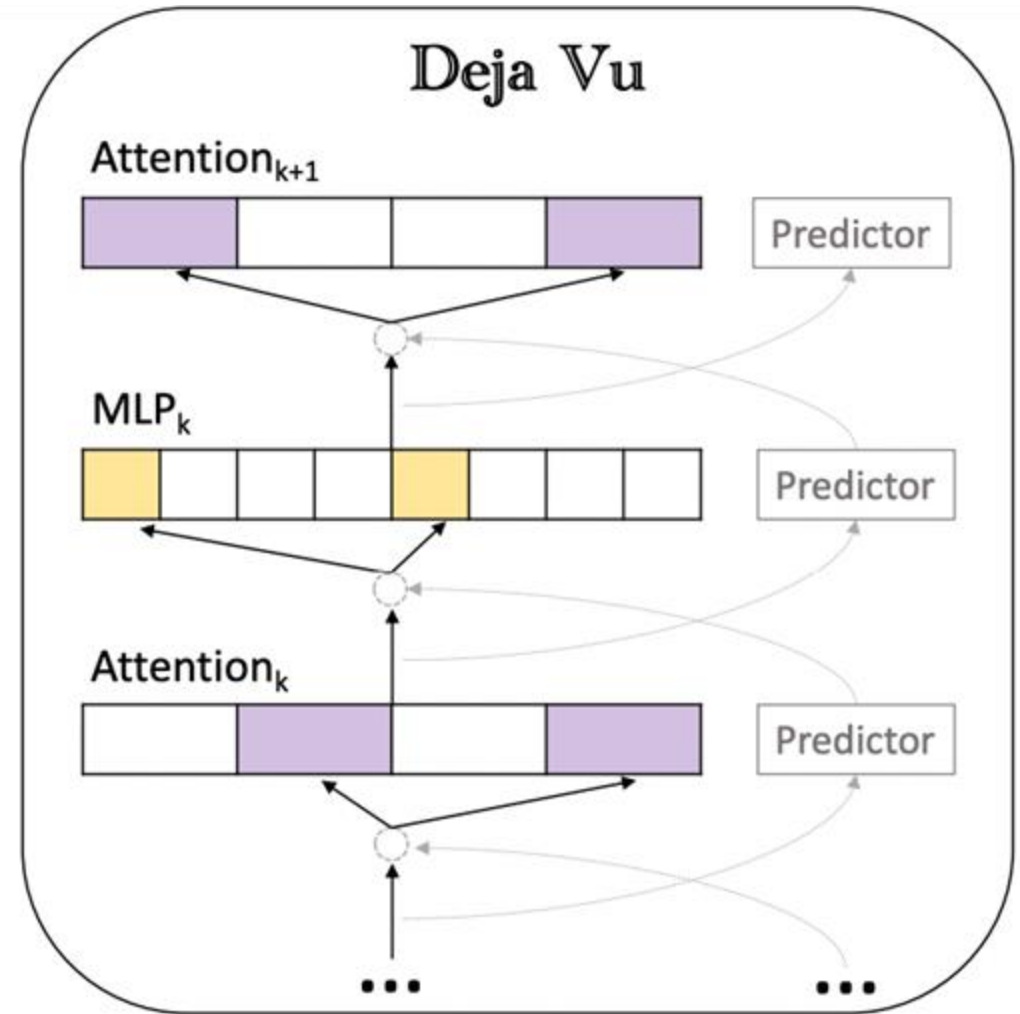


Deja Vu: Contextual Sparsity for Efficient LLMs at Inference Time [Liu et al., 2023]



# Deja Vu: Contextual Sparsity

Sparsification: Use predictors in each layer to determine which neurons to activate and which attention heads to use – ignore all unpredicted heads/neurons



Deja Vu: Contextual Sparsity for Efficient LLMs at Inference Time [Liu et al., 2023]



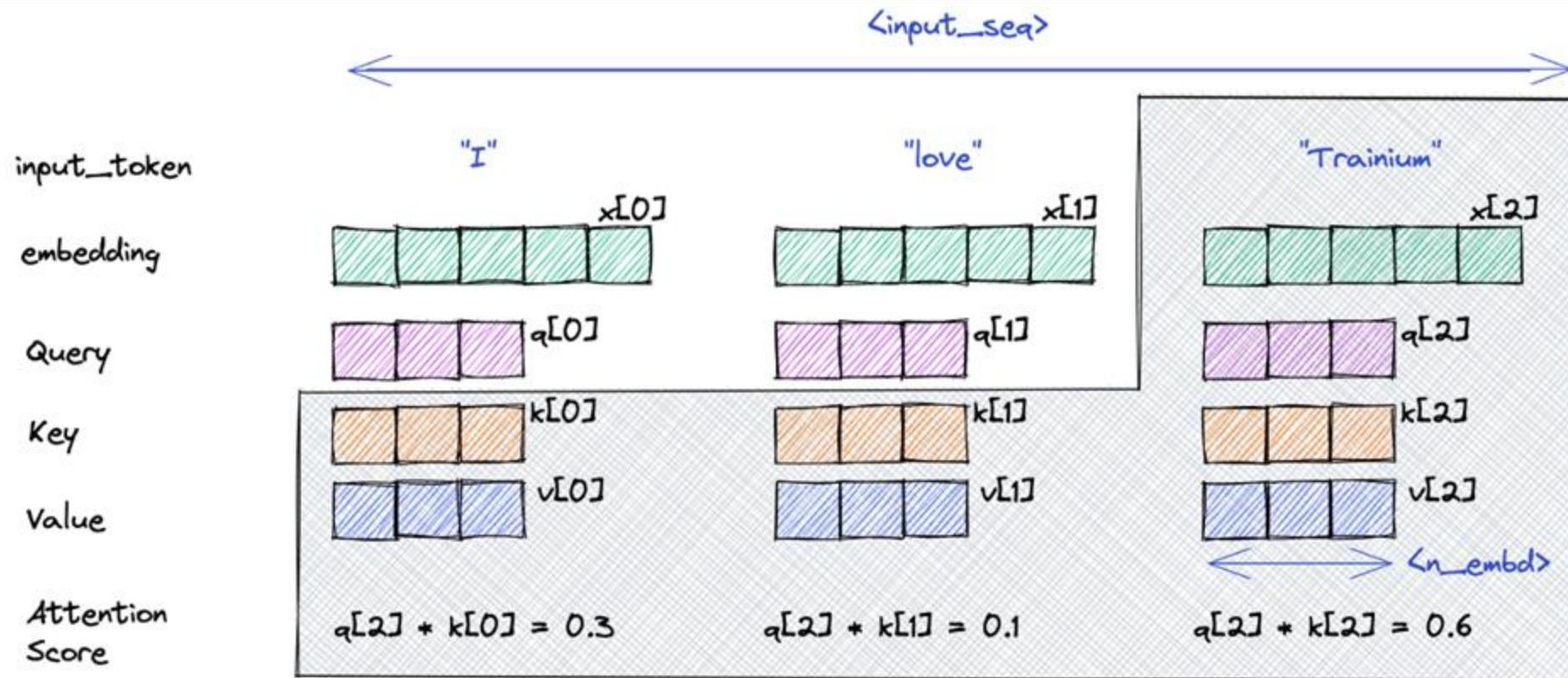
# Efficient LLMs

- ❑ Quantization
  - Background
  - K-Means vs. Linear Quantization
  - Quantization Granularity
  - Quantization Aware Training (QAT) vs Post-Training Quantization (PTQ)
  - LLM Quantization (LLM.int8(), SmoothQuant, AWQ, 1-bit LLMs)
- ❑ Sparsity (Mixture of Experts, Deja Vu: Contextual Sparsity)
- ❑ **Long Context (PagedAttention, StreamingLLM, MHA/GQA/MQA, H<sub>2</sub>O)**
- ❑ Speculative Decoding
- ❑ Parameter Efficient Fine-Tuning (BitFit, Adapter, Prompt Tuning, LoRA)



# The KV-Cache

The transformer needs to have access to the keys and values for all previous tokens in all layers for all heads when



<https://awsdocs-neuron.readthedocs-hosted.com/en/latest/general/appnotes/transformers-neuronx/generative-llm-inference-with-neuron.html>



# The KV-Cache

In total, we must store

$$\text{Batch\_size} * \text{seq\_len} * \text{num\_heads} * \text{num\_layers} * \text{emb\_dim} * 2$$

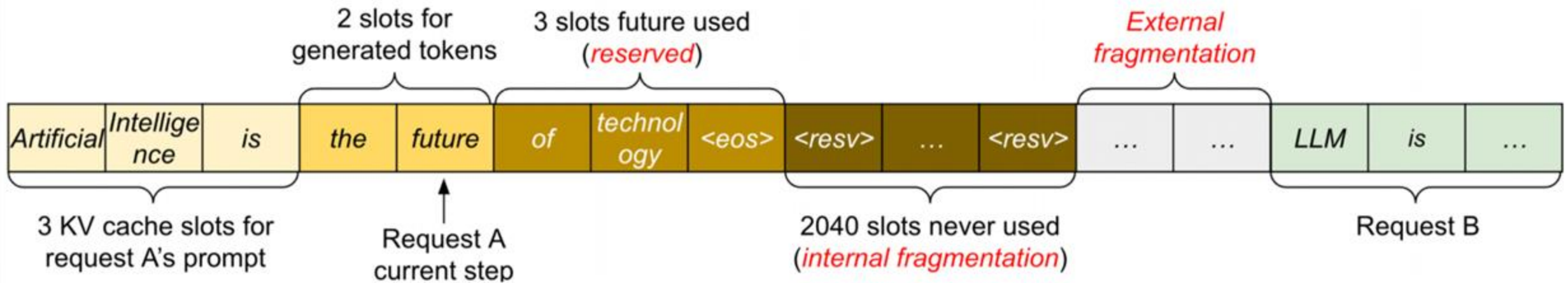
separate values in the kv cache



# PagedAttention

How does a large LLM service (large ChatGPT) handle multiple incoming requests with respect to the KV-cache?

-Originally, most systems just assign fixed sized blocks of memory to each incoming request. How to improve?



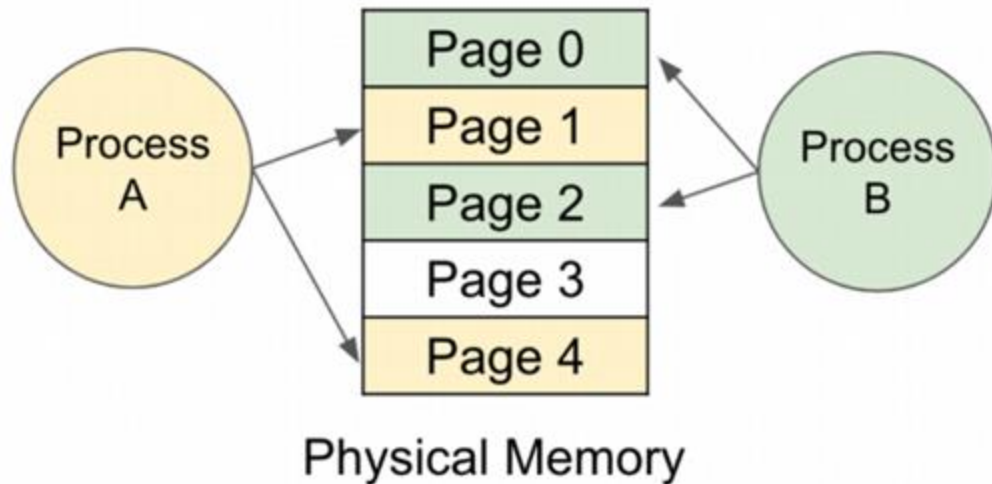
Efficient Memory Management for Large Language Model Serving with PagedAttention (Kwon et al., 2023)



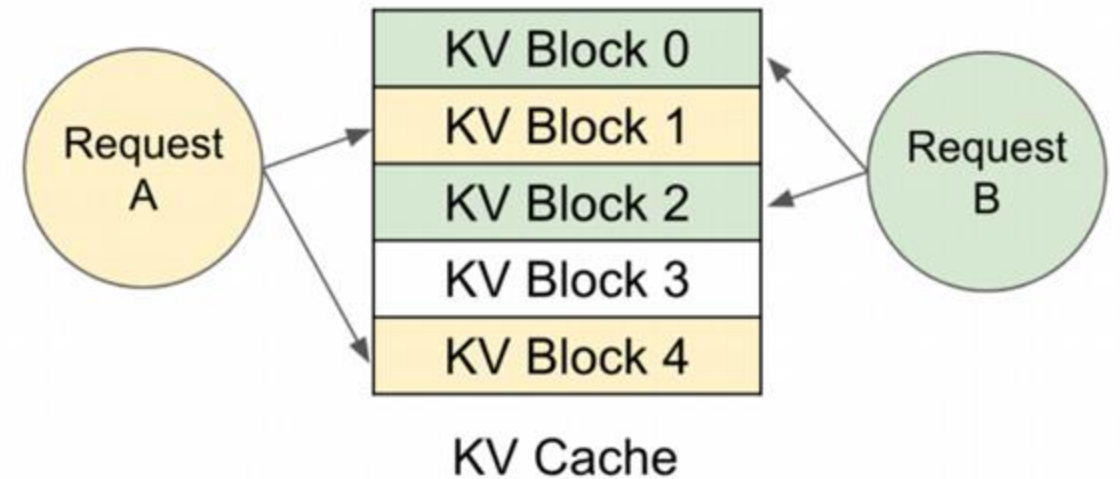
# PagedAttention

Let's adopt a similar approach to that found in virtual memory!

## Memory management in OS



## Memory management in vLLM



Efficient Memory Management for Large Language Model Serving with PagedAttention (Kwon et al., 2023)



# PagedAttention



Block Table

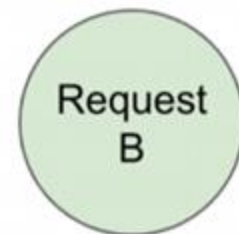

Logical KV blocks

Alan	Turing	is	a
computer	scientist	and	mathematician
renowned			

Physical KV blocks

computer	scientist	and	mathematician
Artificial	Intelligence	is	the
renowned			
future	of	technology	
Alan	Turing	is	a

Block Table

Logical KV blocks

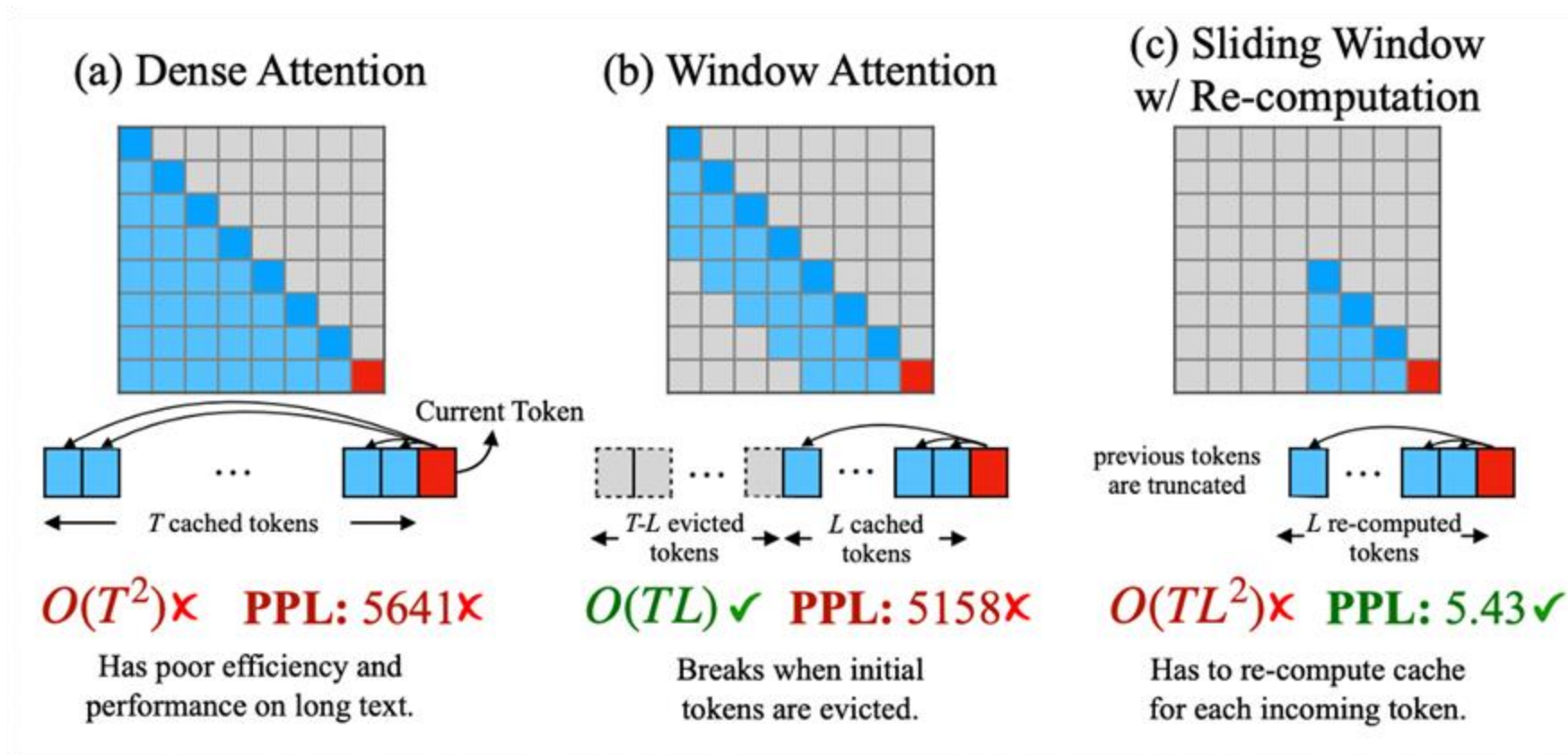
Artificial	Intelligence	is	the
future	of	technology	

Efficient Memory Management for Large Language Model Serving with PagedAttention (Kwon et al., 2023)



# StreamingLLM

How can we extend models to have much longer context length at minimal cost?

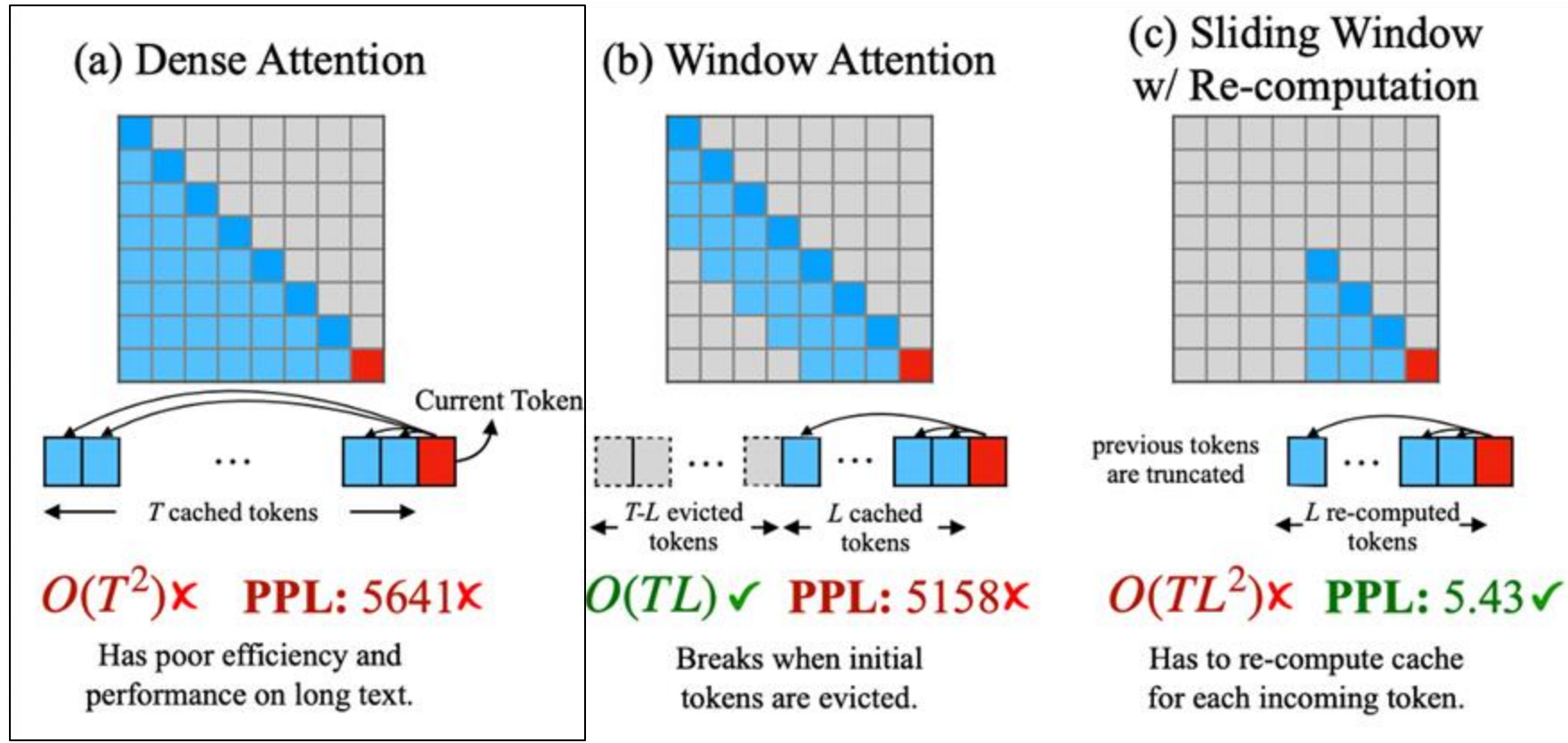


Efficient Streaming Language Models with Attention Sinks [Xiao et al., 2023]

# StreamingLLM

How can we extend models to have much longer context length at minimal cost?

Too much storage



Efficient Streaming Language Models with Attention Sinks [Xiao et al., 2023]

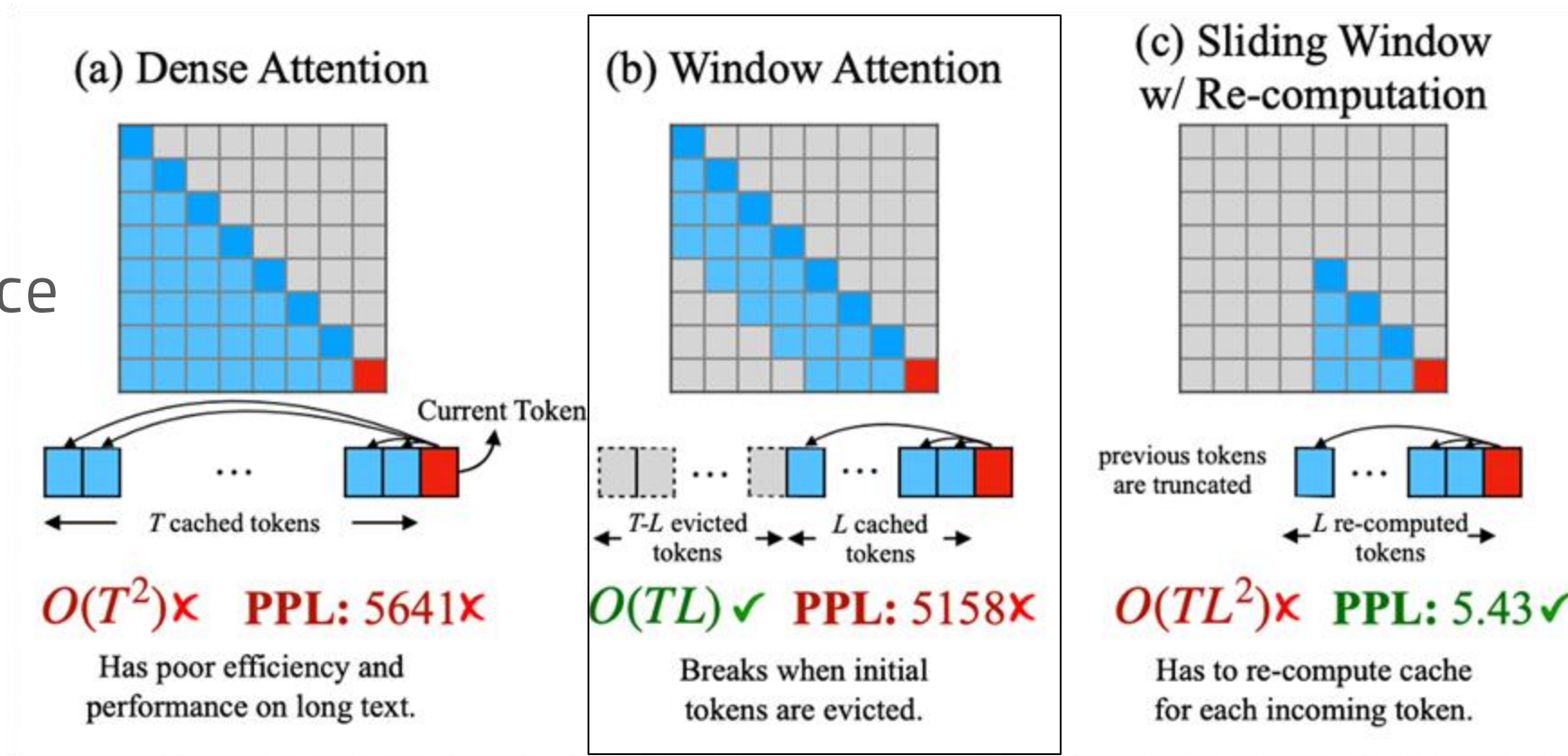




# StreamingLLM

How can we extend models to have much longer context length at minimal cost?

Bad performance



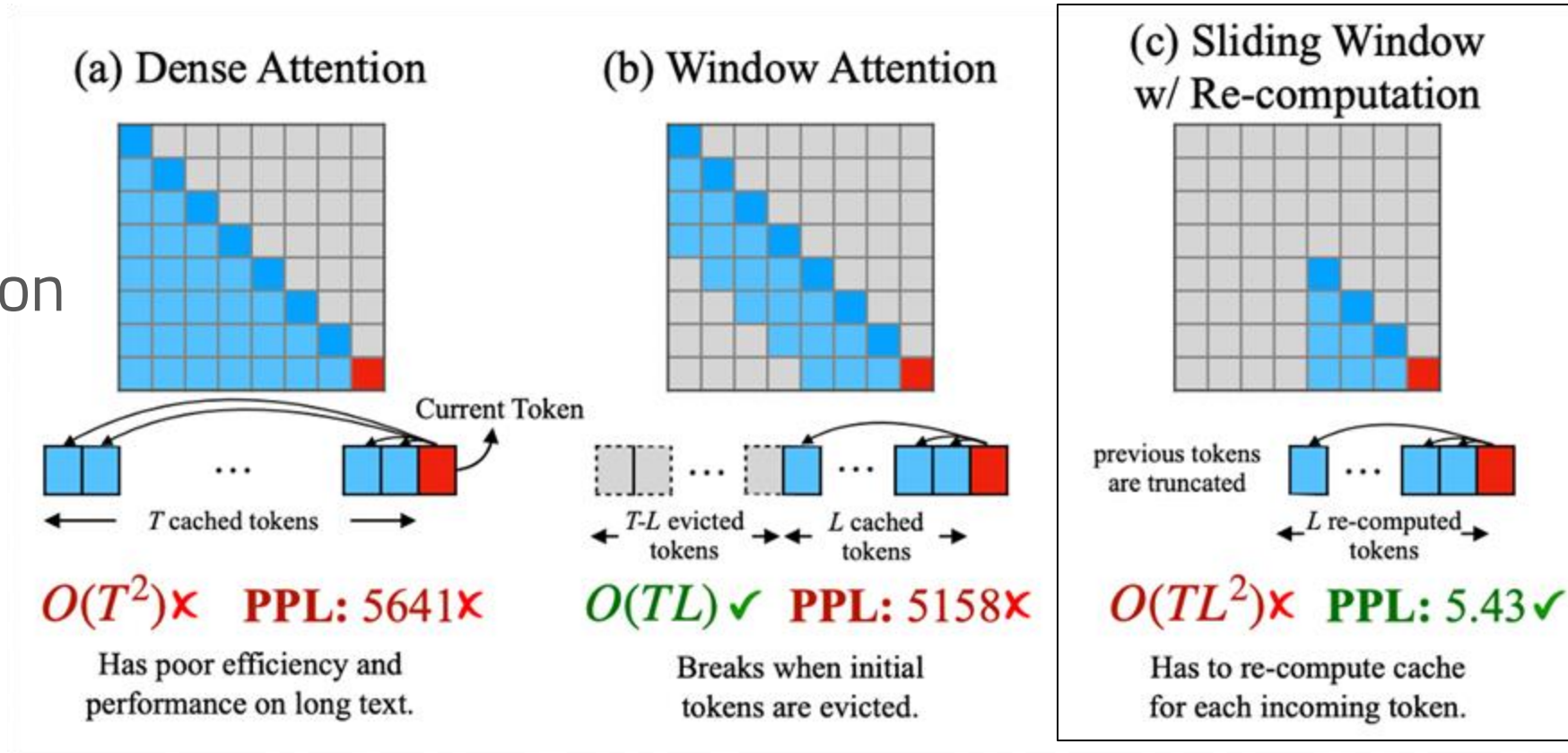
Efficient Streaming Language Models with Attention Sinks [Xiao et al., 2023]



# StreamingLLM

How can we extend models to have much longer context length at minimal cost?

Too much  
recomputation

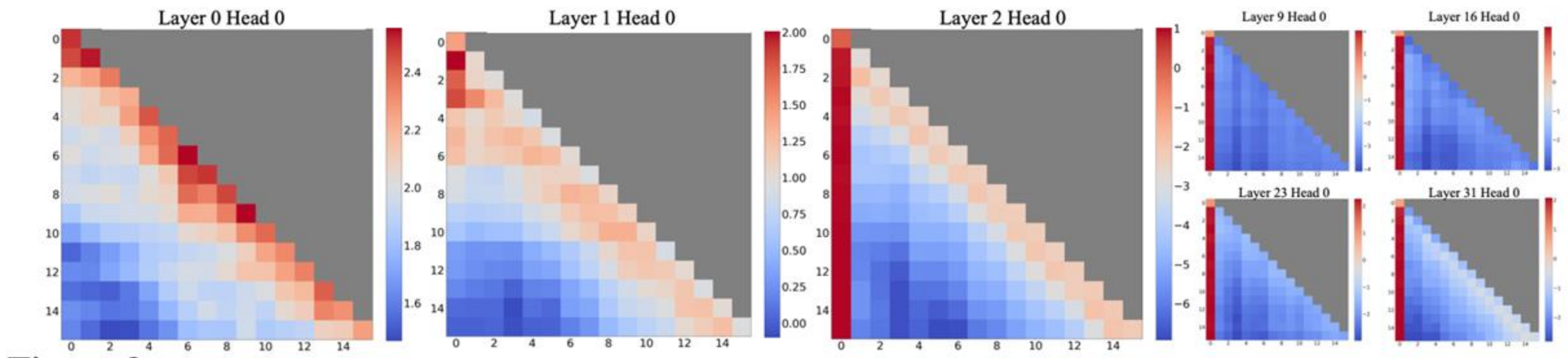


Efficient Streaming Language Models with Attention Sinks [Xiao et al., 2023]



# StreamingLLM

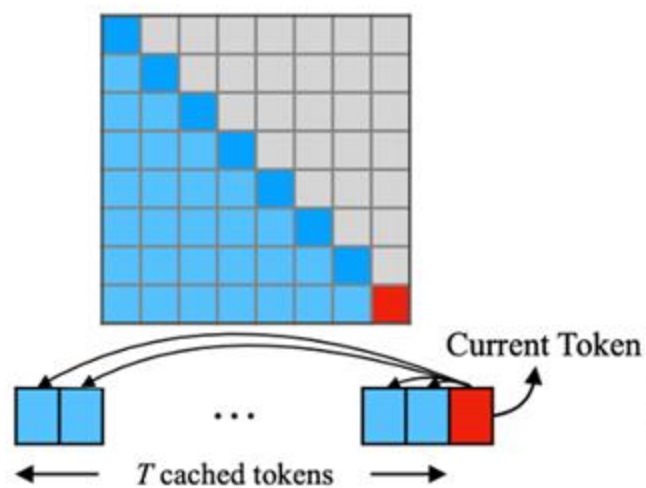
Observation: Most attention is either placed on the first token or to tokens that the model has recently seen.



Efficient Streaming Language Models with Attention Sinks [Xiao et al., 2023]

# StreamingLLM

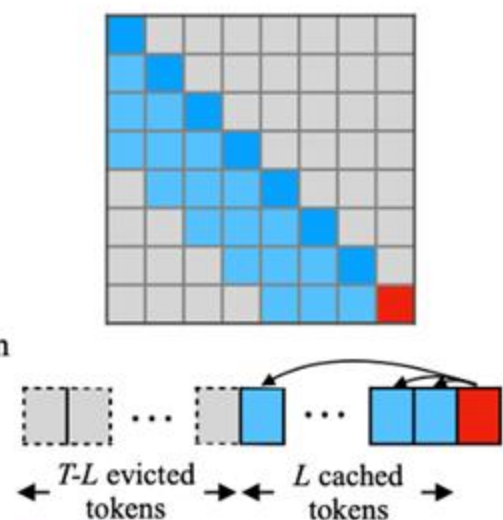
(a) Dense Attention



$O(T^2)$  ✗ PPL: 5641 ✗

Has poor efficiency and performance on long text.

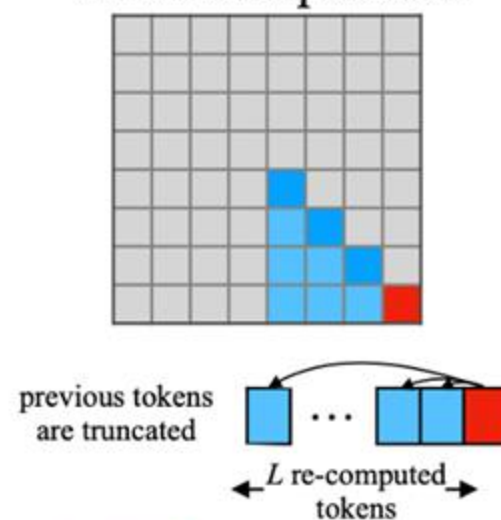
(b) Window Attention



$O(TL)$  ✓ PPL: 5158 ✗

Breaks when initial tokens are evicted.

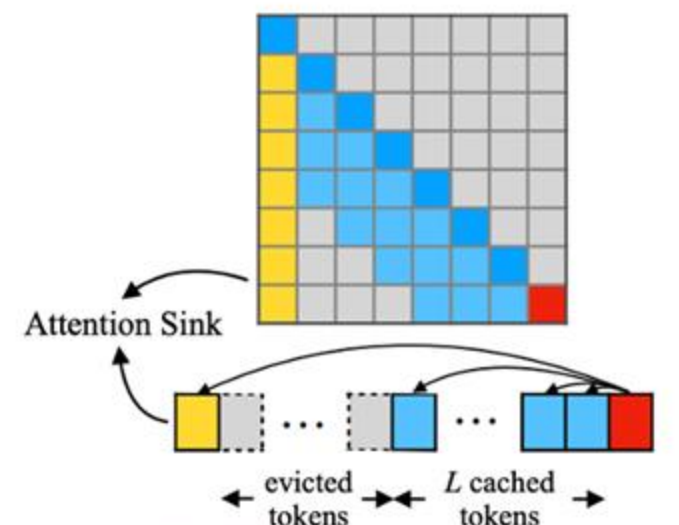
(c) Sliding Window w/ Re-computation



$O(TL^2)$  ✗ PPL: 5.43 ✓

Has to re-compute cache for each incoming token.

(d) StreamingLLM (ours)



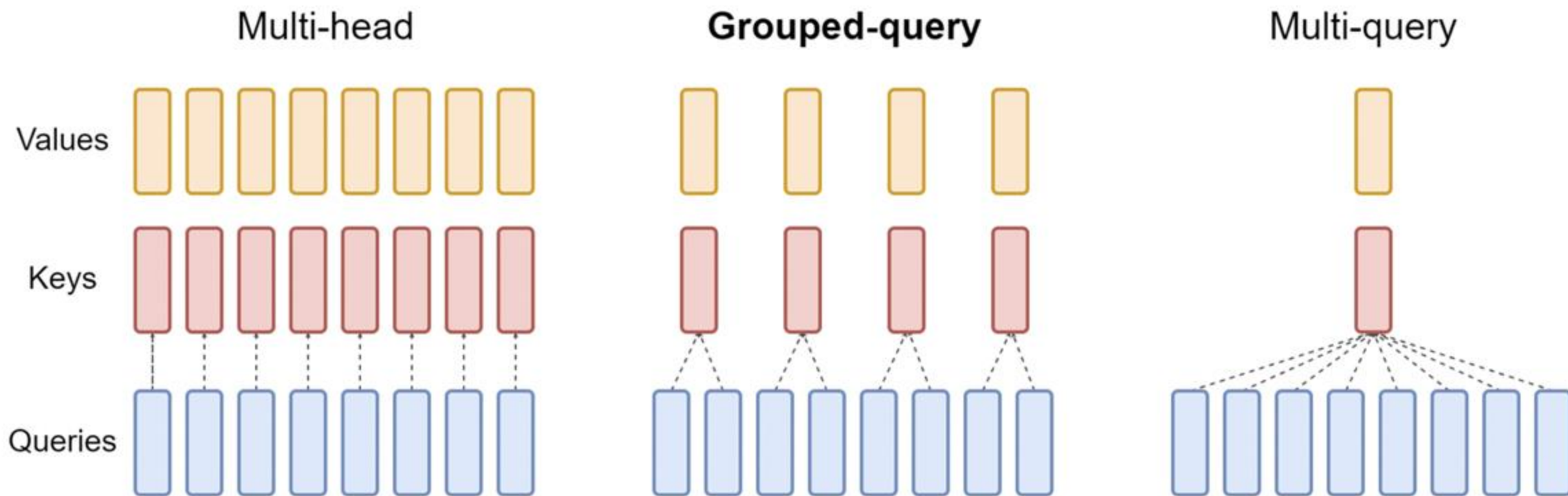
$O(TL)$  ✓ PPL: 5.40 ✓

Can perform efficient and stable language modeling on long texts.

Efficient Streaming Language Models with Attention Sinks [Xiao et al., 2023]



# MHA/GQA/MQA

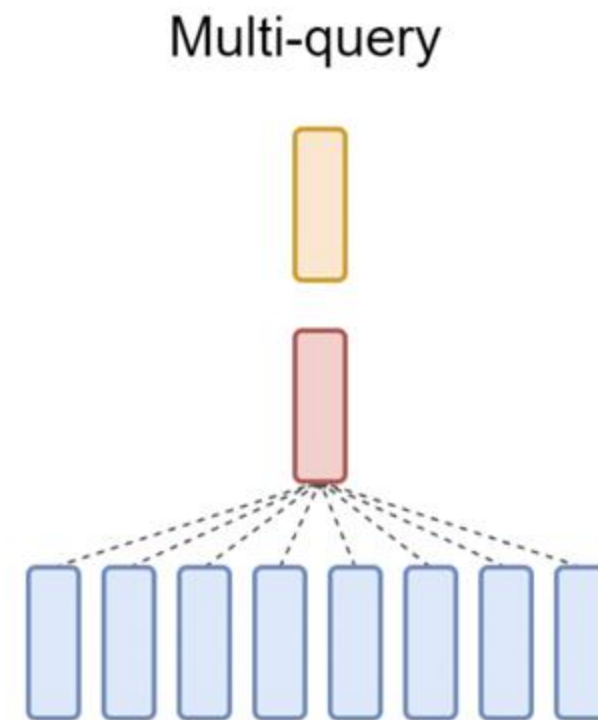
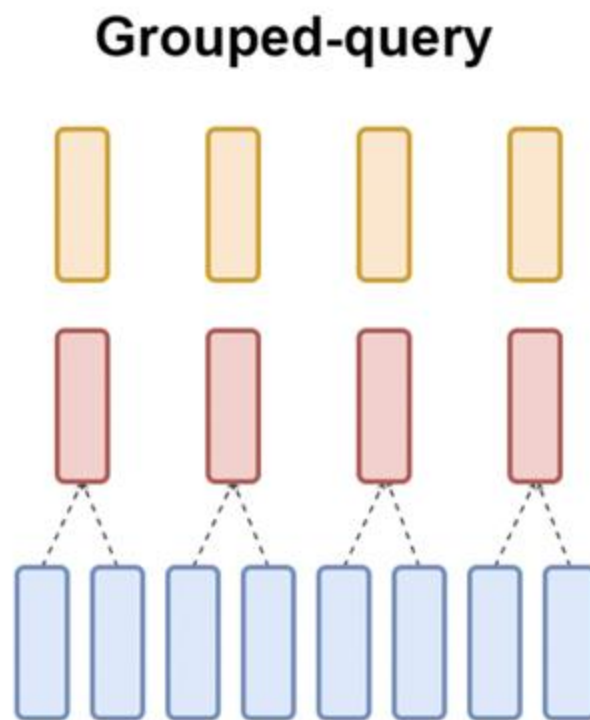
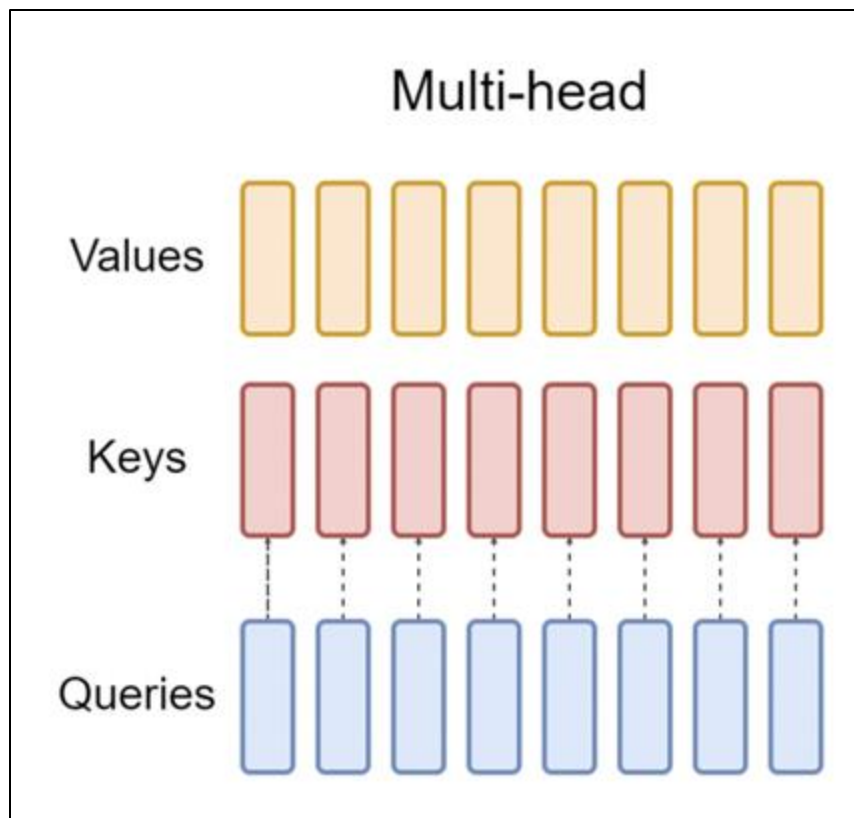


GQA: Training Generalized Multi-Query Transformer Models from Multi-Head Checkpoints



# MHA/GQA/MQA

Each attention head calculates separate keys and values for each token

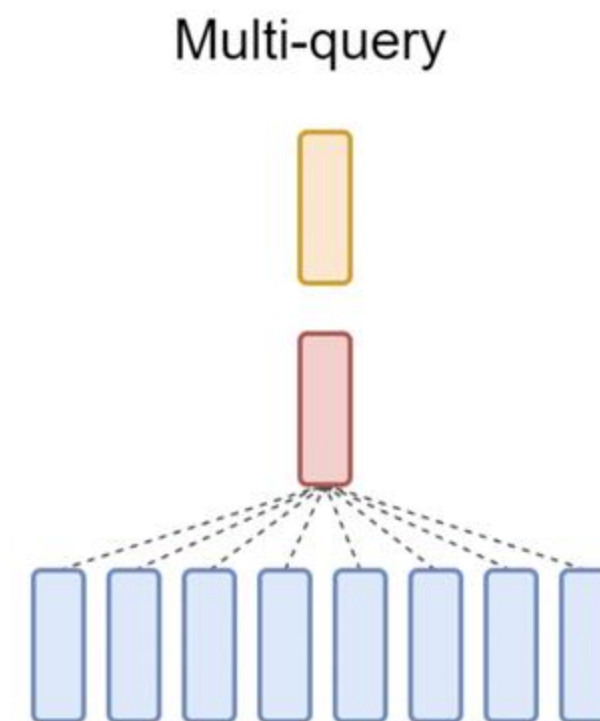
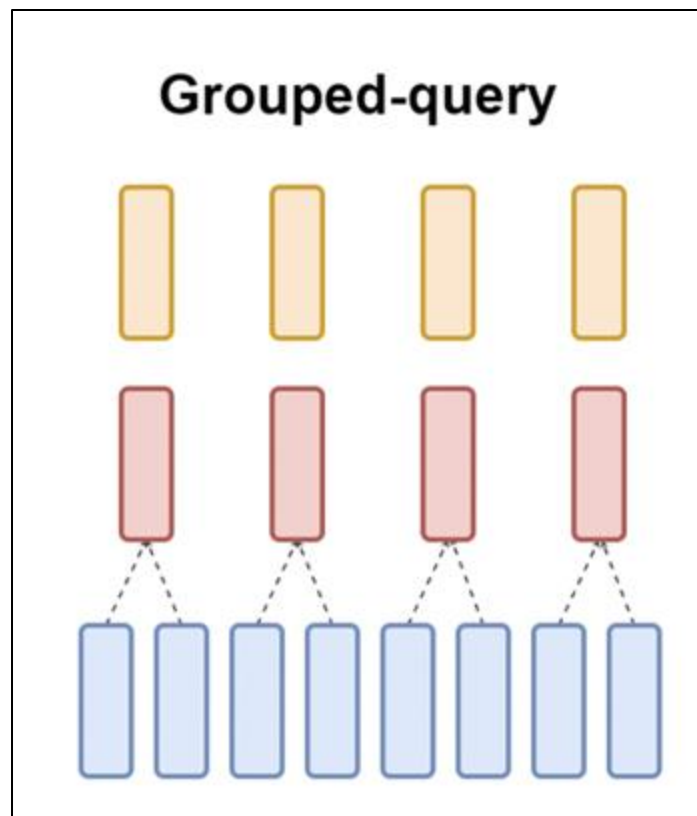
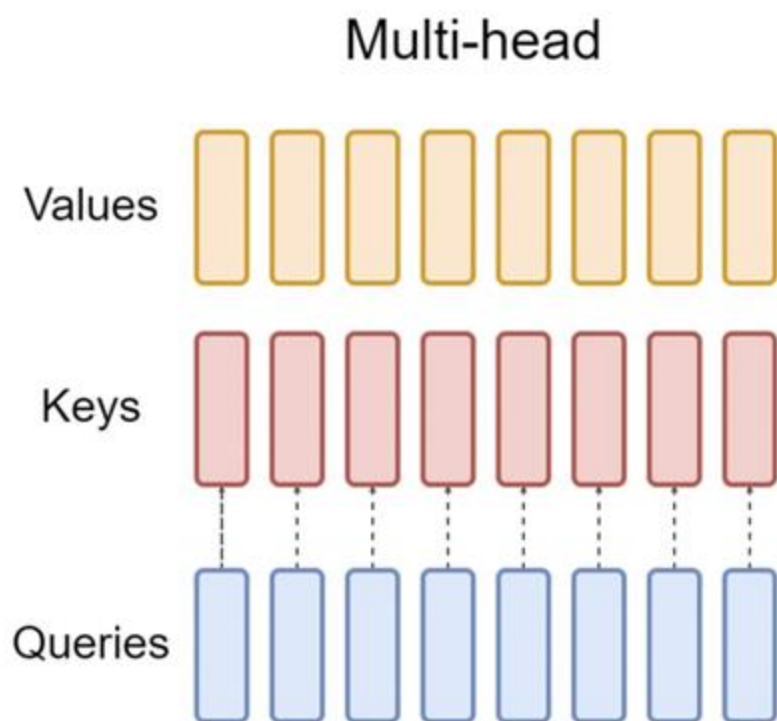


GQA: Training Generalized Multi-Query Transformer Models from Multi-Head Checkpoints



# MHA/GQA/MQA

Attention heads are split into groups.  
Each group has one key/value per token.

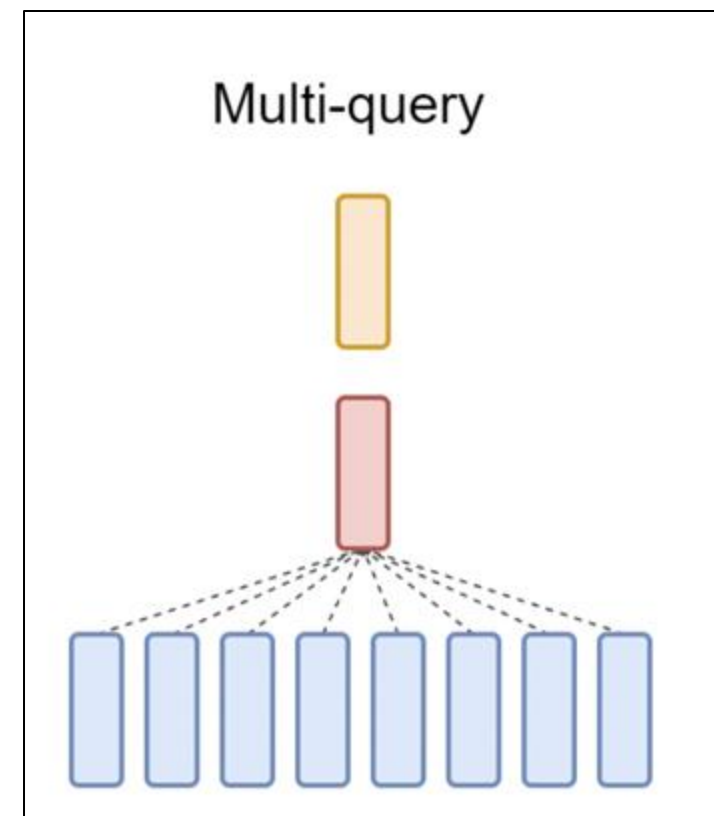
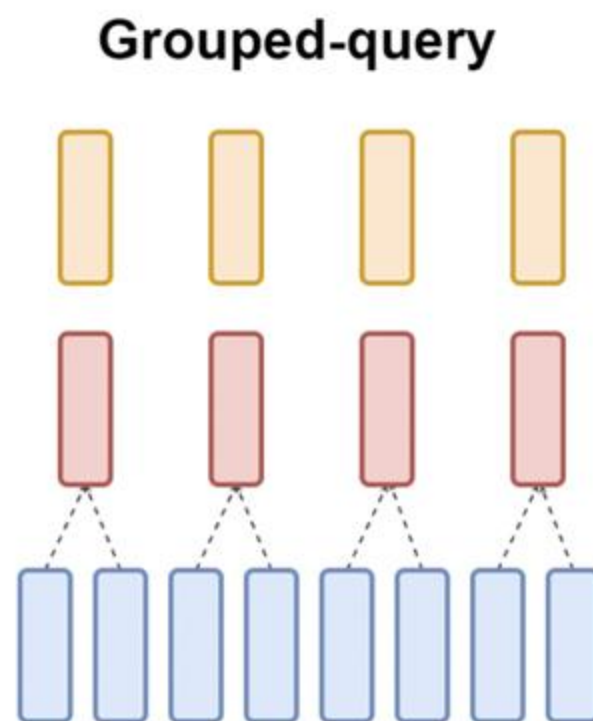
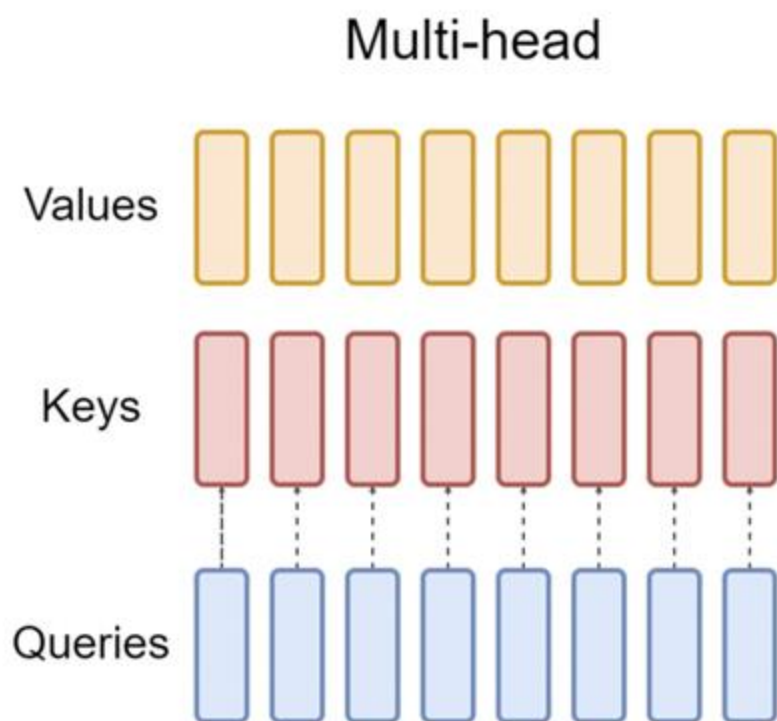


GQA: Training Generalized Multi-Query Transformer Models from Multi-Head Checkpoints



# MHA/GQA/MQA

Attention heads share the same keys and values for each token

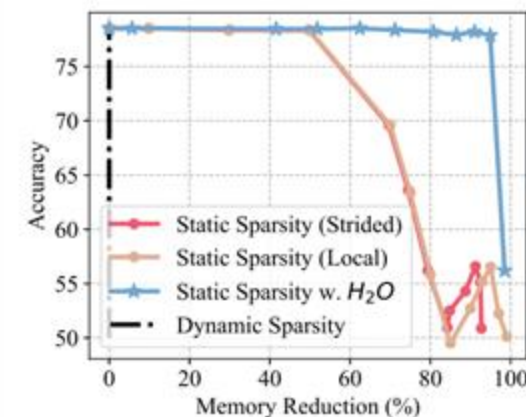
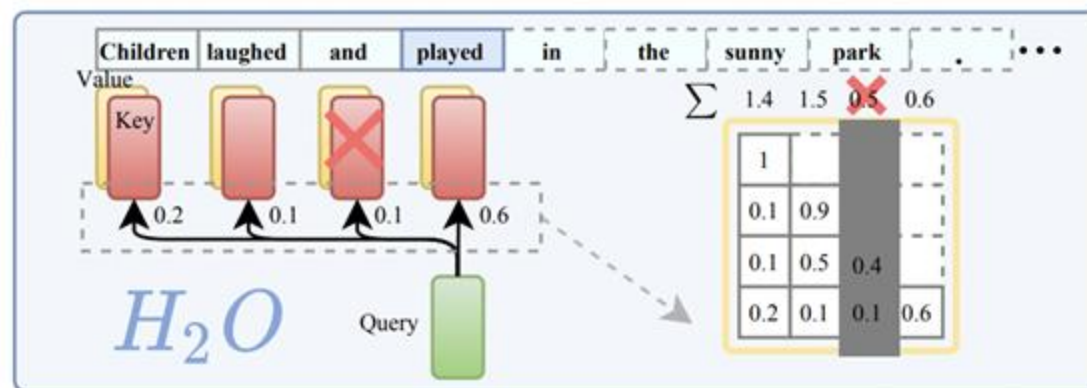
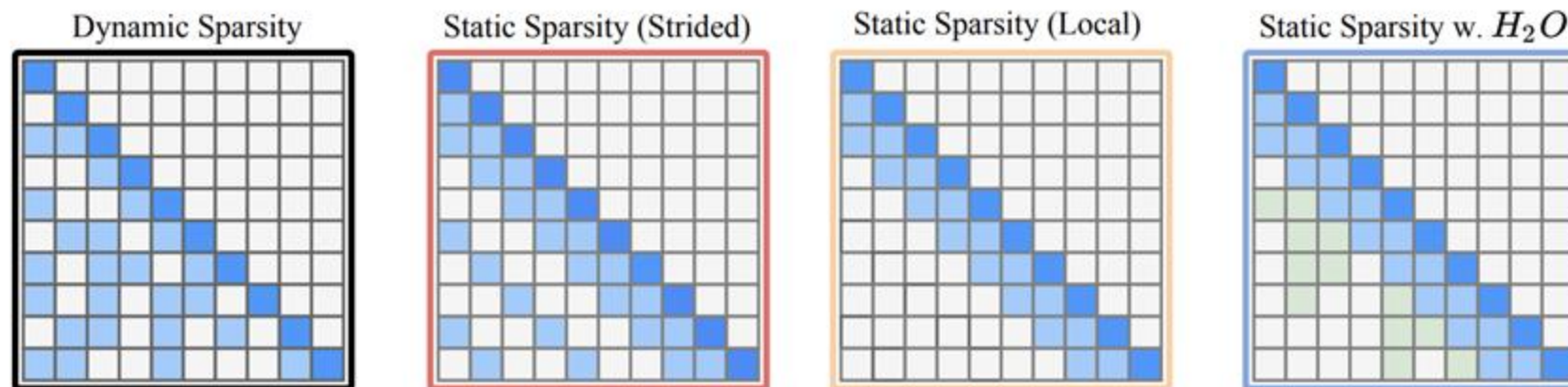


GQA: Training Generalized Multi-Query Transformer Models from Multi-Head Checkpoints



# H<sub>2</sub>O: Heavy Hitter Oracle

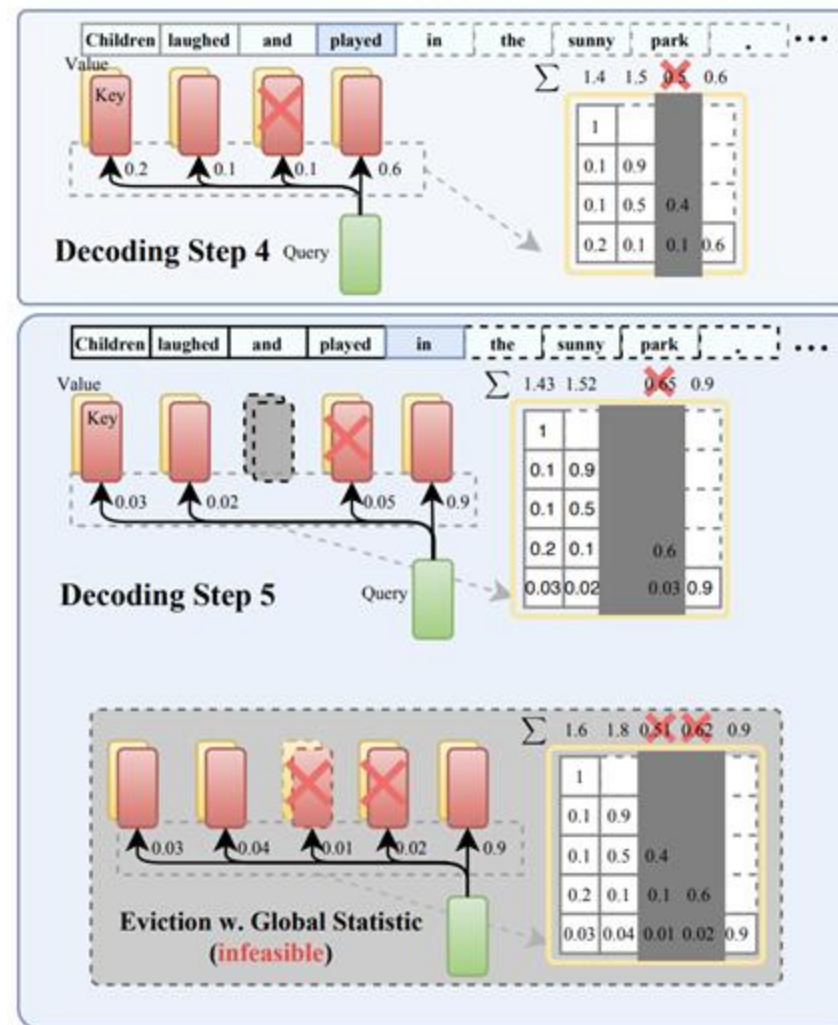
Evict all but k-  
highest  
cumulative  
attention tokens  
from cache



H2O: Heavy-Hitter Oracle for Efficient Generative Inference of Large Language Models [Zhang et. al, 2023]

# H<sub>2</sub>O: Heavy Hitter Oracle

Evict all but k-  
highest  
cumulative  
attention tokens  
from cache



H2O: Heavy-Hitter Oracle for Efficient Generative Inference of Large Language Models [Zhang et. al, 2023]



# Efficient LLMs

## □ Quantization

- Background
- K-Means vs. Linear Quantization
- Quantization Granularity
- Quantization Aware Training (QAT) vs Post-Training Quantization (PTQ)
- LLM Quantization (LLM.int8(), SmoothQuant, AWQ, 1-bit LLMs)

## □ Sparsity (Mixture of Experts, Deja Vu: Contextual Sparsity)

## □ Long Context (PagedAttention, StreamingLLM, MHA/GQA/MQA, H<sub>2</sub>O)

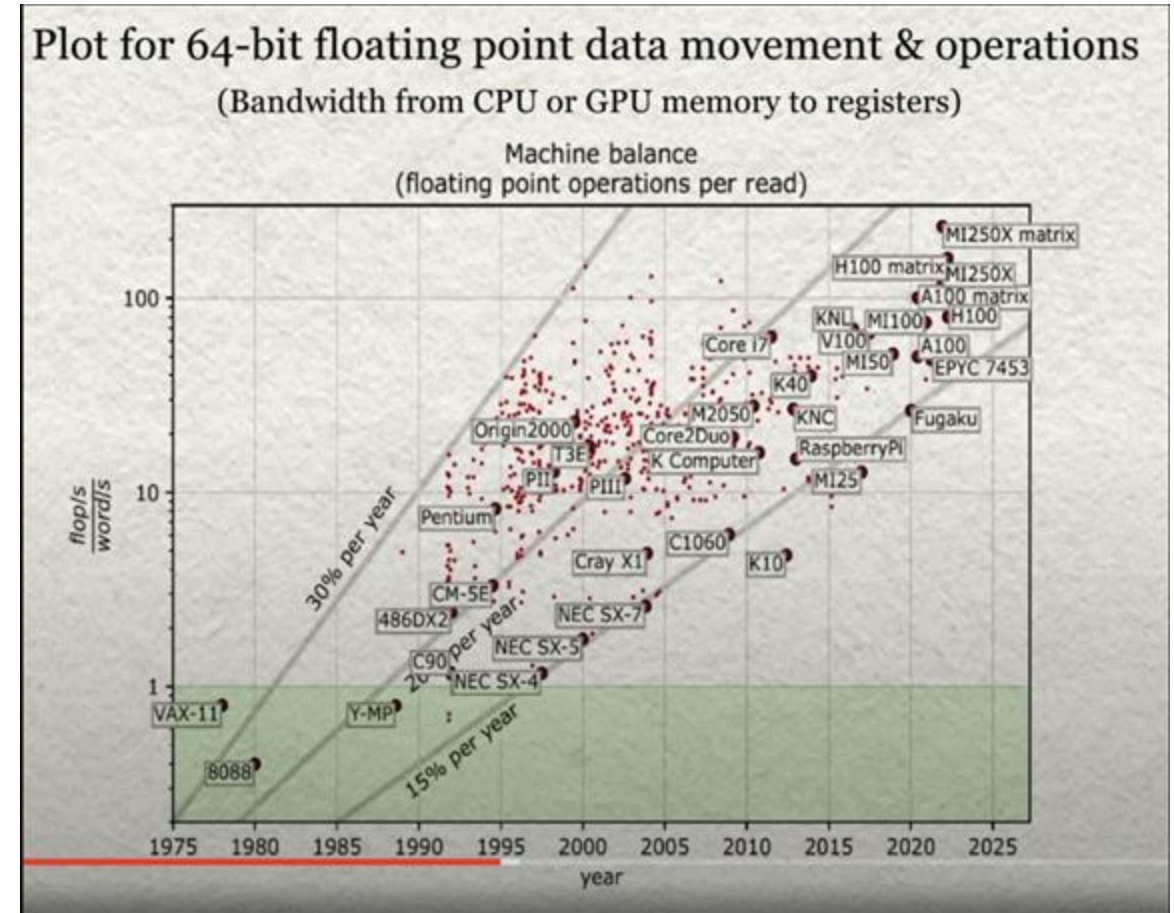
## □ **Speculative Decoding**

## □ Parameter Efficient Fine-Tuning (BitFit, Adapter, Prompt Tuning, LoRA)



# The Memory Bandwidth Bottleneck

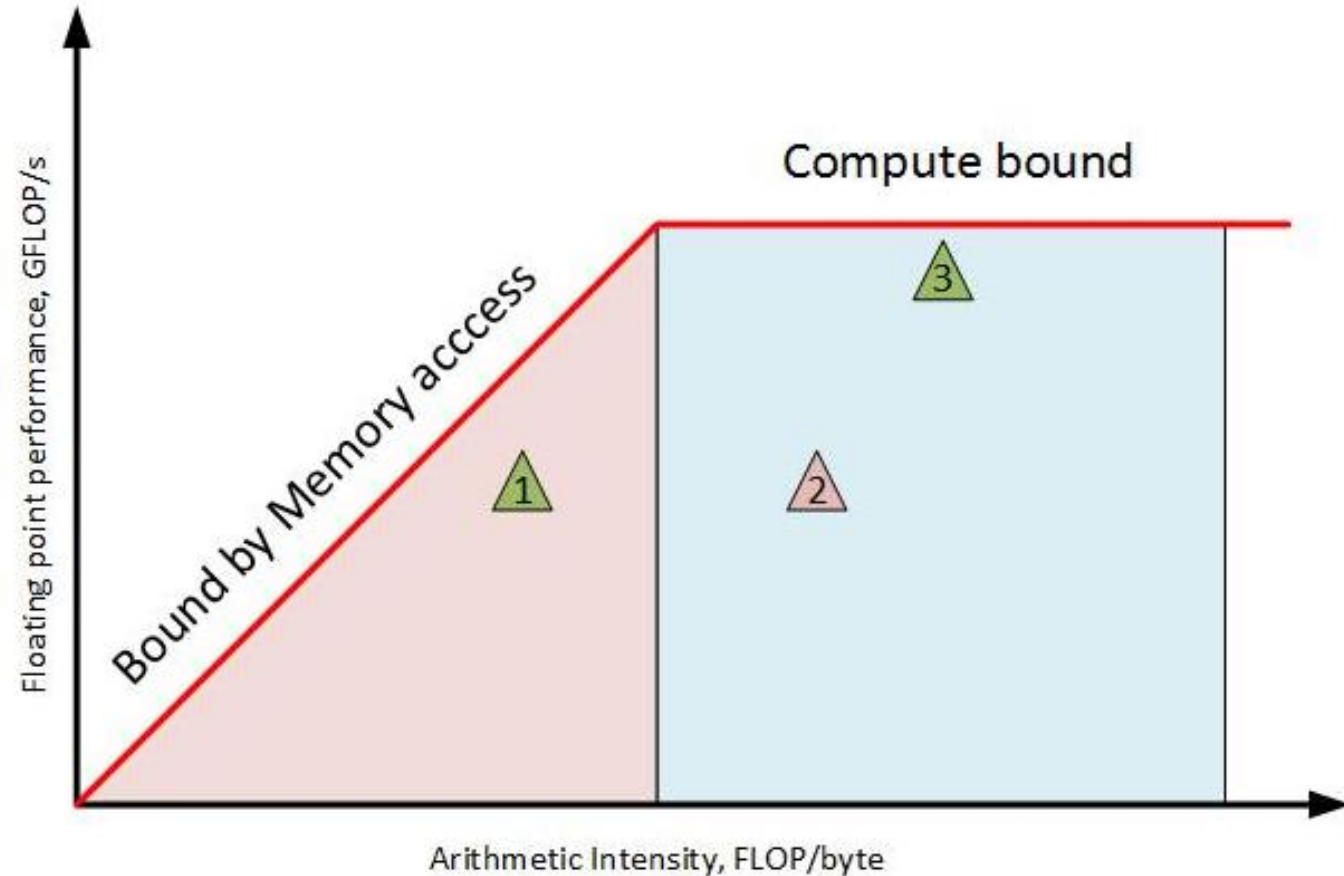
- Two parts of computing each add time to any given task →
  - Memory loading (Gb/s)
  - Computation (FLOPS)
- Over time, memory loading has gotten slower relative to computation
- This means memory loading can be more of a bottleneck if we are only using things we load from memory one time



[Turing Award Presentation, 2021](#) [Dongarra]

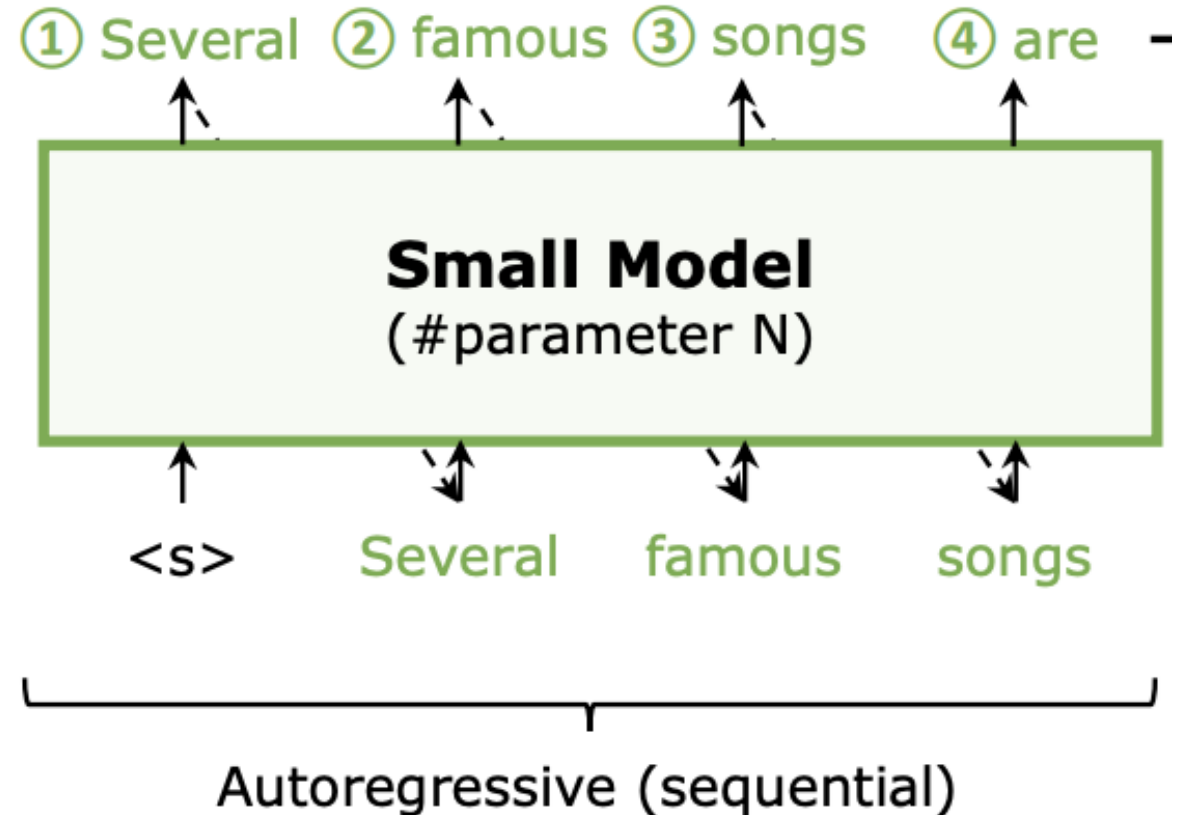
# Compute vs. IO

- One way to alleviate this is to increase the amount of computation we perform for each byte we load from memory
- This is called the *arithmetic intensity* of a given program/function
- Generally, we would prefer to be in the compute bound region more often as this is where work is being done



# Speculative Decoding

- In standard autoregressive decoding, we are only using each parameter **one time** when the batch size is 1
- This means standard decoding has a **low** arithmetic intensity and is memory bound
- We have a bunch more compute we could be getting for free given how massively parallel GPUs are

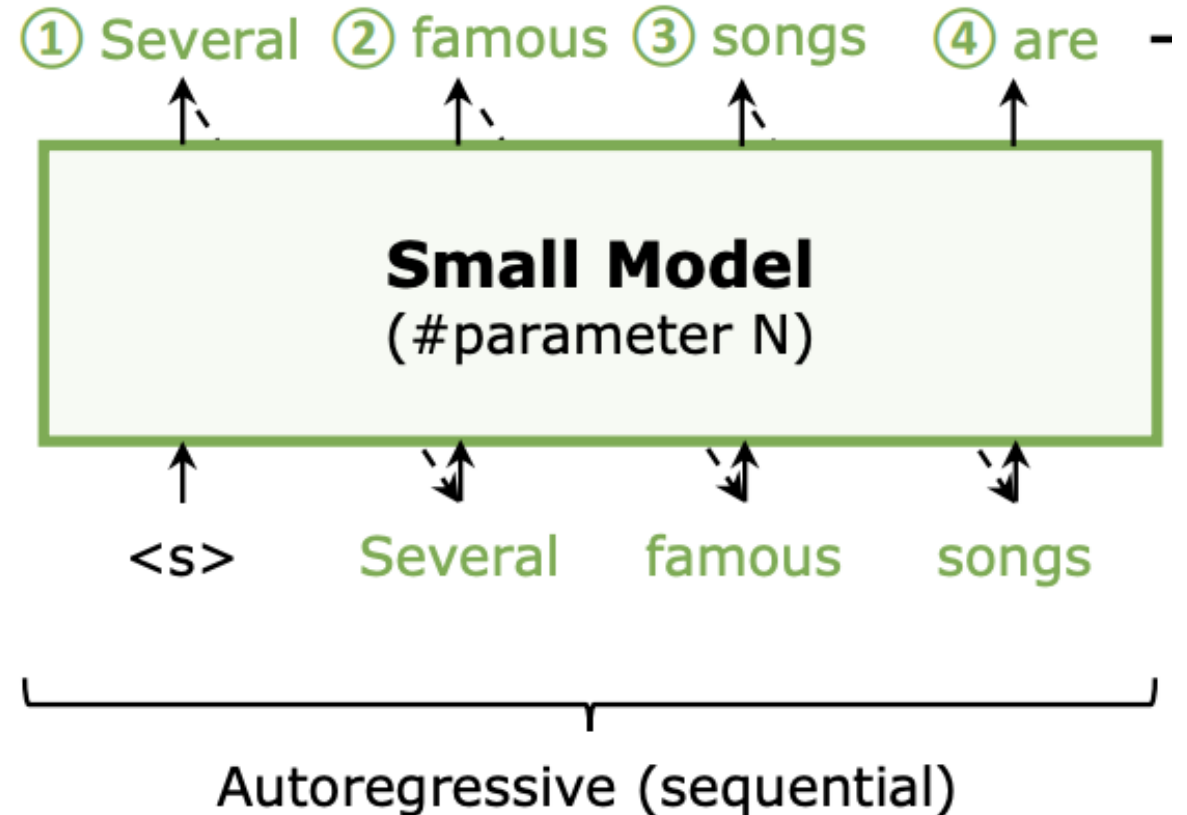


Big Little Decoder [Kim, et. al]



# Speculative Decoding

- Speculative decoding resolves this by 'speculating' multiple tokens into the future with a smaller, cheaper model

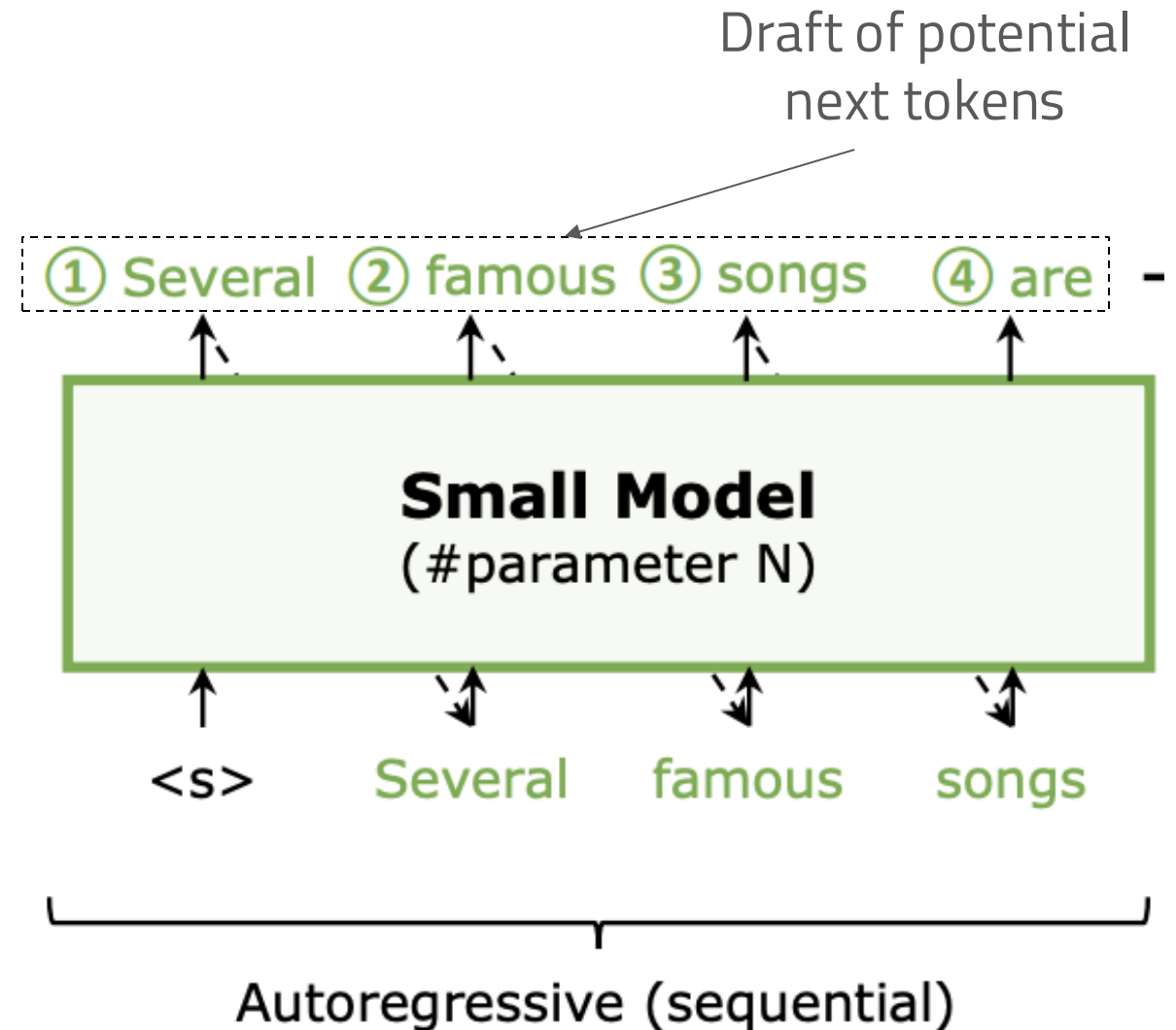


Big Little Decoder [Kim, et. al]



# Speculative Decoding

- Speculative decoding resolves this by 'speculating' multiple tokens into the future with a smaller, cheaper model



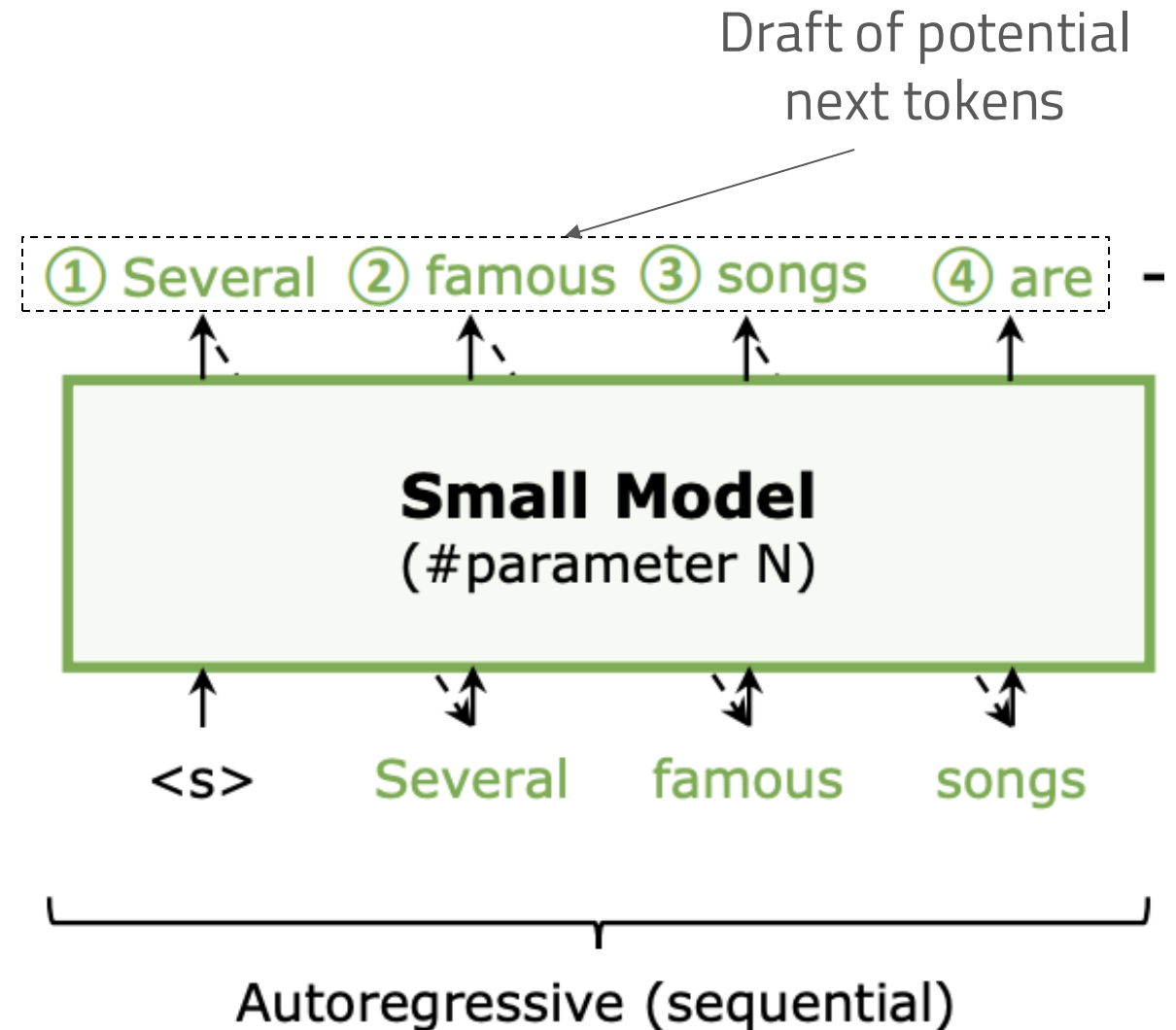
Big Little Decoder [Kim, et. al]





# Speculative Decoding

- Speculative decoding resolves this by 'speculating' multiple tokens into the future with a smaller, cheaper model
- We can now send this set of tokens on to a much larger model to verify the sequence

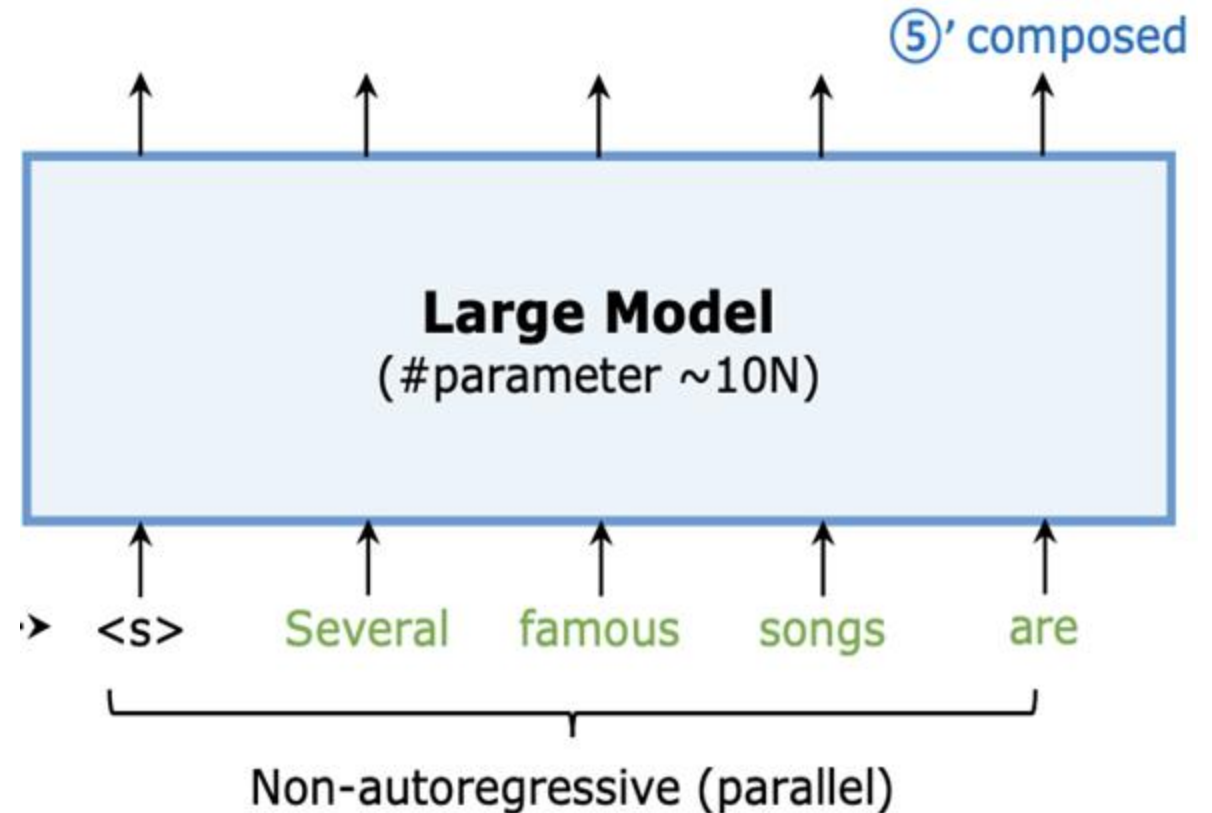


Big Little Decoder [Kim, et. al]



# Speculative Decoding

- Because the sequence will be run in parallel the arithmetic intensity will be proportional to the number of draft tokens
- We run each token through and see if the output of the large model matches that of the smaller, draft model

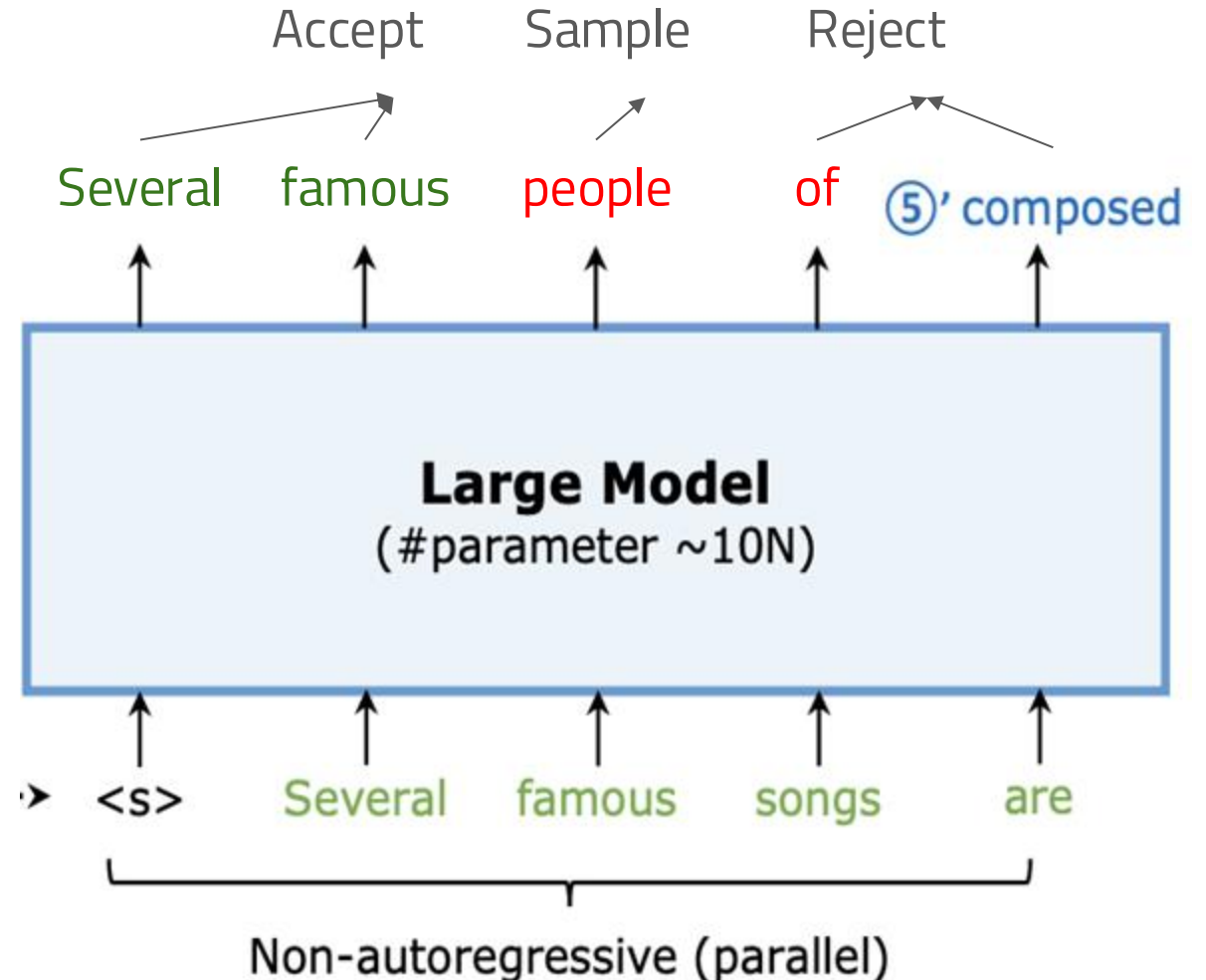


Big Little Decoder [Kim, et. al]



# Speculative Decoding

- Because the sequence will be run in parallel the arithmetic intensity will be proportional to the number of draft tokens
- We run each token through and see if the output of the large model matches that of the smaller, draft model
- We accept the matching tokens

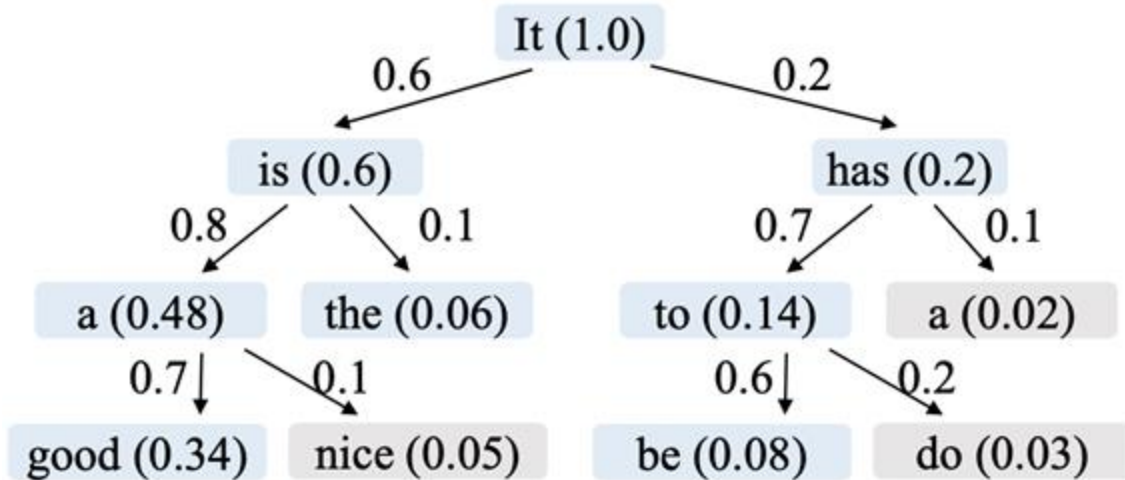


Big Little Decoder [Kim, et. al]



# Advanced Approaches

More advanced approaches will use draft trees, rather than draft sequences



Attention mask

	It	is	has	a	the	to	good	be
It	✓							
is	✓	✓						
has	✓		✓					
a	✓	✓		✓				
the	✓	✓			✓			
to	✓		✓			✓		
good	✓	✓		✓			✓	
be	✓		✓			✓		✓

Big Little Decoder [Kim, et. al]

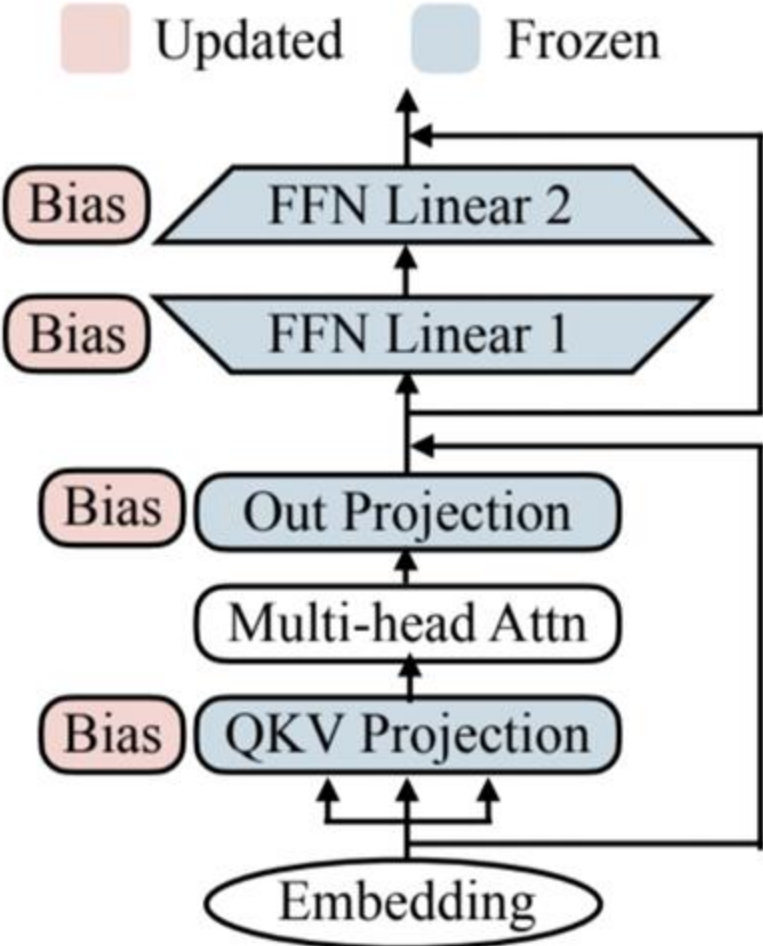


# Efficient LLMs

- ❑ Quantization
  - Background
  - K-Means vs. Linear Quantization
  - Quantization Granularity
  - Quantization Aware Training (QAT) vs Post-Training Quantization (PTQ)
  - LLM Quantization (LLM.int8(), SmoothQuant, AWQ, 1-bit LLMs)
- ❑ Sparsity (Mixture of Experts, Deja Vu: Contextual Sparsity)
- ❑ Long Context (PagedAttention, StreamingLLM, MHA/GQA/MQA, H<sub>2</sub>O)
- ❑ Speculative Decoding
- ❑ **Parameter Efficient Fine-Tuning (BitFit, Adapter, Prompt Tuning, LoRA)**



# BitFit



Update only the bias parameters

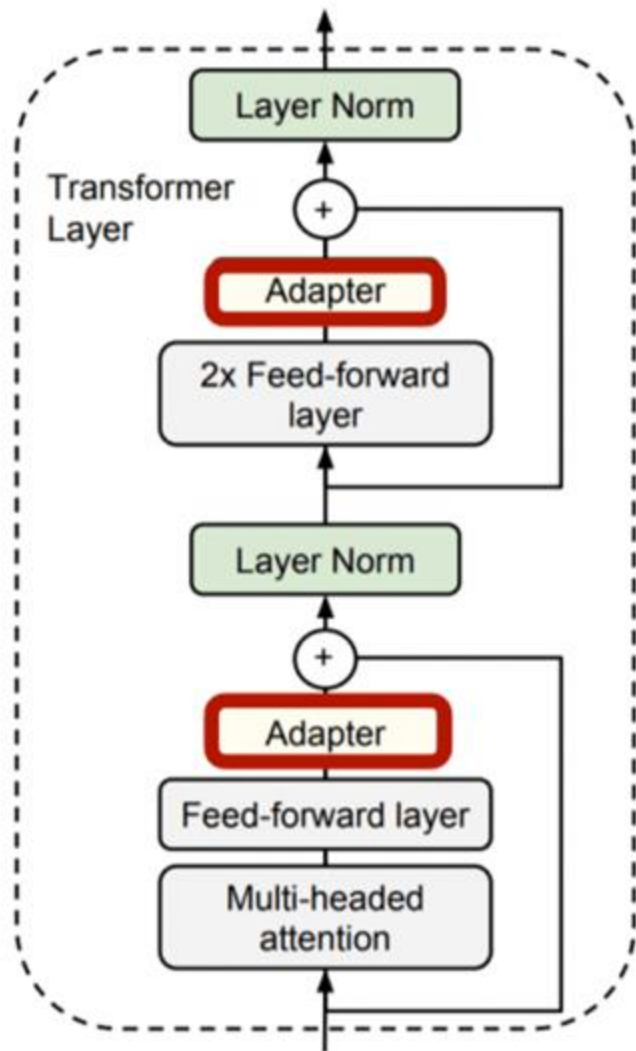
	Train size	%Param	QNLI 105k	SST-2 67k	MNLI <sub>m</sub> 393k	MNLI <sub>mm</sub> 393k	CoLA 8.5k	MRPC 3.7k	STS-B 7k	RTE 2.5k	QQP 364k	Avg.
(V)	Full-FT <sup>†</sup>	100%	<b>93.5</b>	<b>94.1</b>	<b>86.5</b>	<b>87.1</b>	<b>62.8</b>	<b>91.9</b>	89.8	71.8	<b>87.6</b>	<b>84.8</b>
(V)	Full-FT	100%	91.7±0.1	93.4±0.2	85.5±0.4	85.7±0.4	62.2±1.2	90.7±0.3	<b>90.0±0.4</b>	<b>71.9±1.3</b>	87.5±0.4	84.1
(V)	Diff-Prune <sup>‡</sup>	0.5%	<b>93.4</b>	<b>94.2</b>	<b>86.4</b>	<b>86.9</b>	63.5	91.3	89.5	71.5	<b>86.6</b>	<b>84.6</b>
(V)	BitFit	0.08%	91.4±2.4	93.2±0.4	84.4±0.2	84.8±0.1	<b>63.6±0.7</b>	<b>91.7±0.5</b>	<b>90.3±0.1</b>	<b>73.2±3.7</b>	85.4±0.1	84.2
(T)	Full-FT <sup>‡</sup>	100%	91.1	<b>94.9</b>	86.7	85.9	<b>60.5</b>	<b>89.3</b>	<b>87.6</b>	70.1	<b>72.1</b>	<b>81.8</b>
(T)	Full-FT <sup>†</sup>	100%	<b>93.4</b>	94.1	86.7	<b>86.0</b>	59.6	88.9	86.6	<b>71.2</b>	71.7	81.5
(T)	Adapters <sup>‡</sup>	3.6%	90.7	94.0	84.9	85.1	59.5	89.5	<b>86.9</b>	71.5	<b>71.8</b>	81.1
(T)	Diff-Prune <sup>‡</sup>	0.5%	<b>93.3</b>	94.1	<b>86.4</b>	<b>86.0</b>	<b>61.1</b>	<b>89.7</b>	86.0	70.6	71.1	<b>81.5</b>
(T)	BitFit	0.08%	92.0	<b>94.2</b>	84.5	84.8	59.7	88.9	85.5	<b>72.0</b>	70.5	80.9

Table 1: BERT<sub>LARGE</sub> model performance on the GLUE benchmark validation set (V) and test set (T). Lines with <sup>†</sup> and <sup>‡</sup> indicate results taken from Guo et al. (2020) and Houlsby et al. (2019) (respectively).

BitFit: Simple Parameter-efficient Fine-tuning for Transformer-based Masked Language-models [Zeken et al, ACL 2021]



# Adapter



Add trainable layers after each feedforward layer

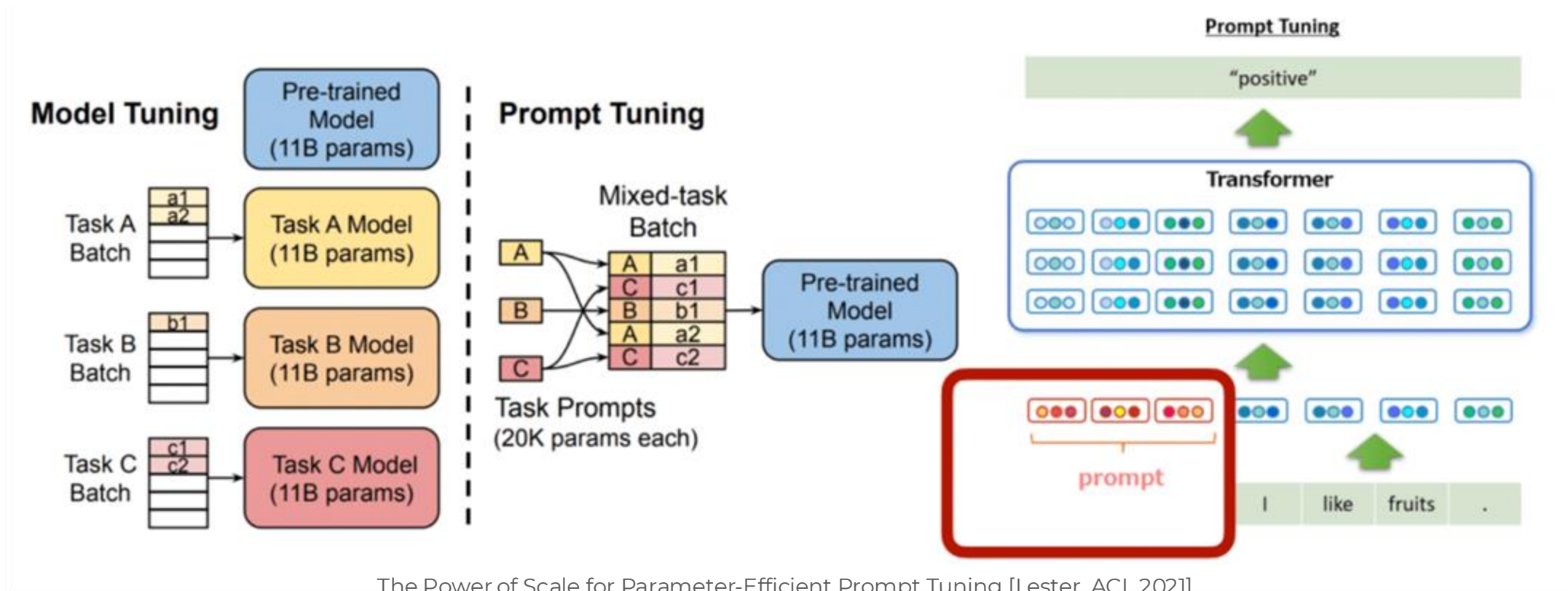
	Total num params	Trained params / task	CoLA	SST	MRPC	STS-B	QQP	MNLI <sub>m</sub>	MNLI <sub>mm</sub>	QNLI	RTE	Total
BERT <sub>LARGE</sub>	9.0×	100%	60.5	94.9	89.3	87.6	72.1	86.7	85.9	91.1	70.1	80.4
Adapters (8-256)	1.3×	3.6%	59.5	94.0	89.5	86.9	71.8	84.9	85.1	90.7	71.5	80.0
Adapters (64)	1.2×	2.1%	56.9	94.2	89.6	87.3	71.8	85.3	84.6	91.4	68.8	79.6

Parameter-Efficient Transfer Learning for NLP [Houlsby et al, ICML 2019]



# Prompt Tuning (Soft Prompting)

Train a continuous, learnable prompt in embedding space for each task we are training on

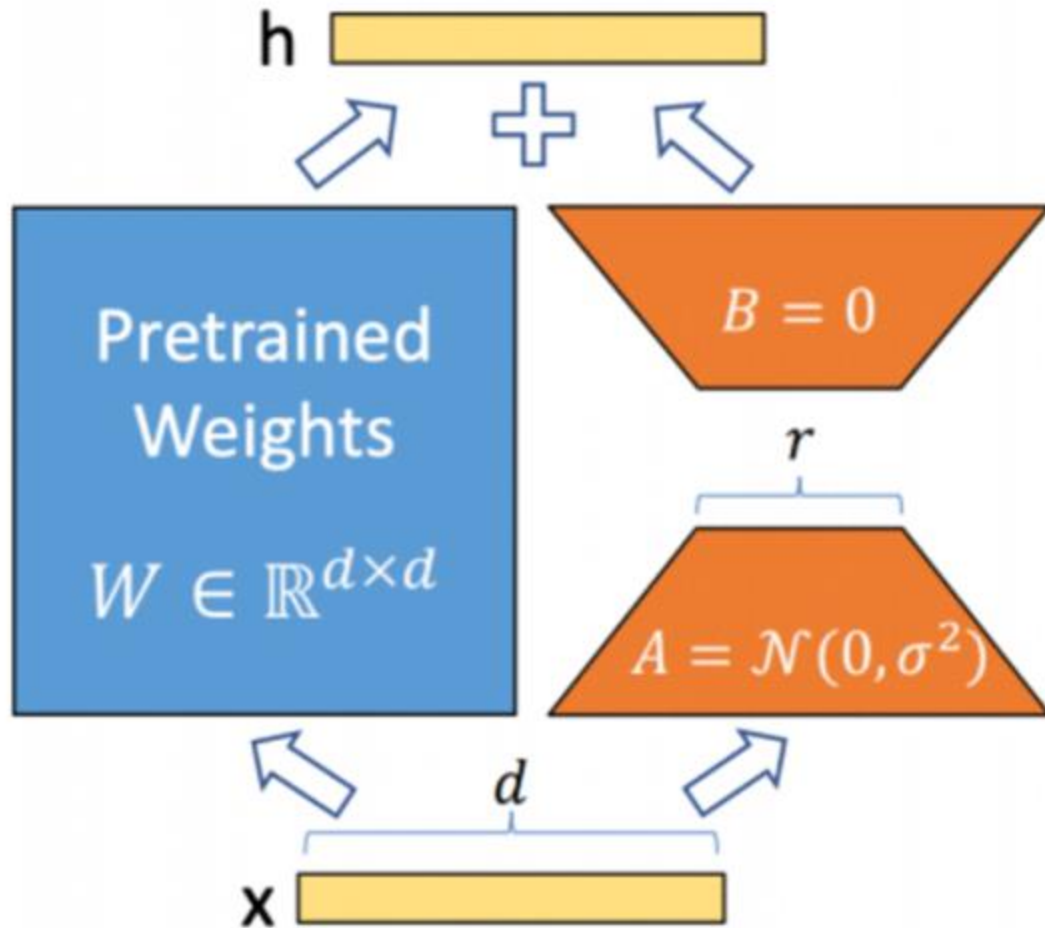


The Power of Scale for Parameter-Efficient Prompt Tuning [Lester, ACL 2021]





# LoRA

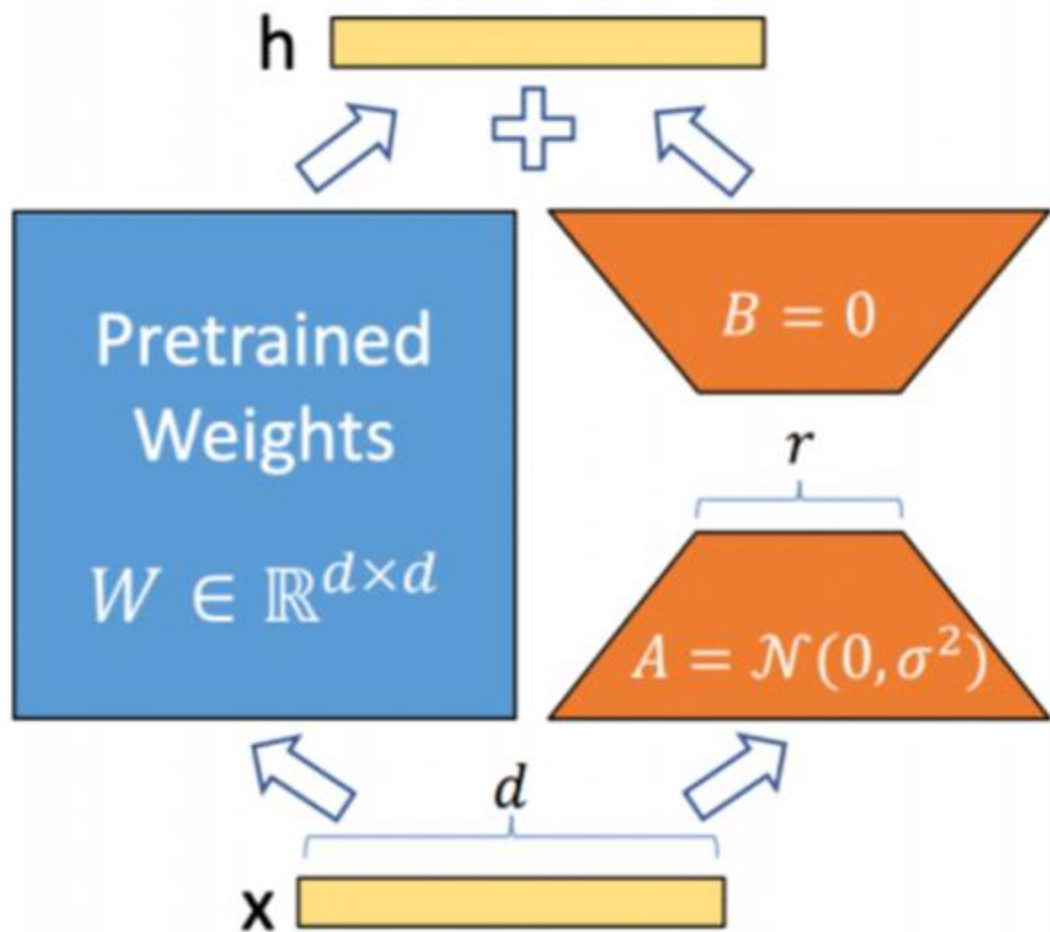


The Power of Scale for Parameter-Efficient Prompt Tuning [Lester, ACL 2021]

- Hypothesizes that fine-tuning results in only low rank updates
- Thus, we may approximate the updates themselves as low-rank and train on this low-rank approximation directly



# LoRA



The Power of Scale for Parameter-Efficient Prompt Tuning [Lester, ACL 2021]

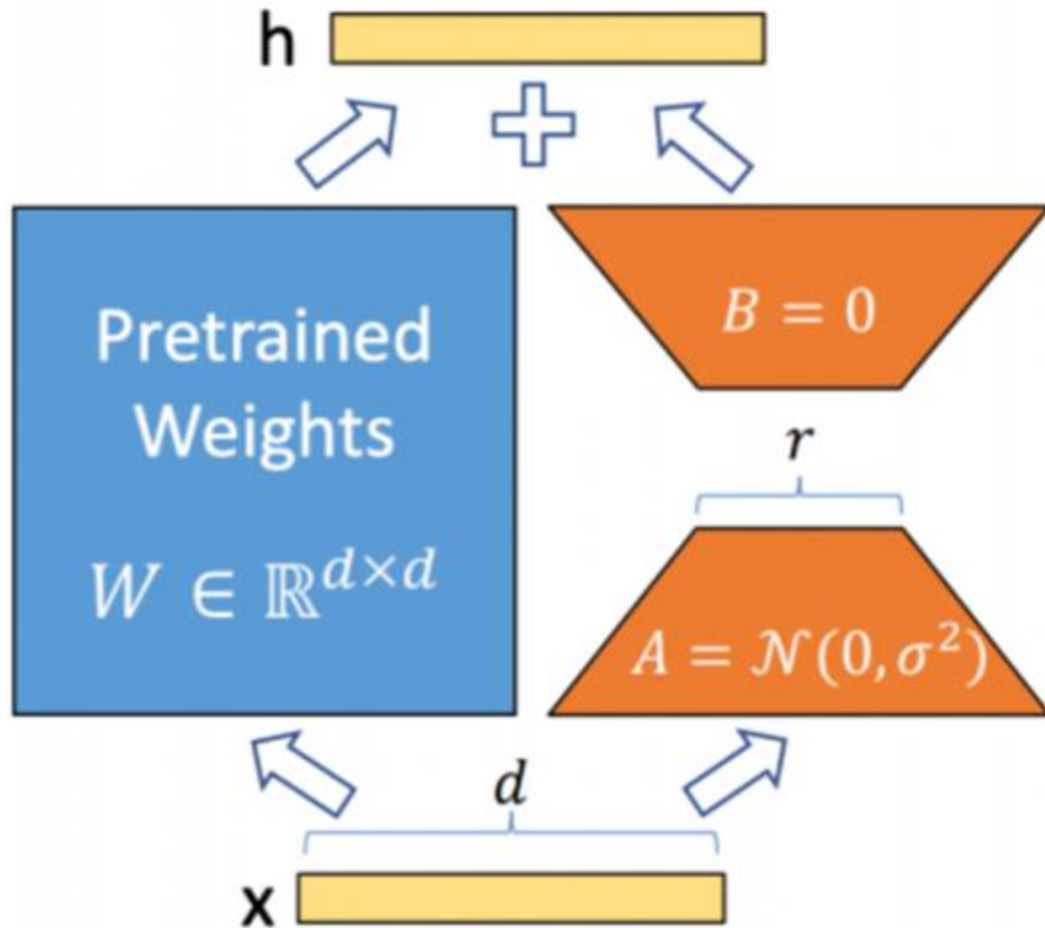
Normal Finetuning:

$$h = Wx$$

Update  $W$



# LoRA



The Power of Scale for Parameter-Efficient Prompt Tuning [Lester, ACL 2021]

Normal Finetuning:

$$h = Wx$$

Update  $W$

LoRA Finetuning:

$$h = Wx + BAx$$

Update  $B, A$   
Leave  $W$  unchanged



# LoRA

Model & Method	# Trainable Parameters	MNLI	SST-2	MRPC	CoLA	QNLI	QQP	RTE	STS-B	Avg.
RoB <sub>base</sub> (FT)*	125.0M	<b>87.6</b>	94.8	90.2	<b>63.6</b>	92.8	<b>91.9</b>	78.7	91.2	86.4
RoB <sub>base</sub> (BitFit)*	0.1M	84.7	93.7	<b>92.7</b>	62.0	91.8	84.0	81.5	90.8	85.2
RoB <sub>base</sub> (Adpt <sup>D</sup> )*	0.3M	87.1 $\pm$ .0	94.2 $\pm$ .1	88.5 $\pm$ 1.1	60.8 $\pm$ .4	93.1 $\pm$ .1	90.2 $\pm$ .0	71.5 $\pm$ 2.7	89.7 $\pm$ .3	84.4
RoB <sub>base</sub> (Adpt <sup>D</sup> )*	0.9M	87.3 $\pm$ .1	94.7 $\pm$ .3	88.4 $\pm$ .1	62.6 $\pm$ .9	93.0 $\pm$ .2	90.6 $\pm$ .0	75.9 $\pm$ 2.2	90.3 $\pm$ .1	85.4
RoB <sub>base</sub> (LoRA)	0.3M	87.5 $\pm$ .3	<b>95.1<math>\pm</math>.2</b>	89.7 $\pm$ .7	63.4 $\pm$ 1.2	<b>93.3<math>\pm</math>.3</b>	90.8 $\pm$ .1	<b>86.6<math>\pm</math>.7</b>	<b>91.5<math>\pm</math>.2</b>	<b>87.2</b>
RoB <sub>large</sub> (FT)*	355.0M	90.2	<b>96.4</b>	<b>90.9</b>	68.0	94.7	<b>92.2</b>	86.6	92.4	88.9
RoB <sub>large</sub> (LoRA)	0.8M	<b>90.6<math>\pm</math>.2</b>	96.2 $\pm$ .5	<b>90.9<math>\pm</math>1.2</b>	<b>68.2<math>\pm</math>1.9</b>	<b>94.9<math>\pm</math>.3</b>	91.6 $\pm$ .1	<b>87.4<math>\pm</math>2.5</b>	<b>92.6<math>\pm</math>.2</b>	<b>89.0</b>
RoB <sub>large</sub> (Adpt <sup>P</sup> )†	3.0M	90.2 $\pm$ .3	96.1 $\pm$ .3	90.2 $\pm$ .7	<b>68.3<math>\pm</math>1.0</b>	<b>94.8<math>\pm</math>.2</b>	<b>91.9<math>\pm</math>.1</b>	83.8 $\pm$ 2.9	92.1 $\pm$ .7	88.4
RoB <sub>large</sub> (Adpt <sup>P</sup> )†	0.8M	<b>90.5<math>\pm</math>.3</b>	<b>96.6<math>\pm</math>.2</b>	89.7 $\pm$ 1.2	67.8 $\pm$ 2.5	<b>94.8<math>\pm</math>.3</b>	91.7 $\pm$ .2	80.1 $\pm$ 2.9	91.9 $\pm$ .4	87.9
RoB <sub>large</sub> (Adpt <sup>H</sup> )†	6.0M	89.9 $\pm$ .5	96.2 $\pm$ .3	88.7 $\pm$ 2.9	66.5 $\pm$ 4.4	94.7 $\pm$ .2	92.1 $\pm$ .1	83.4 $\pm$ 1.1	91.0 $\pm$ 1.7	87.8
RoB <sub>large</sub> (Adpt <sup>H</sup> )†	0.8M	90.3 $\pm$ .3	96.3 $\pm$ .5	87.7 $\pm$ 1.7	66.3 $\pm$ 2.0	94.7 $\pm$ .2	91.5 $\pm$ .1	72.9 $\pm$ 2.9	91.5 $\pm$ .5	86.4
RoB <sub>large</sub> (LoRA)†	0.8M	<b>90.6<math>\pm</math>.2</b>	96.2 $\pm$ .5	<b>90.2<math>\pm</math>1.0</b>	68.2 $\pm$ 1.9	<b>94.8<math>\pm</math>.3</b>	91.6 $\pm$ .2	<b>85.2<math>\pm</math>1.1</b>	<b>92.3<math>\pm</math>.5</b>	<b>88.6</b>
DeB <sub>XXL</sub> (FT)*	1500.0M	91.8	<b>97.2</b>	92.0	72.0	<b>96.0</b>	92.7	93.9	92.9	91.1
DeB <sub>XXL</sub> (LoRA)	4.7M	<b>91.9<math>\pm</math>.2</b>	96.9 $\pm$ .2	<b>92.6<math>\pm</math>.6</b>	<b>72.4<math>\pm</math>1.1</b>	<b>96.0<math>\pm</math>.1</b>	<b>92.9<math>\pm</math>.1</b>	<b>94.9<math>\pm</math>.4</b>	<b>93.0<math>\pm</math>.2</b>	<b>91.3</b>

Table 2: RoBERTa<sub>base</sub>, RoBERTa<sub>large</sub>, and DeBERTa<sub>XXL</sub> with different adaptation methods on the GLUE benchmark. We report the overall (matched and mismatched) accuracy for MNLI, Matthew’s correlation for CoLA, Pearson correlation for STS-B, and accuracy for other tasks. Higher is better for all metrics. \* indicates numbers published in prior works. † indicates runs configured in a setup similar to Houlsby et al. (2019) for a fair comparison.

The Power of Scale for Parameter-Efficient Prompt Tuning [Lester, ACL 2021]



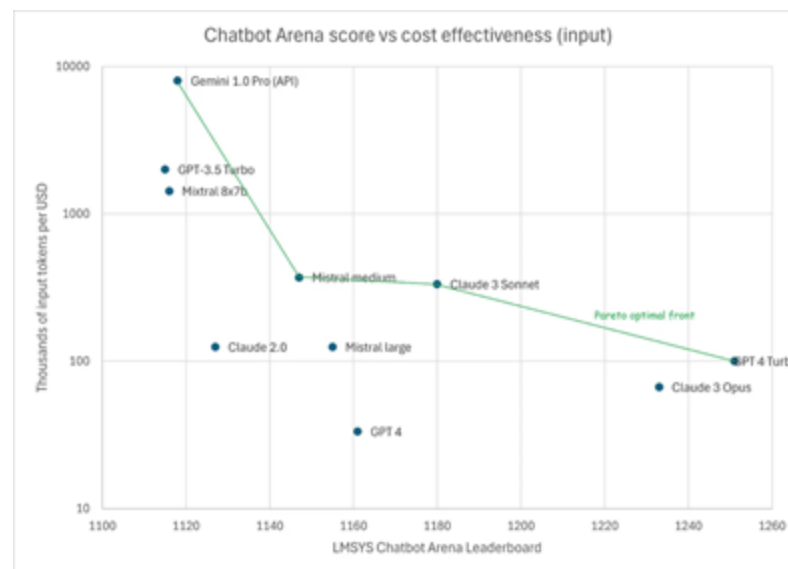
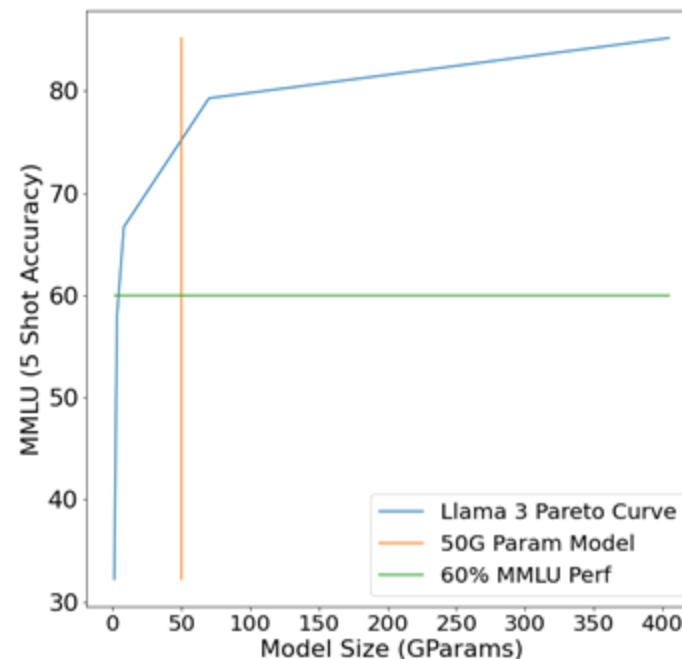
# Summary

- ❑ Efficient inference algorithms in LLMs lead to lower cost, faster inference, and smaller models
- ❑ Quantization and sparsity are the primary techniques for realizing these efficiencies
- ❑ PEFT techniques allow for faster fine-tuning with smaller storage requirements



# Future Directions

- Better, more adaptive inference systems
  - Adaptive speculative decoding
  - Variable Model Serving
- Improved efficiency benchmarking
- More efficient architectures



# Open Source Models/Inference Systems

## □ Models

- [Llama3.2](#)
- [Qwen2.5](#)
- [Mixtral](#)

## □ Quantization

- [AWQ](#)
- [LLM.int8\(\)](#)
- [QLoRA](#)
- [GGUF](#)

## □ Inference Systems

- [vLLM](#)
- [SGLang](#)
- [Tensor-RT LLM](#)
- [Llama.cpp](#)
- [oLLama](#)
- [Huggingface TGI](#)

