



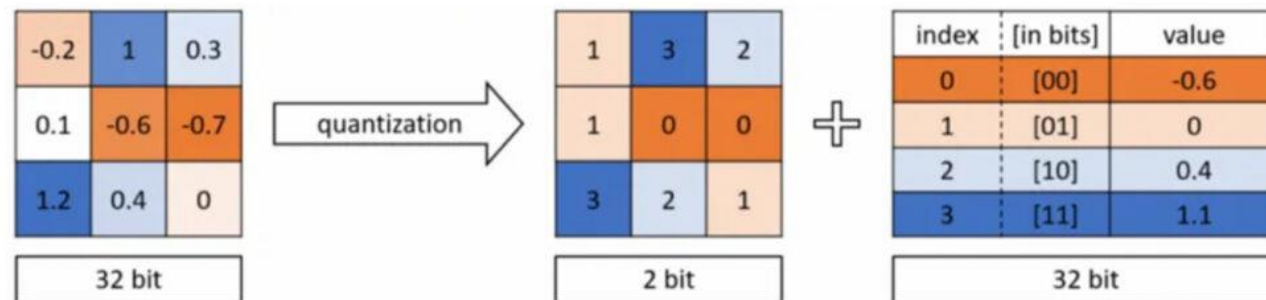
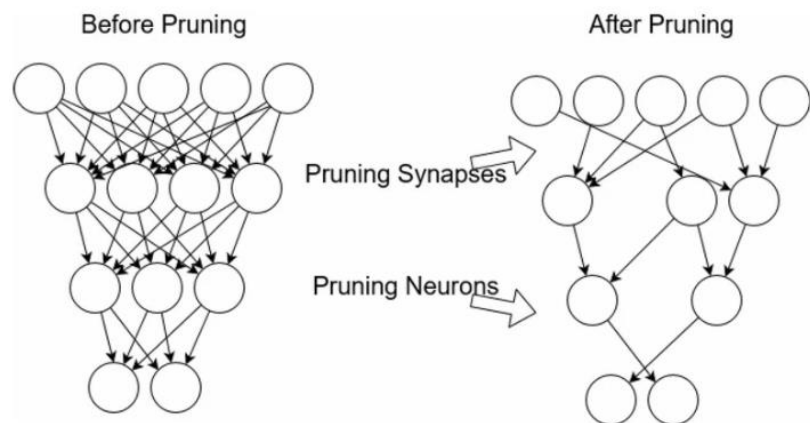
**CS 598**

**AI Efficiency: Systems and Algorithms**  
**Overview of Efficient and Effective Algorithms**

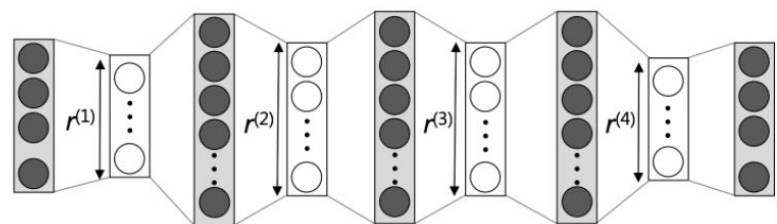
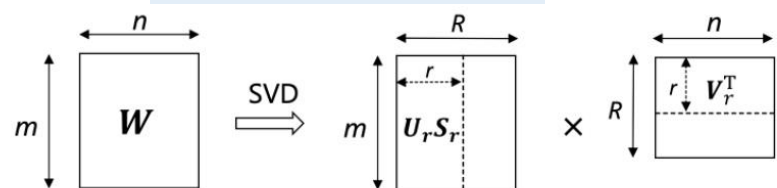
Minjia Zhang

Computer Science Department

# Compression Strategies

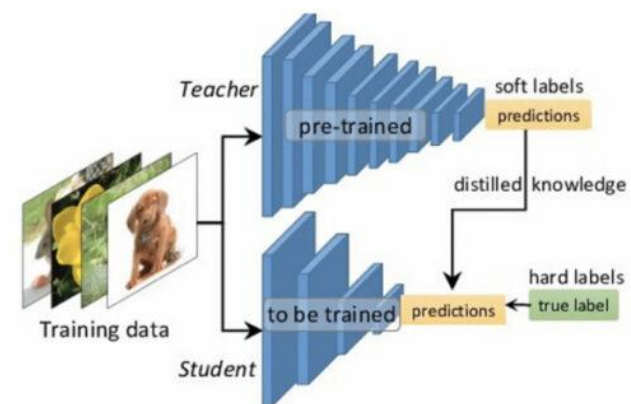
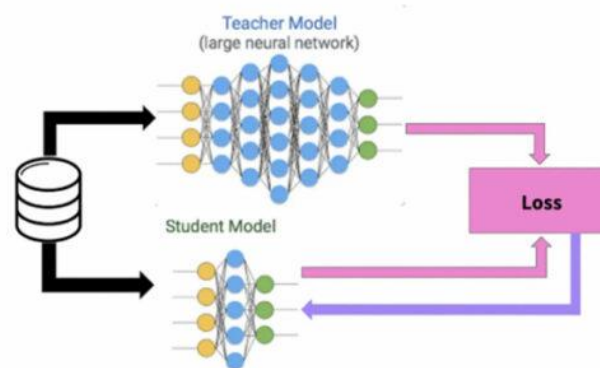


## Quantization



## Low-rank decomposition

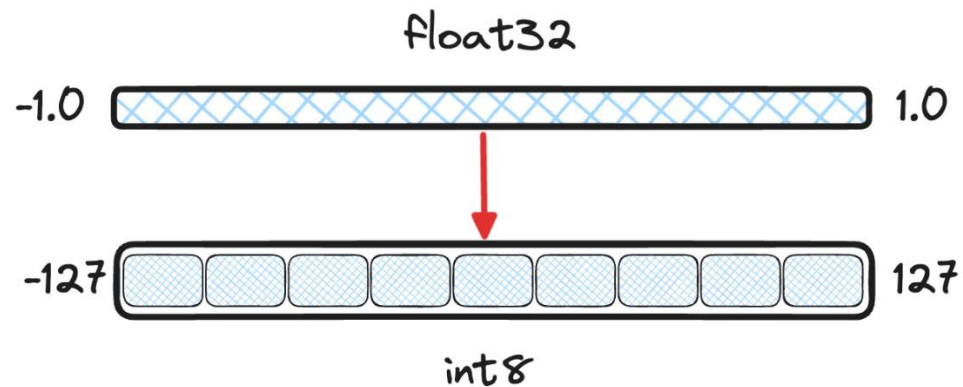
## Quantization



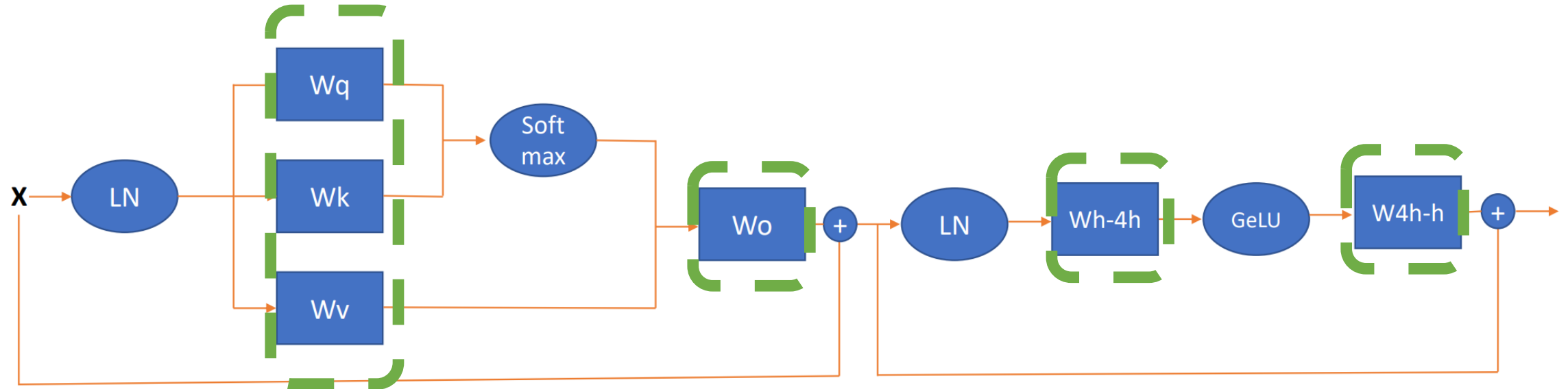
## Distillation

# Quantization: Quick Recap

- Reduce the bits per weight, saving memory consumption
- Accelerate inference speed on supporting hardware



# 8-bit Weight Quantization



- 8-bit weight quantization

$$\mathbf{x}_{quantize} = \text{round} \left( \text{clamp} \left( \frac{\mathbf{x}}{S}, -2^{bit-1}, 2^{bit-1} - 1 \right) \right)$$

FP32 weight matrix

1.1	2.2	0.1	-0.1	-5.5	-6.6
...					
...					
...					
...					
1.1	2.1	0.1	-0.1	-4.8	-6.6

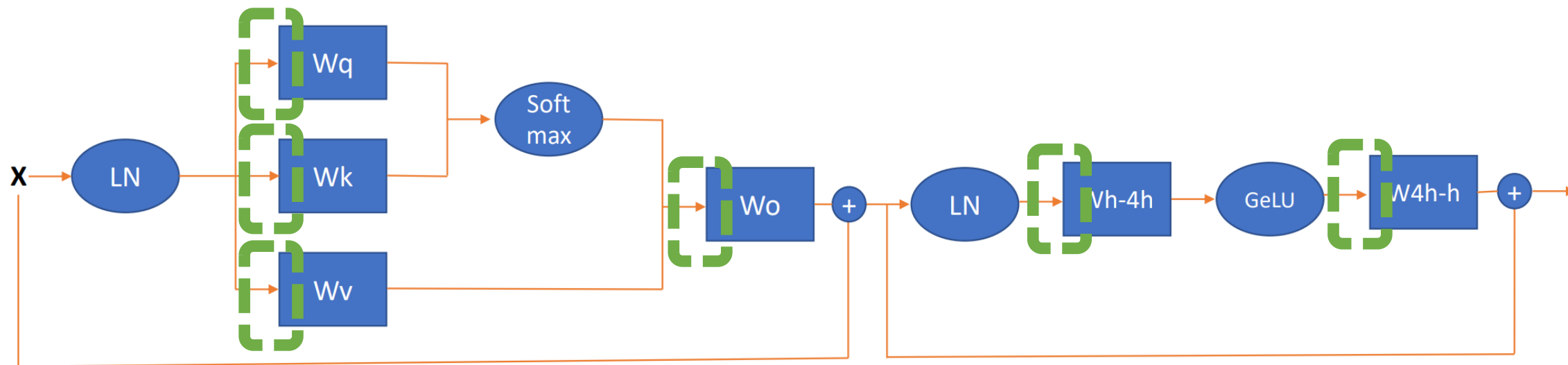
Scaling  
Factor  
 $1/S$

$\approx 0.05^*$

8-bit quantization

21	42	2	-2	-106	-127
...					
...					
...					
...					
21	40	2	-2	-92	-127

# 8-bit Activation Quantization



- 8-bit activation  
(Input to the linear layer)

$$x_{quantize} = \text{round} \left( \text{clamp} \left( \frac{x}{S}, -2^{bit-1}, 2^{bit-1} - 1 \right) \right)$$

FP32 input matrix

1.1	2.2	0.1	-0.1	-5.5	-6.6
...					
...					
...					
...					
1.1	2.1	0.1	-0.1	-4.8	-6.6

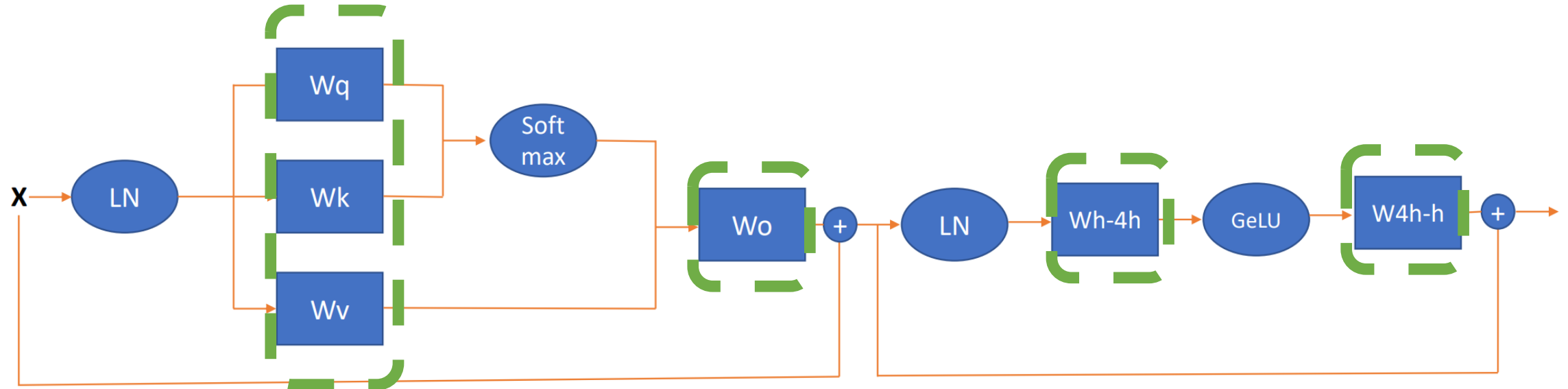
Scaling  
Factor  
 $1/S$

$\approx 0.05^*$

8-bit quantization

21	42	2	-2	-106	-127
...					
...					
...					
...					
21	40	2	-2	-92	-127

# Weight Ternarization



- Ternarization (weight)

$W$ : weight matrix, FP32.

$Q(W)$ : Quantization mapping, 2-bit.

With  $\alpha = \|W\|_1/n$ , for some scalar  $s$

$$Q(W_{ij}) = \begin{cases} \alpha \cdot \text{sign}(W_{ij}) & \text{when } |W_{ij}| > s \\ 0 & \text{when } |W_{ij}| < s \end{cases}$$

FP32 weight matrix

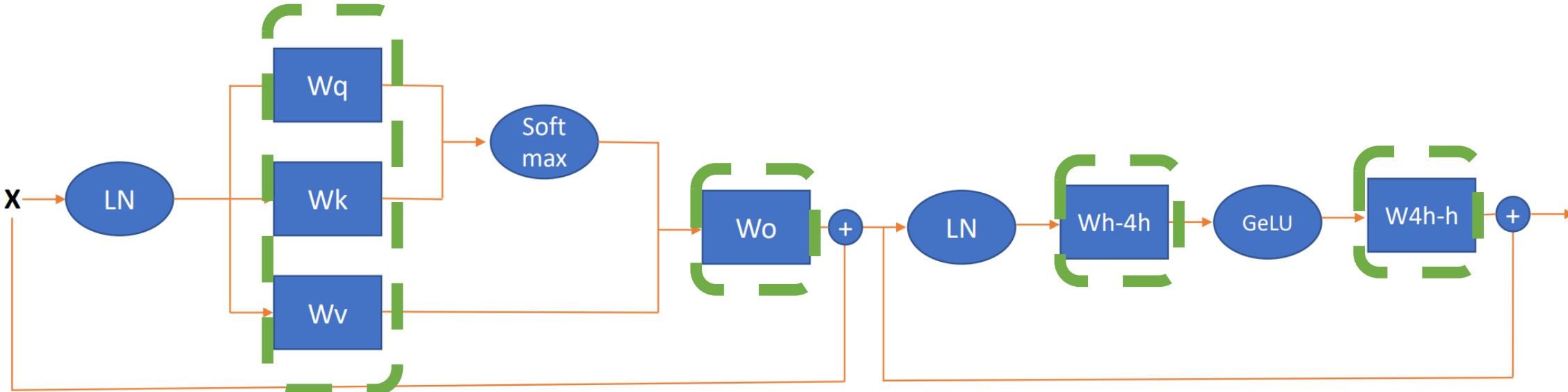
1.1	2.2	0.1	-0.1	-5.5	-6.6
...					
...					
...					
...					
1.1	2.1	0.1	-0.1	-4.8	-6.0

Scaling  
Factor  $\alpha$   
 $\approx 2.06^*$

2-bit quantization

1	1	0	0	-1	-1
...					
...					
...					
...					
1	1	0	0	-1	-1

# Weight Binarization



- Binarization (weight)
  - $W$ : weight matrix, FP32.
  - $Q(W)$ : Quantization mapping, 1-bit.
  - With  $\alpha = \|W\|_1/n$
  - $Q(W_{ij}) = \alpha \cdot \text{sign}(W_{ij})$

FP32 weight matrix

1.1	2.2	0.1	-0.1	-5.5	-6.6
...					
...					
...					
...					
1.1	2.1	0.1	-0.1	-4.8	-6.0

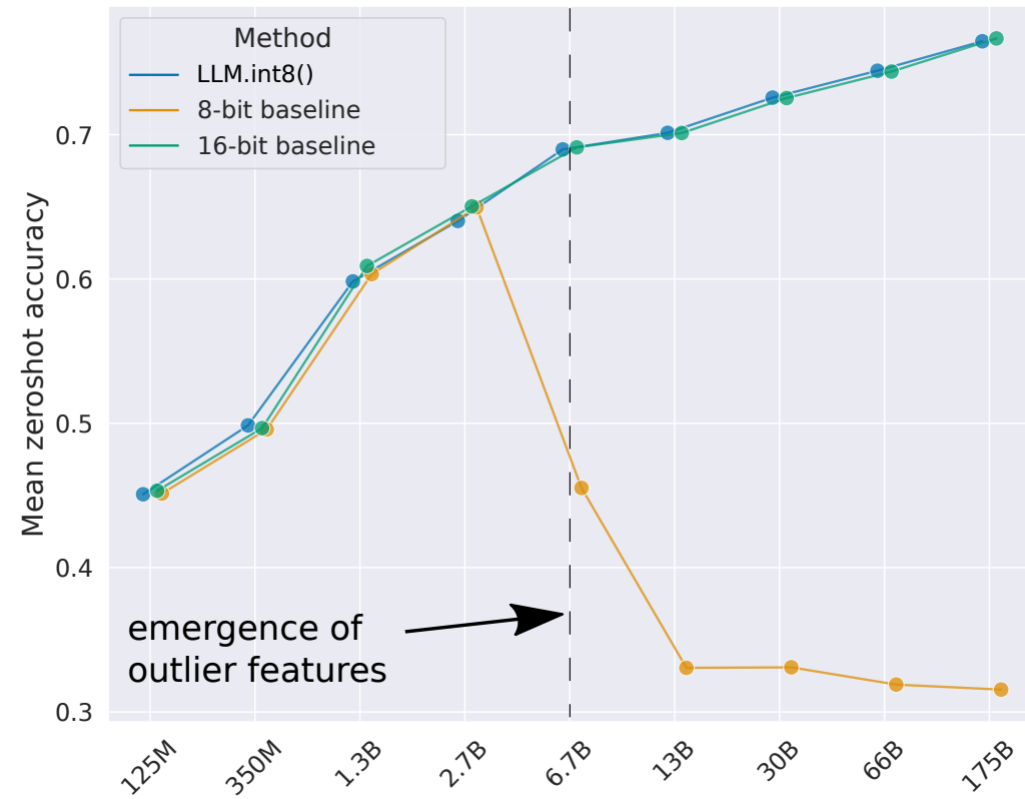
1-bit quantization

$\approx 1.4 \cdot$   
Scaling Factor  $\alpha$

1	1	1	-1	-1	-1
...					
...					
...					
...					
1	1	1	-1	-1	-1

# Challenges to Quantize LLMs

- Standard quantization strategy leads to catastrophic accuracy drop

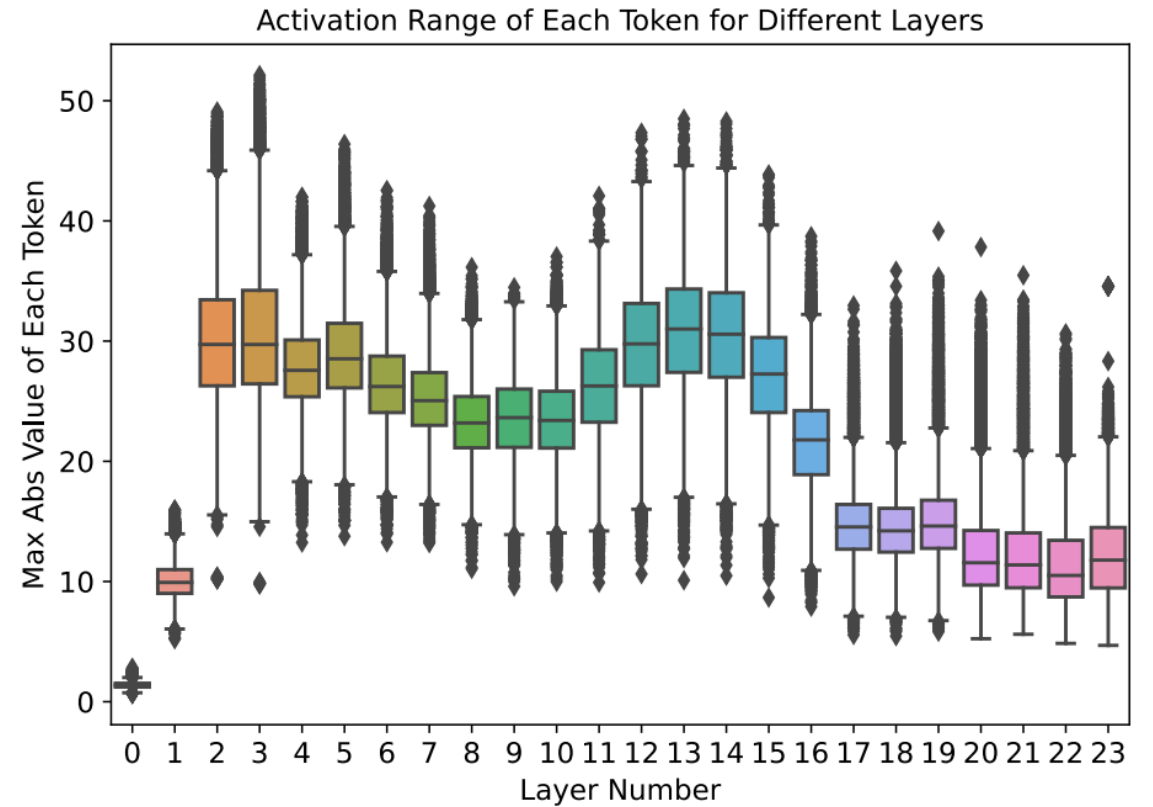


LLM.int8(): 8-bit Matrix Multiplication for Transformers at Scale, 2023



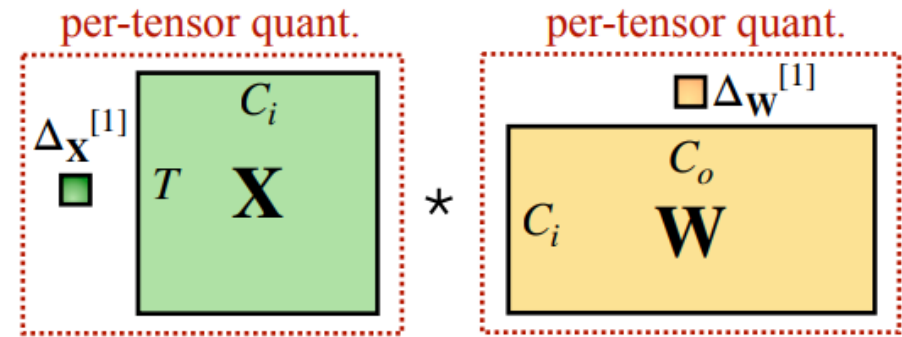
# Challenges to Quantize LLMs

- High dynamic ranges of activation, leading to large quantization errors

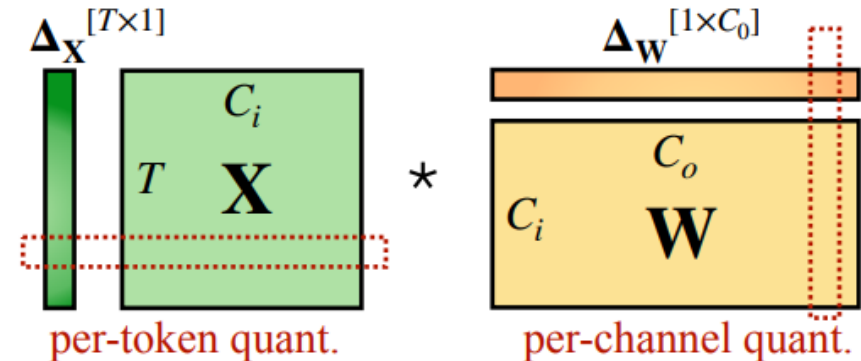


# Fine-grained Quantization

- Per-tensor quantization
  - Low accuracy
  - Fast to quantize/dequantize
- Per-token/channel quantization
  - High accuracy
  - Slower to quantize/dequantize
  - Custom kernels required



(a) per-tensor quantization

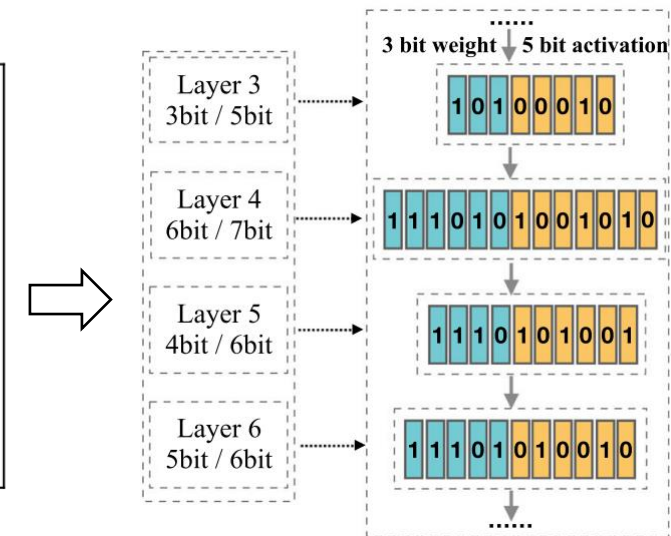
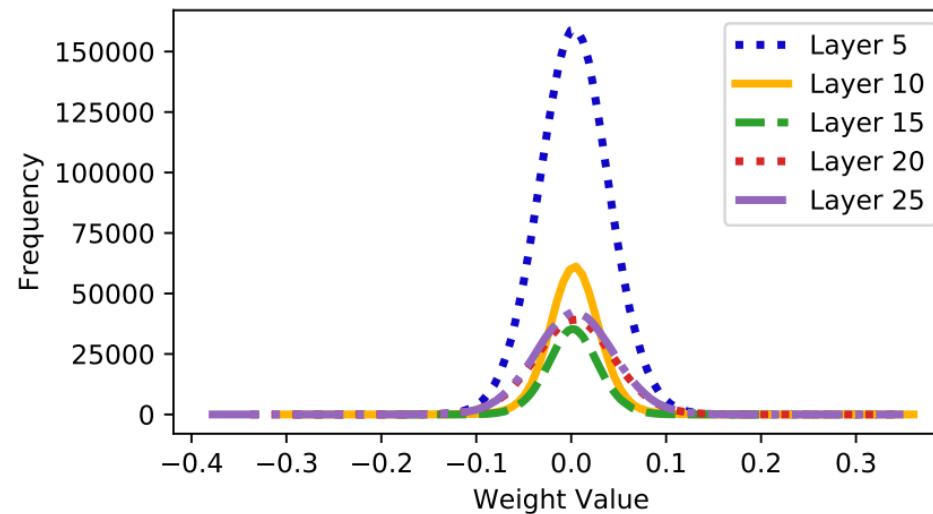


(b) per-token + per-channel quantization

ZeroQuant: Efficient and Affordable Post-Training Quantization for Large-Scale Transformers, NeurIPS 2022

# Mixed Precision Quantization

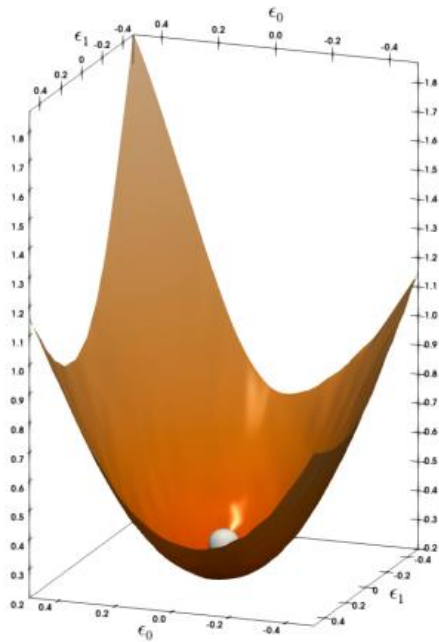
- Weights follow Gaussian distribution
- Outliers remain in original form, quantize the rest of the values
- Different bits for different layers



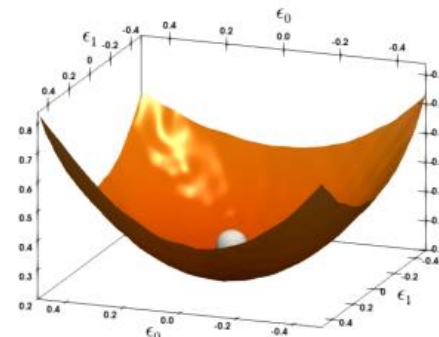
GOBO: Quantizing Attention-Based NLP Models for Low Latency and Energy Efficient Inference, MICRO 2020

# Second Order Information

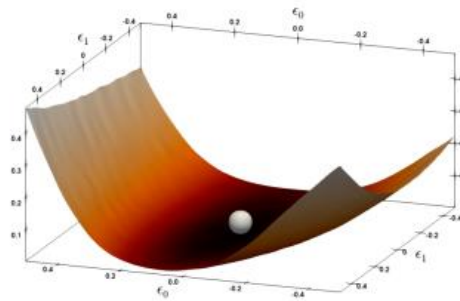
- Analyze the loss curvature (Hessian matrices) to help identify layer sensitivity



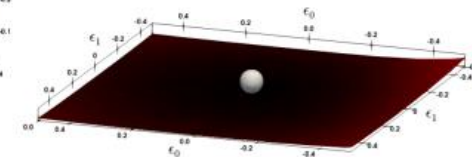
(a) MNLI 4<sup>th</sup> layer



(b) MNLI 10<sup>th</sup> layer



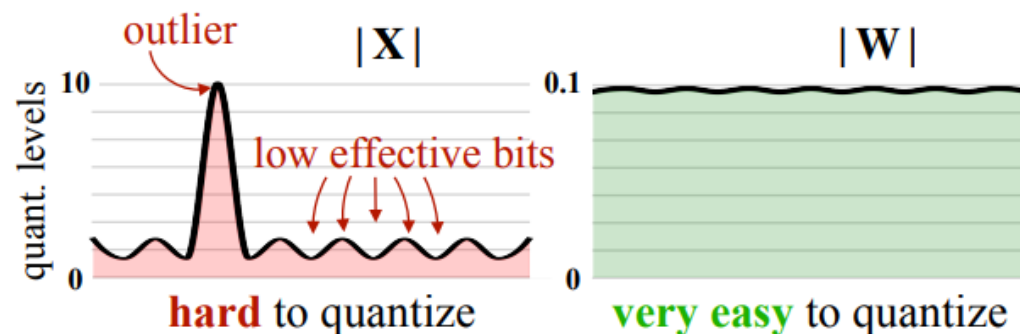
(c) CoNLL-03 4<sup>th</sup> layer



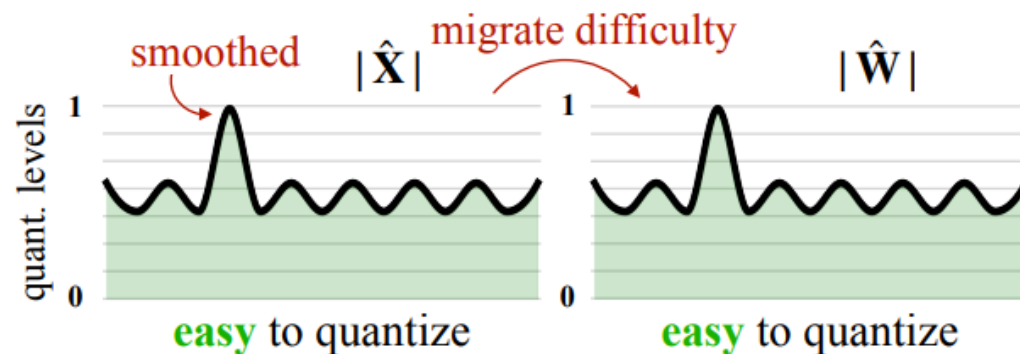
(d) CoNLL-03 11<sup>th</sup> layer

GPTQ: Accurate Post-Training Quantization for Generative Pre-trained Transformers, ICLR 2023

# Outlier Smoothing



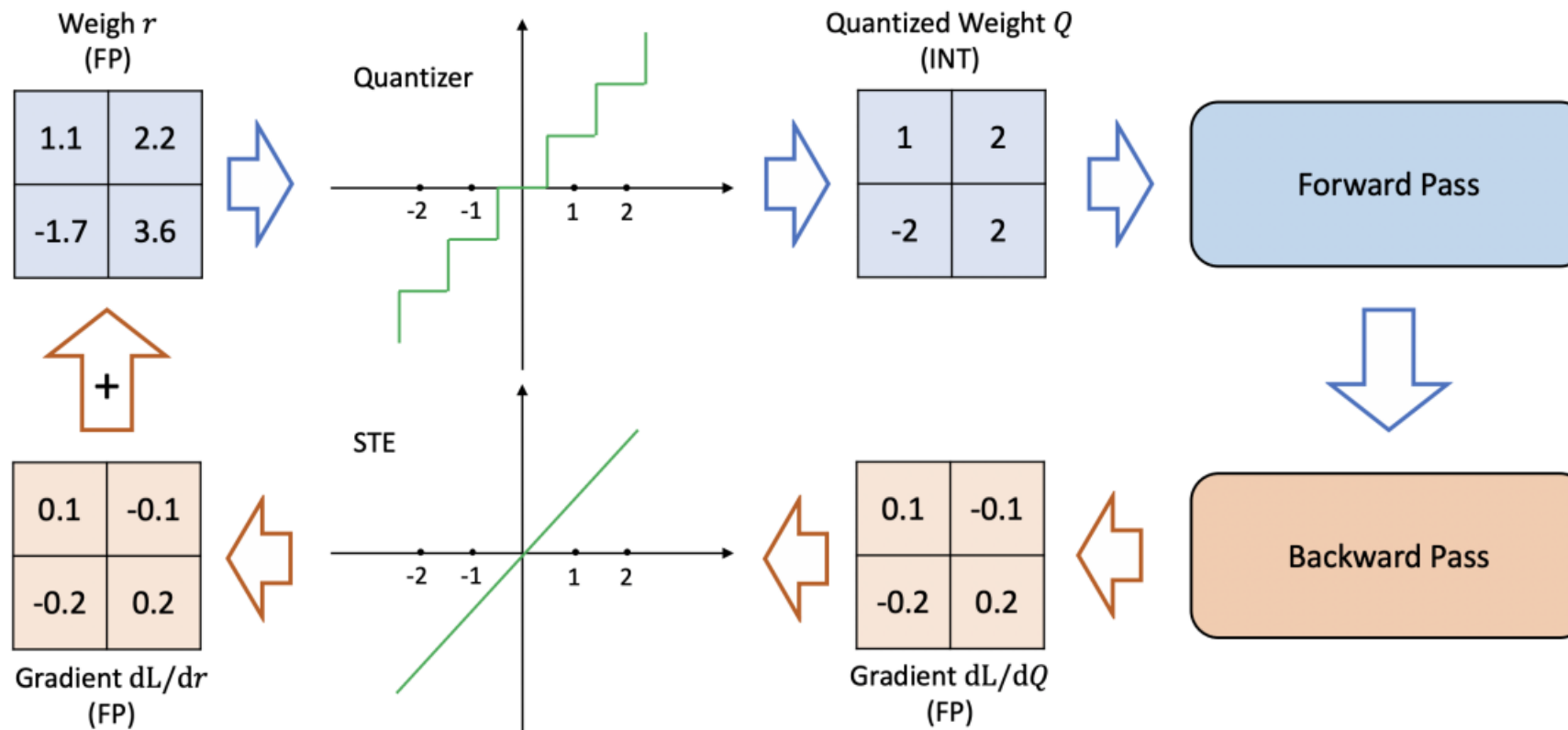
(a) Original



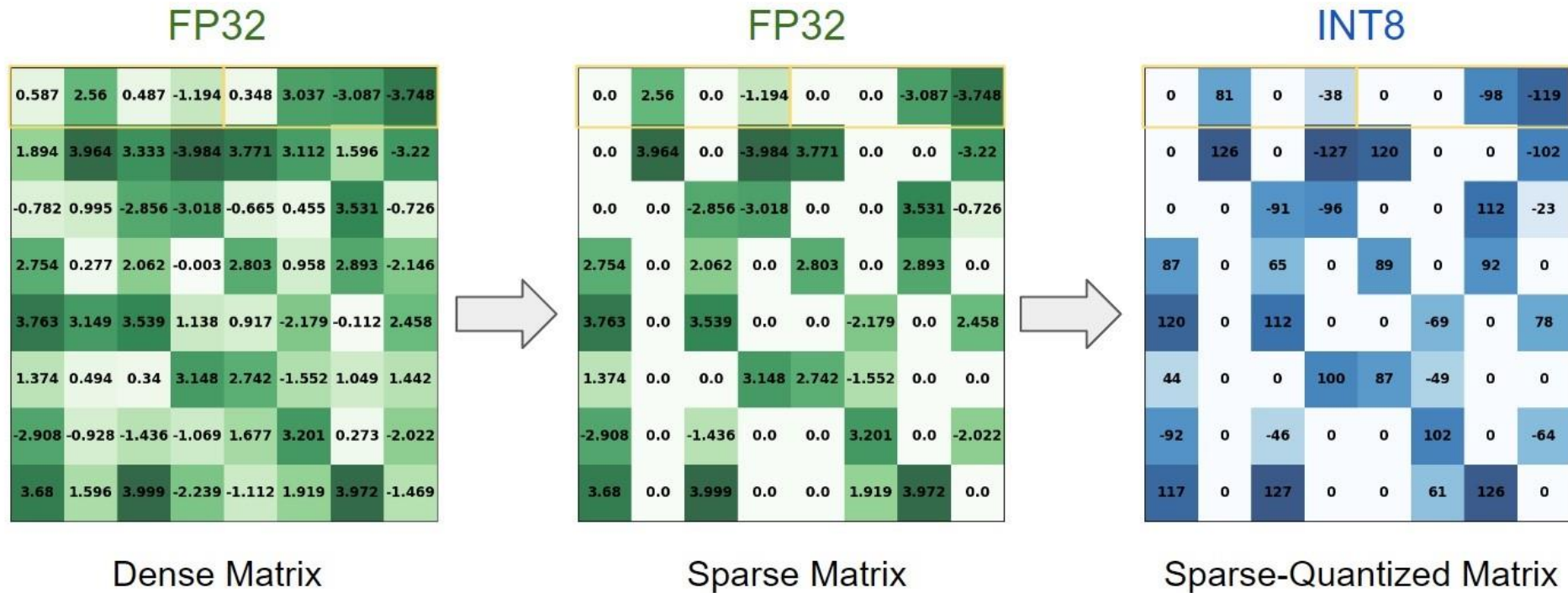
(b) SmoothQuant

SmoothQuant: Accurate and Efficient Post-Training Quantization for Large Language Models, ICML 2023

# Quantization-Aware Training



# Sparsification



# Model Pruning

- **Unstructured (connection) Sparsity:**
- High accuracy
- No performance improvement or performance regression

1	7	-3	4	2	-1	-3	3
5	3	6	2	0	2	1	7
-8	-2	1	3	5	-6	2	0
3	9	1	4	5	-3	0	1
9	0	-1	3	6	2	-1	3
4	-5	2	8	7	6	-3	0
-1	-1	4	7	0	7	6	1
9	1	2	4	6	8	9	0

Unstructured Sparsity

- **N:M Semi-Structured Sparsity:**
- High accuracy
- High performance improvement

1	7	-3	4	2	-1	-3	3
5	3	6	2	0	2	1	7
-8	-2	1	3	5	-6	2	0
3	9	1	4	5	-3	0	1
9	0	-1	3	6	2	-1	3
4	-5	2	8	7	6	-3	0
-1	-1	4	7	0	7	6	1
9	1	2	4	6	8	9	0

Semi-Structured Sparsity (4:2 N:M)

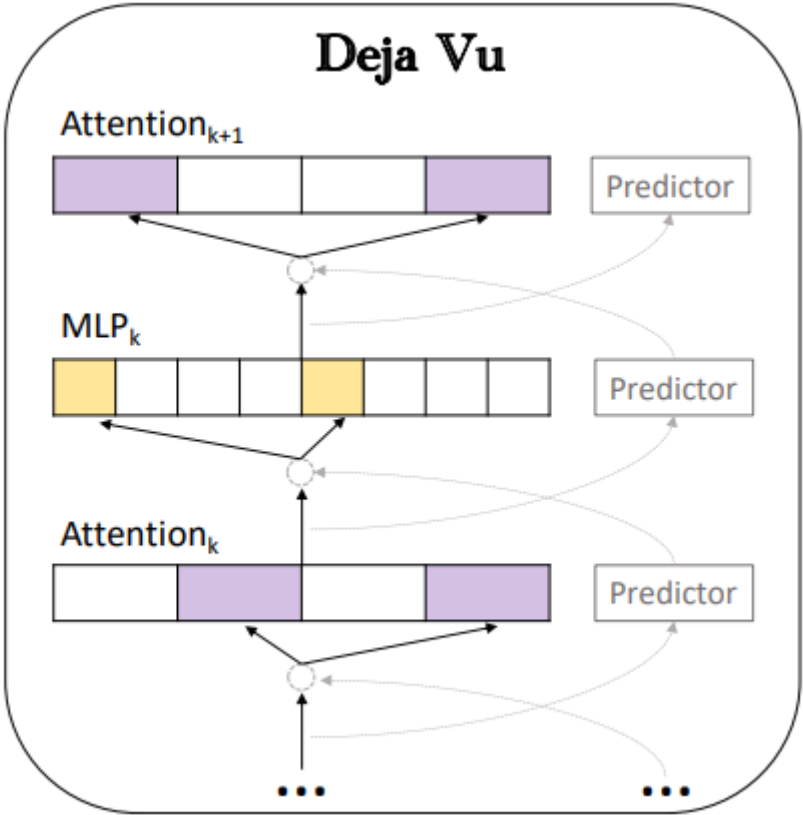
- **Structured Sparsity:**
- Large accuracy degradation
- High performance scalability

1	7	-3	4	2	-1	-3	3
5	3	6	2	0	2	1	7
-8	-2	1	3	5	-6	2	0
3	9	1	4	5	-3	0	1
9	0	-1	3	6	2	-1	3
4	-5	2	8	7	6	-3	0
-1	-1	4	7	0	7	6	1
9	1	2	4	6	8	9	0

Structured Sparsity (Column-wise Sparsity)

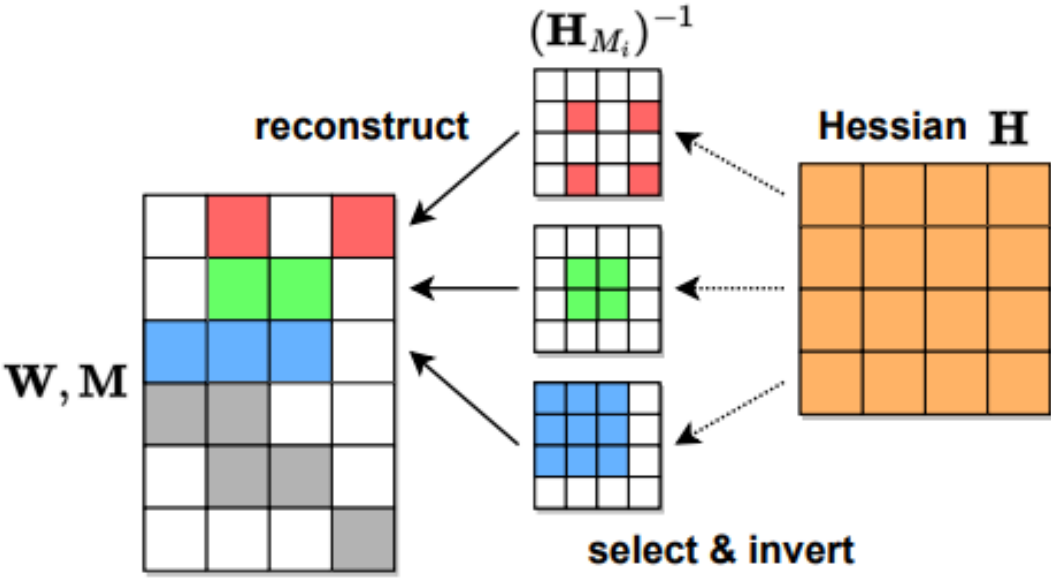


# Model Pruning

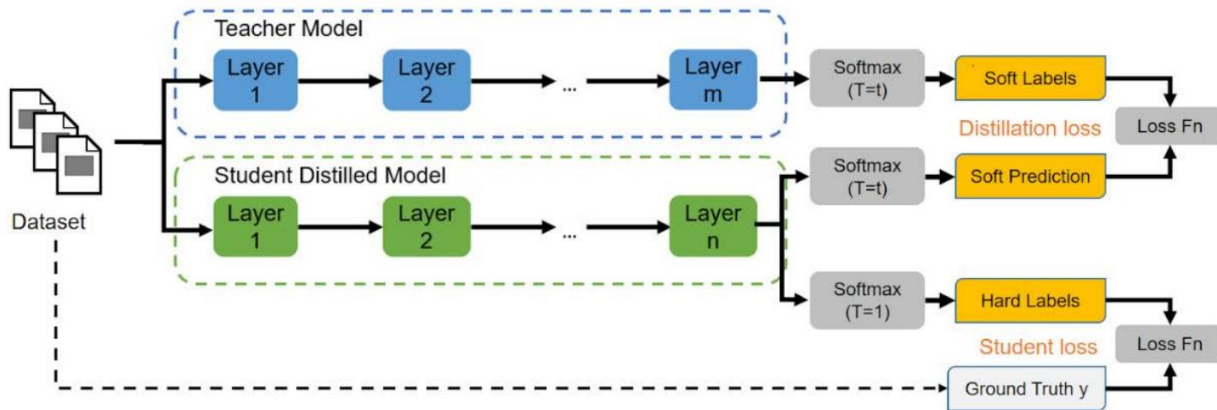


Deja Vu: Contextual Sparsity for Efficient LLMs at Inference Time, 2023

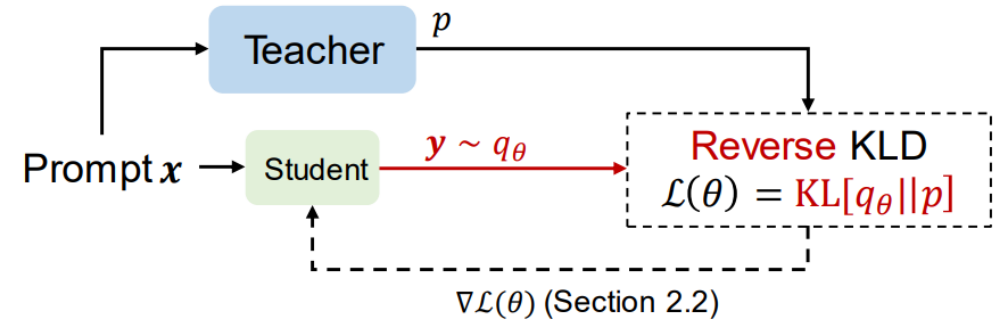
SparseGPT: Massive Language Models Can be Accurately Pruned in One-Shot, 2023



# Knowledge Distillation

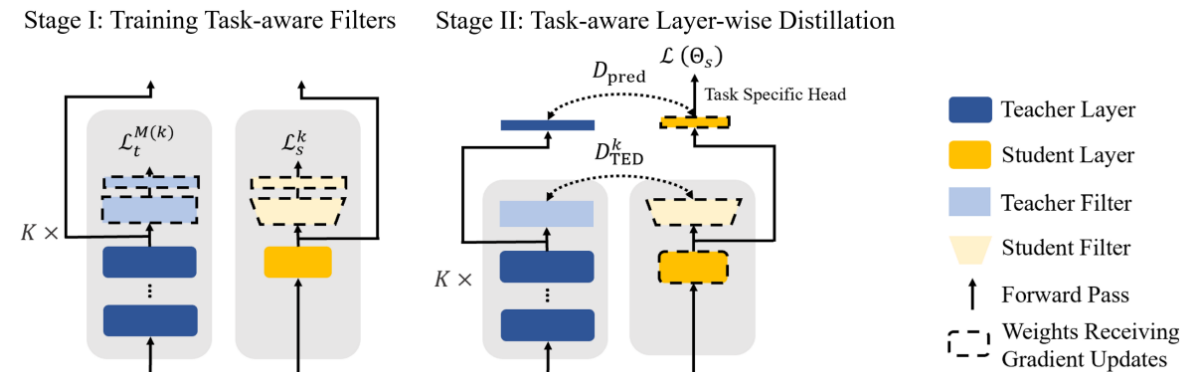


Distilling the Knowledge in a Neural Network, 2015



MiniLLM (Ours)

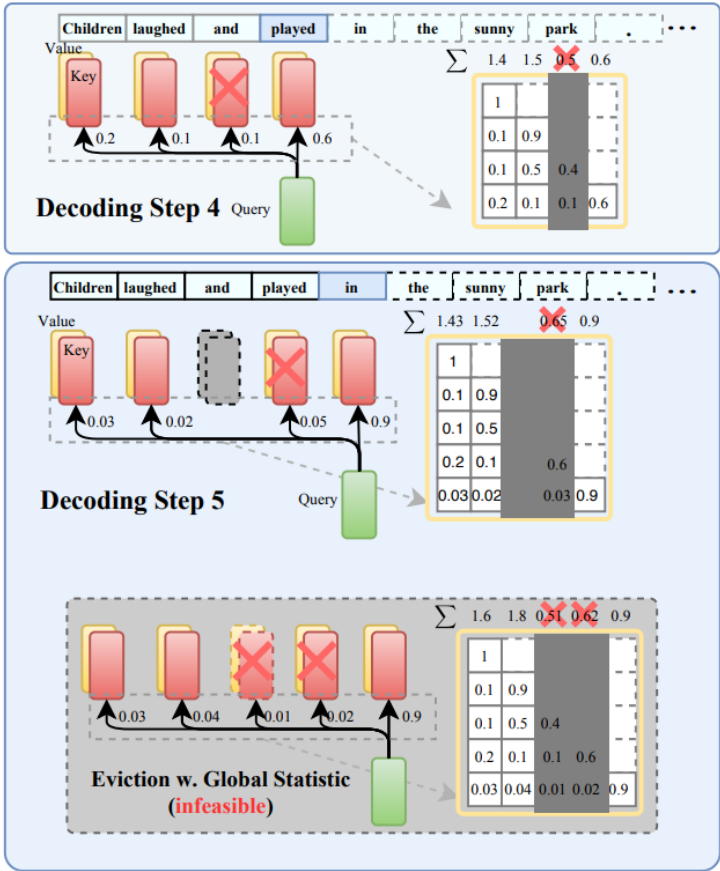
MiniLLM: Knowledge distillation of large language models, 2024



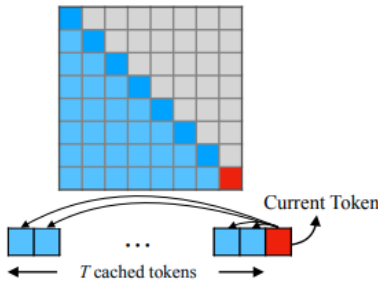
Less is More: Task-aware Layer-wise Distillation for Language Model Compression, 2024

# KV Cache Compression

Efficient Streaming Language Models with Attention Sinks, ICL 2024



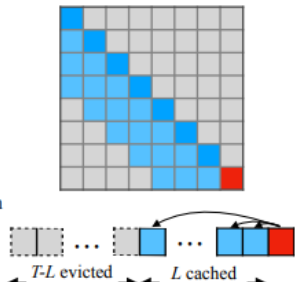
(a) Dense Attention



$O(T^2)$  **PPL: 5641**

Has poor efficiency and performance on long text.

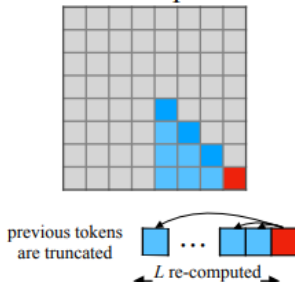
(b) Window Attention



$O(TL)$  **PPL: 5158**

Breaks when initial tokens are evicted.

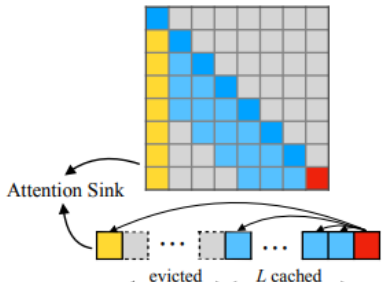
(c) Sliding Window w/ Re-computation



$O(TL^2)$  **PPL: 5.43**

Has to re-compute cache for each incoming token.

(d) StreamingLLM (ours)

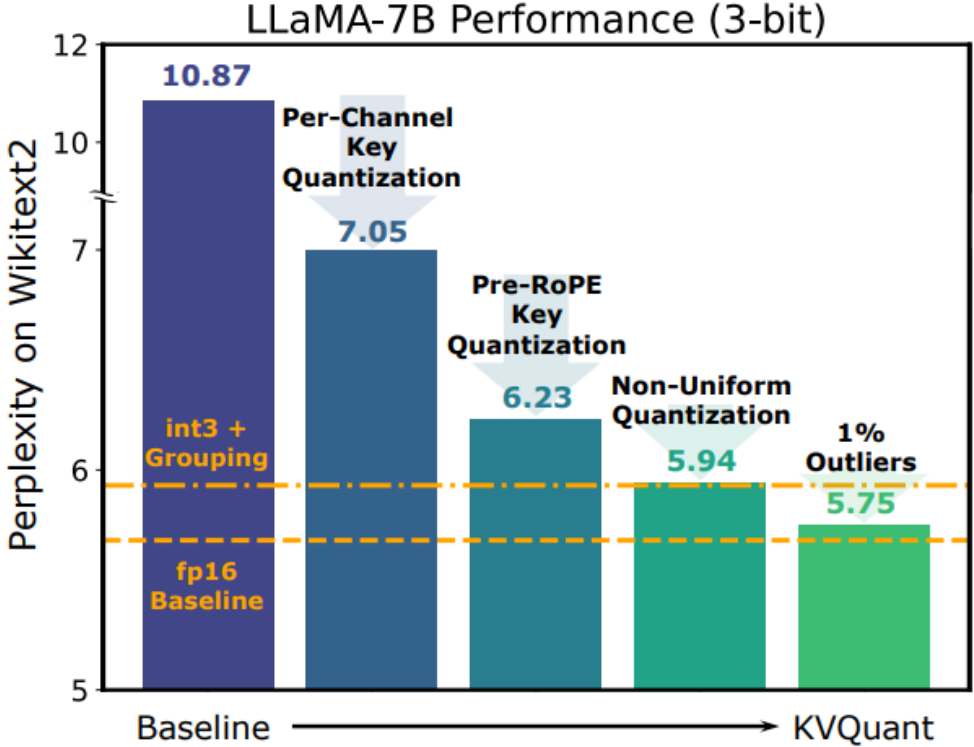


$O(TL)$  **PPL: 5.40**

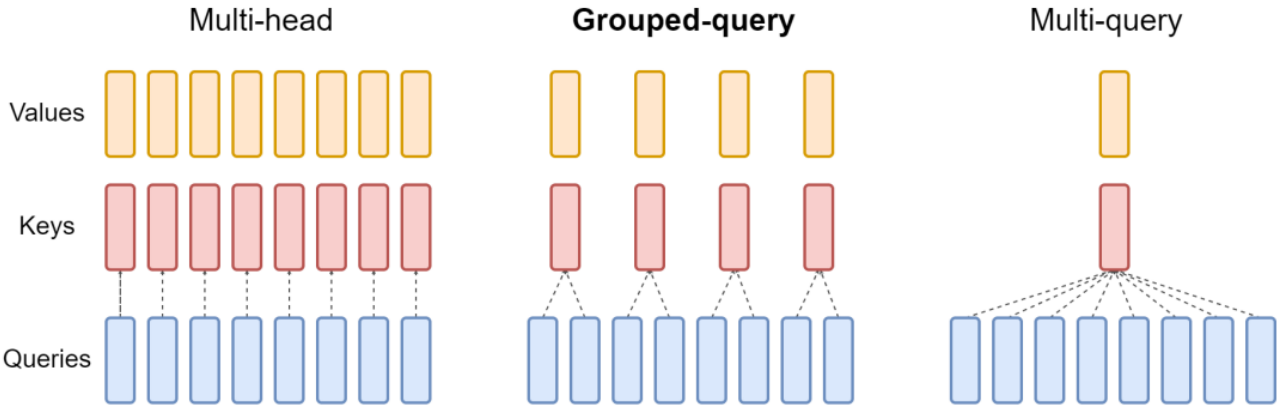
Can perform efficient and stable language modeling on long texts.

H2O: Heavy-Hitter Oracle for Efficient Generative Inference of Large Language Models, 2023

# KV Cache Compression



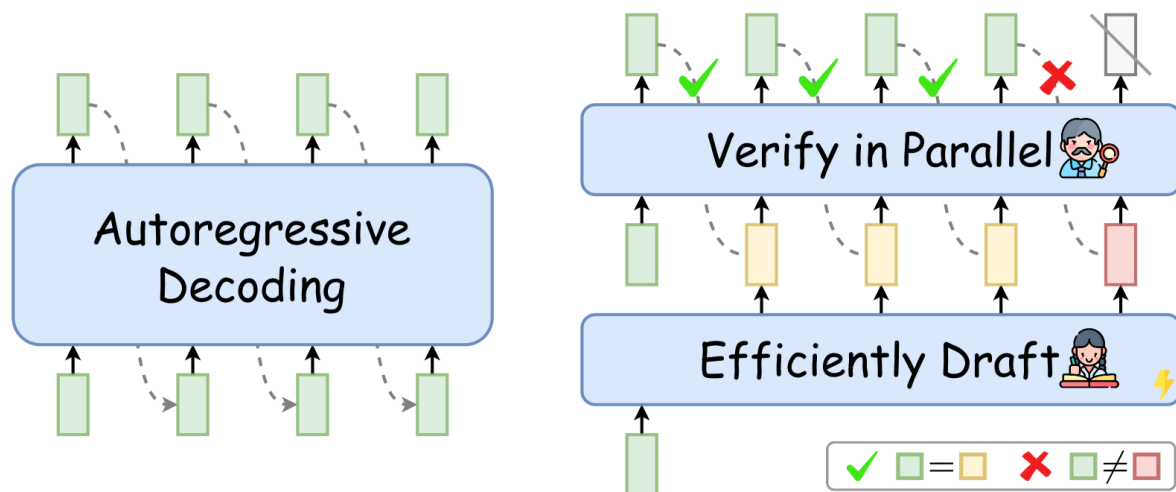
GQA: Training Generalized Multi-Query Transformer Models from Multi-Head Checkpoints, 2023



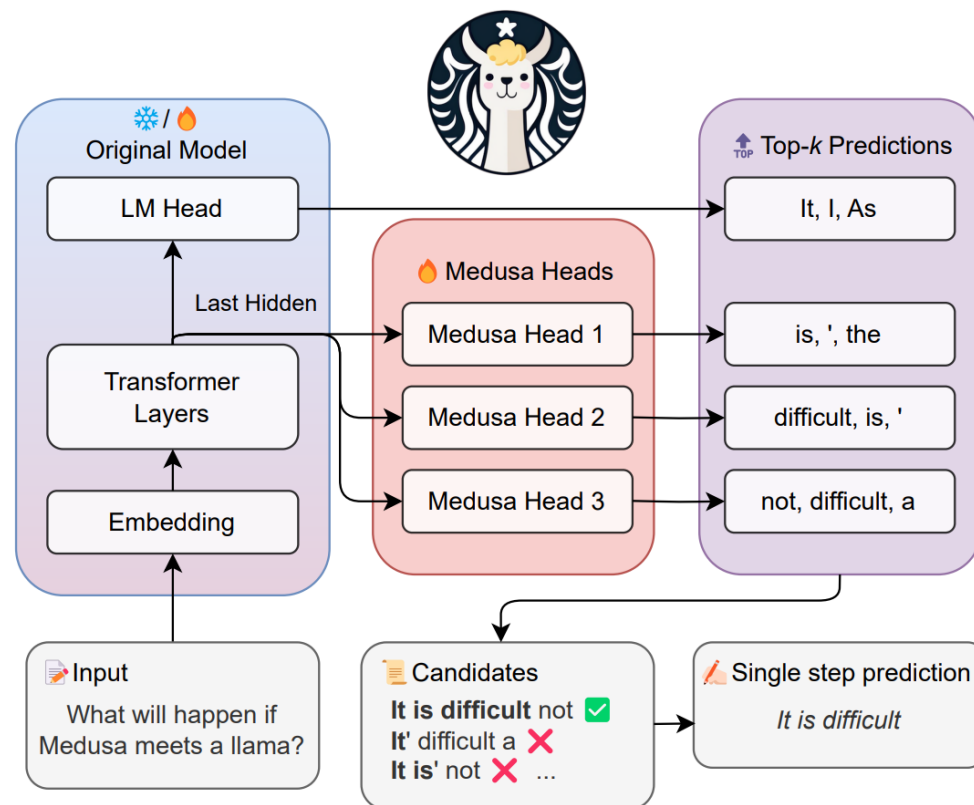
KVQuant: Towards 10 Million Context Length LLM Inference with KV Cache Quantization, 2024

# Speculative/Parallel Decoding

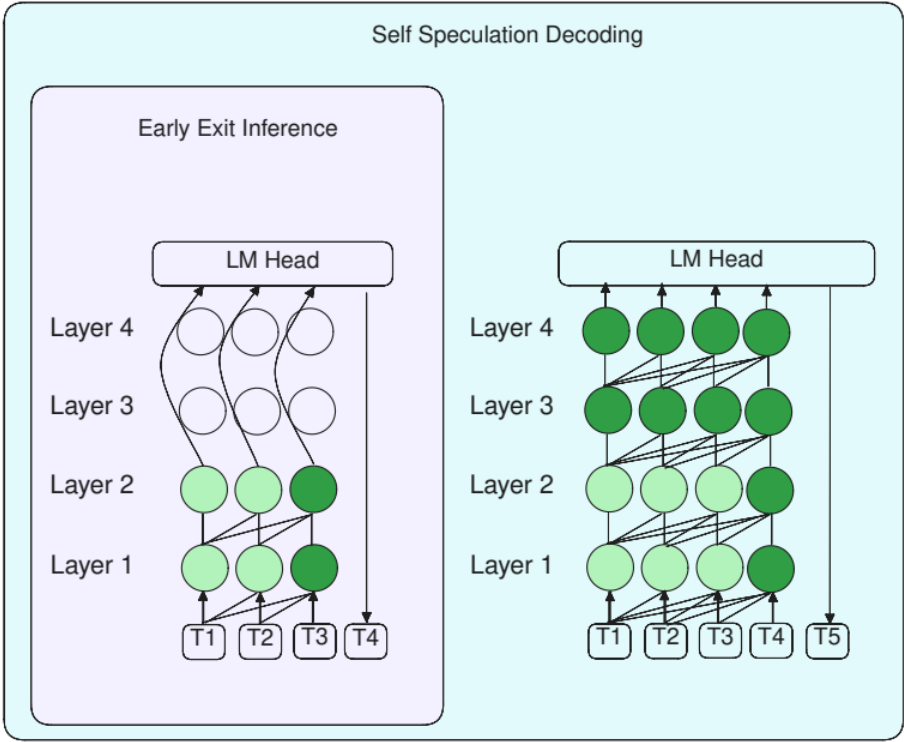
## MEDUSA: Simple LLM Inference Acceleration Framework with Multiple, 2024



Fast Inference from Transformers via Speculative Decoding, 2023



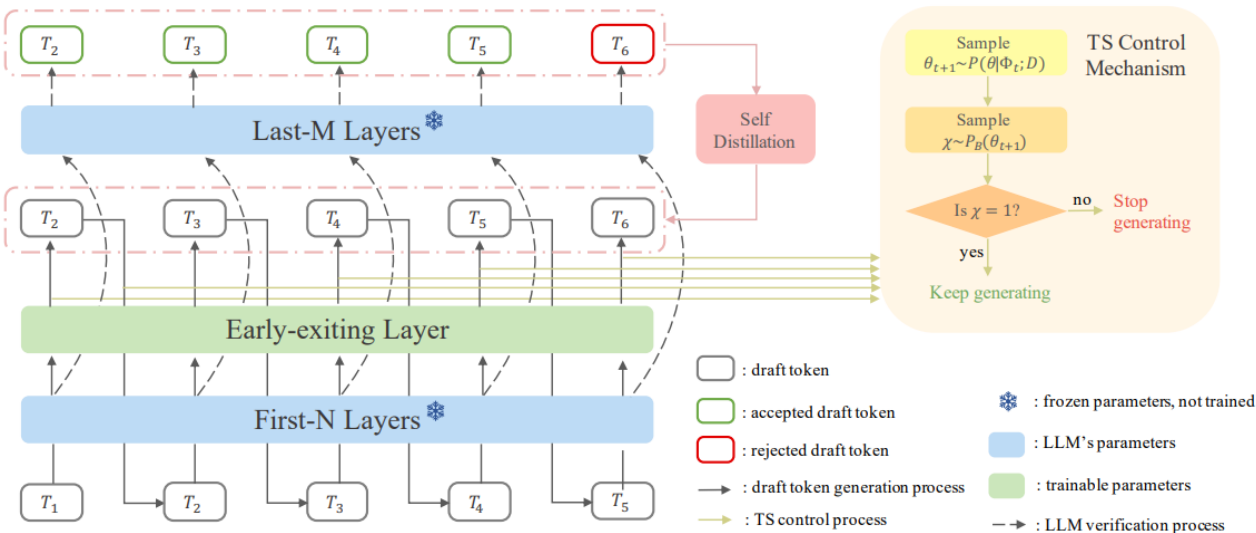
# Early-Exit Inference



... enables inference with subset of layers with higher accuracy...

... and we can improve accuracy by verifying and correcting with remaining layers

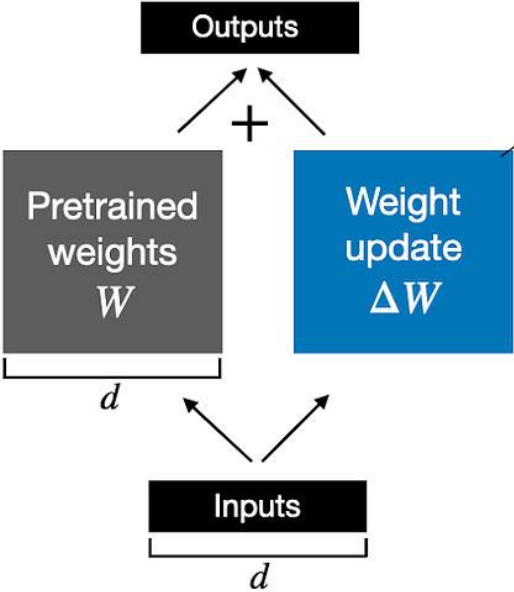
## Speculative Decoding via Early-exiting for Faster LLM Inference with Thompson Sampling Control Mechanism, 2024



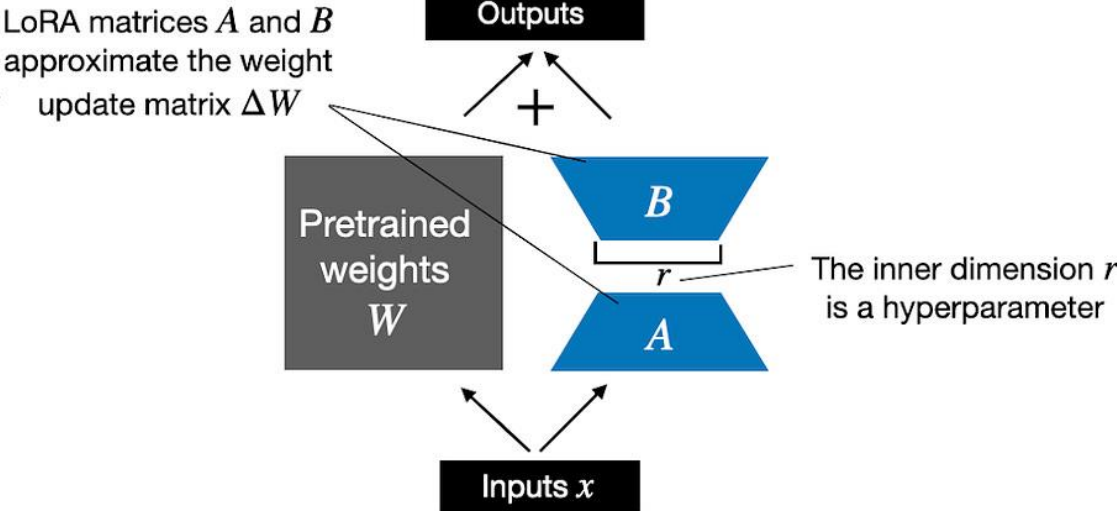
## LayerSkip: Enabling Early Exit Inference and Self-Speculative Decoding, 2024

# Efficient Fine-Tuning

Weight update in **regular finetuning**



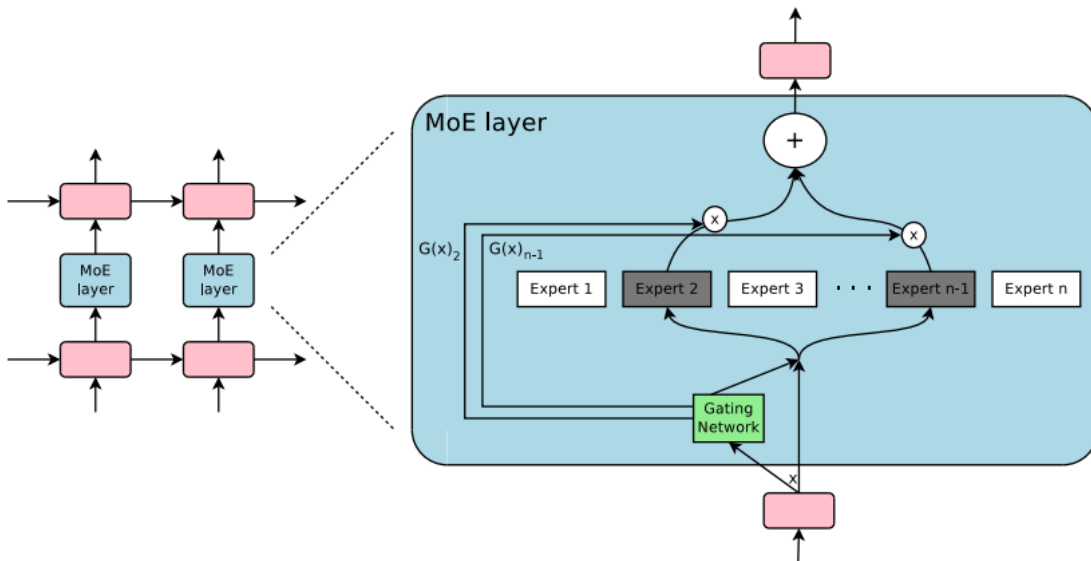
Weight update in **LoRA**



LoRA: Low-Rank Adaptation of Large Language Models, 2021

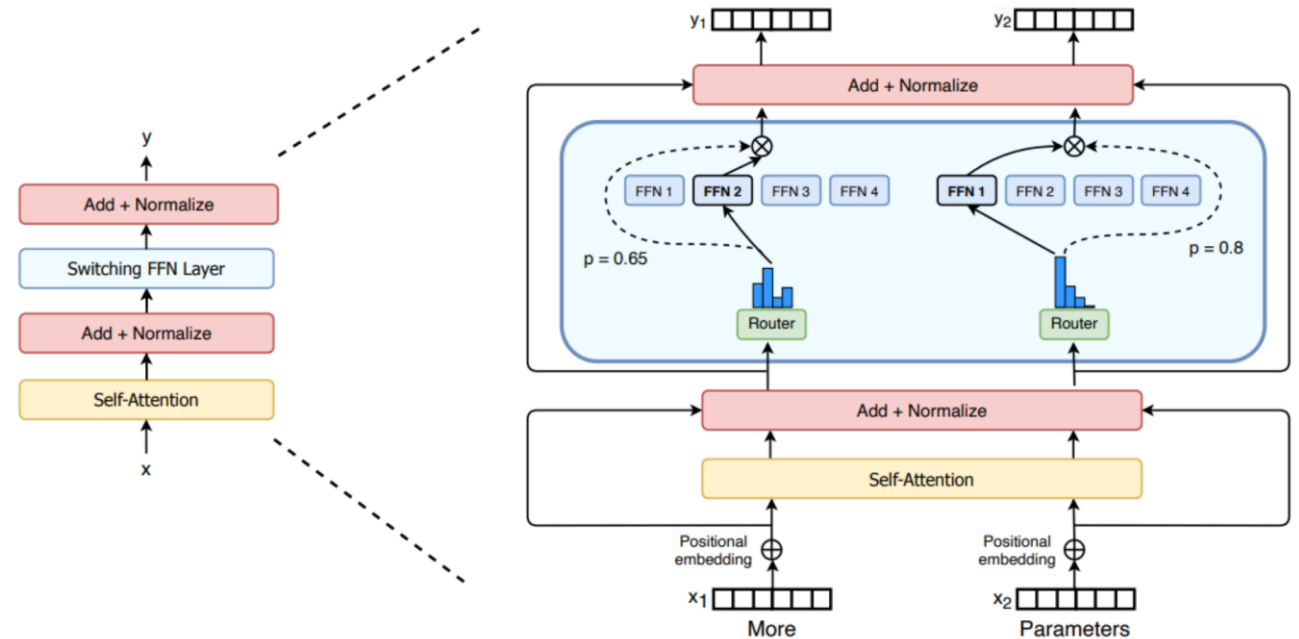
QLoRA: Efficient Finetuning of Quantized LLMs, 2024

# Mixture-of-Expert Models are Sparse and Need Less Compute



Outrageously Large Neural Networks: The Sparsely-Gated Mixture-of-Experts Layer, 2017

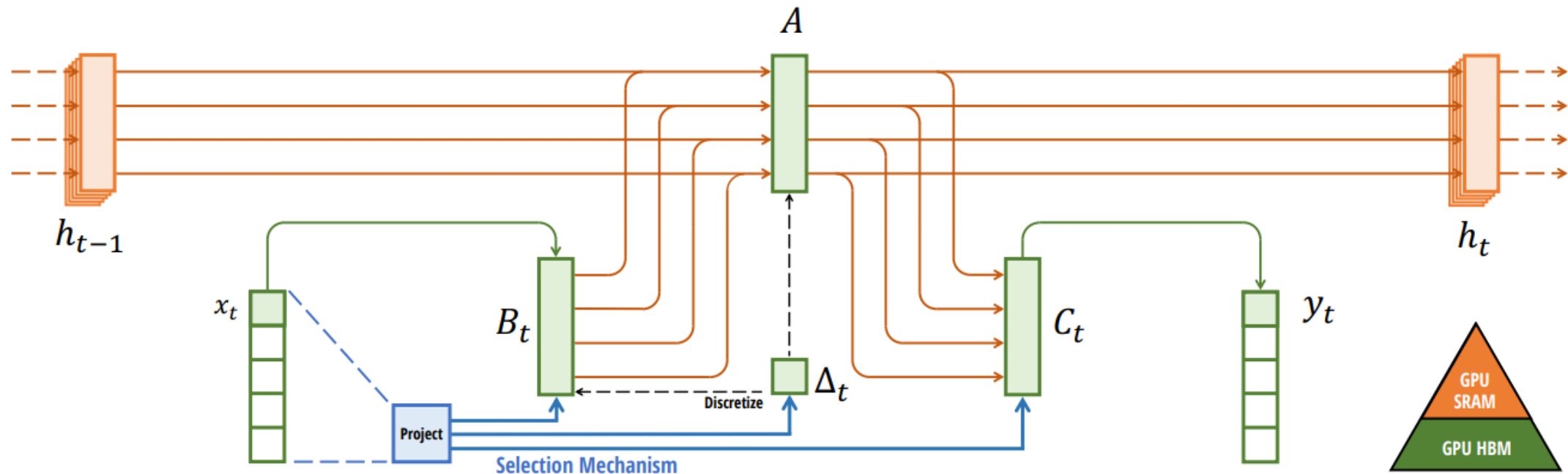
Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity, 2021





# Mamba – Linear Time Sequence Model

## Selective State Space Model *with Hardware-aware State Expansion*



Mamba: Linear-Time Sequence Modeling with Selective State Spaces, 2024

QA