

CSCI 5541: Natural Language Processing

Lecture XX: 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?

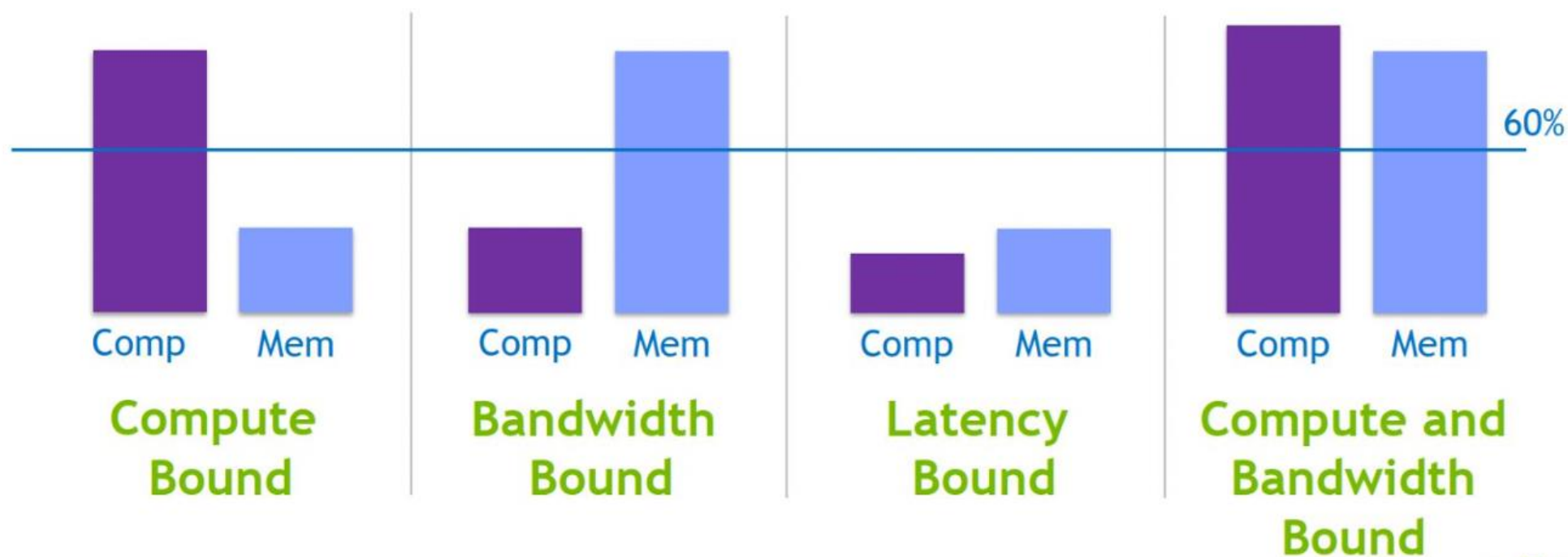
- ❑ Efficiency for NLP is concerned with delivering **faster, cheaper, smaller, less energy** intensive solutions to problems involving natural language
- ❑ **Faster** models means LLM model services (GPT3.5, Claude 2.0, etc.) can meet the demands of many clients more quickly
- ❑ **Cheaper** models reduce 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 Speed

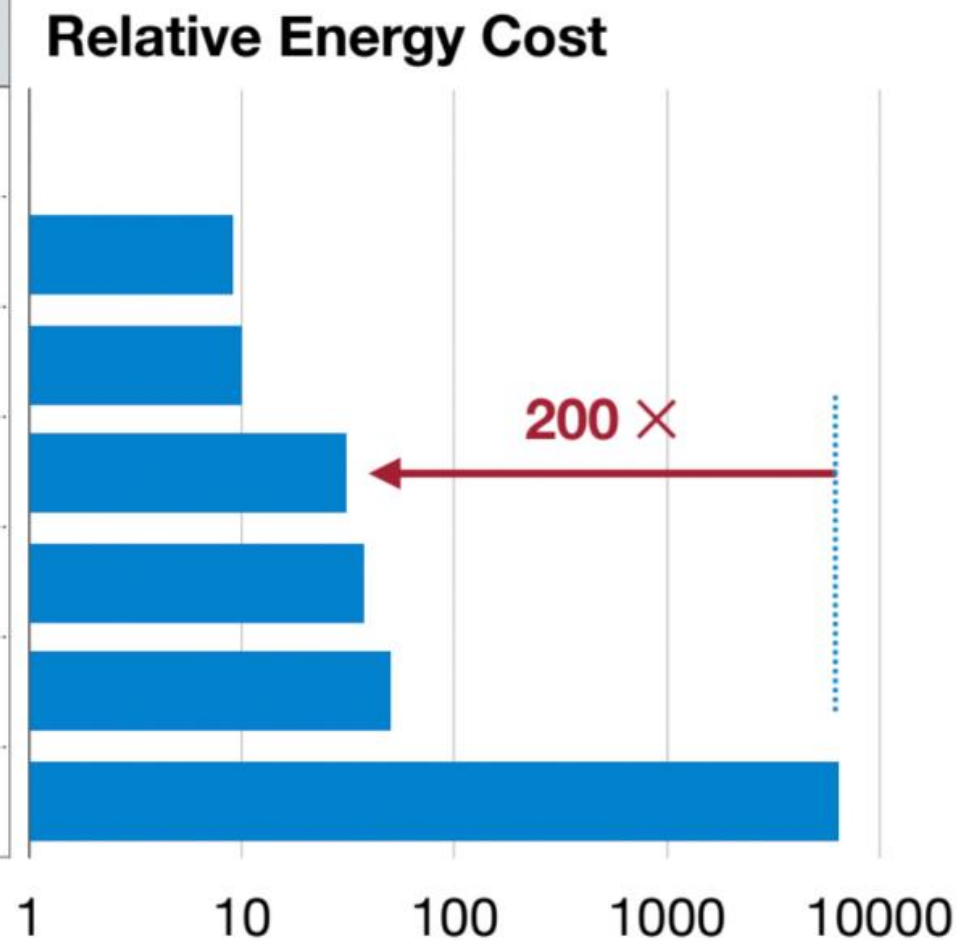
Memory Utilization vs Compute Utilization

Four possible combinations:



Model Energy Use

Operation	Energy [pJ]
32 bit int ADD	0.1
32 bit float ADD	0.9
32 bit Register File	1
32 bit int MULT	3.1
32 bit float MULT	3.7
32 bit SRAM Cache	5
32 bit DRAM Memory	640

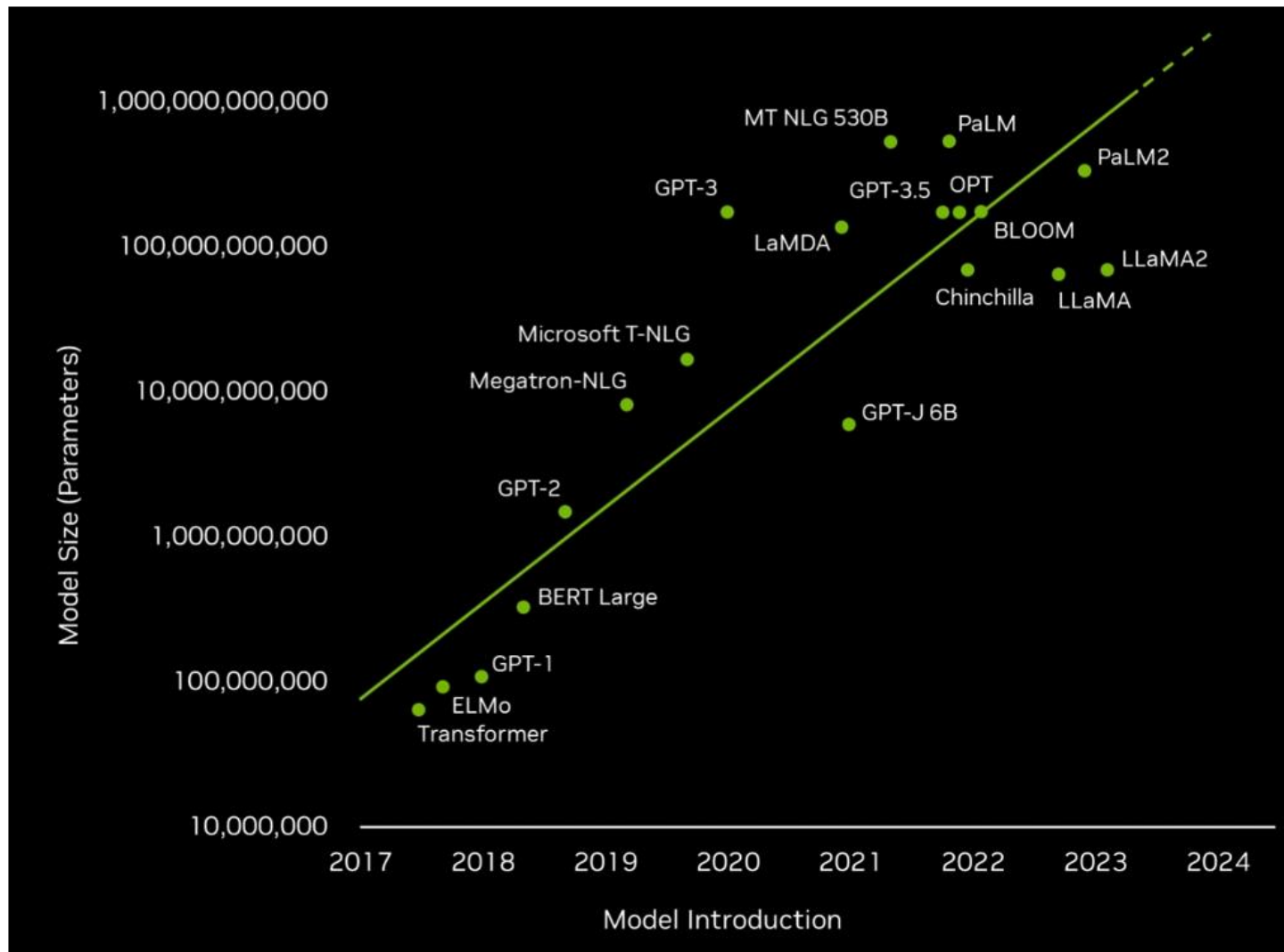


Rough Energy Cost For Various Operations in 45nm 0.9V

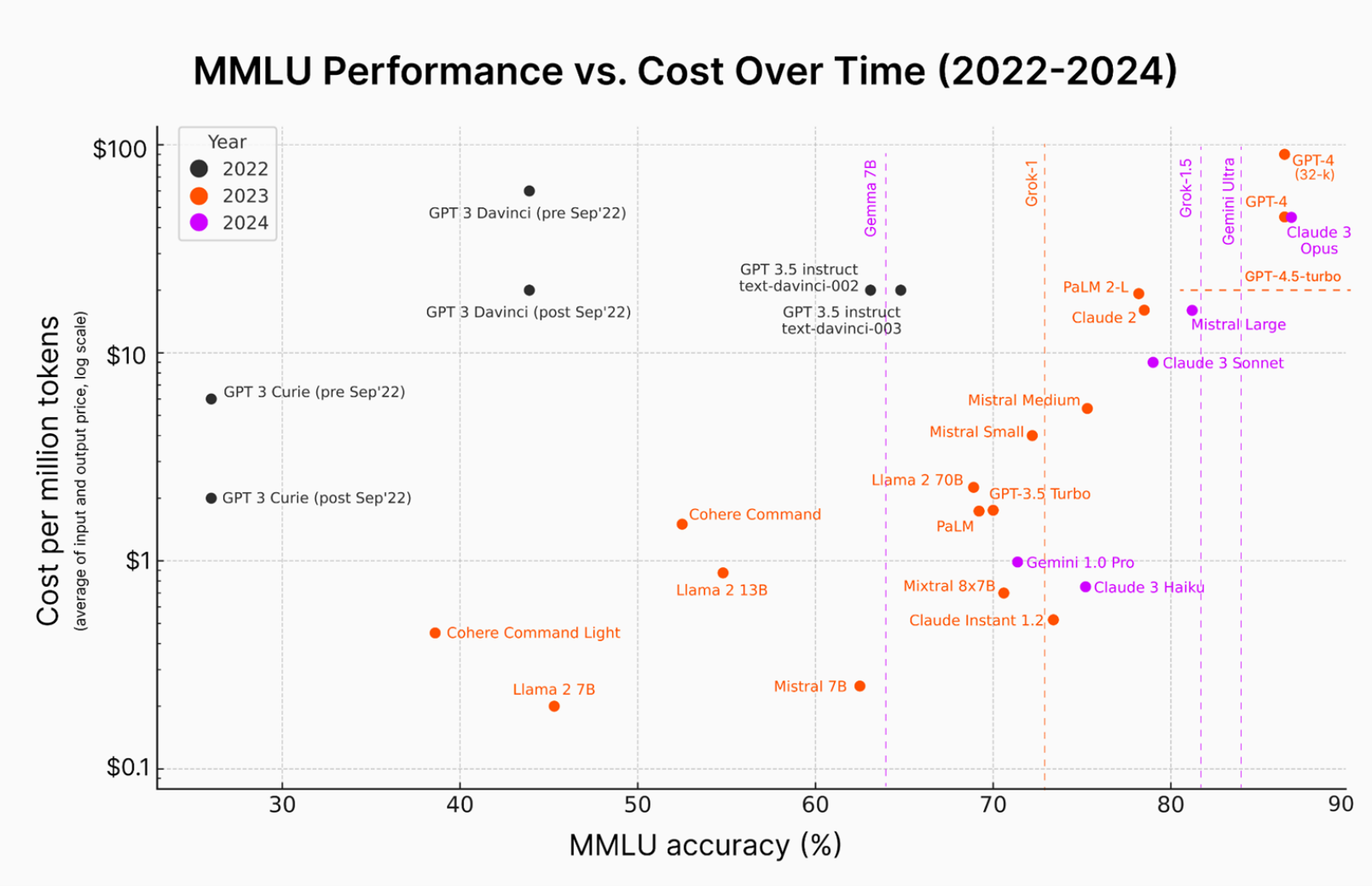
Computing's Energy Problem (and What We Can Do About it) [Horowitz, M., IEEE ISSCC 2014]



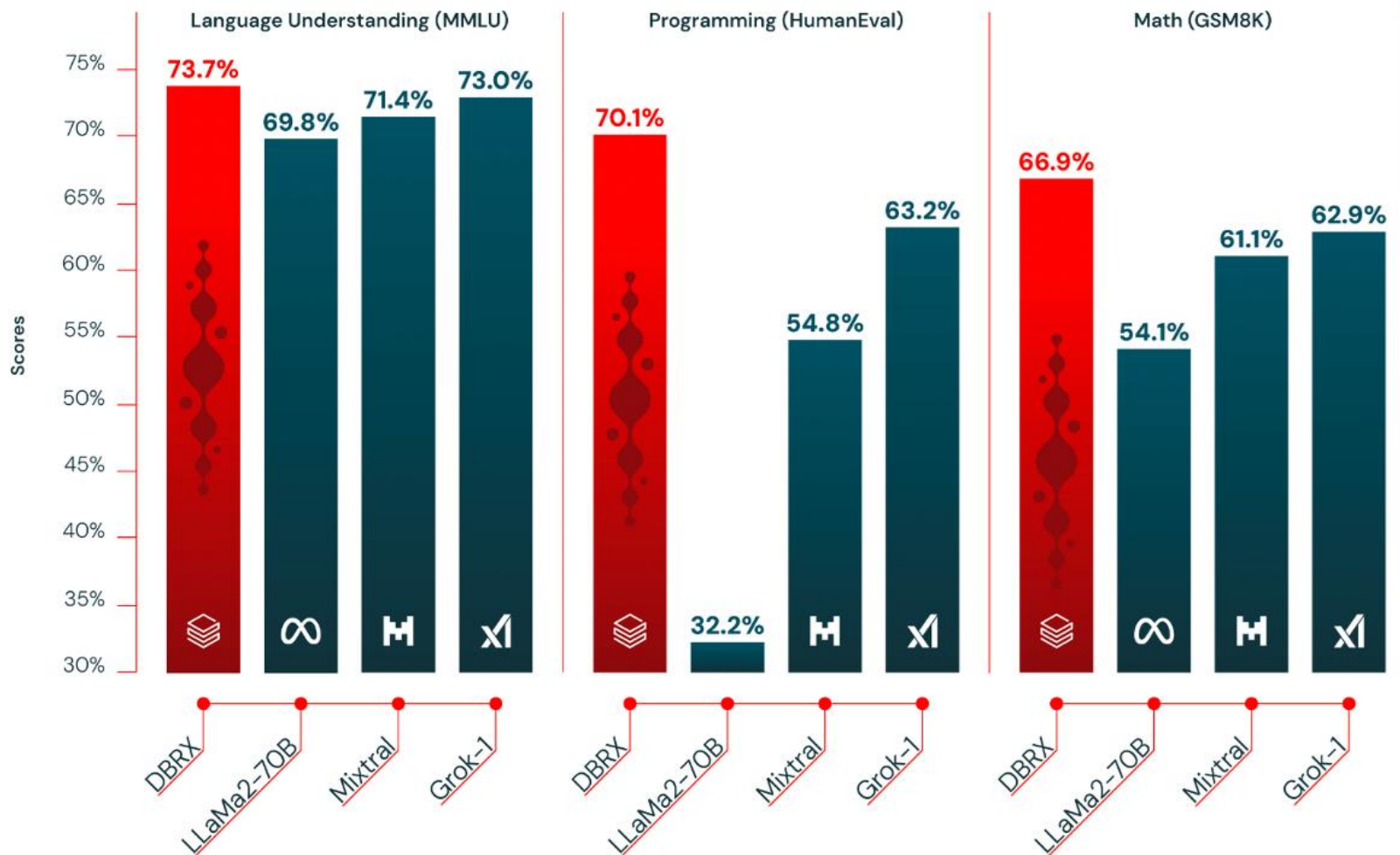
Model Size



Model Cost



Development Speed

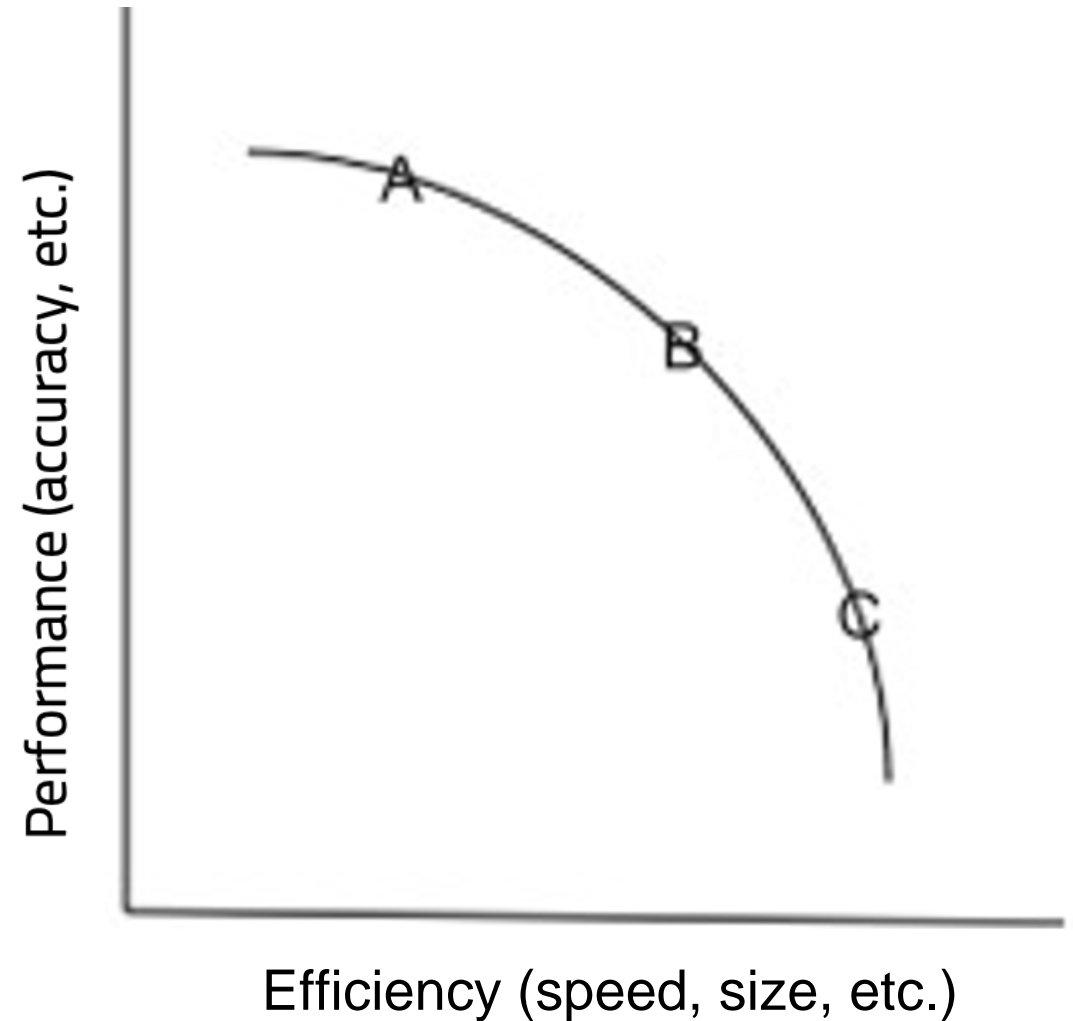


<https://www.databricks.com/blog/introducing-dbrx-new-state-art-open-llm>



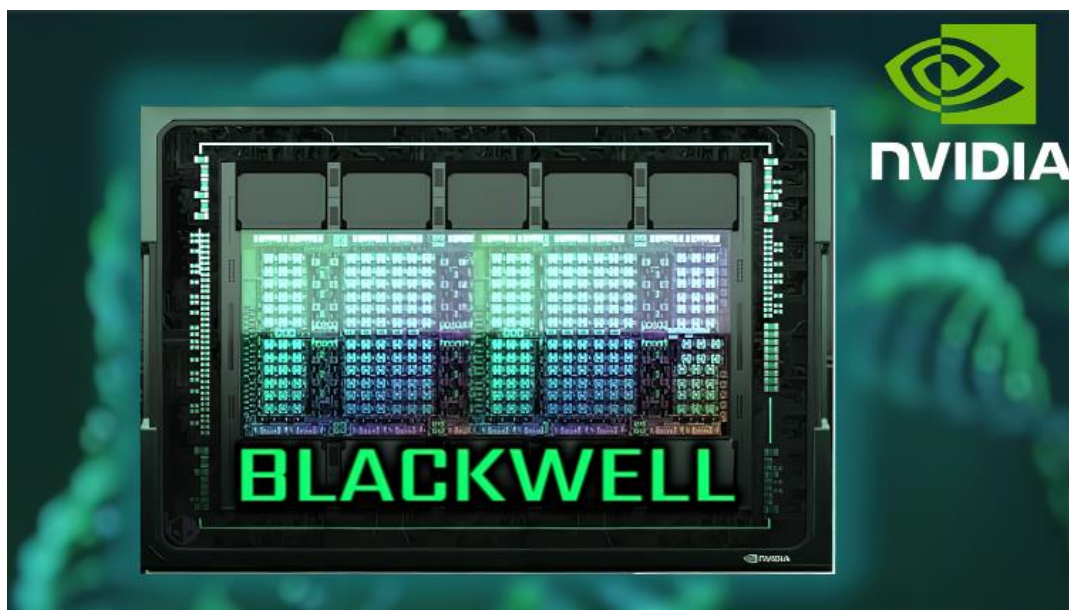
Efficiency Tradeoff

- ❑ 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

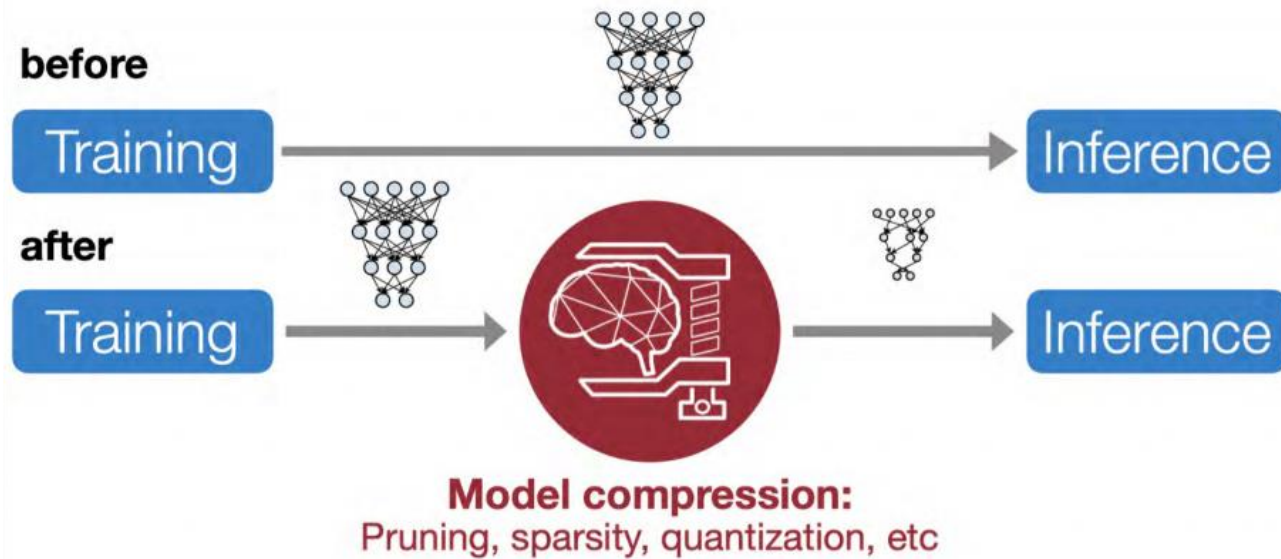


How to Improve Model Efficiency?

Hardware



Software



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)

□ Efficient Inference Systems (vLLM, StreamingLLM, MHA/GQA/MQA)

□ 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)

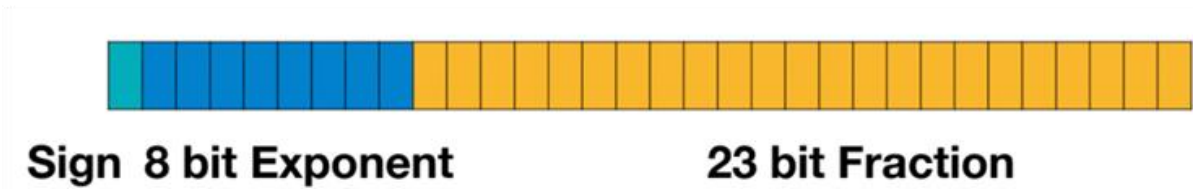
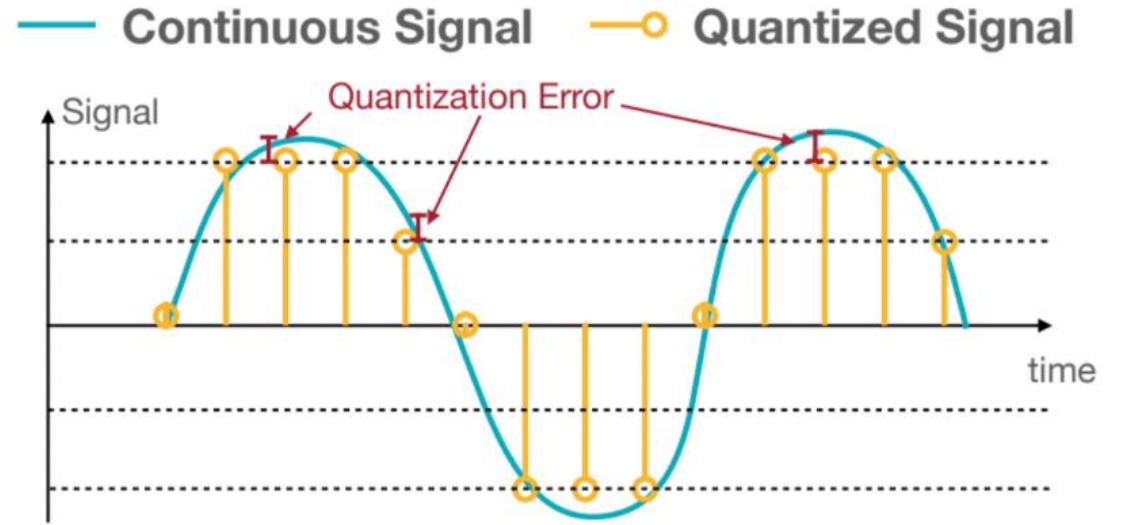
□ Efficient Inference Systems (vLLM, StreamingLLM, MHA/GQA/MQA)

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

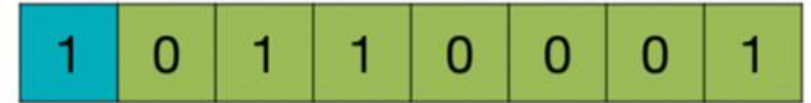


Quantization

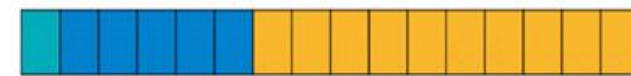
Reduce model size by replacing high bit-width representations with low bit-width representations



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)

□ Efficient Inference Systems (vLLM, StreamingLLM, MHA/GQA/MQA)

□ 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**

Storage	Floating-Point Weights	Integer Weights; Floating-Point Codebook	Integer Weights
Computation	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

$(-1) \times 1.07$

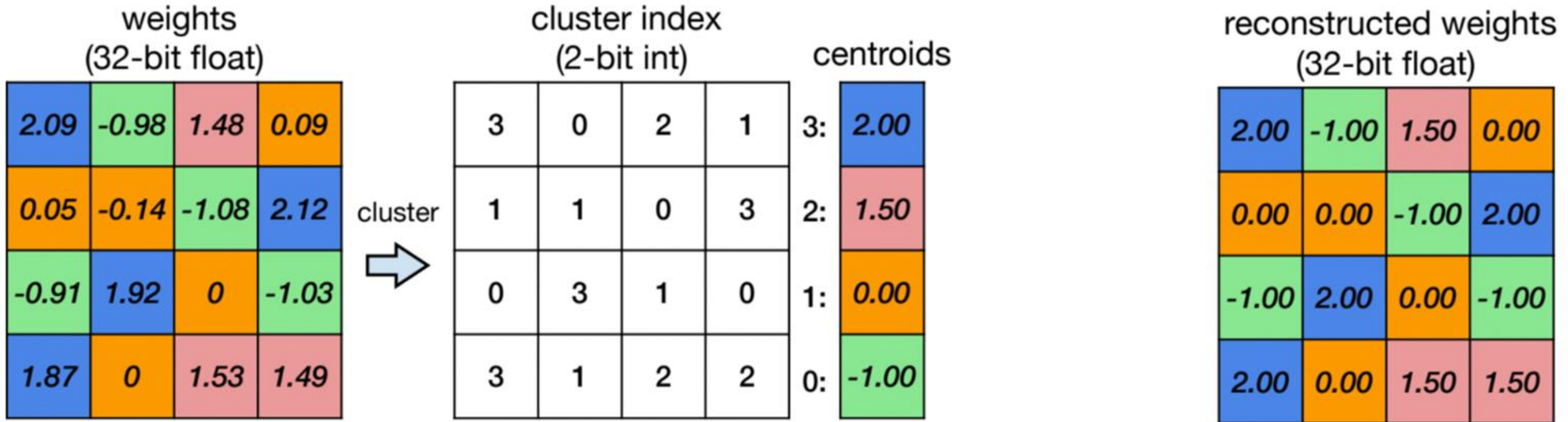
**K-Means-based
Quantization**

**Linear
Quantization**

Storage	Floating-Point Weights	Integer Weights; Floating-Point Codebook	Integer Weights
Computation	Floating-Point Arithmetic	Floating-Point Arithmetic	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)

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

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



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) × 1.07

**K-Means-based
Quantization**

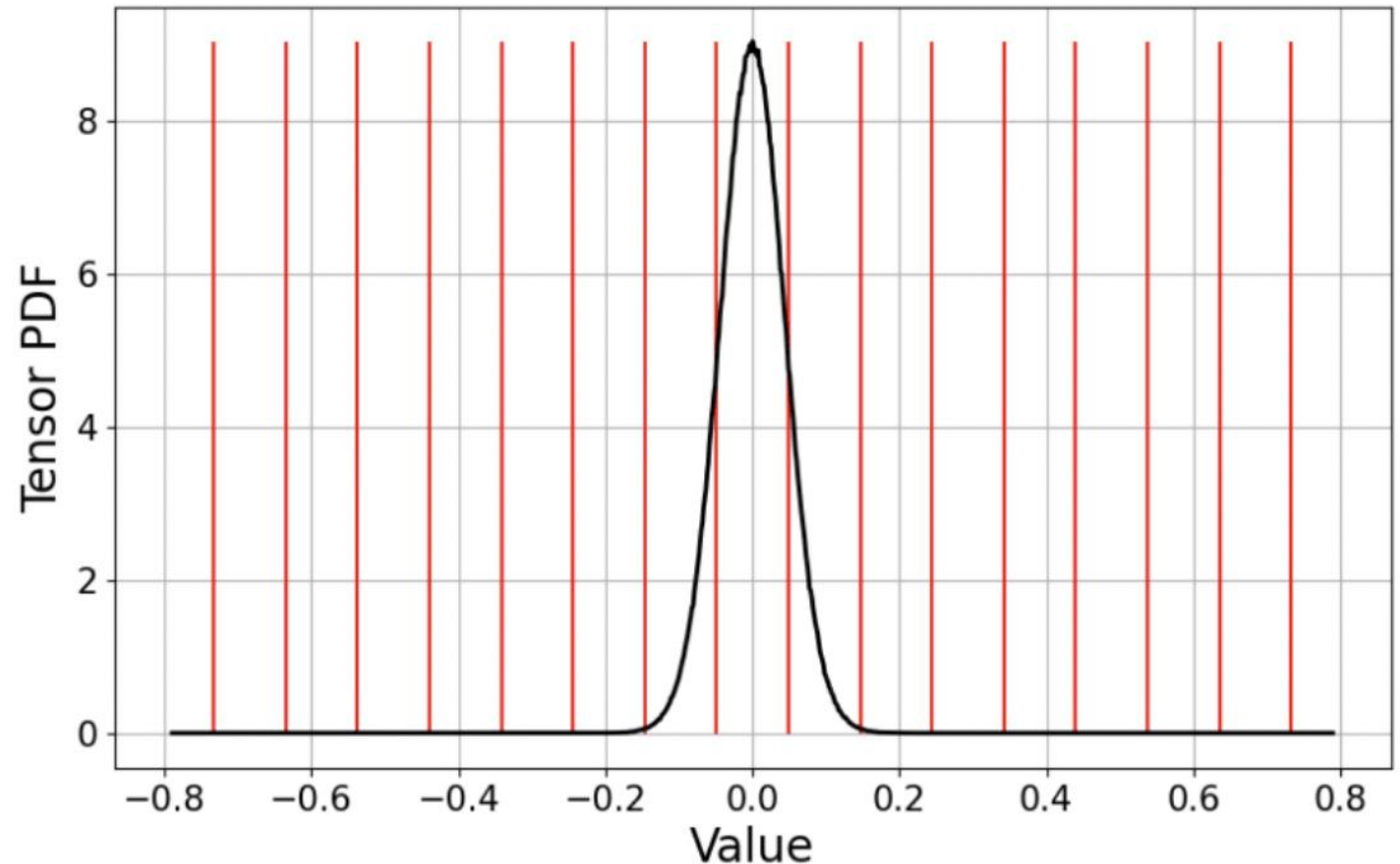
**Linear
Quantization**

Storage	Floating-Point Weights	Integer Weights; Floating-Point Codebook	Integer Weights
Computation	Floating-Point Arithmetic	Floating-Point Arithmetic	Integer Arithmetic

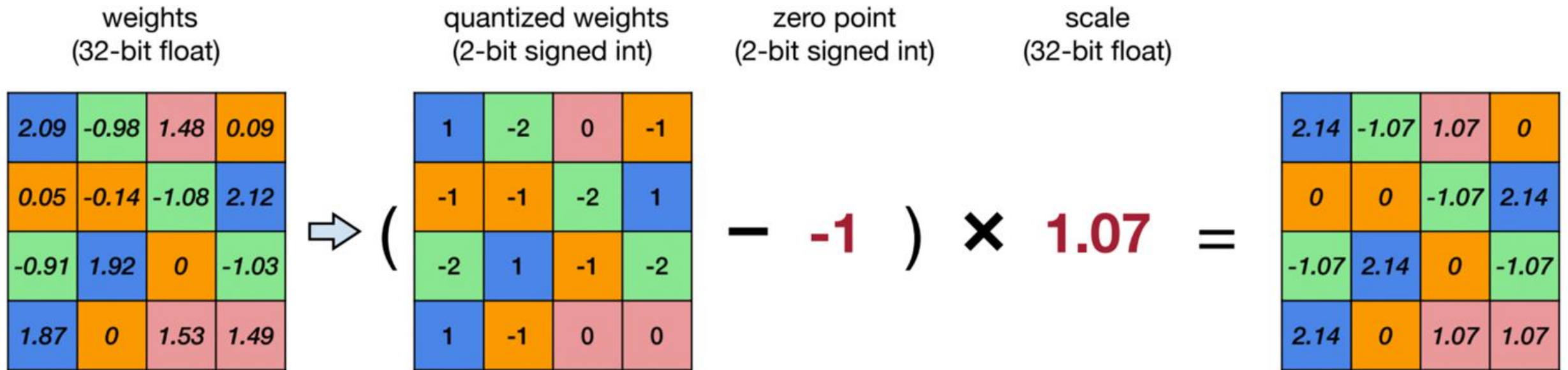


Linear Quantization

- Apply linear function on weights and hidden state activations from floating point values (r) to integer values (q)
- Original weights (black), Quantized bins (red)



Linear Quantization

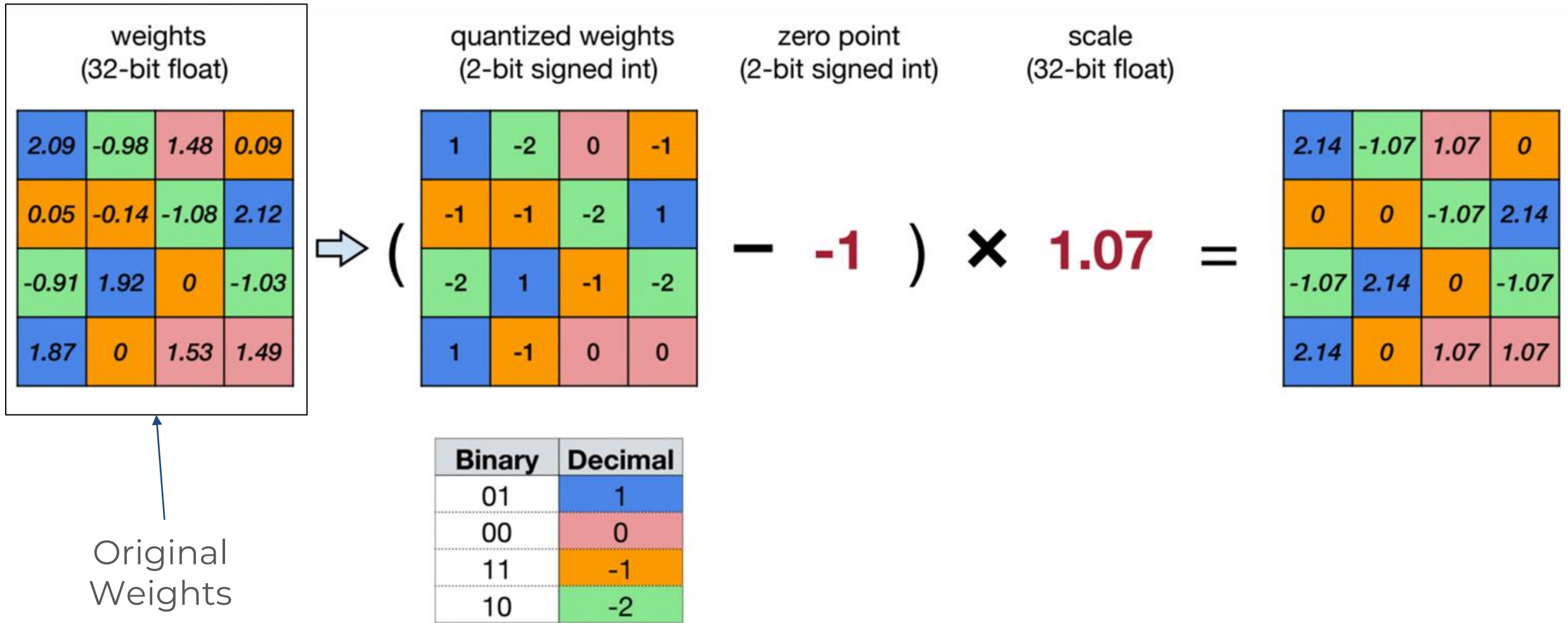


Binary	Decimal
01	1
00	0
11	-1
10	-2

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



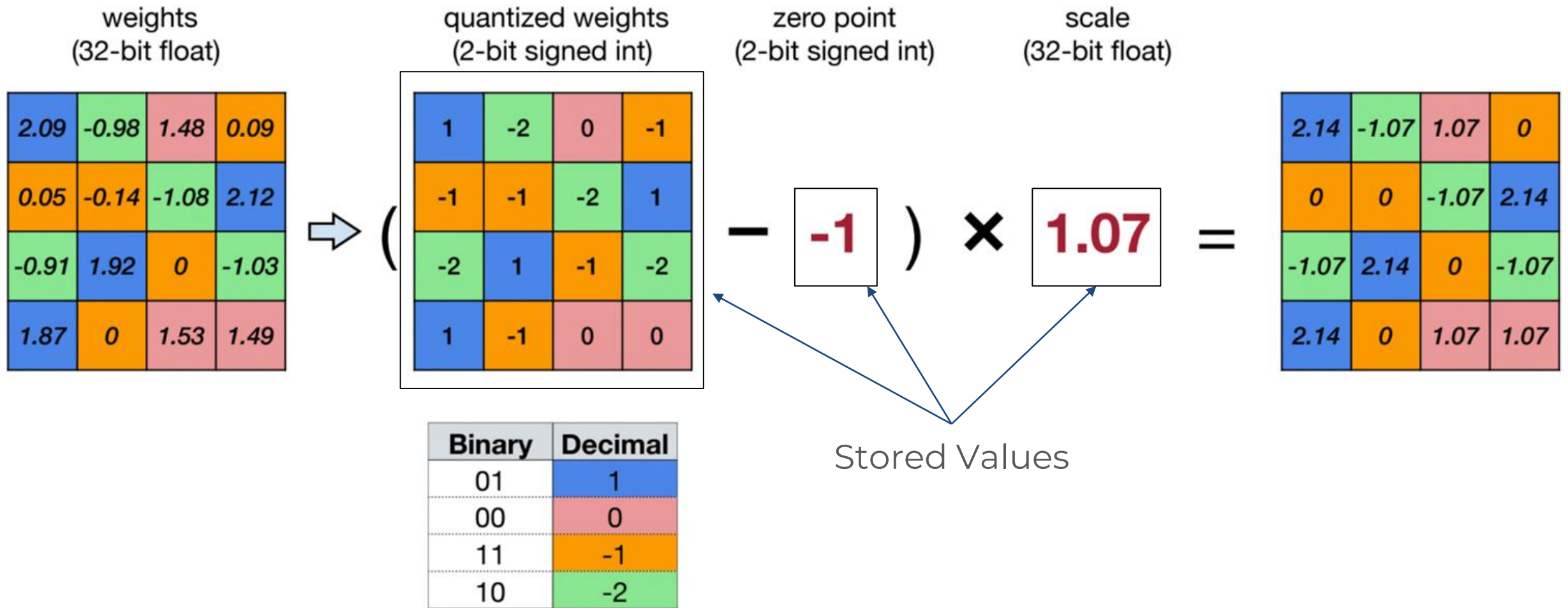
Linear Quantization



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



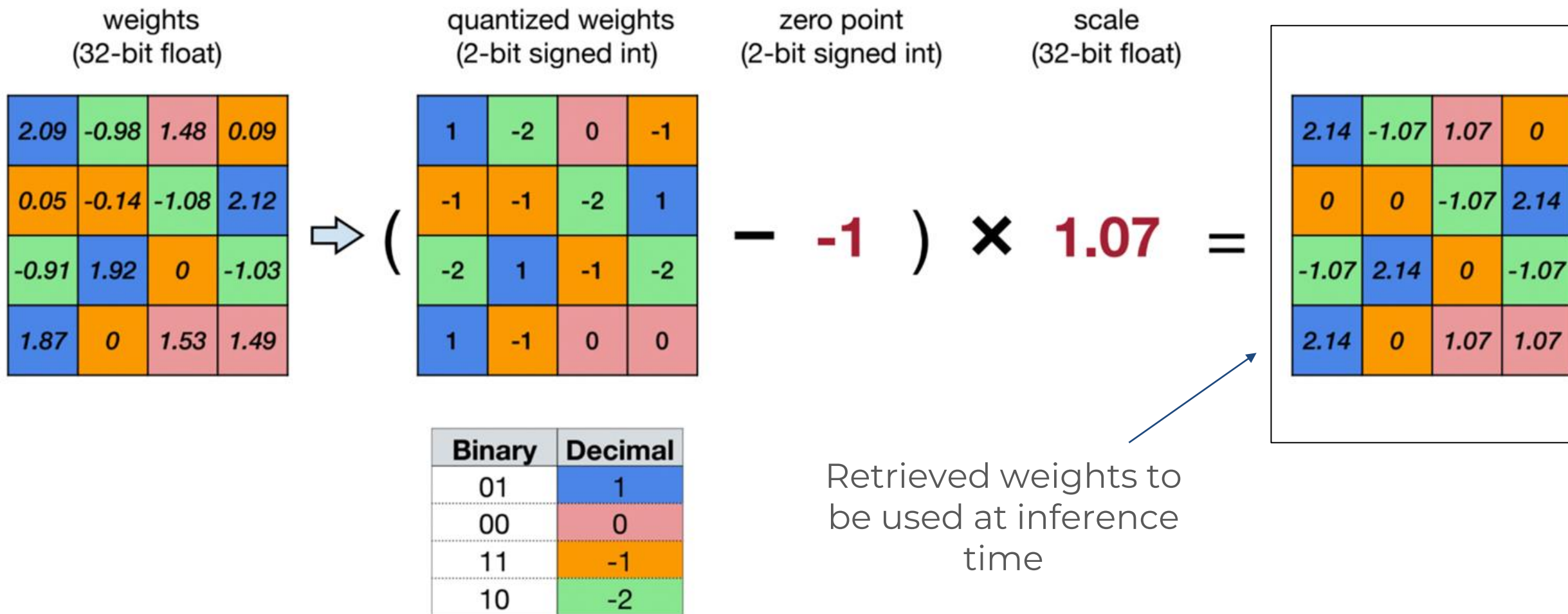
Linear Quantization



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 [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)

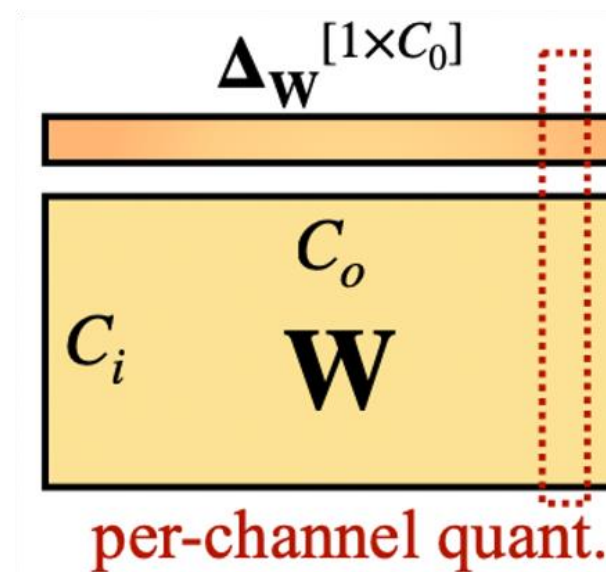
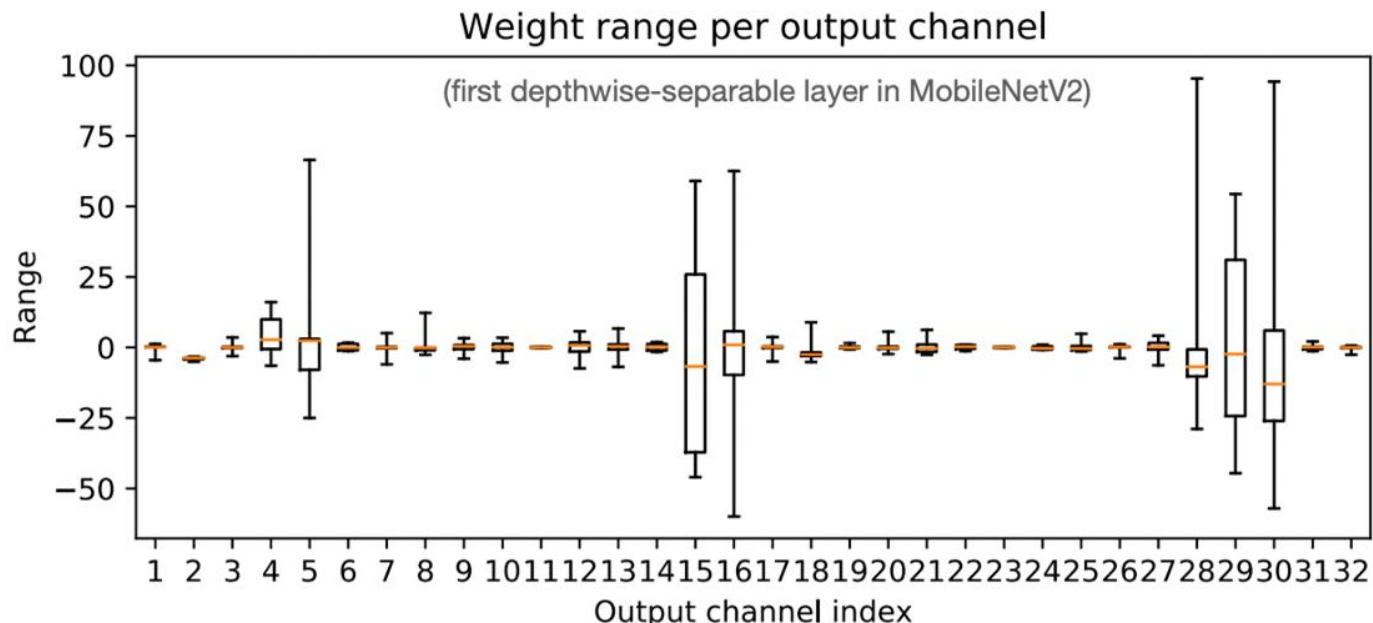
□ Efficient Inference Systems (vLLM, StreamingLLM, MHA/GQA/MQA)

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



Weight Granularity

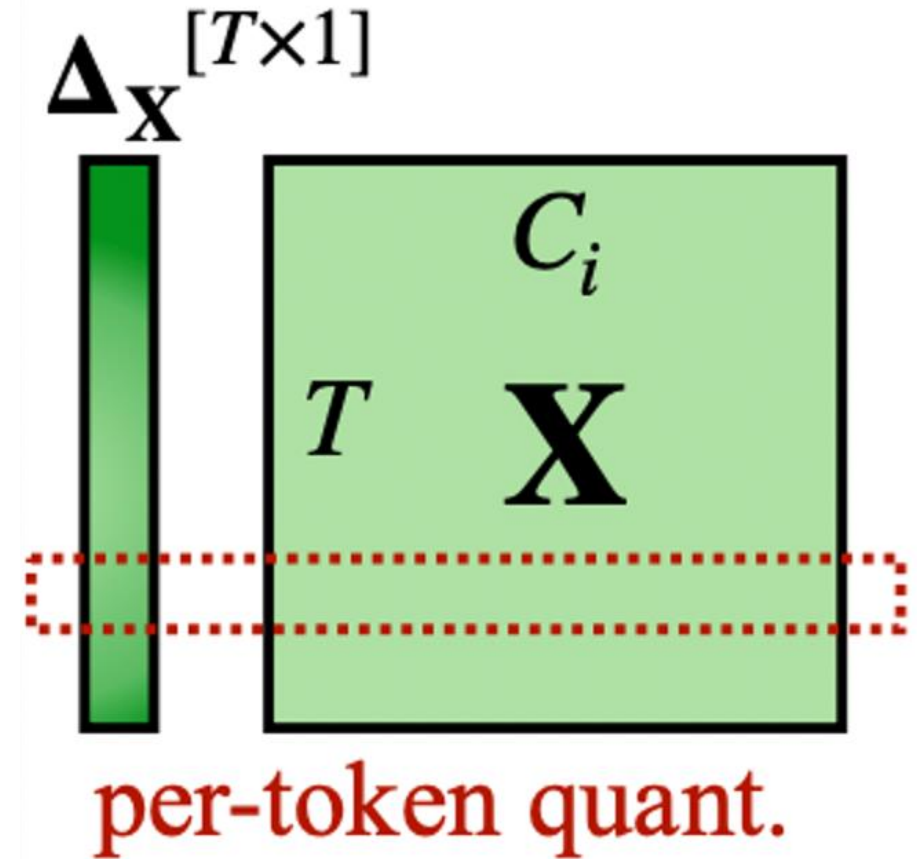
- ❑ Weight matrices will often have different variances along each output channel
- ❑ High variance in weights means that applying linear quantization will result in large performance degradation
- ❑ To fix this, we can perform linear quantization along each channel of the weight tensor separately



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)

□ Efficient Inference Systems (vLLM, StreamingLLM, MHA/GQA/MQA)

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



Quantization Aware Training (QAT)

Quantize while training

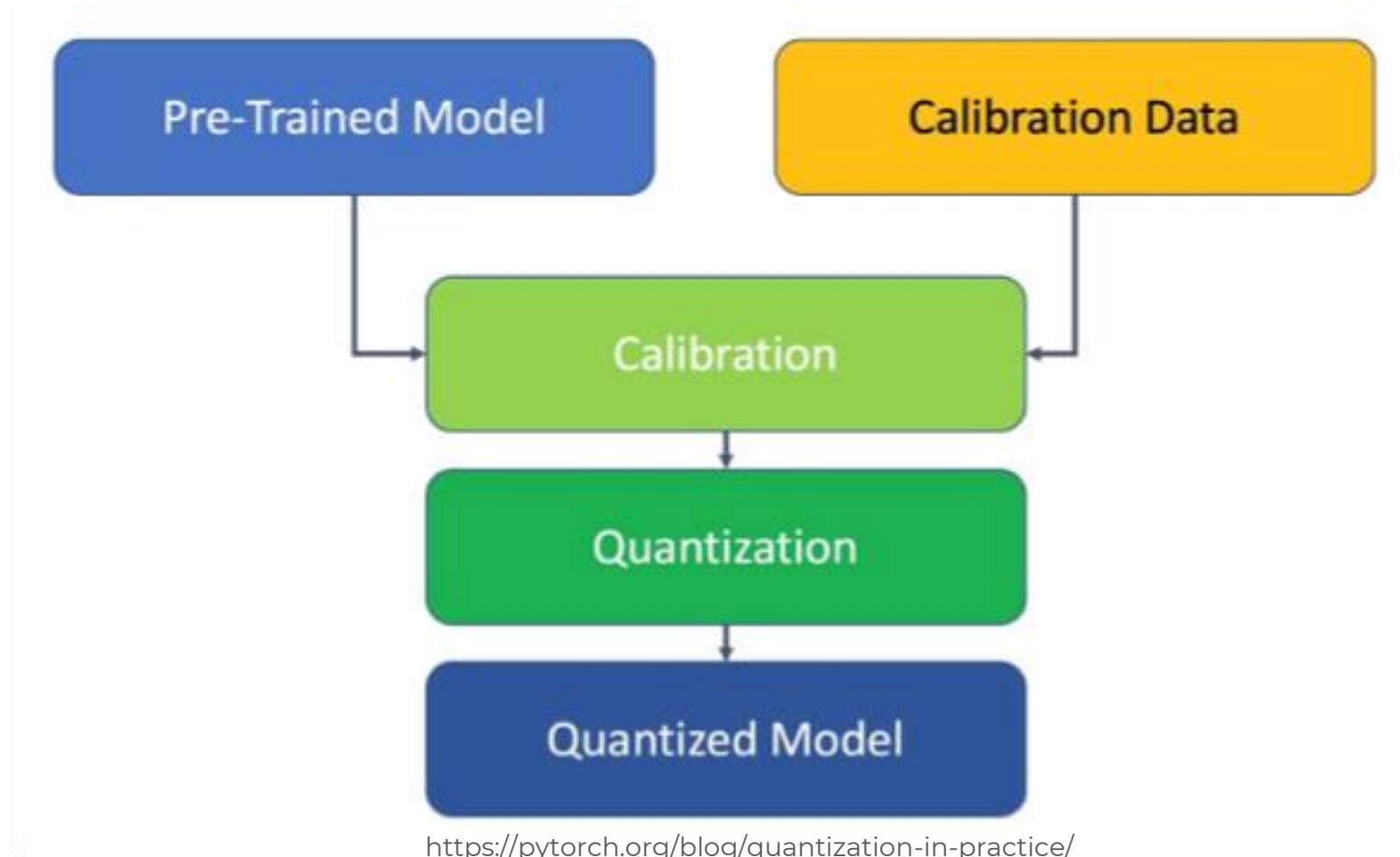


<https://pytorch.org/blog/quantization-in-practice/>



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)

□ Efficient Inference Systems (vLLM, StreamingLLM, MHA/GQA/MQA)

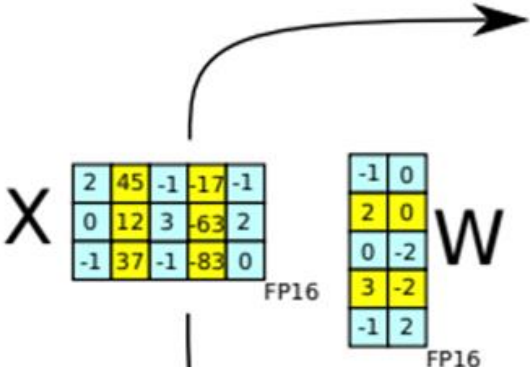
□ Parameter Efficient Fine-Tuning (BitFit, Adapter, Prompt Tuning, LoRA)



LLM.int8() (W8A8)

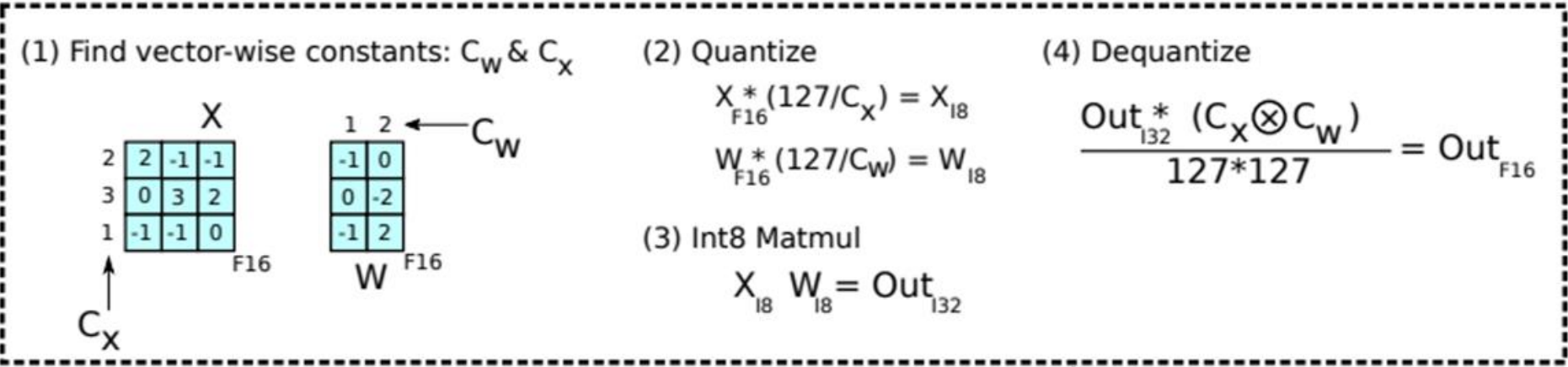
Stored 8-bit weights

LLM.int8()

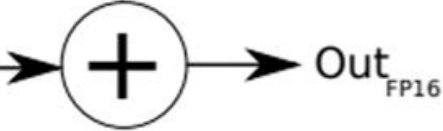
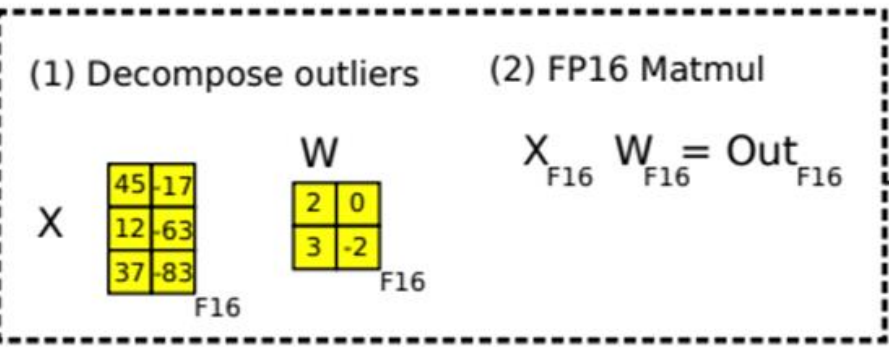


Regular values
Outliers

8-bit Vector-wise Quantization



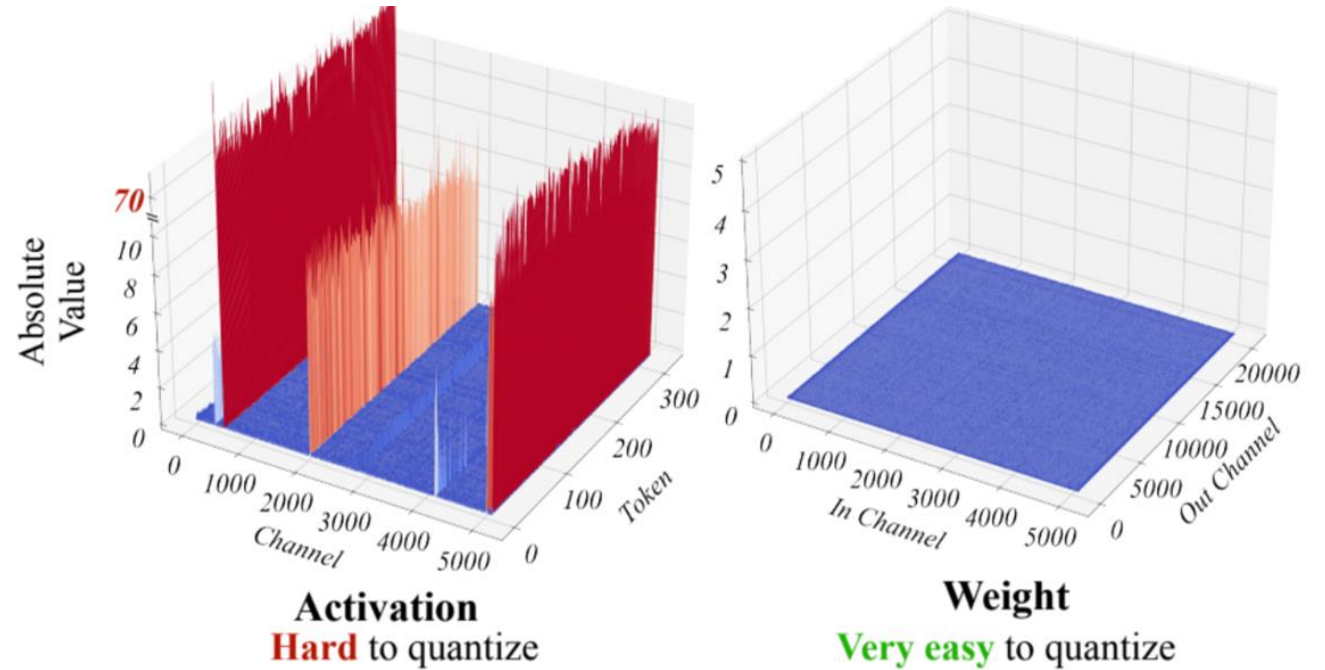
16-bit Decomposition



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

SmoothQuant (W8A8)

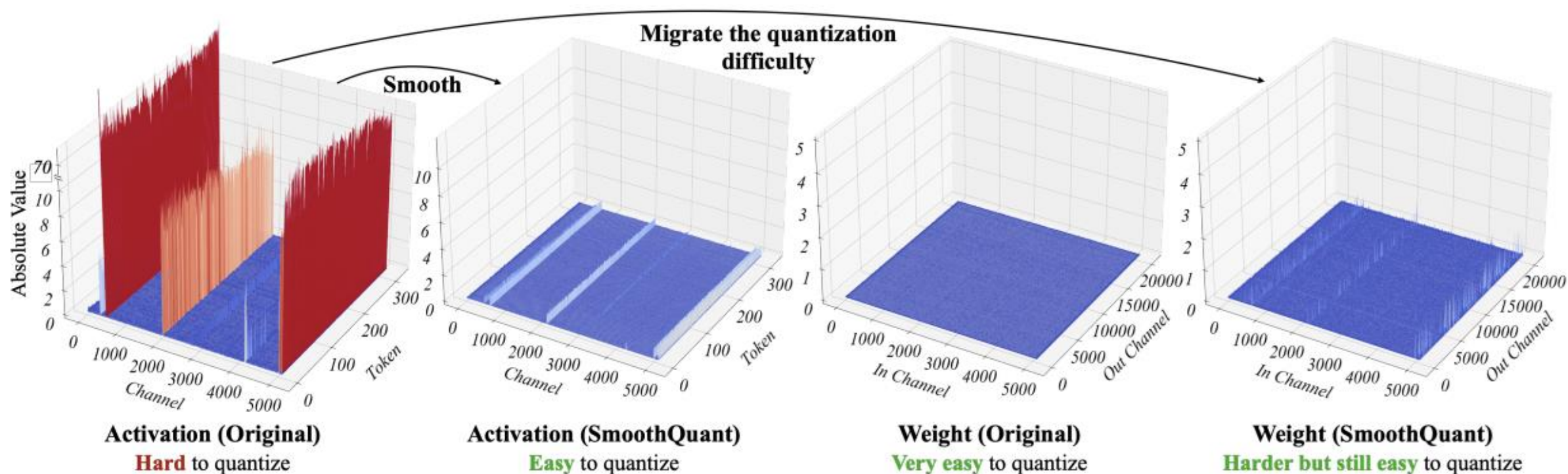
Observation: High variance channels are fixed in activations in LLM FFN layers-weights have relatively little difference in variance



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

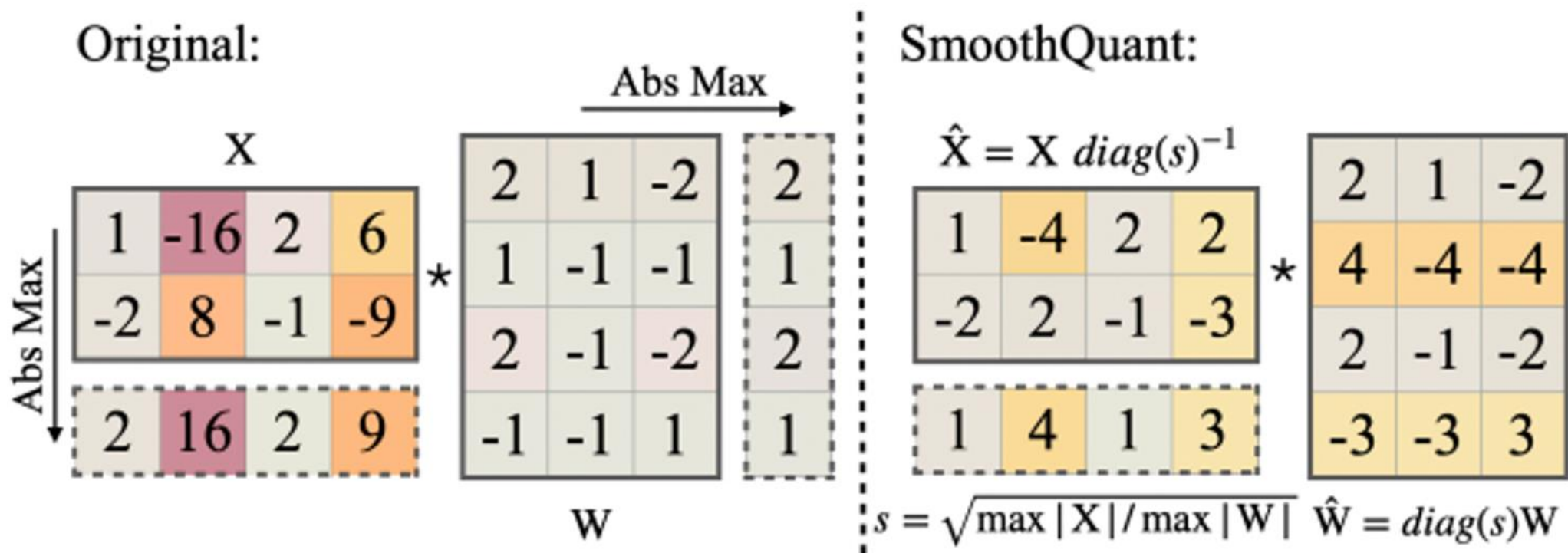


SmoothQuant (W8A8)



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

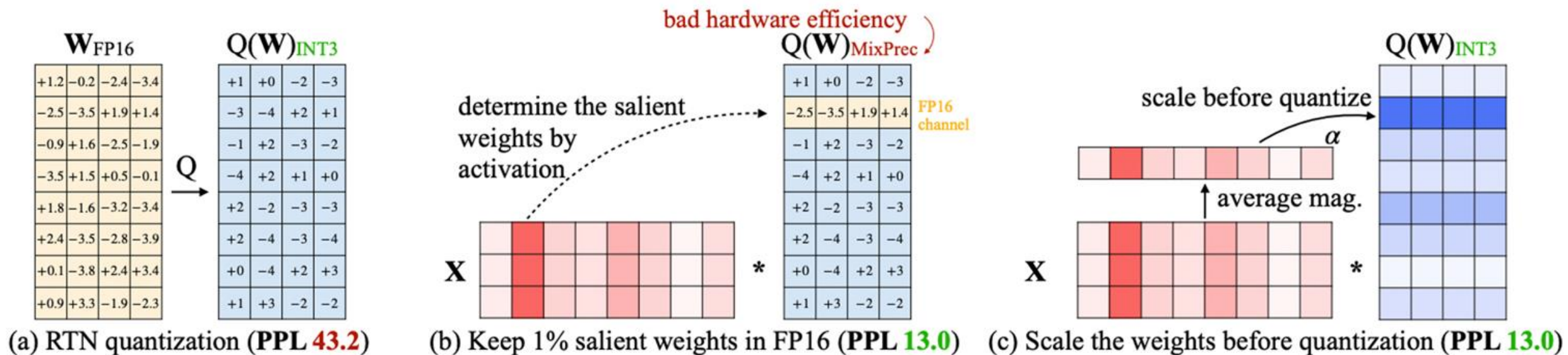
SmoothQuant (W8A8)



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



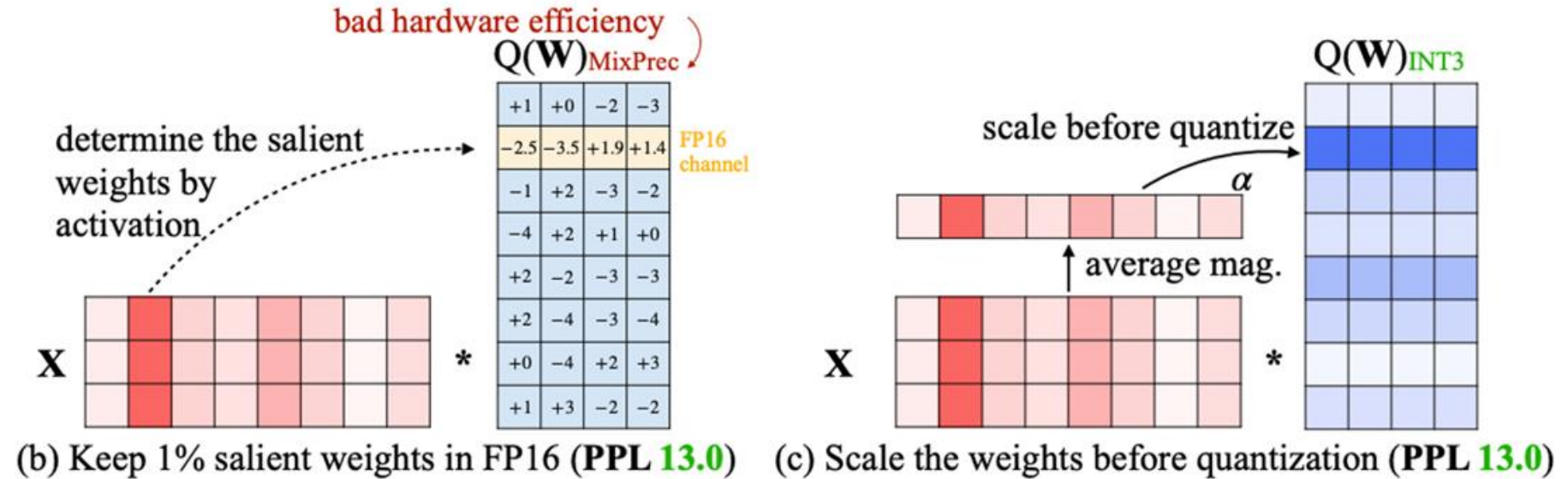
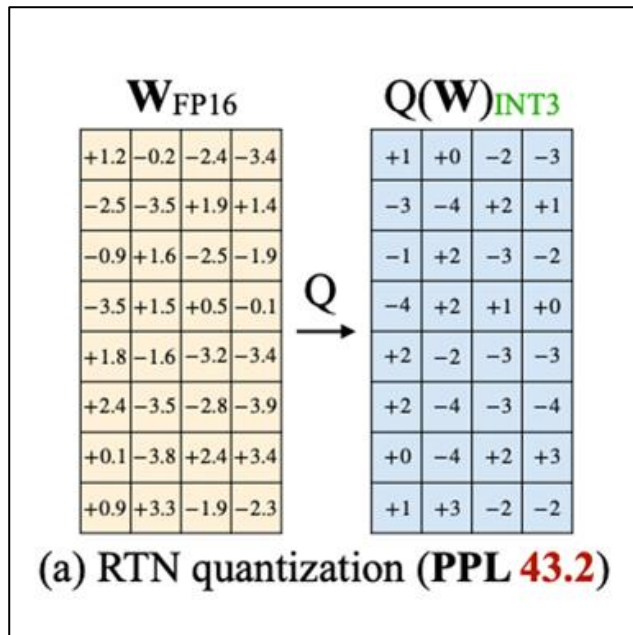
AWQ (W4A16)



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

AWQ (W4A16)

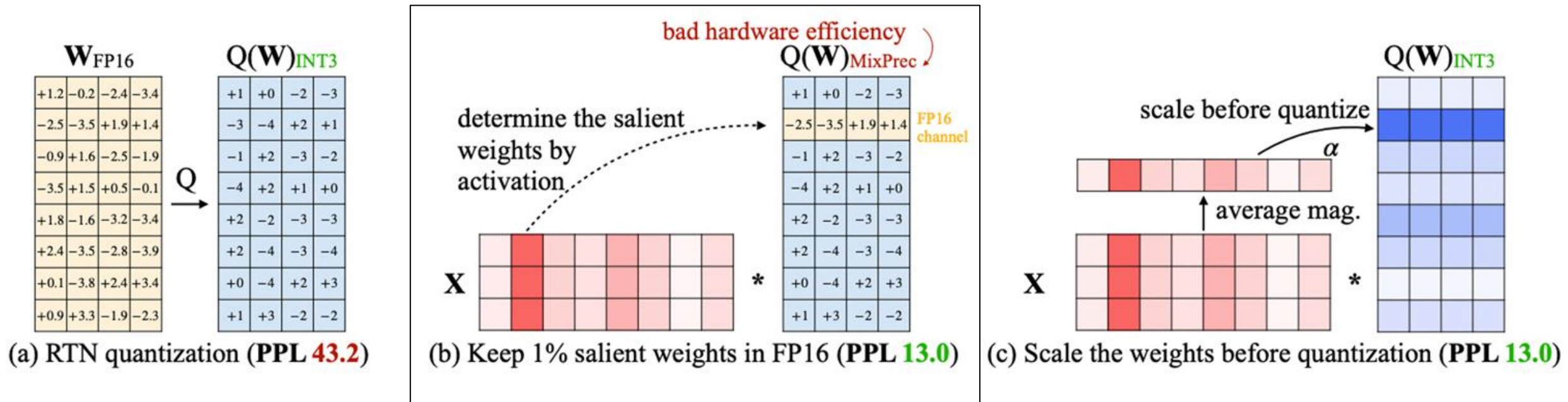
Normal quantization on LLMs performs poorly due to outliers in the model's hidden state



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

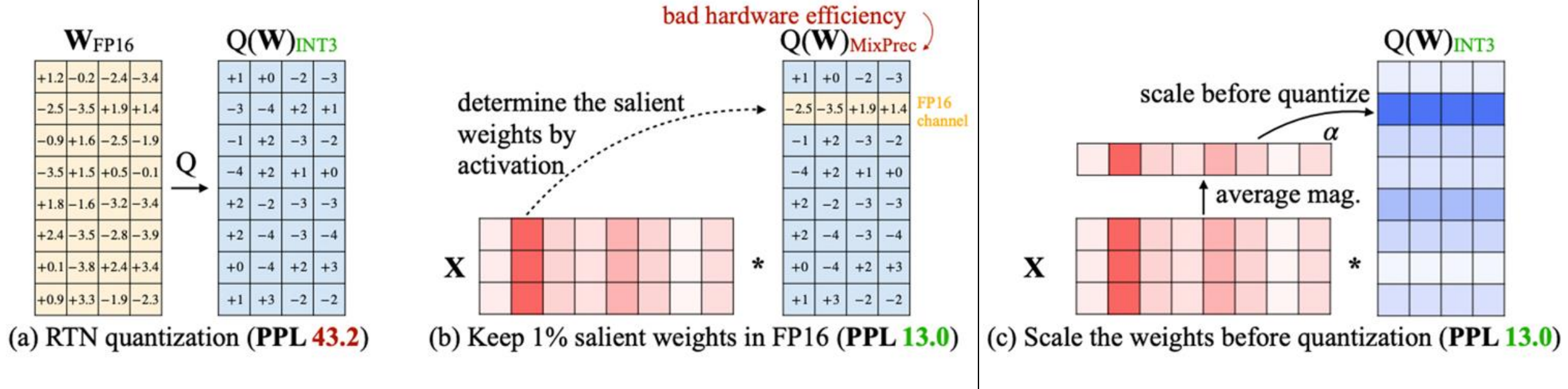


AWQ: Activation-aware Weight Quantization for LLM Compression and Acceleration [Lin et. al., arxiv 2023]



AWQ (W4A16)

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

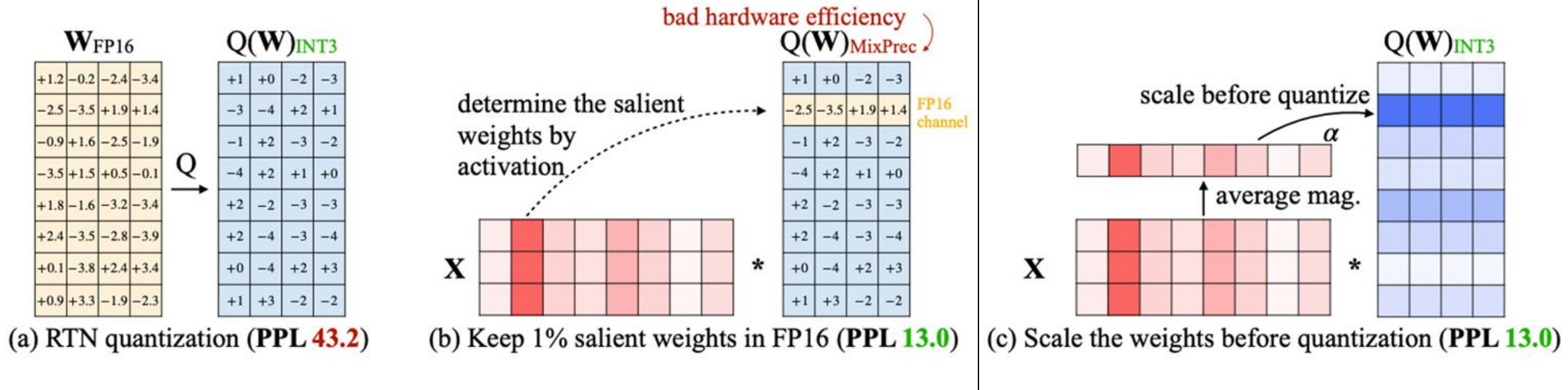


AWQ: Activation-aware Weight Quantization for LLM Compression and Acceleration [Lin et. al., arxiv 2023]



AWQ (W4A16)

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

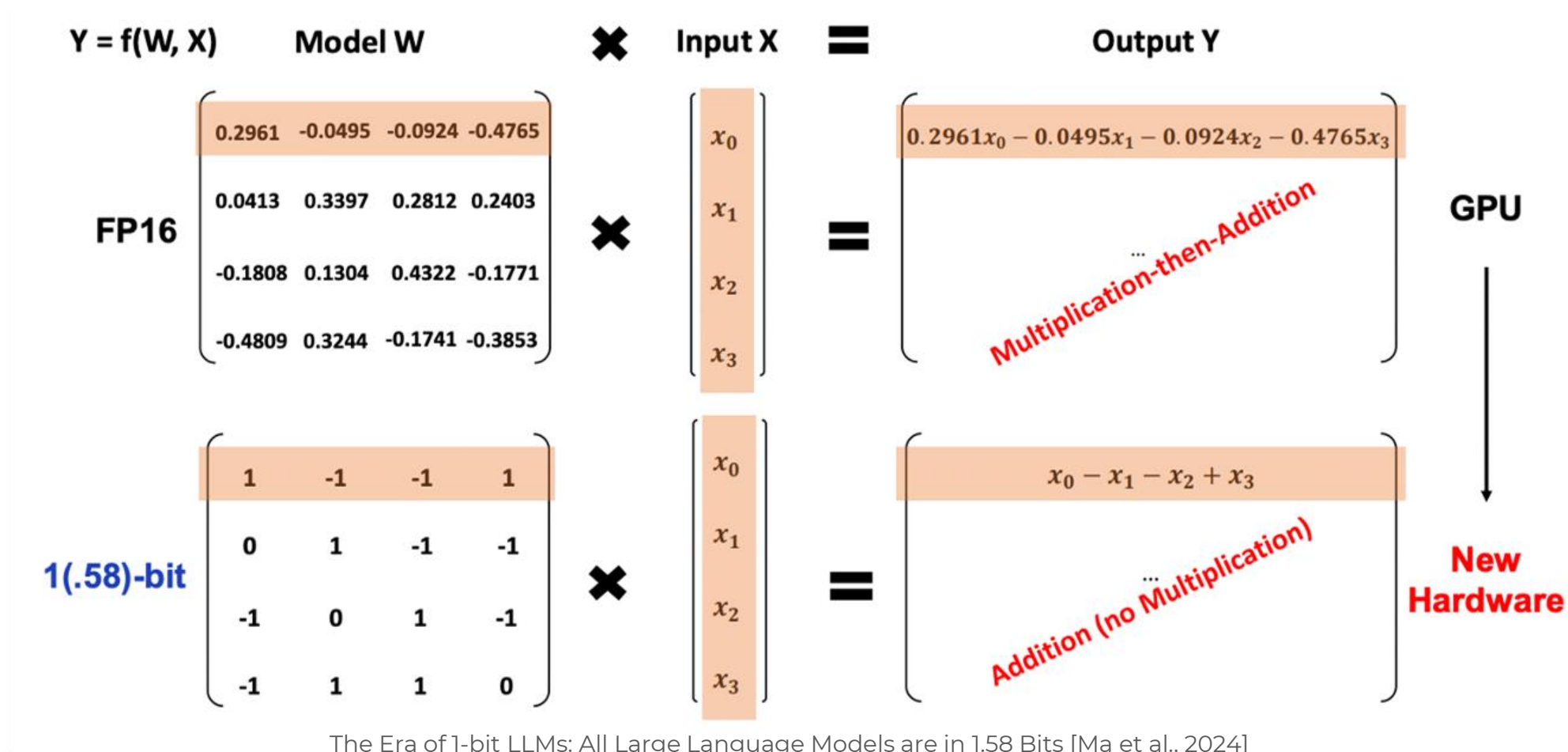


AWQ: Activation-aware Weight Quantization for LLM Compression and Acceleration [Lin et. al., arxiv 2023]



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)

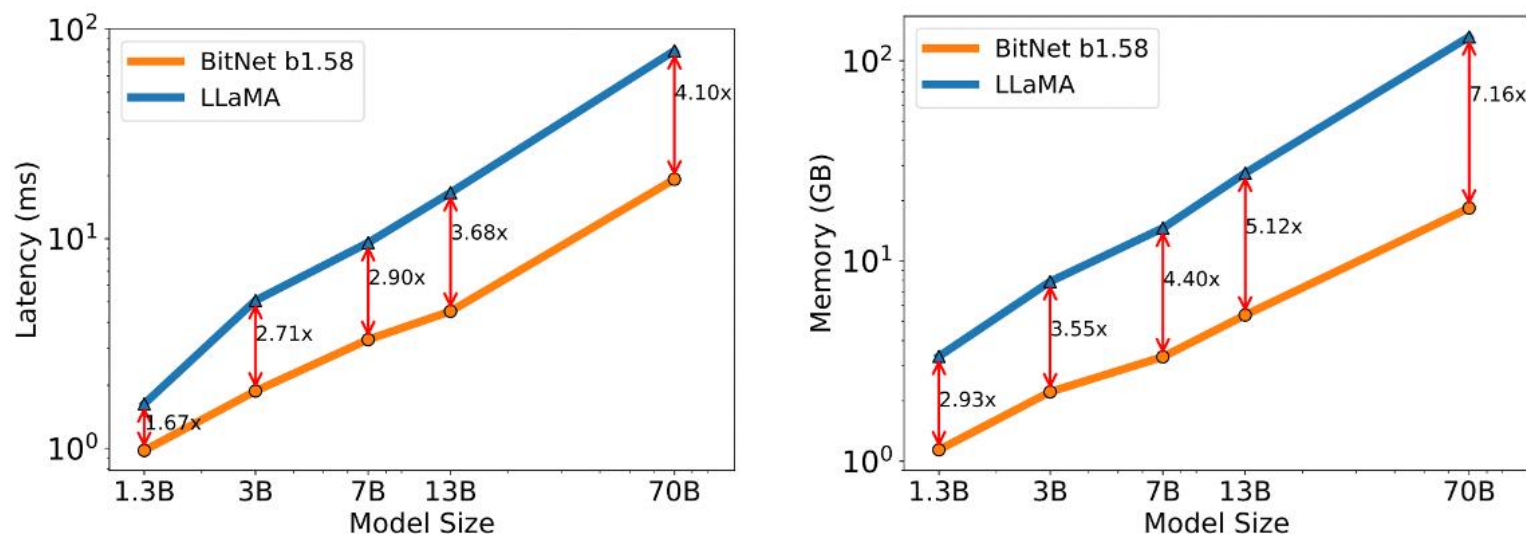


Figure 2: Decoding latency (Left) and memory consumption (Right) of BitNet b1.58 varying the model size.

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

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)**

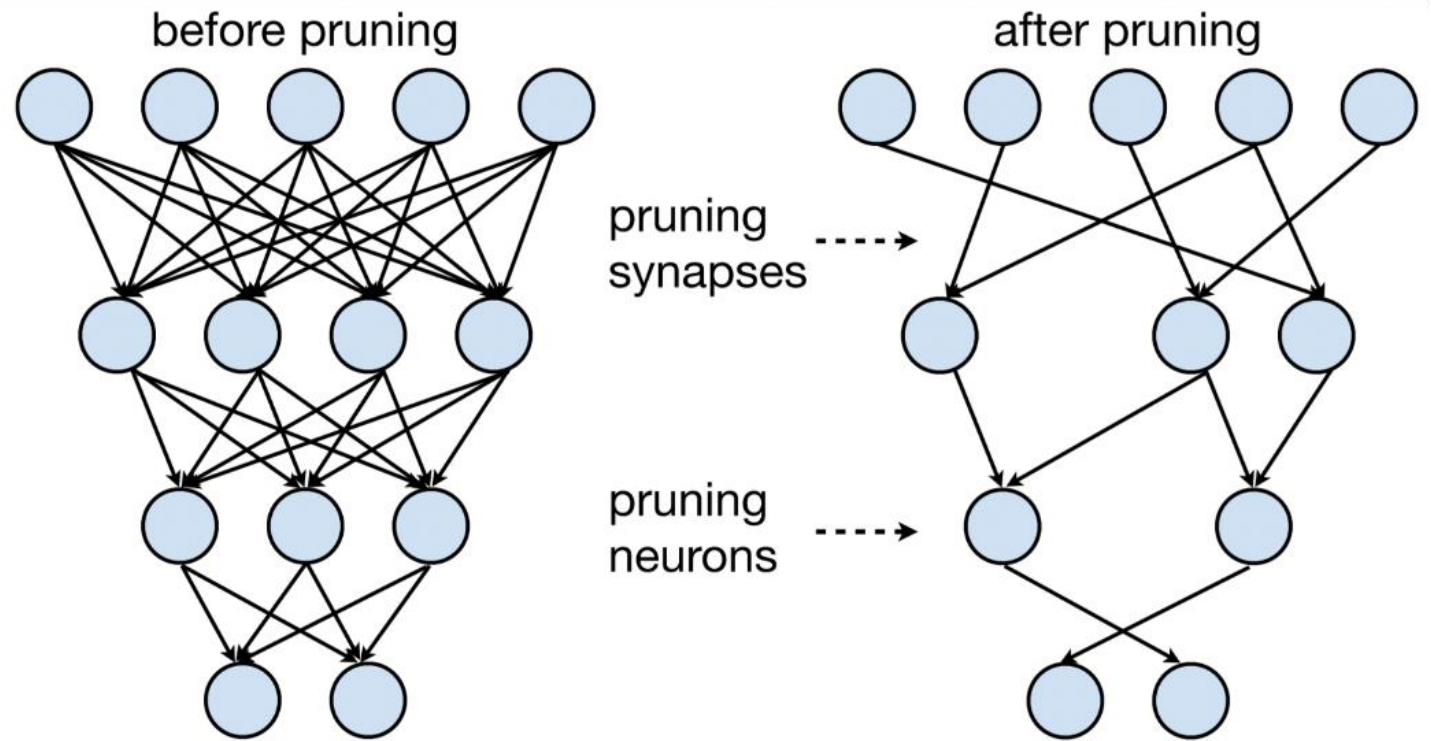
□ Efficient Inference Systems (vLLM, StreamingLLM, MHA/GQA/MQA)

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



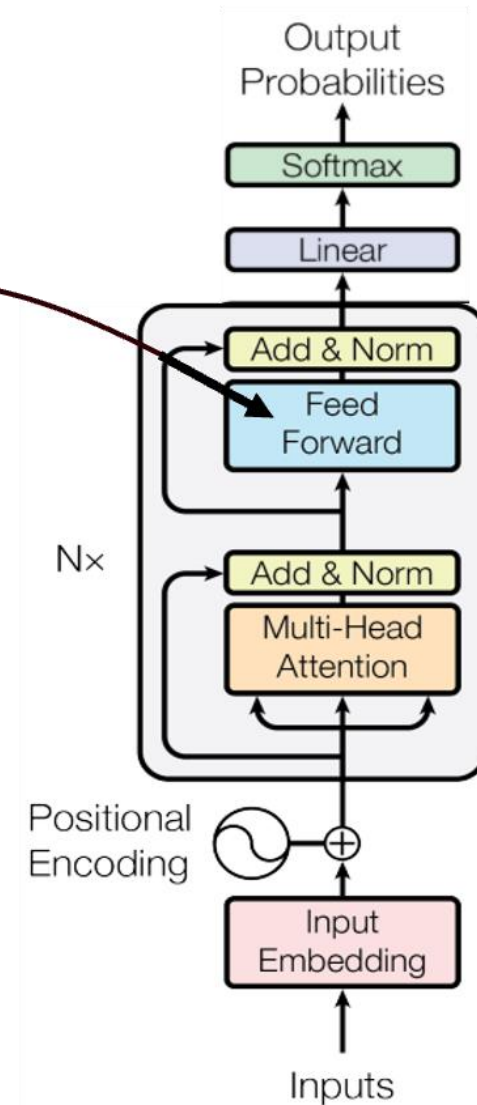
Sparsity

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

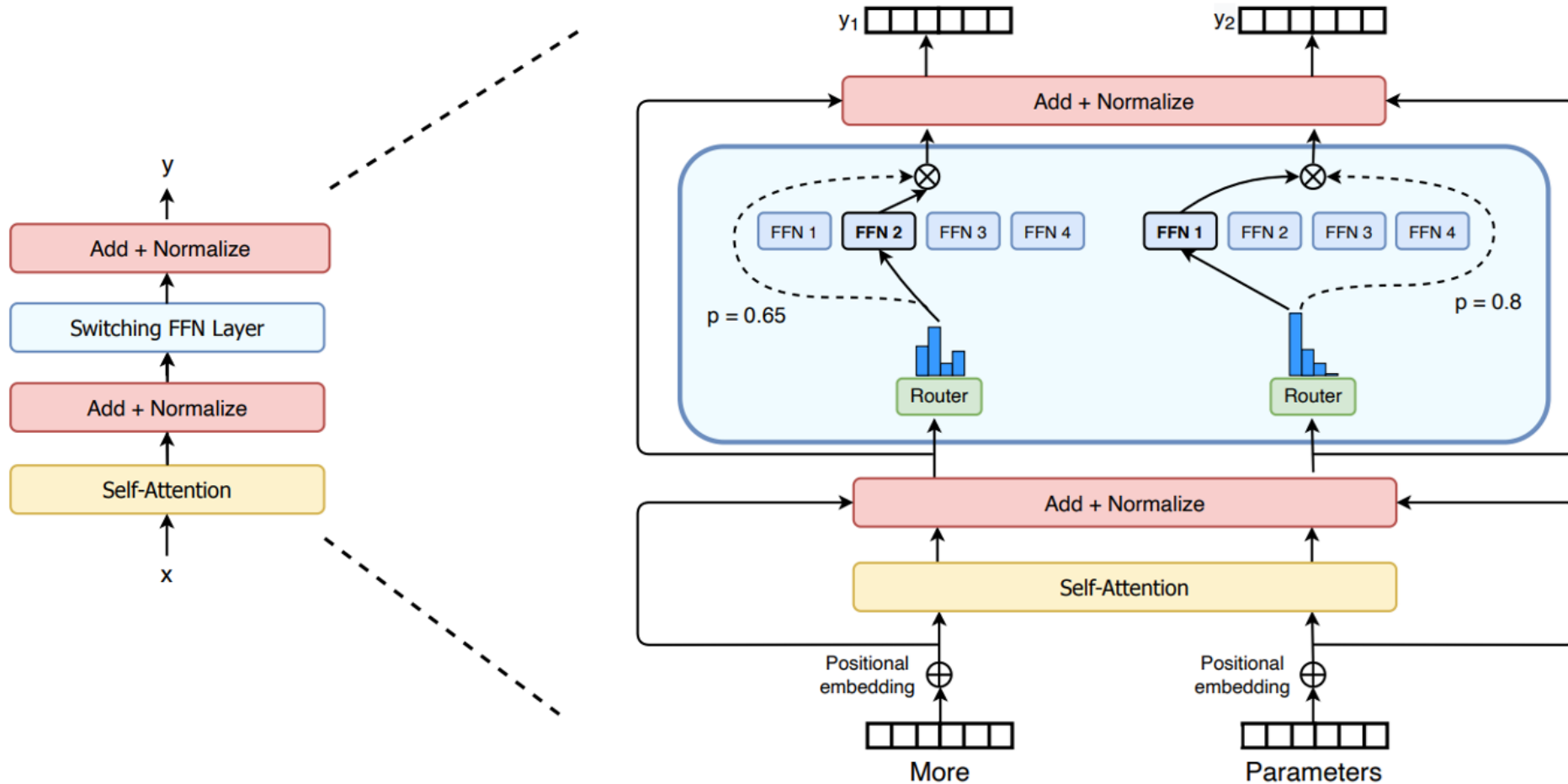


Mixture of Experts (MoE)

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



Mixture of Experts (MoE)

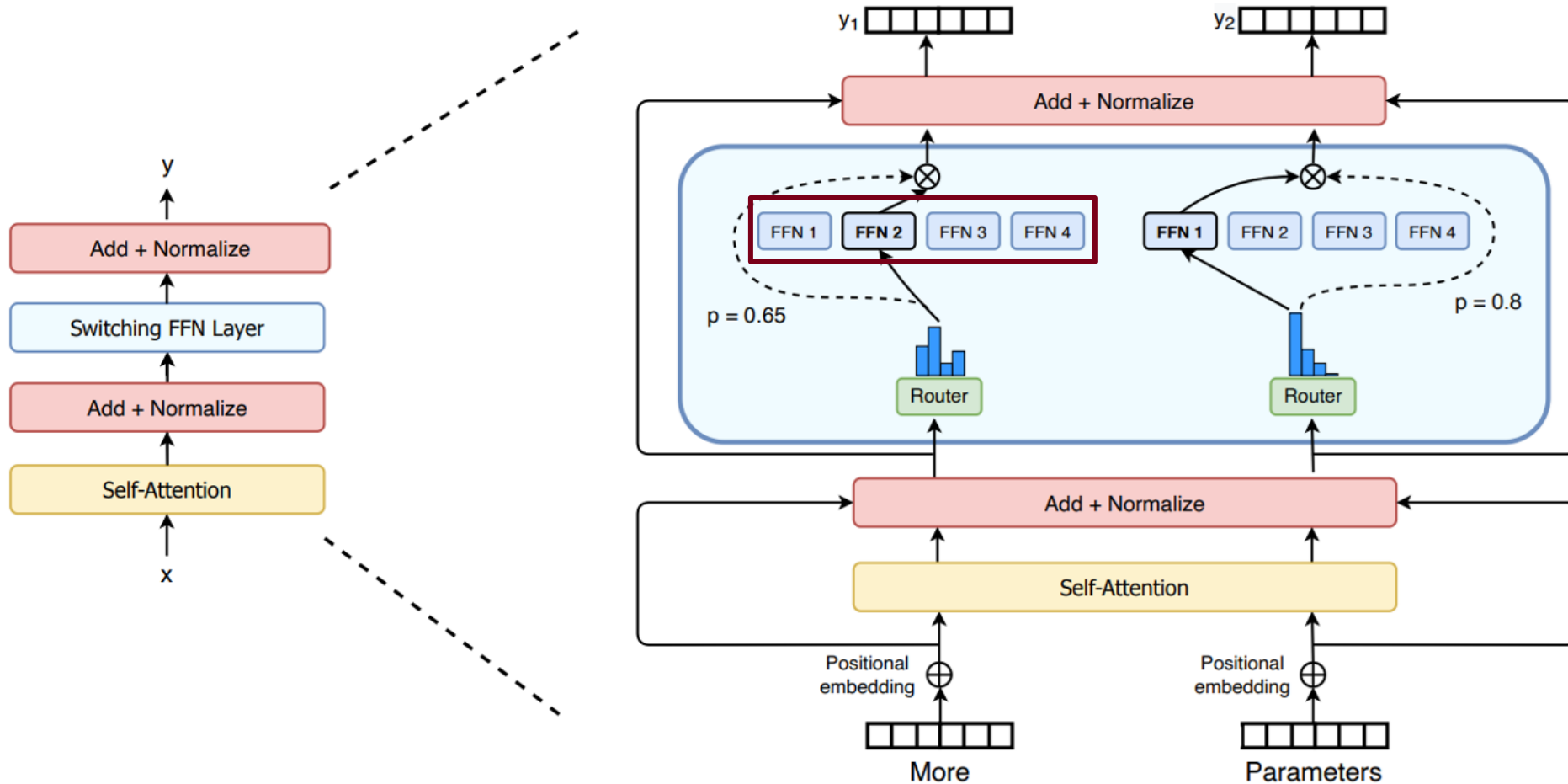


Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity [Fedus et al., CoRR 2021]



Mixture of Experts (MoE)

Four FFN layers

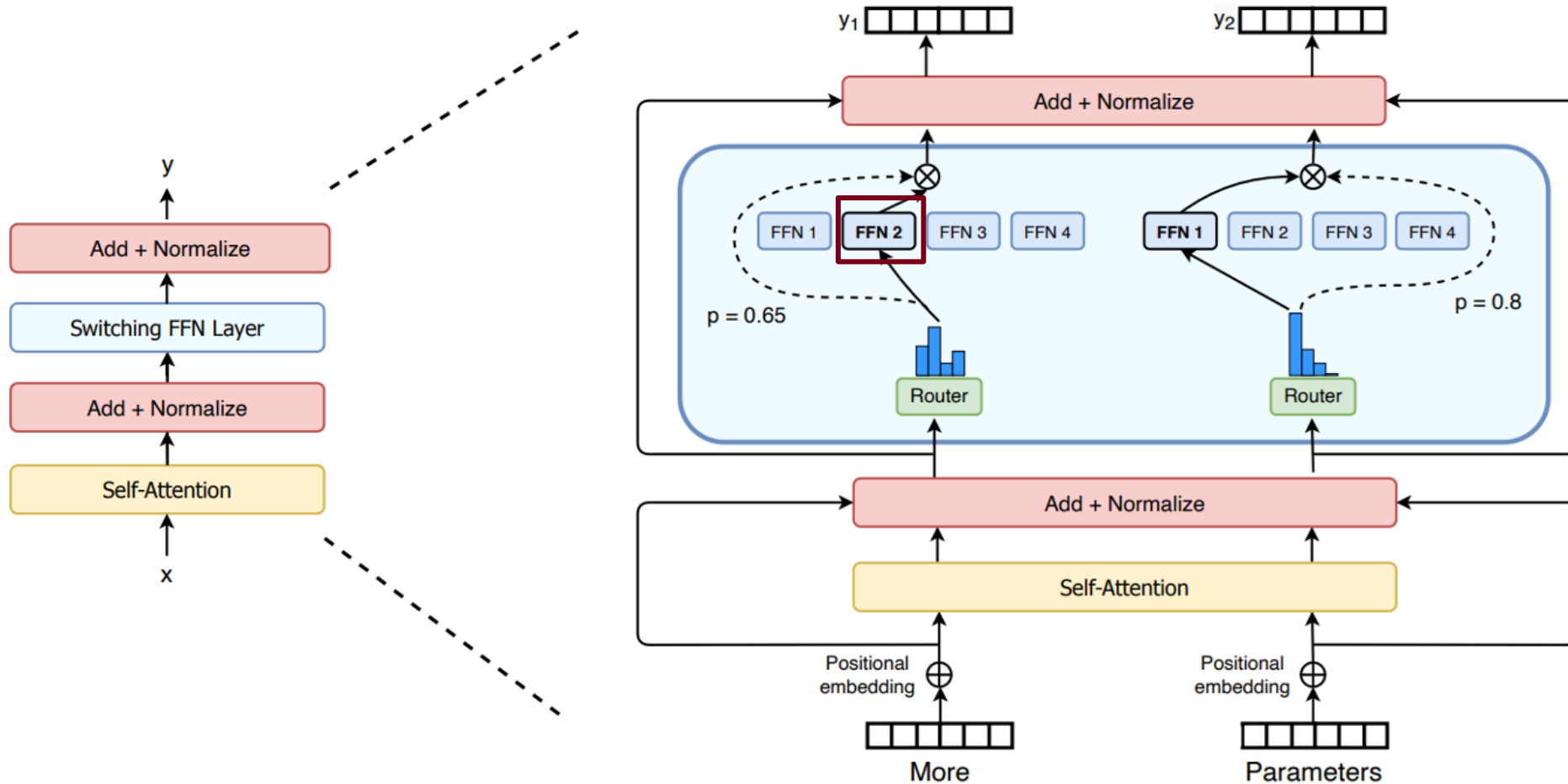


Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity [Fedus et al., CoRR 2021]



Mixture of Experts (MoE)

Only one is used per token

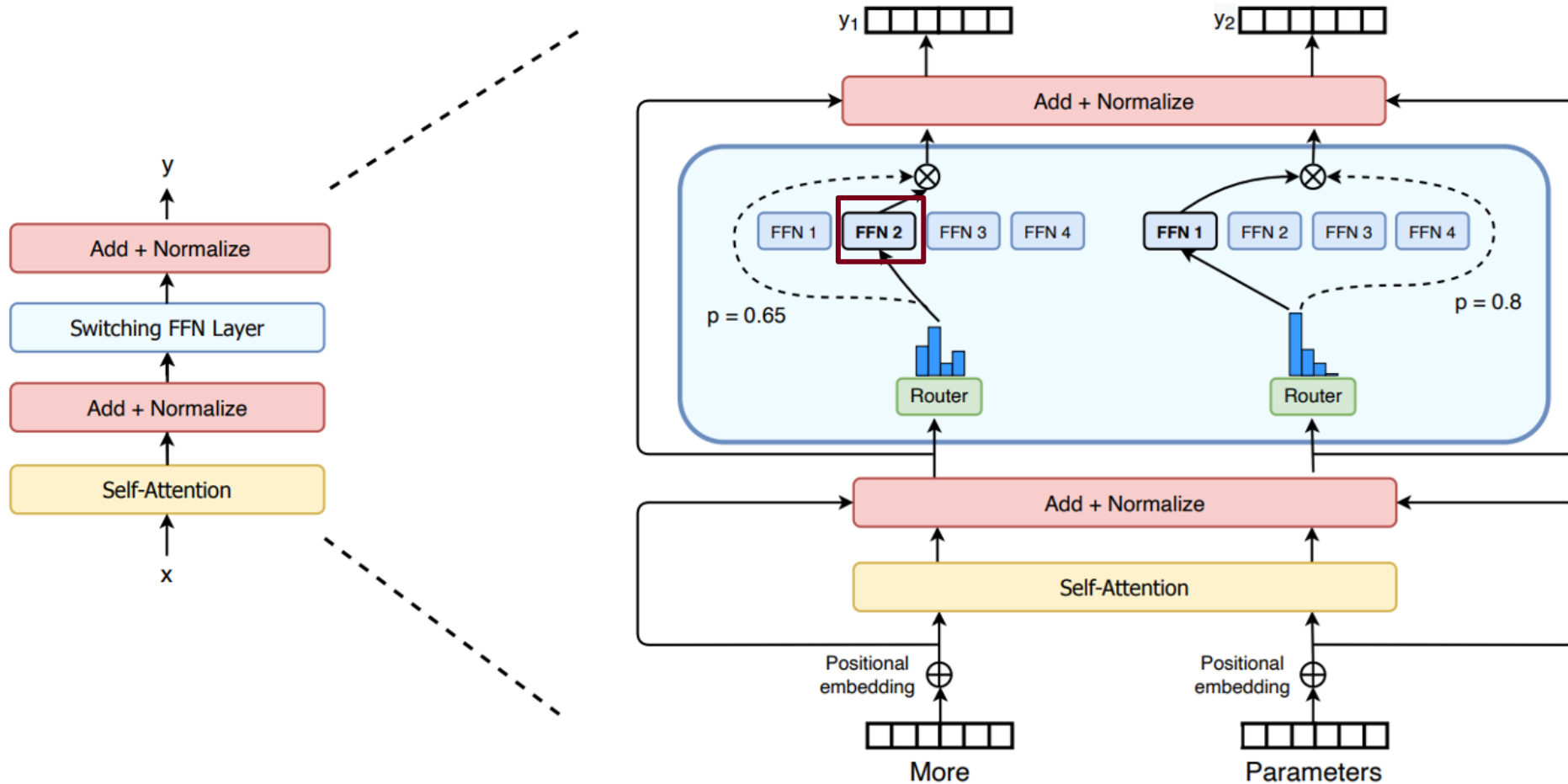


Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity [Fedus et al., CoRR 2021]



Mixture of Experts (MoE)

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

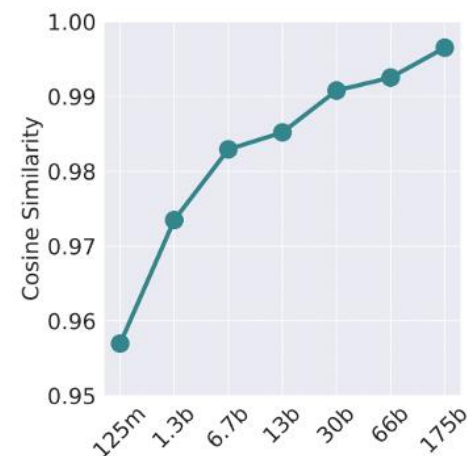


Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity [Fedus et al., CoRR 2021]

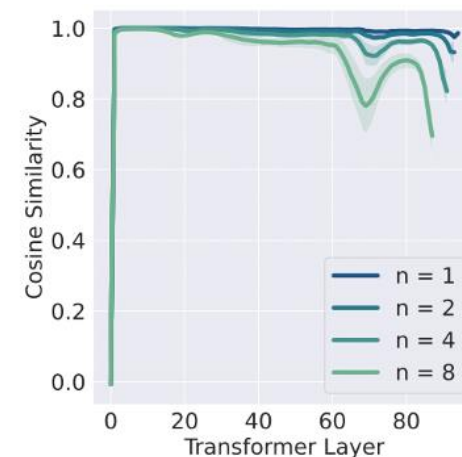


Deja Vu: Contextual Sparsity

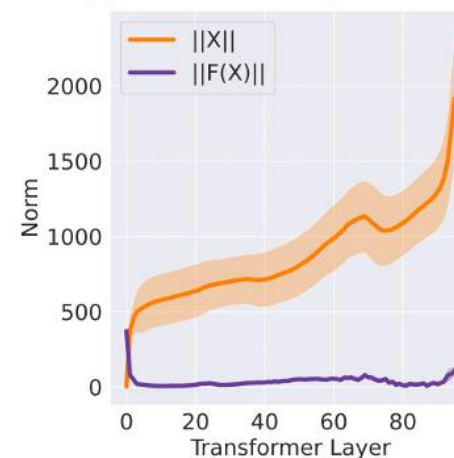
Observation 1: Model activations change very little between consecutive layers of a network



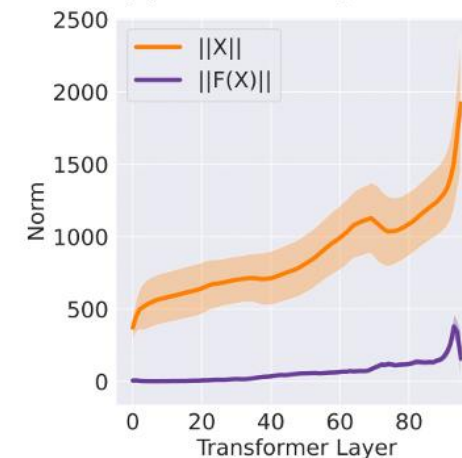
(a) Model Comparison



(b) Across Layer



(c) Residual Around Attention



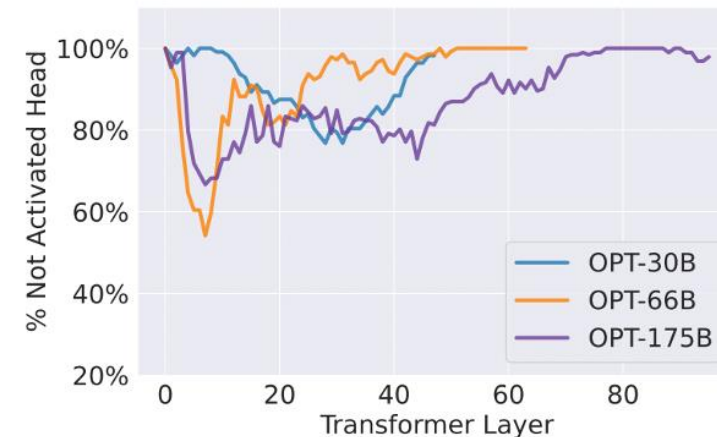
(d) Residual Around MLP

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

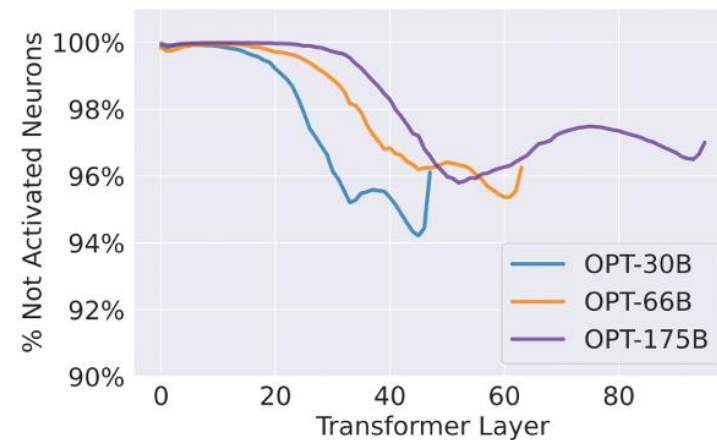


Deja Vu: Contextual Sparsity

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



(a) Contextual sparsity in Attention Head

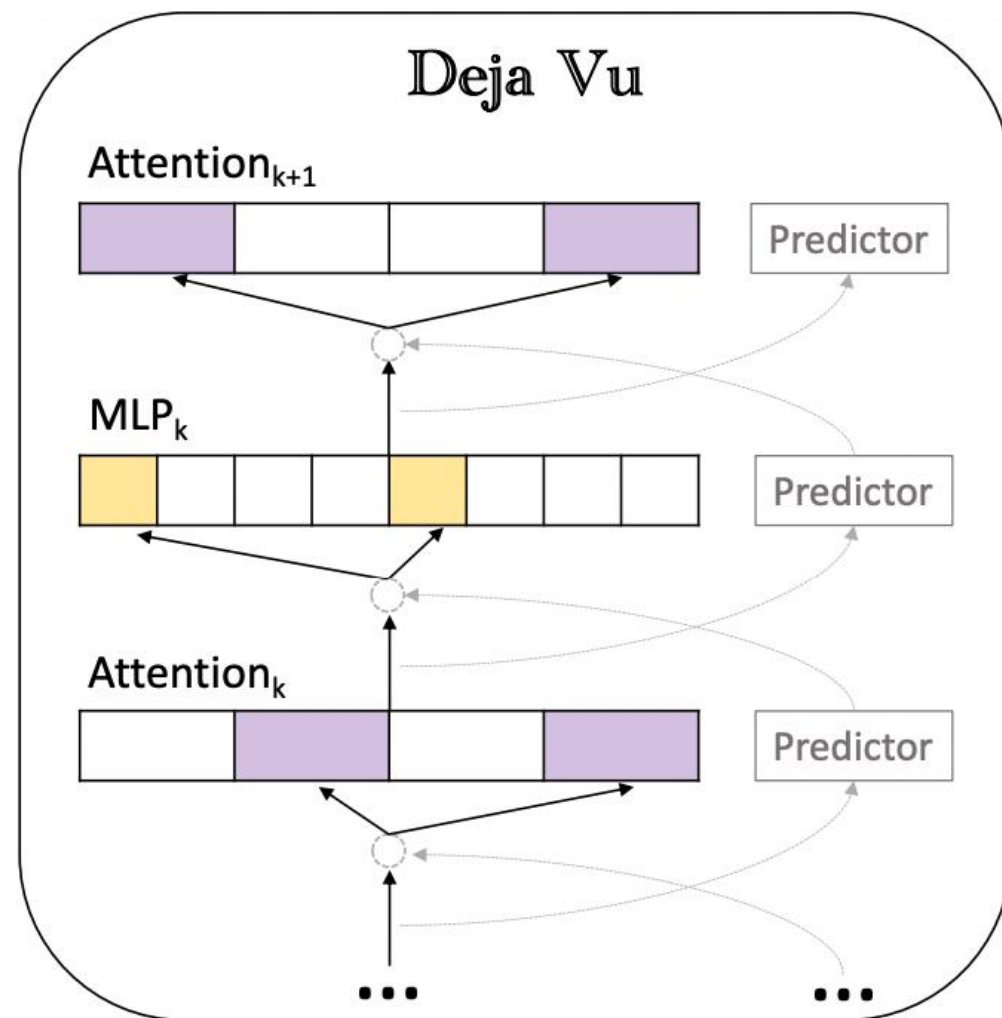


(b) Contextual sparsity in MLP Block



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)

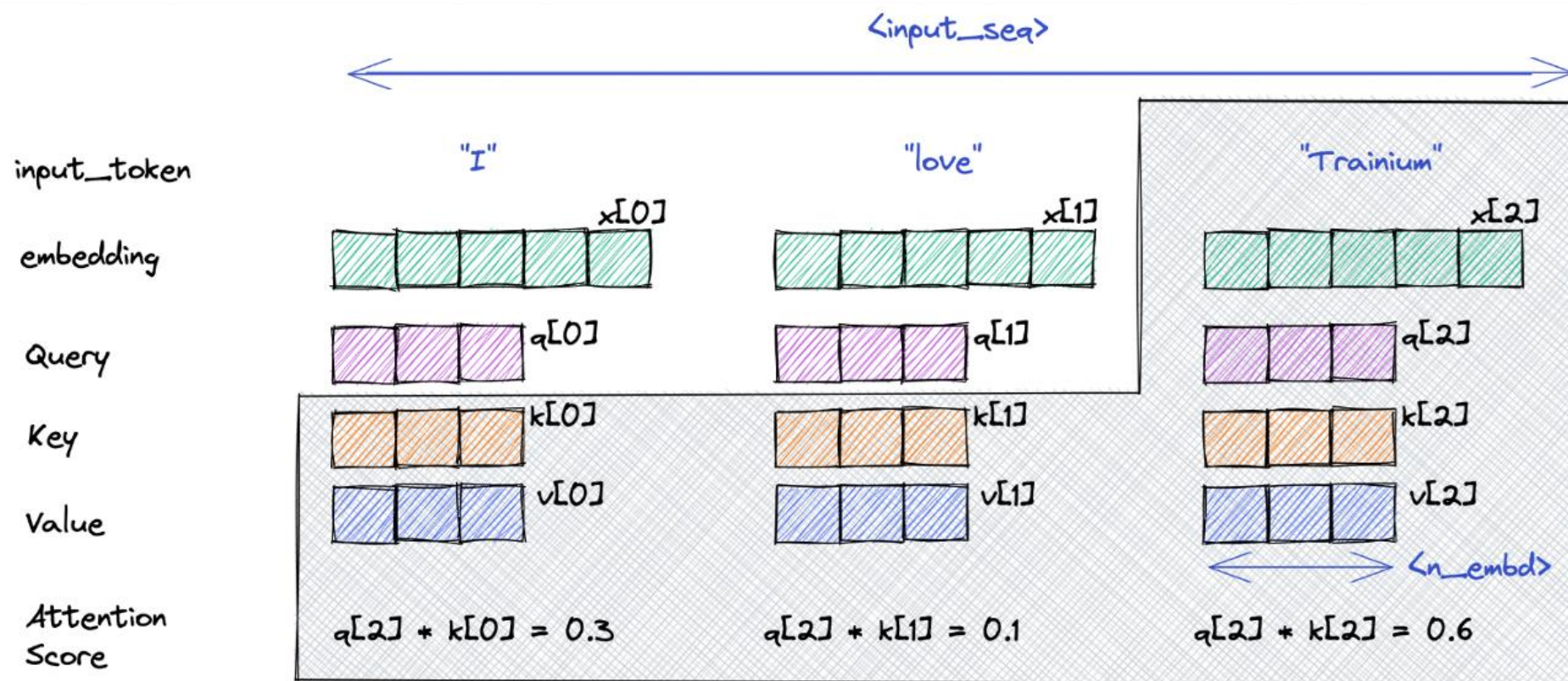
□ **Efficient Inference Systems (vLLM, StreamingLLM, MHA/GQA/MQA)**

□ 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

$$\text{Batch_size} * \text{seq_len} * \text{num_heads} * \text{num_layers} * \text{emb_dim} * 2$$

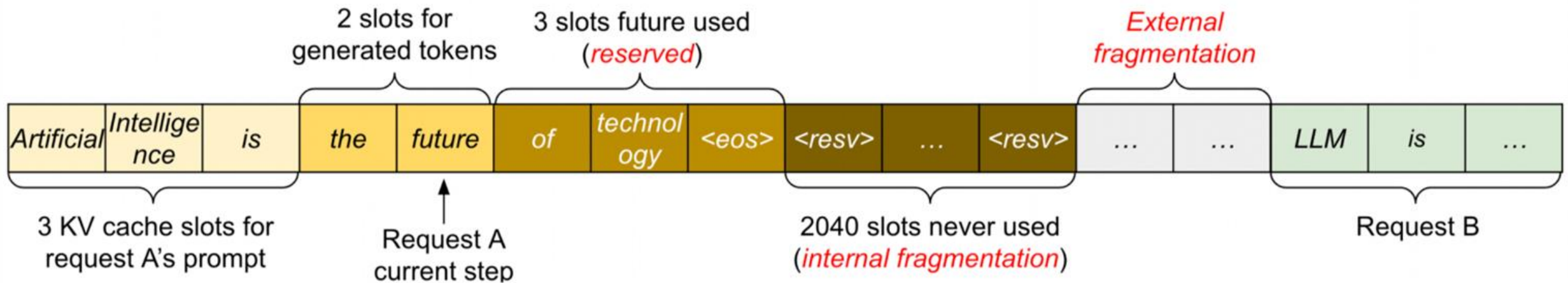
separate values in the kv cache



vLLM

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?



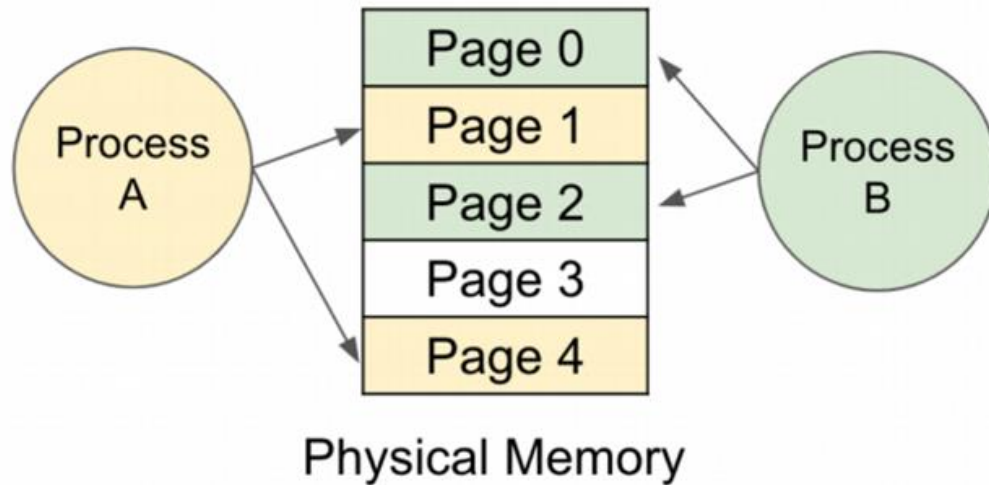
Efficient Memory Management for Large Language Model Serving with PagedAttention (Kwon et al., 2023)



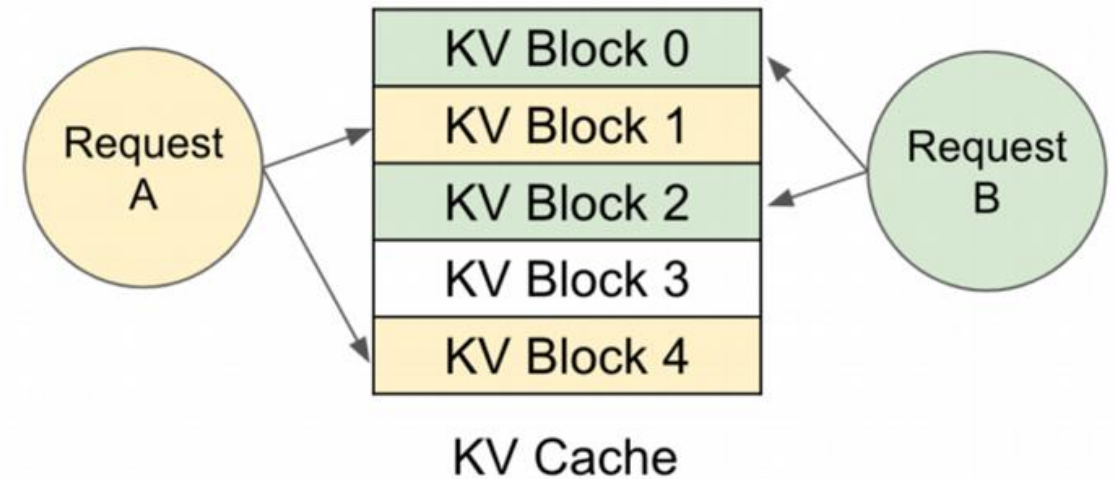
vLLM

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

Memory management in OS

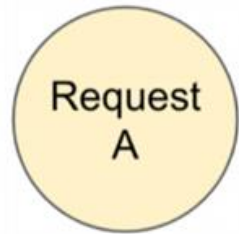


Memory management in vLLM



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

vLLM



Block Table

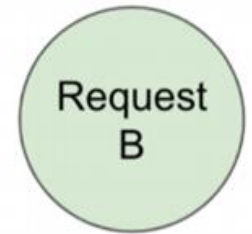
Logical KV blocks

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

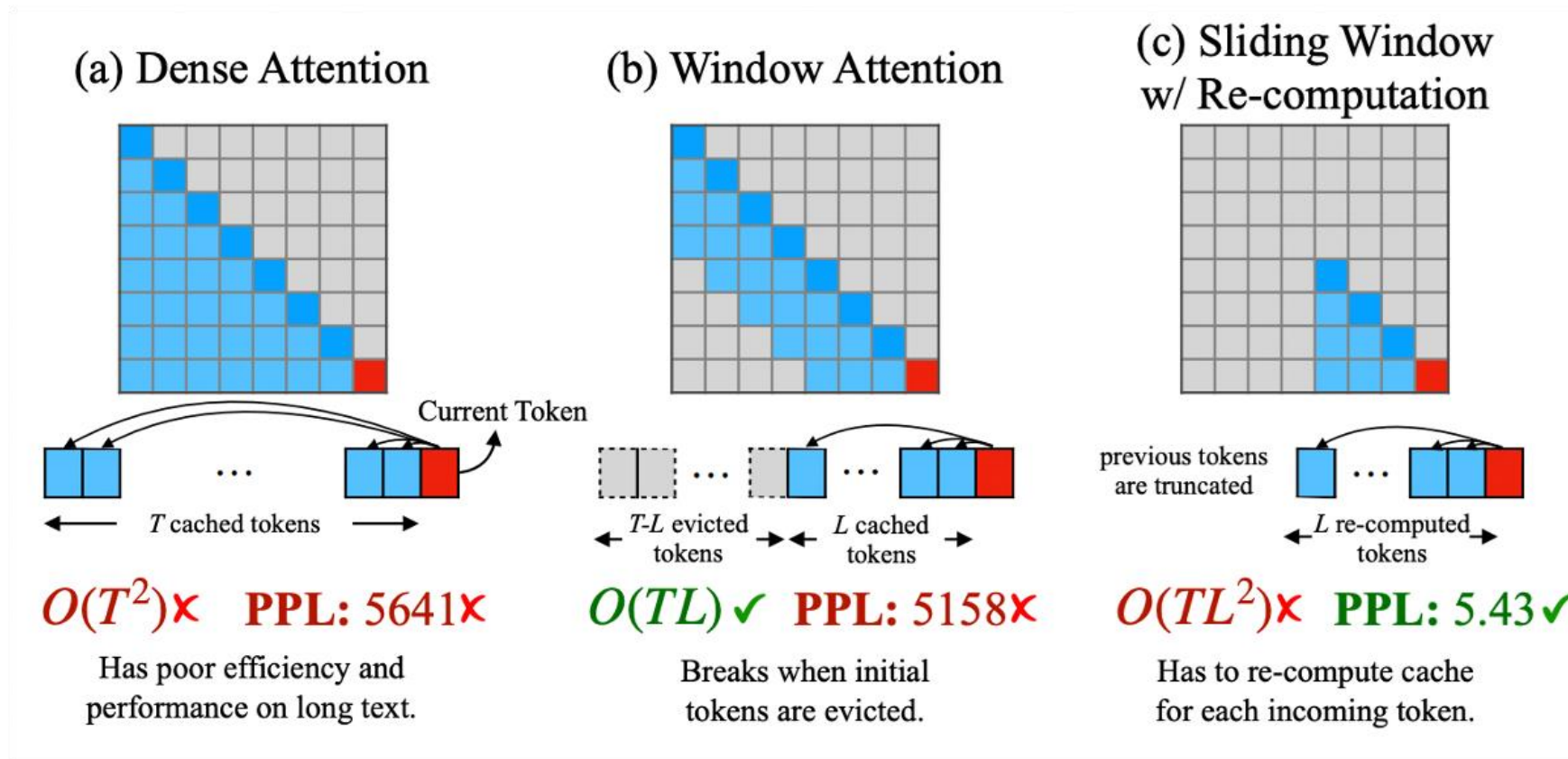
Artificial	Intelligence	is	the
future	of	technology	

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?

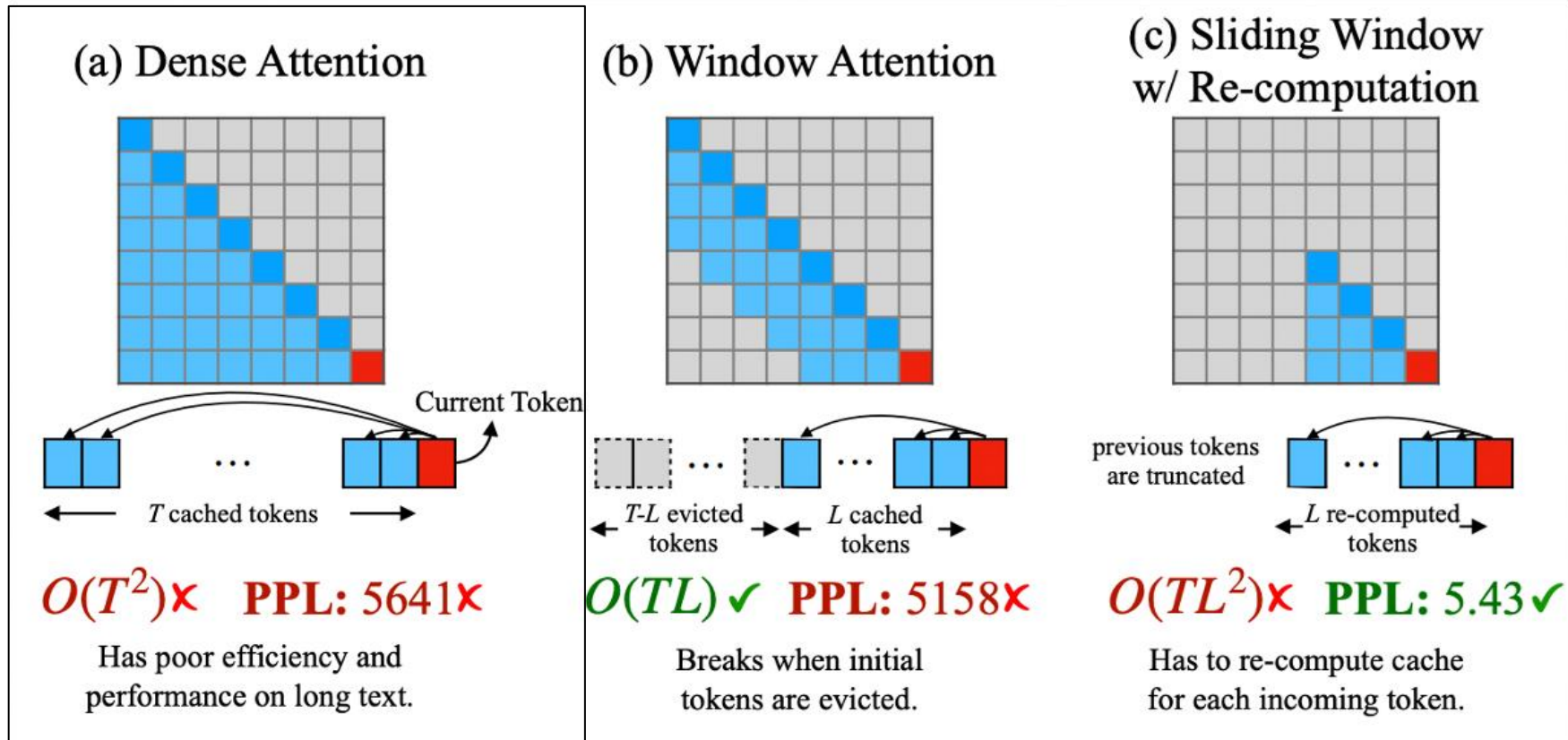


Efficient Streaming Language Models with Attention Sinks [Xiao et al., 2023]

StreamingLLM

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

Too much storage



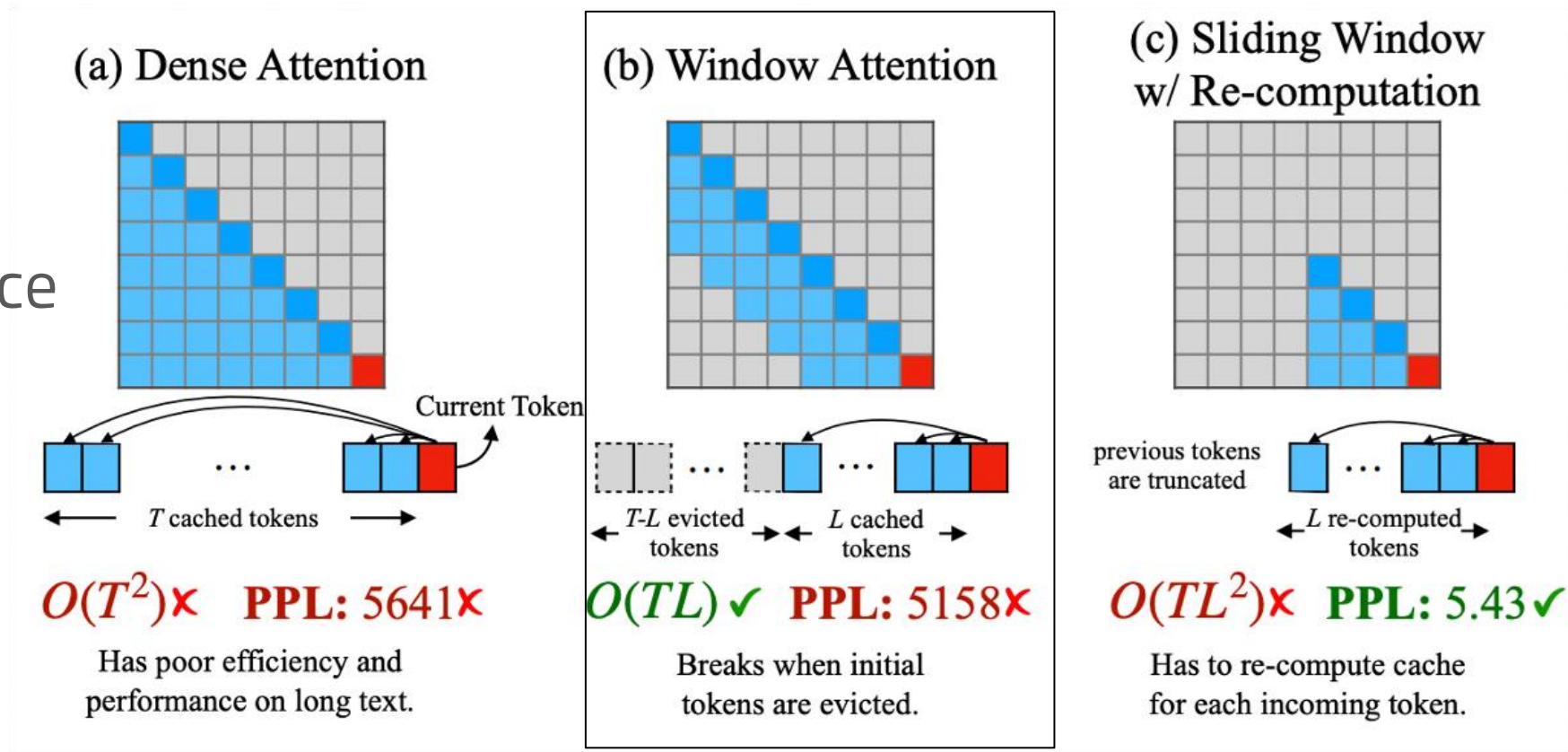
Efficient Streaming Language Models with Attention Sinks [Xiao et al., 2023]



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]

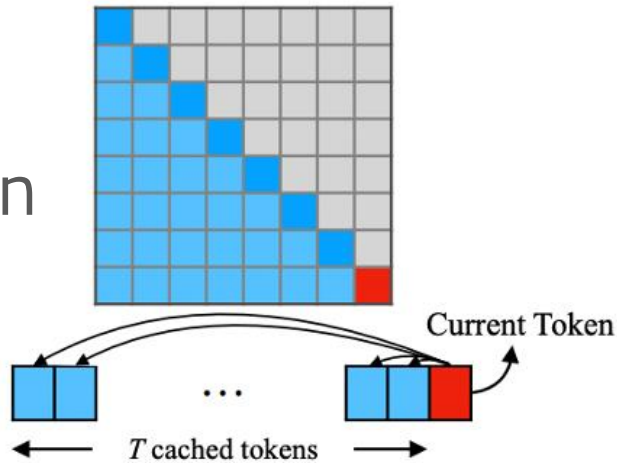


StreamingLLM

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

Too much
recomputation

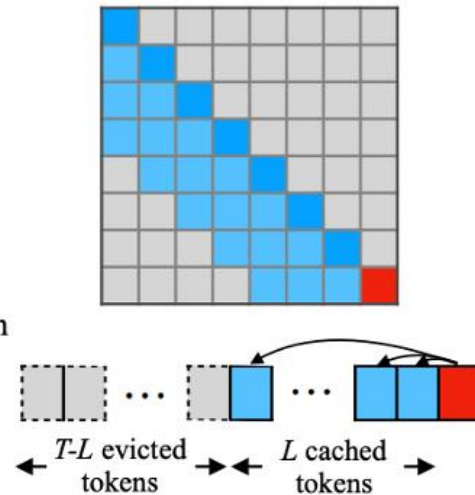
(a) Dense Attention



$O(T^2)$ ✗ PPL: 5641 ✗

Has poor efficiency and performance on long text.

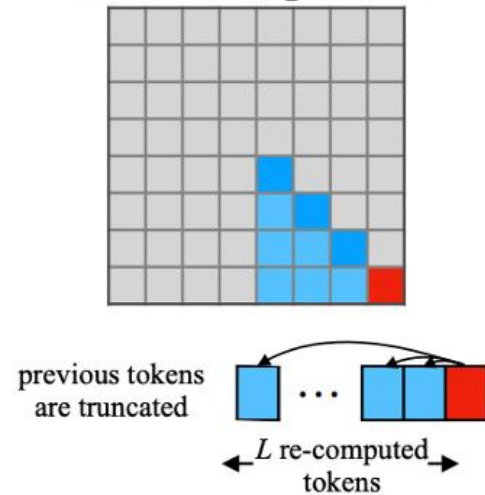
(b) Window Attention



$O(TL)$ ✓ PPL: 5158 ✗

Breaks when initial tokens are evicted.

(c) Sliding Window w/ Re-computation



$O(TL^2)$ ✗ PPL: 5.43 ✓

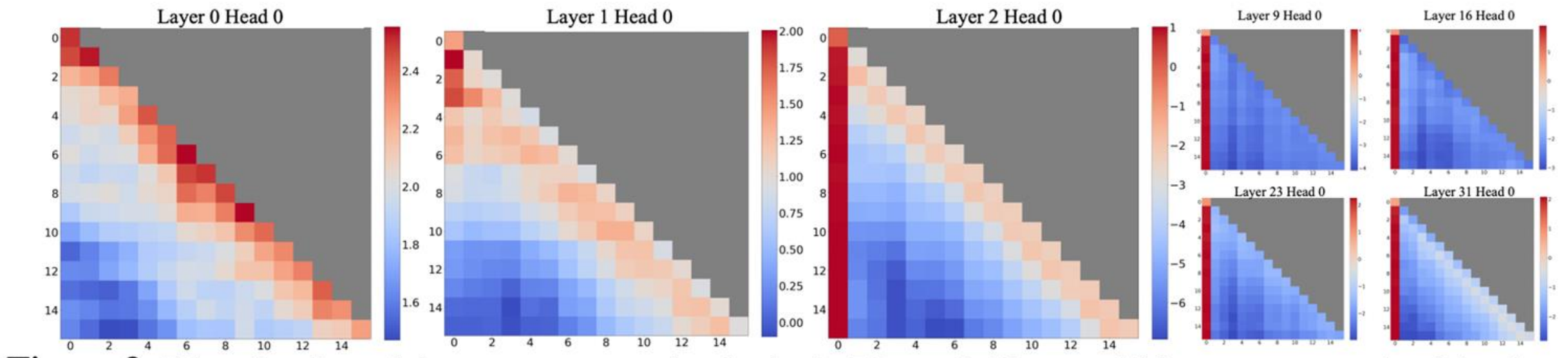
Has to re-compute cache for each incoming token.

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.

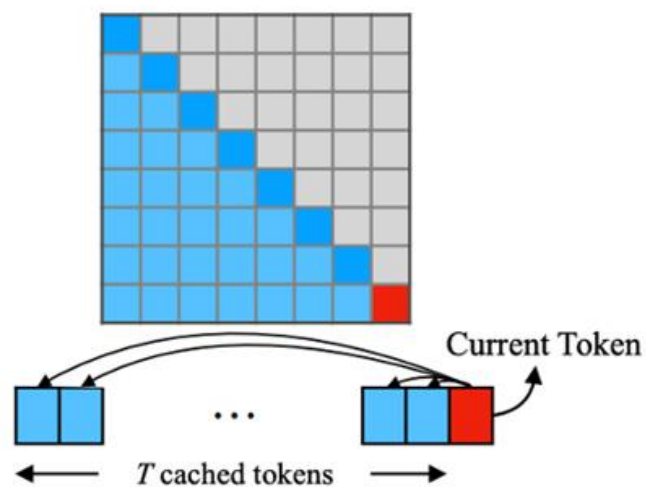


Efficient Streaming Language Models with Attention Sinks [Xiao et al., 2023]



StreamingLLM

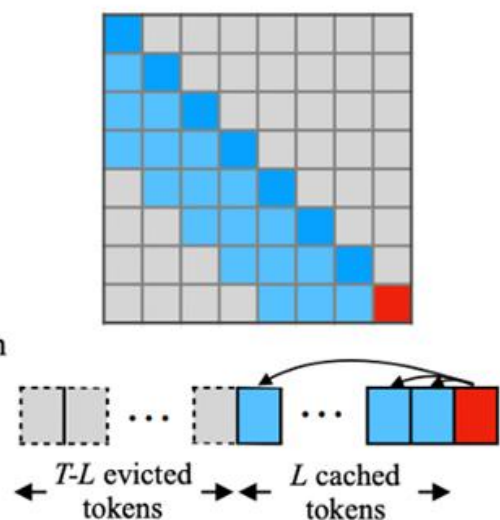
(a) Dense Attention



$O(T^2)$ ✗ PPL: 5641 ✗

Has poor efficiency and performance on long text.

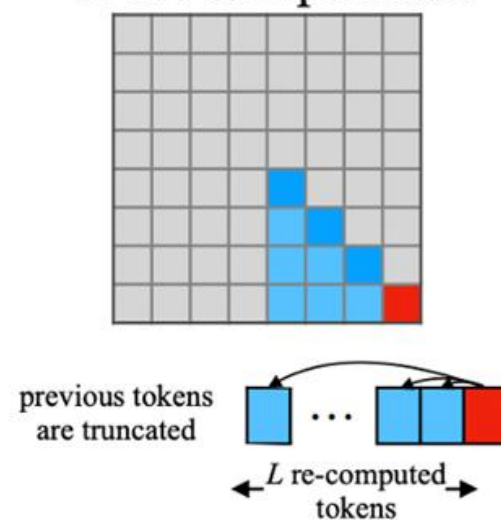
(b) Window Attention



$O(TL)$ ✓ PPL: 5158 ✗

Breaks when initial tokens are evicted.

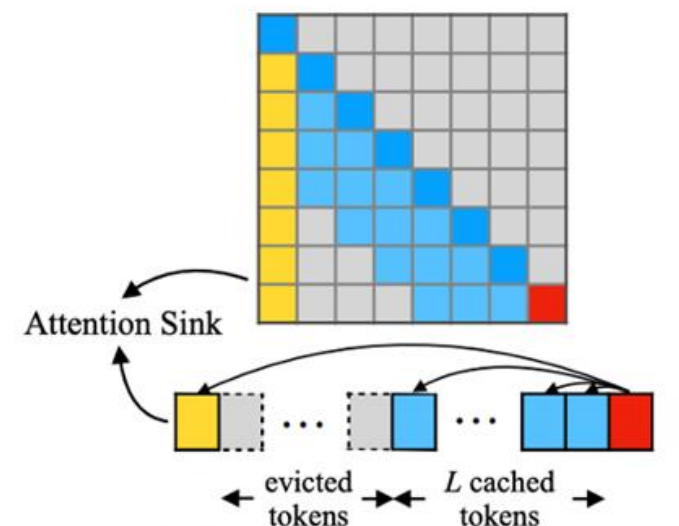
(c) Sliding Window w/ Re-computation



$O(TL^2)$ ✗ PPL: 5.43 ✓

Has to re-compute cache for each incoming token.

(d) StreamingLLM (ours)



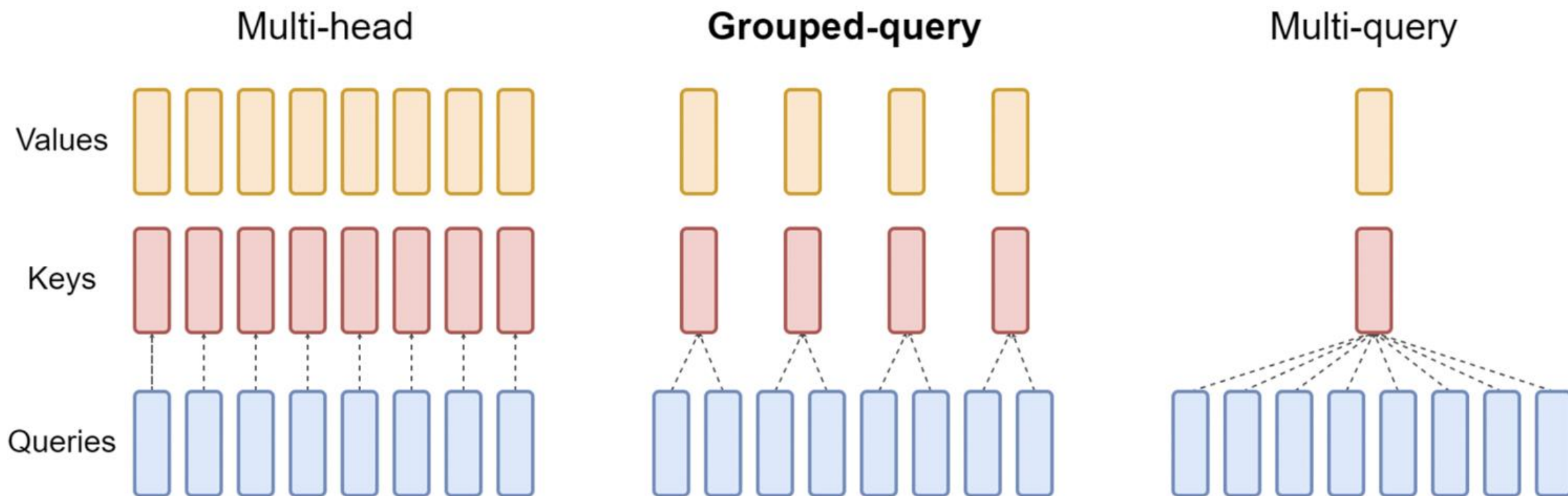
$O(TL)$ ✓ PPL: 5.40 ✓

Can perform efficient and stable language modeling on long texts.

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



MHA/GQA/MQA

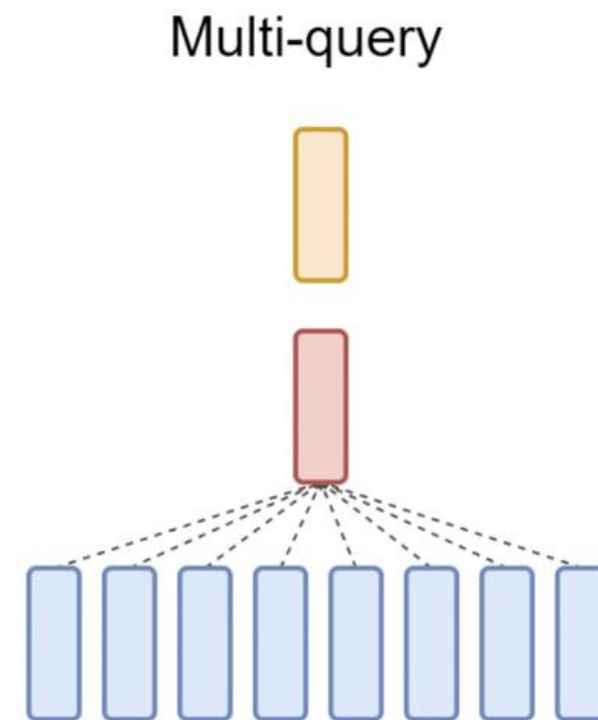
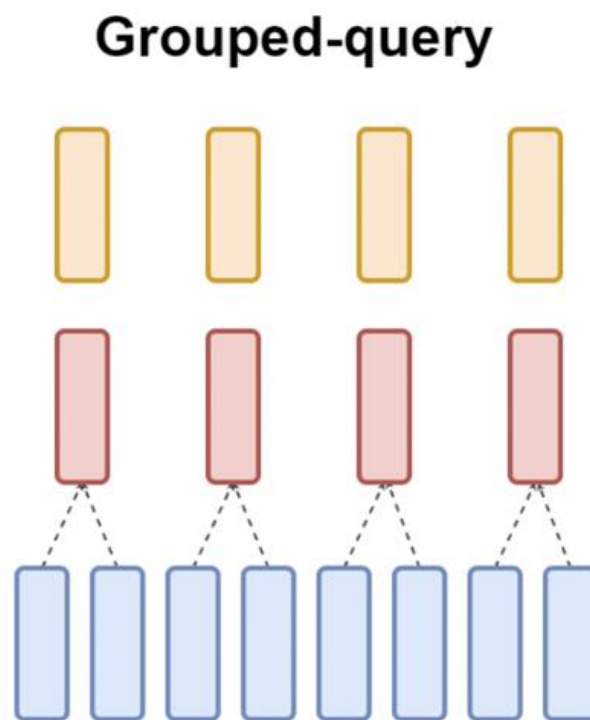
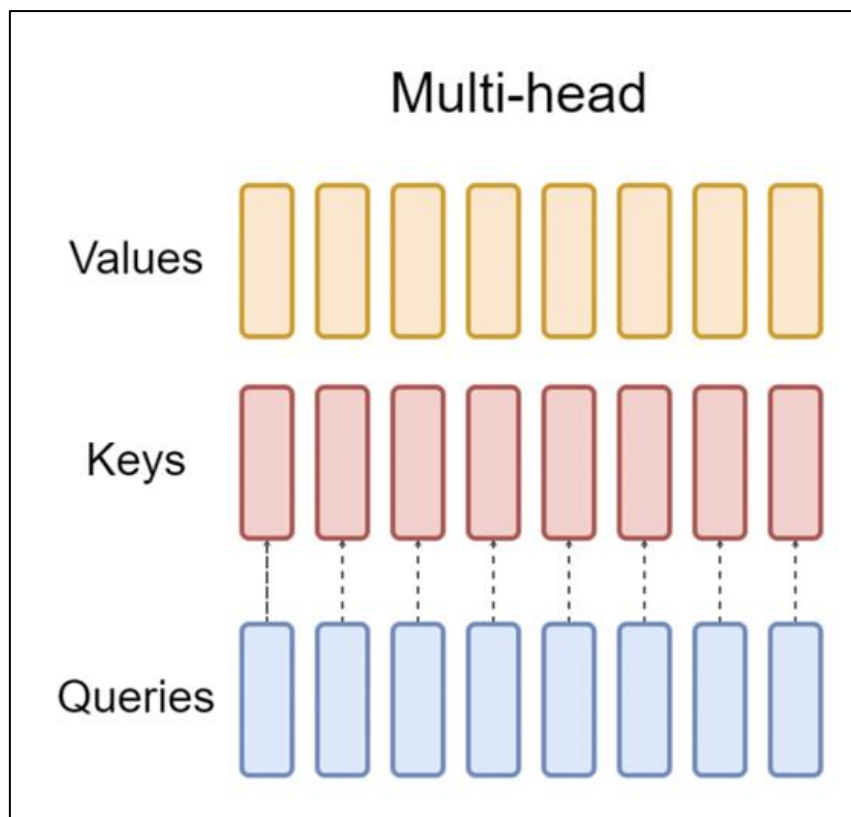


GQA: Training Generalized Multi-Query Transformer Models from Multi-Head Checkpoints



MHA/GQA/MQA

Each attention head calculates separate keys and values for each token

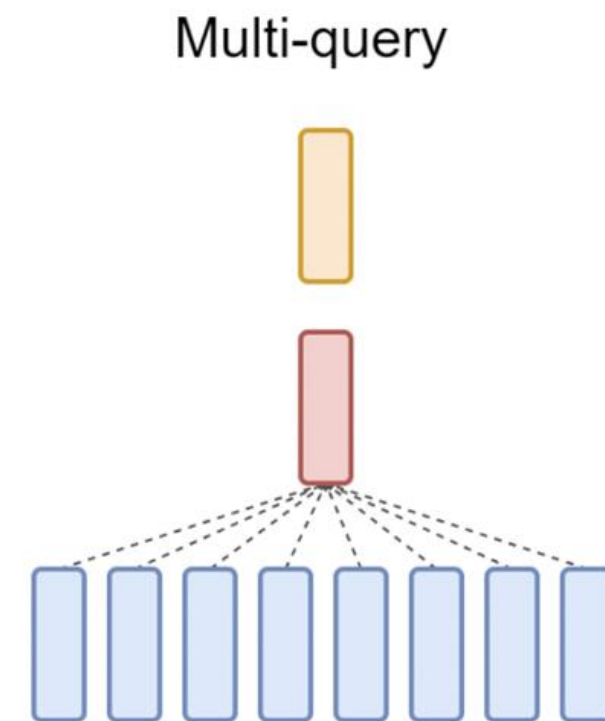
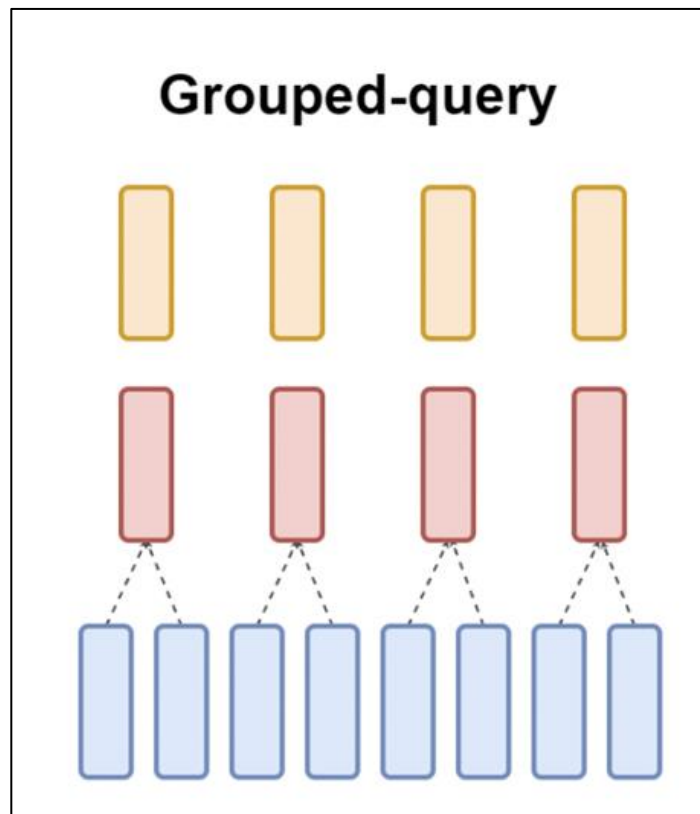
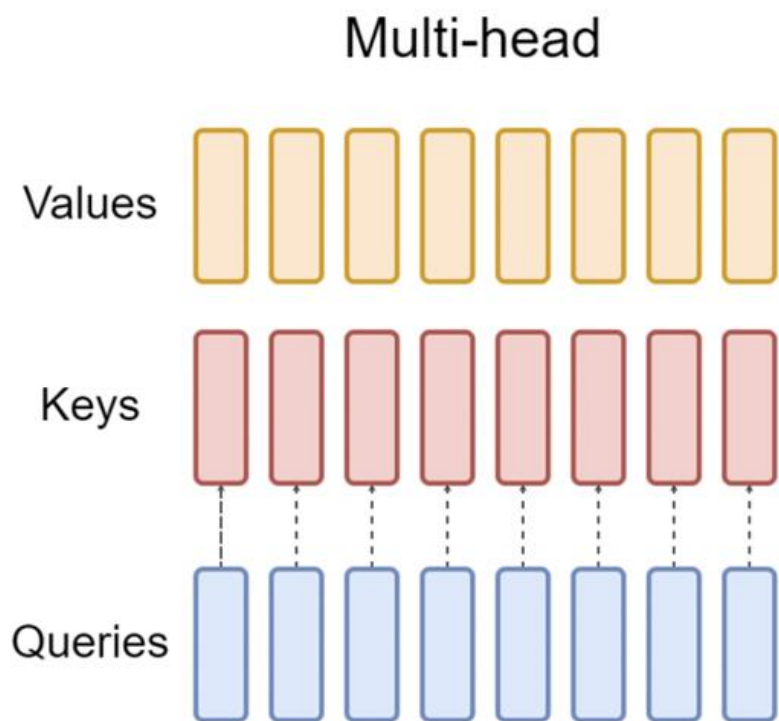


GQA: Training Generalized Multi-Query Transformer Models from Multi-Head Checkpoints



MHA/GQA/MQA

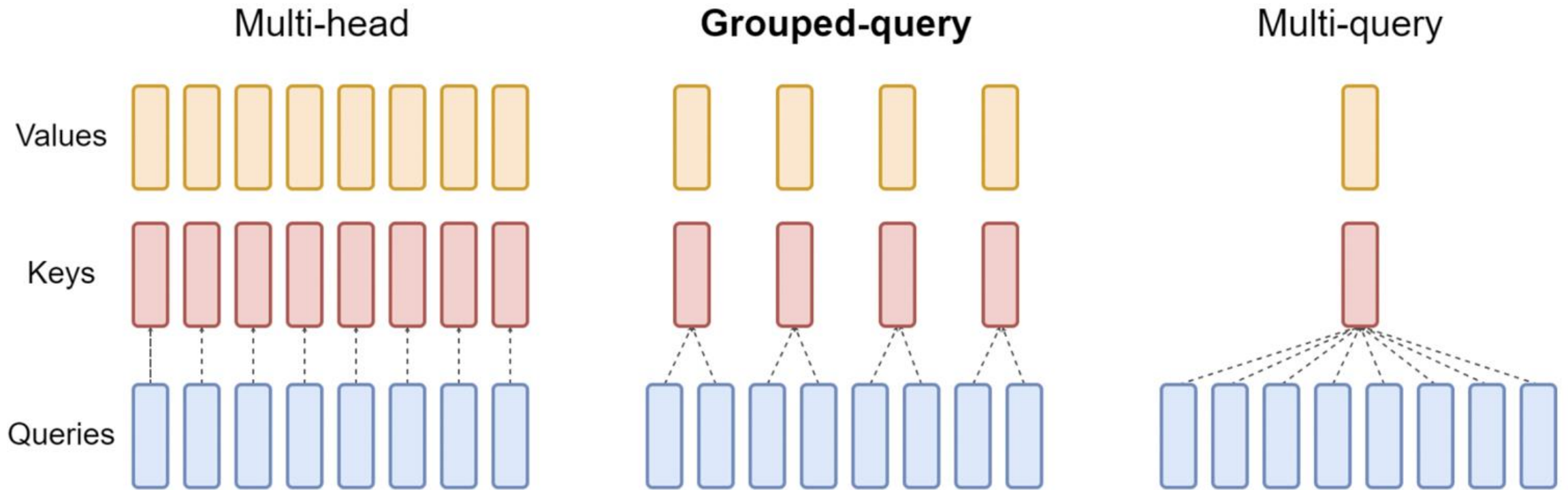
Attention heads are split into groups.
Each group has one key/value per token.



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



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)

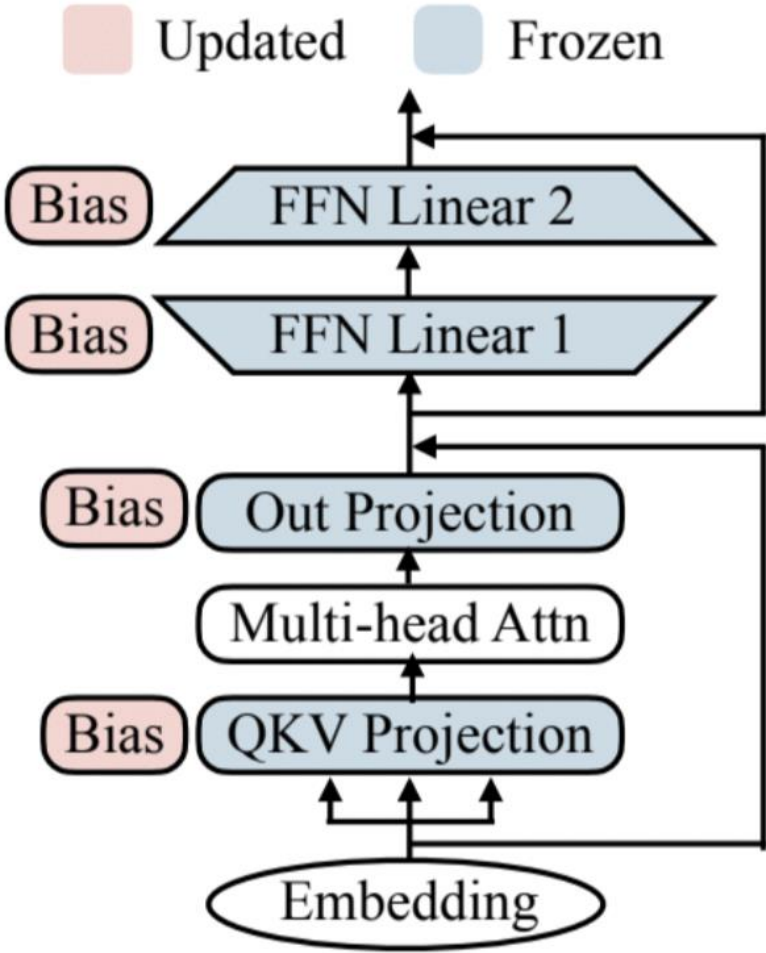
□ Sparsity (Mixture of Experts, Deja Vu: Contextual Sparsity)

□ Efficient Inference Systems (vLLM, StreamingLLM, MHA/GQA/MQA)

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



BitFit



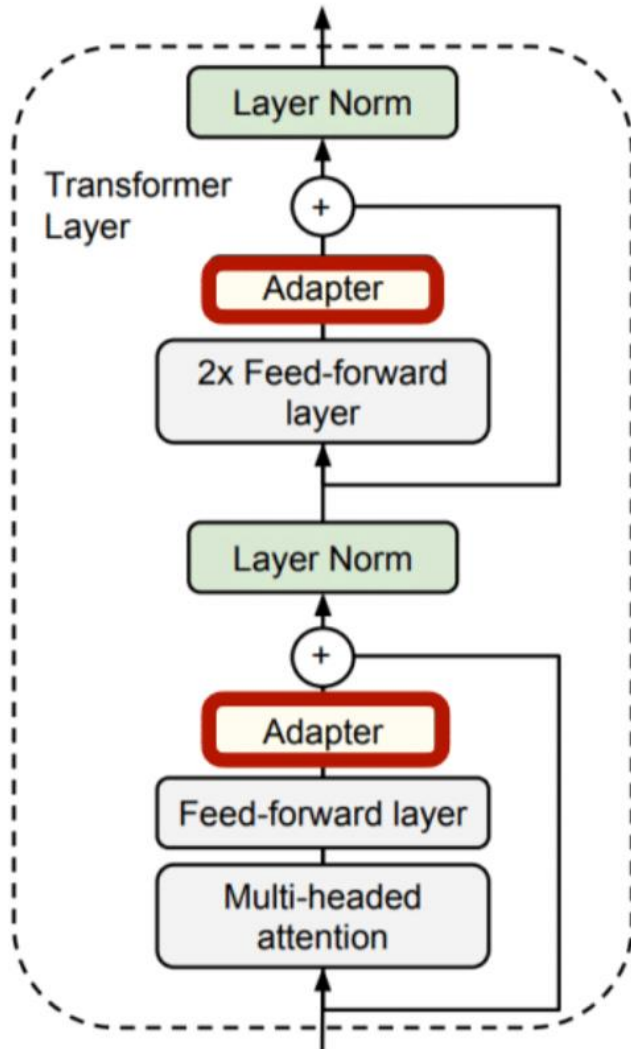
Update only the bias parameters

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

Table 1: BERT_{LARGE} model performance on the GLUE benchmark validation set (V) and test set (T). Lines with † and ‡ 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

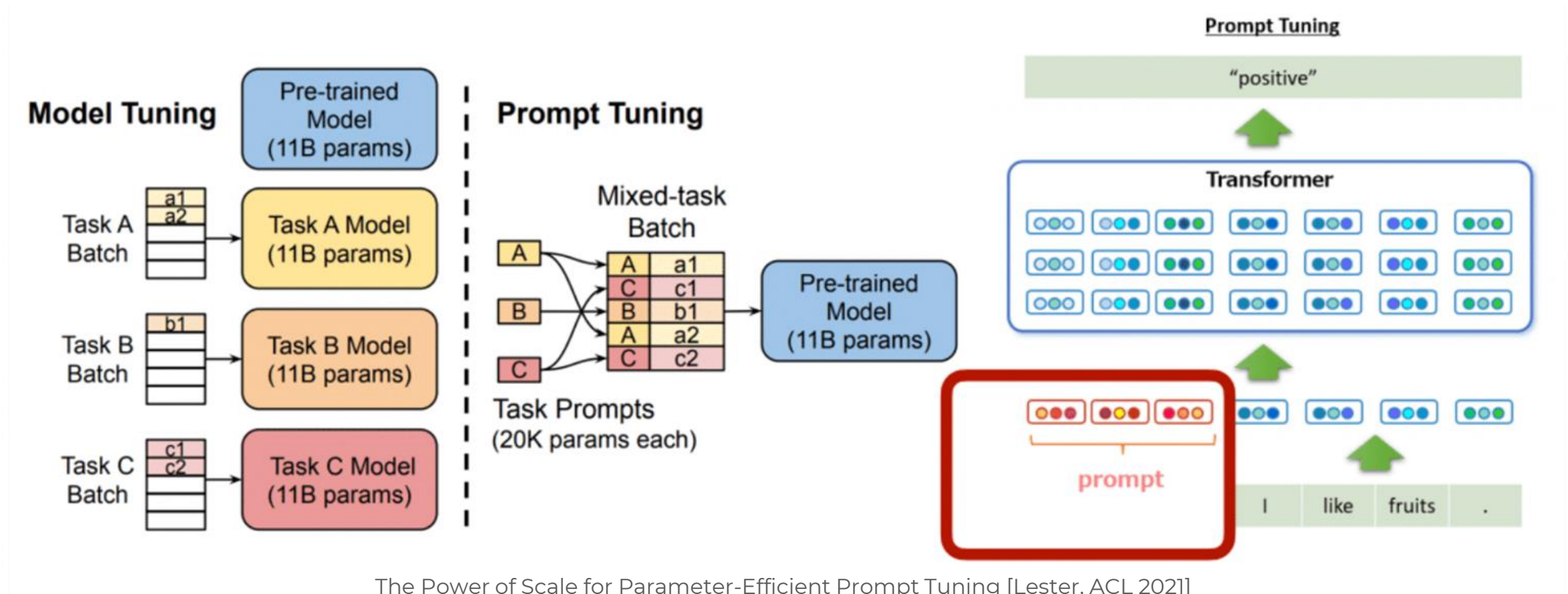
	Total num params	Trained params / task	CoLA	SST	MRPC	STS-B	QQP	MNLI _m	MNLI _{mm}	QNLI	RTE	Total
BERT _{LARGE}	9.0×	100%	60.5	94.9	89.3	87.6	72.1	86.7	85.9	91.1	70.1	80.4
Adapters (8-256)	1.3×	3.6%	59.5	94.0	89.5	86.9	71.8	84.9	85.1	90.7	71.5	80.0
Adapters (64)	1.2×	2.1%	56.9	94.2	89.6	87.3	71.8	85.3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

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



Prompt Tuning (Soft Prompting)

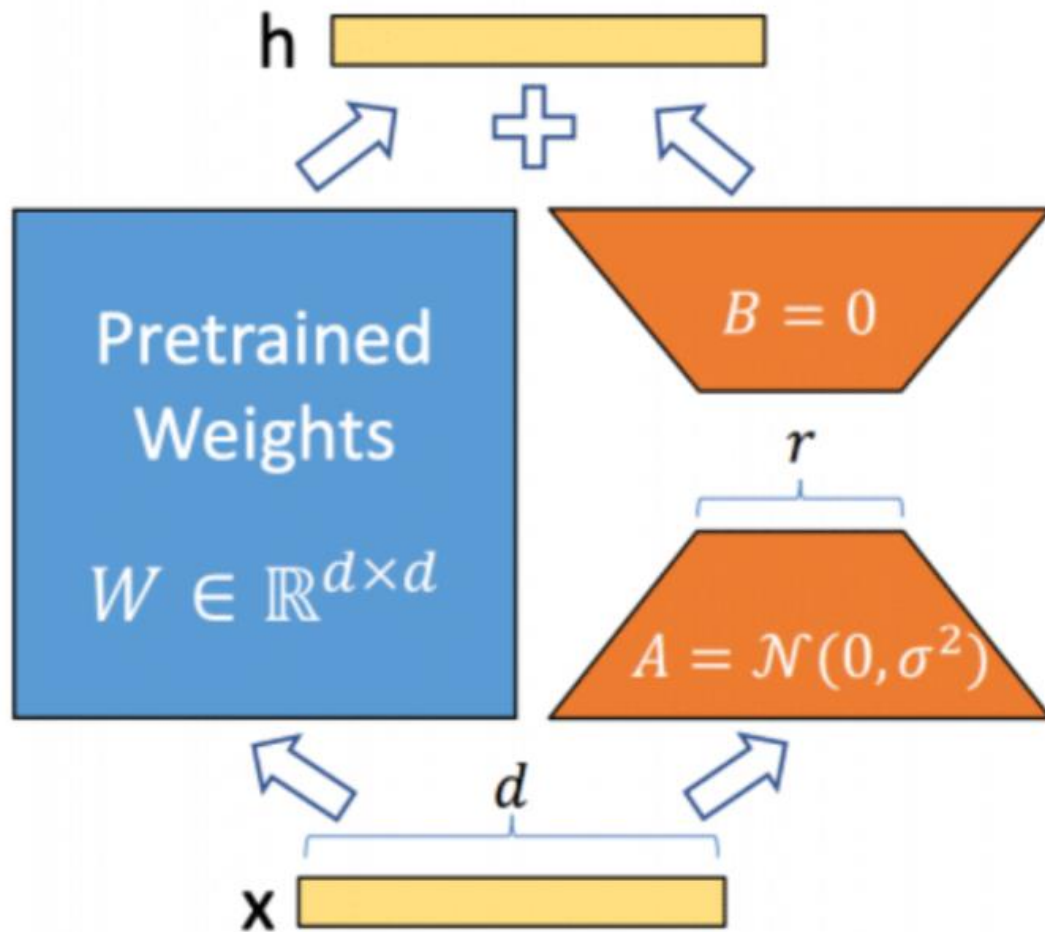
Train a continuous, learnable prompt in embedding space for each task we are training on



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



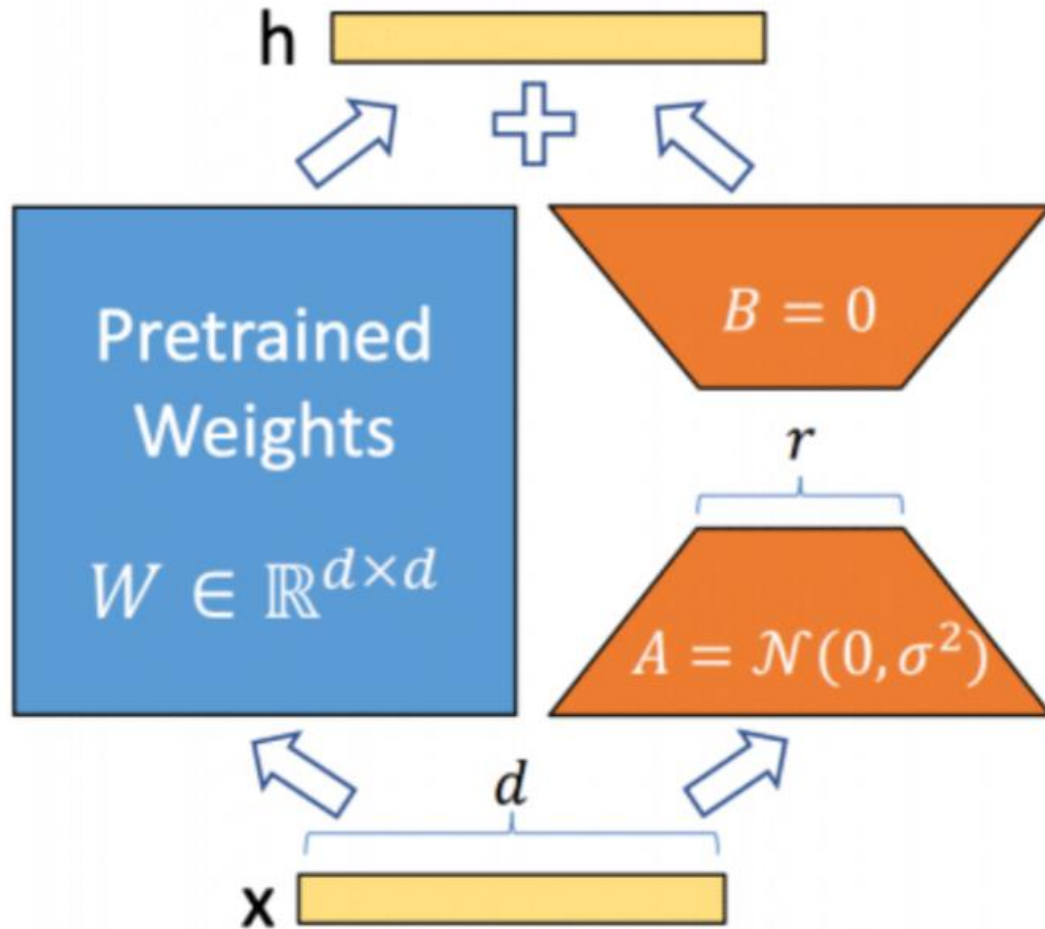
LoRA



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



$$\mathbf{h} = W\mathbf{x}$$

$$\mathbf{h} = W\mathbf{x} + \mathbf{B}\mathbf{A}\mathbf{x}$$

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



LoRA

Model & Method	# Trainable Parameters	MNLI	SST-2	MRPC	CoLA	QNLI	QQP	RTE	STS-B	Avg.
RoB _{base} (FT)*	125.0M	87.6	94.8	90.2	63.6	92.8	91.9	78.7	91.2	86.4
RoB _{base} (BitFit)*	0.1M	84.7	93.7	92.7	62.0	91.8	84.0	81.5	90.8	85.2
RoB _{base} (Adpt ^D)*	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 _{base} (Adpt ^D)*	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 _{base} (LoRA)	0.3M	87.5 \pm .3	95.1\pm.2	89.7 \pm .7	63.4 \pm 1.2	93.3\pm.3	90.8 \pm .1	86.6\pm.7	91.5\pm.2	87.2
RoB _{large} (FT)*	355.0M	90.2	96.4	90.9	68.0	94.7	92.2	86.6	92.4	88.9
RoB _{large} (LoRA)	0.8M	90.6\pm.2	96.2 \pm .5	90.9\pm1.2	68.2\pm1.9	94.9\pm.3	91.6 \pm .1	87.4\pm2.5	92.6\pm.2	89.0
RoB _{large} (Adpt ^P)†	3.0M	90.2 \pm .3	96.1 \pm .3	90.2 \pm .7	68.3\pm1.0	94.8\pm.2	91.9\pm.1	83.8 \pm 2.9	92.1 \pm .7	88.4
RoB _{large} (Adpt ^P)†	0.8M	90.5\pm.3	96.6\pm.2	89.7 \pm 1.2	67.8 \pm 2.5	94.8\pm.3	91.7 \pm .2	80.1 \pm 2.9	91.9 \pm .4	87.9
RoB _{large} (Adpt ^H)†	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 _{large} (Adpt ^H)†	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 _{large} (LoRA)†	0.8M	90.6\pm.2	96.2 \pm .5	90.2\pm1.0	68.2 \pm 1.9	94.8\pm.3	91.6 \pm .2	85.2\pm1.1	92.3\pm.5	88.6
DeB _{XXL} (FT)*	1500.0M	91.8	97.2	92.0	72.0	96.0	92.7	93.9	92.9	91.1
DeB _{XXL} (LoRA)	4.7M	91.9\pm.2	96.9 \pm .2	92.6\pm.6	72.4\pm1.1	96.0\pm.1	92.9\pm.1	94.9\pm.4	93.0\pm.2	91.3

Table 2: RoBERTa_{base}, RoBERTa_{large}, and DeBERTa_{XXL} 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



Open Source Models

❑ Llama (Meta) →

https://huggingface.co/docs/transformers/en/model_doc/llama

❑ Mixtral (Mistral) →

https://huggingface.co/docs/transformers/en/model_doc/mixtral

❑ DBRX (Databricks) → <https://huggingface.co/databricks/dbrx-base>

❑ Grok (xai) → <https://huggingface.co/xai-org/grok-1>

❑ Gemma (Google) → <https://huggingface.co/google/gemma-2b-it>



Repos to Make Models More Efficient

- ❑ Megablocks (MoE library) → <https://github.com/databricks/megablocks>
- ❑ LLM.int8() → <https://huggingface.co/blog/hf-bitsandbytes-integration>
- ❑ AutoAWQ (AWQ integration) → <https://github.com/casper-hansen/AutoAWQ>
- ❑ LoRA → <https://huggingface.co/docs/diffusers/en/training/lora>
- ❑ QLoRA (not covered here) → <https://huggingface.co/blog/4bit-transformers-bitsandbytes>

