CSCI 5541: Natural Language Processing

Lecture 15: LLM Compute efficiency and engineering

James Mooney

With slides borrowed from Song Han (MIT)



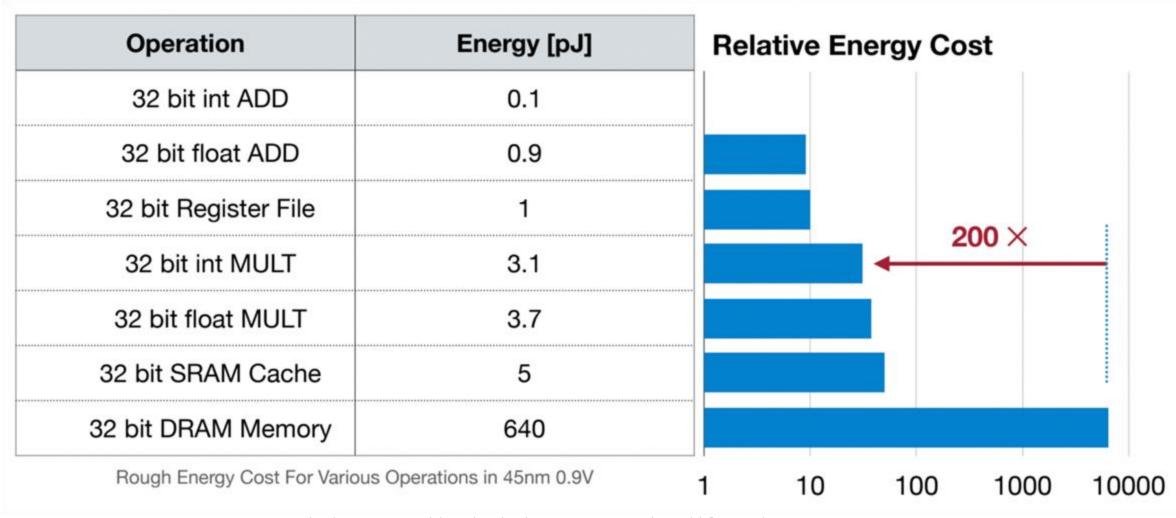
What Is Efficiency and Why Does It Matter?

- ☐ 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 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 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

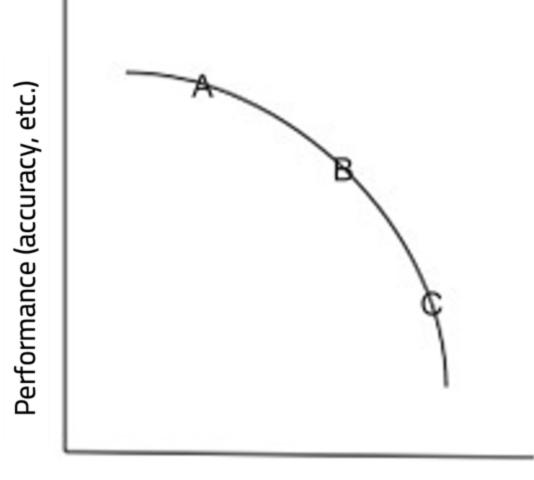
Model Energy Use



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

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



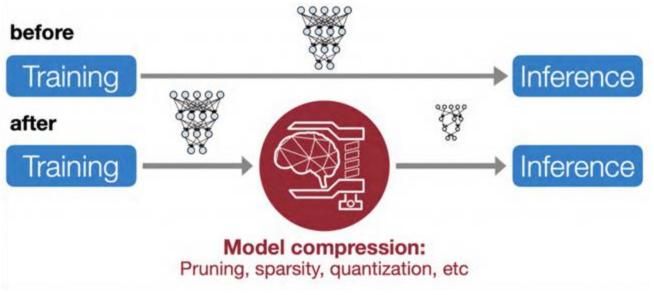
Efficiency (speed, 1/size, etc.)

How to Improve Model Efficiency?

Hardware

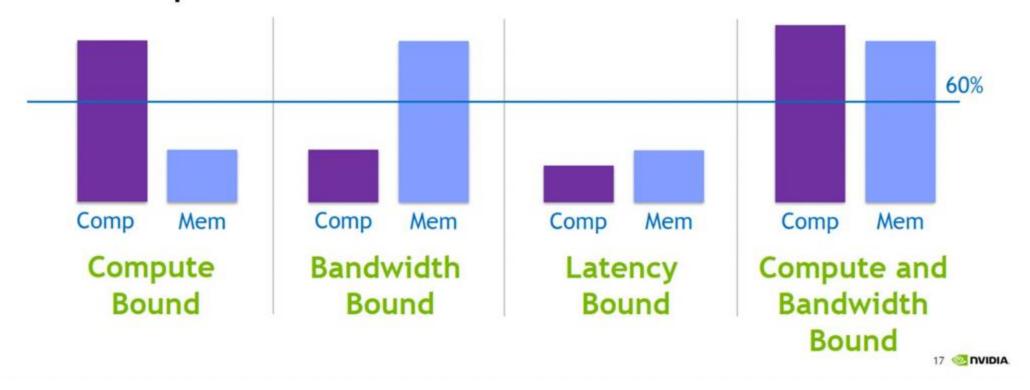
Software





What Makes a Language Model Slow

Memory Utilization vs Compute Utilization Four possible combinations:



Efficient LLMs

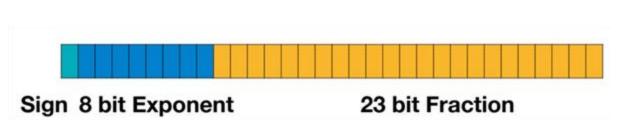
- Quantization
- Sparsity
- ☐ Long Context
- ☐ Serving & Systems
- Distillation
- ☐ Parameter Efficient Fine-Tuning
- Alternative Paradigms

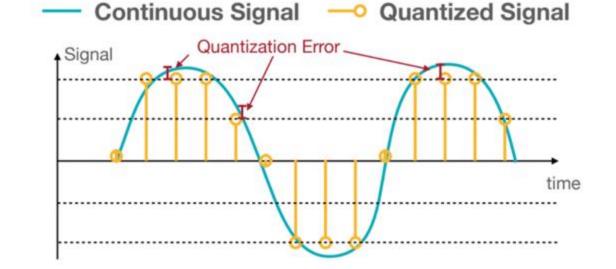
Efficient LLMs

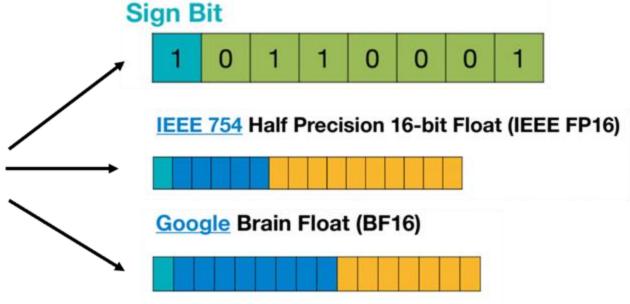
- Quantization
- Sparsity
- ☐ Long Context
- ☐ Serving & Systems
- Distillation
- ☐ Parameter Efficient Fine-Tuning
- ☐ Alternative Paradigms

Quantization

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



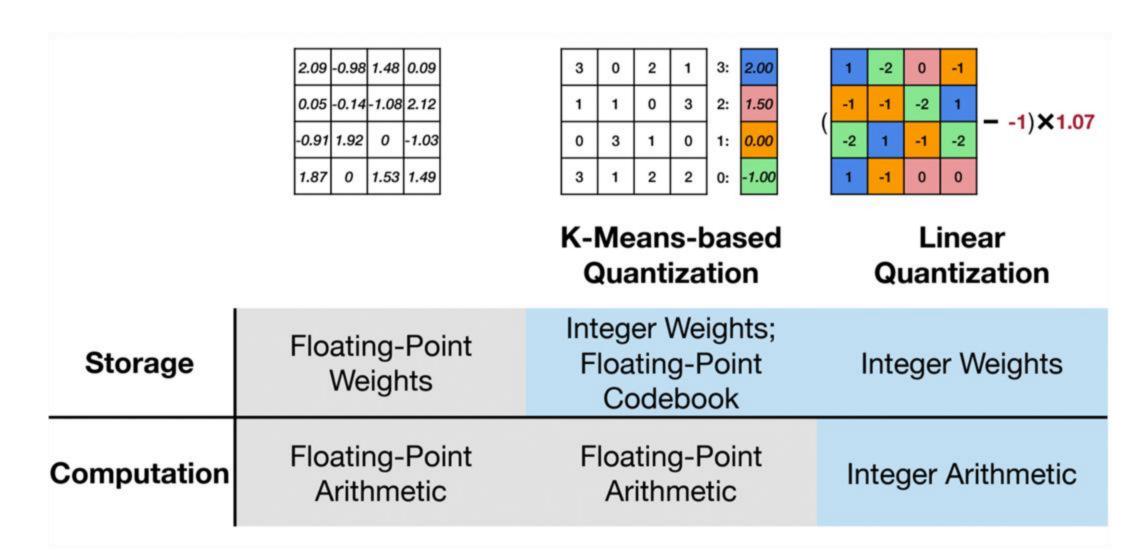




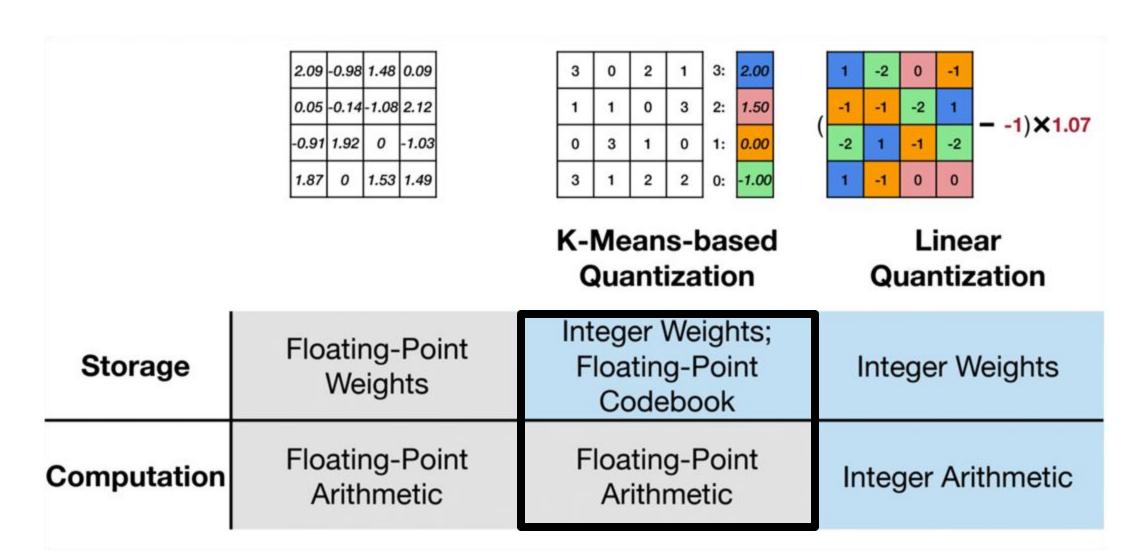
How do we go from a high-bit width data type to a low-bit width data type?

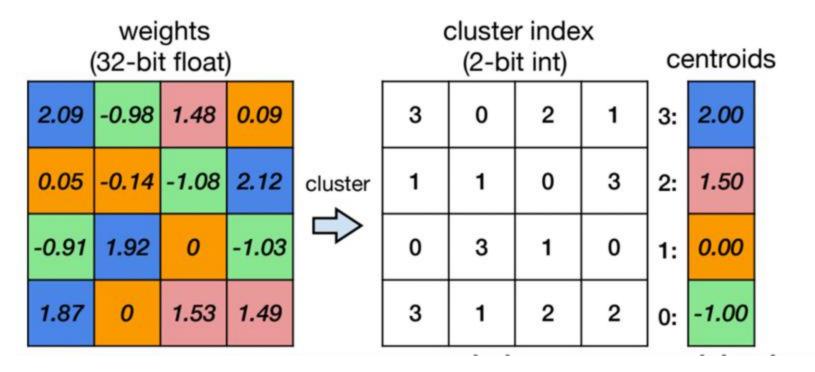


K-Means Quantization vs Linear Quantization

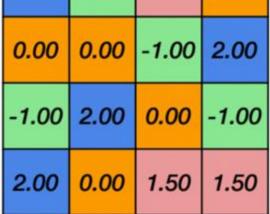


K-Means Quantization vs Linear Quantization





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





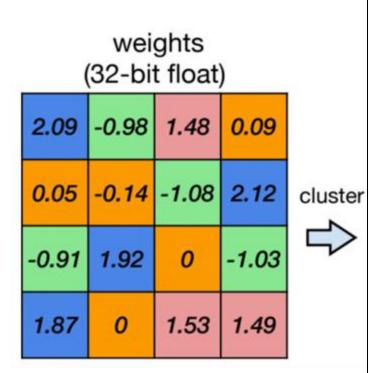
reconstructed weights (32-bit float)

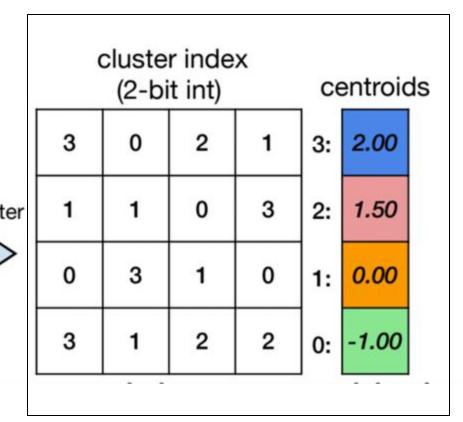
2.00	-1.00	1.50	0.00
0.00	0.00	-1.00	2.00
-1.00	2.00	0.00	-1.00
2.00	0.00	1.50	1.50

Deep Compression [Han et al., ICLR 2016]

Stored weights after clustering







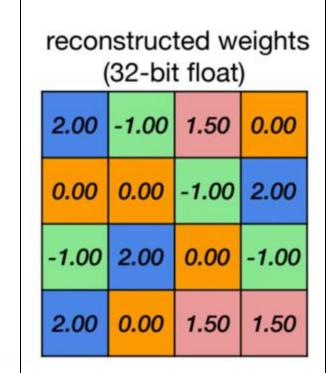
reconstructed weights (32-bit float)

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]

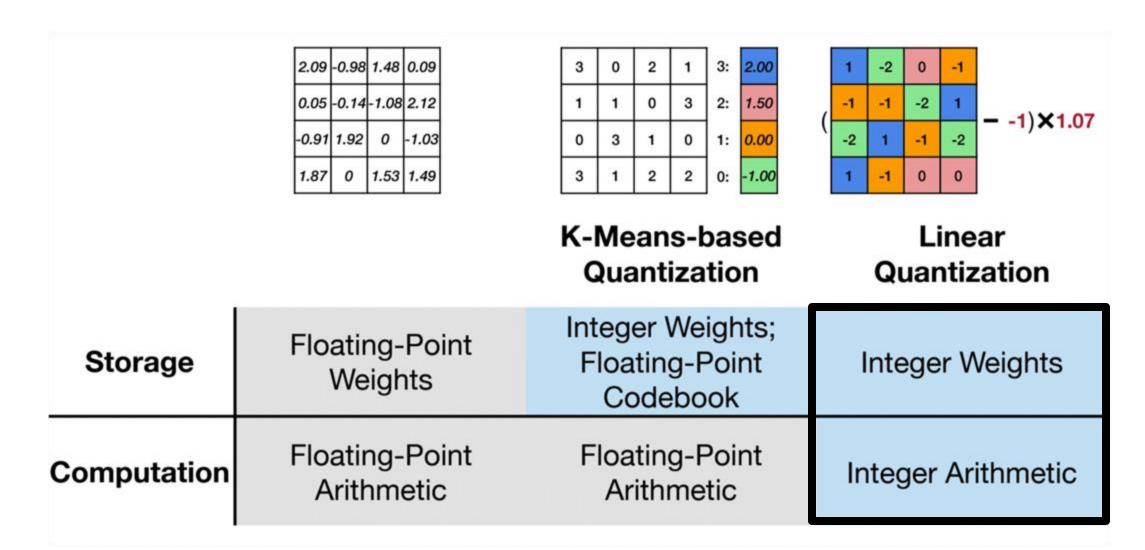
Retrieved weights to be used at inference time



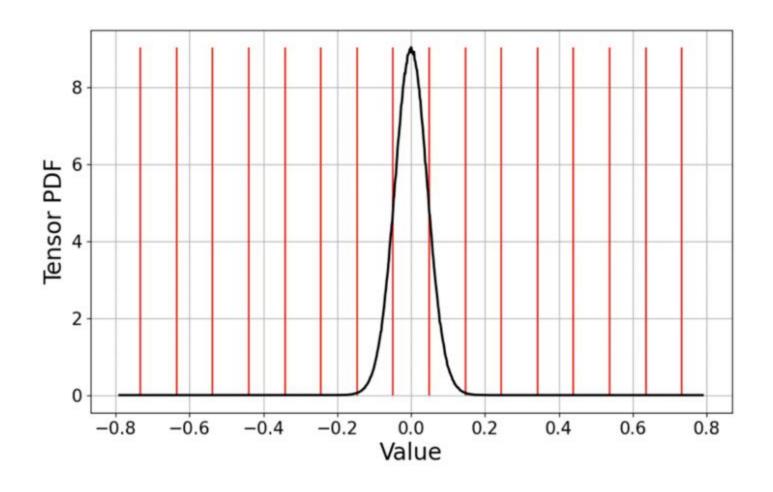


Deep Compression [Han et al., ICLR 2016]

K-Means Quantization vs Linear Quantization



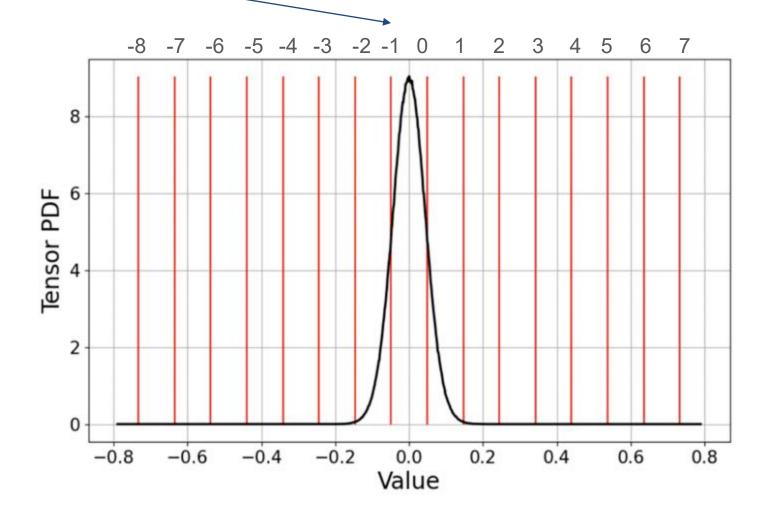
- □ Apply linear function on weights and hidden state activations from floating point values (r) to integer values (q)
- Original weights (black),Quantized bins (red)
- Black weights are mapped to one of the vertical red lines



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

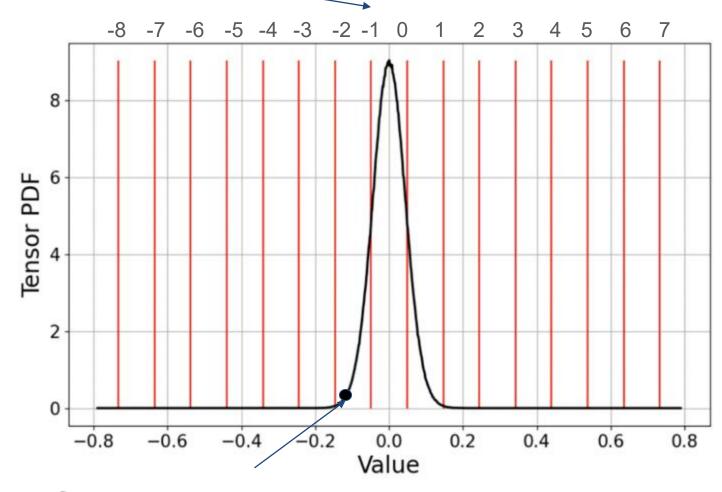
values (q)

□ Black weights are mapped to one of the vertical red lines 32-bit float to 4-bit int



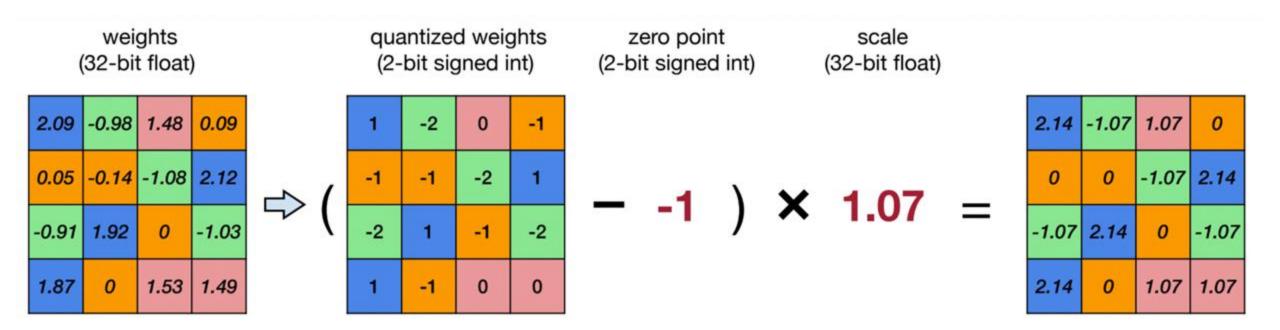
- □ Apply linear function on weights and hidden state activations from floating point values (r) to integer values (q)
- Original weights (black),Quantized bins (red)
- □ Black weights are mapped to one of the vertical red lines



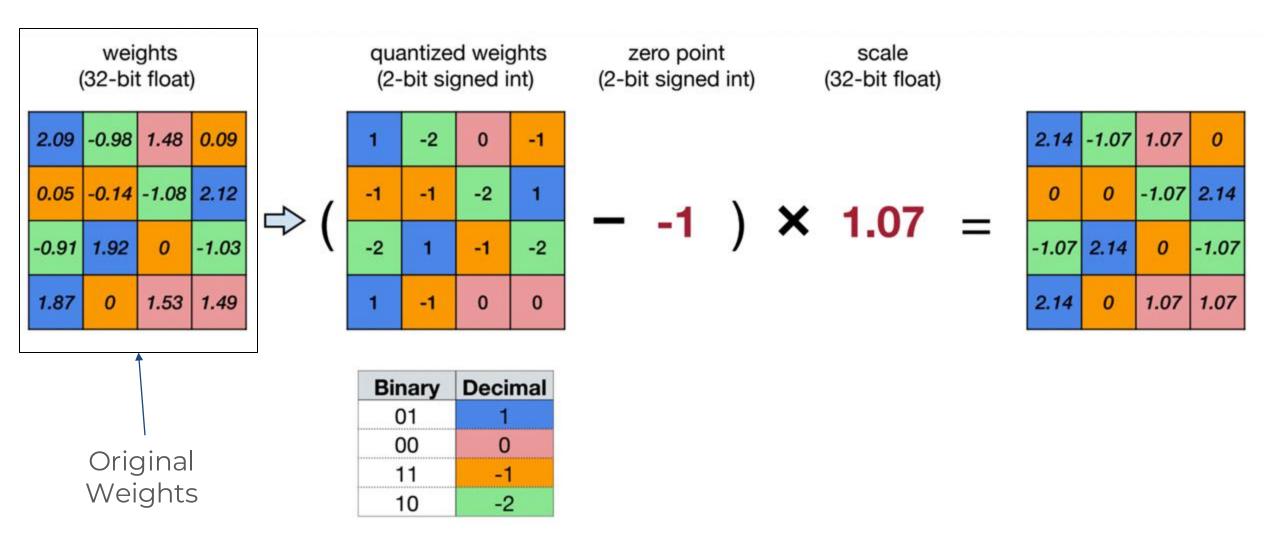


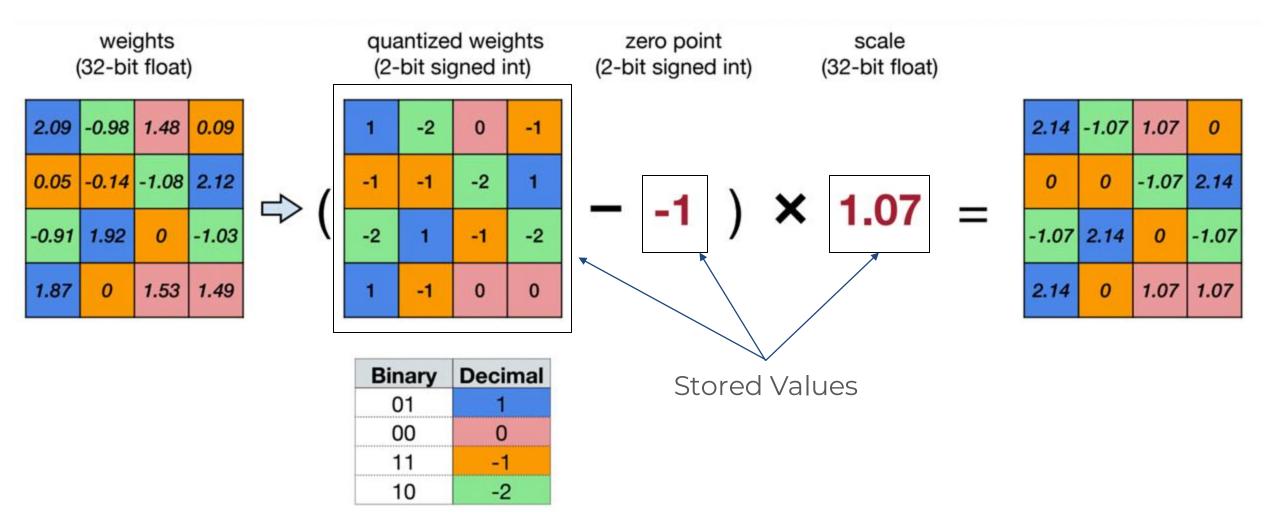
Before Quantization: -.14

After Quantization: -2



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



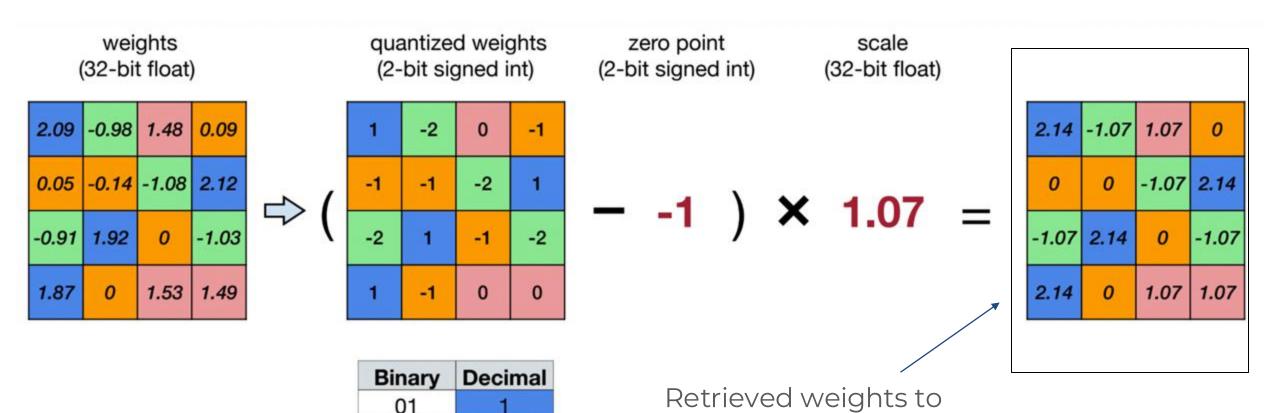


00

11

10

-2



be used at inference time

00

11

10

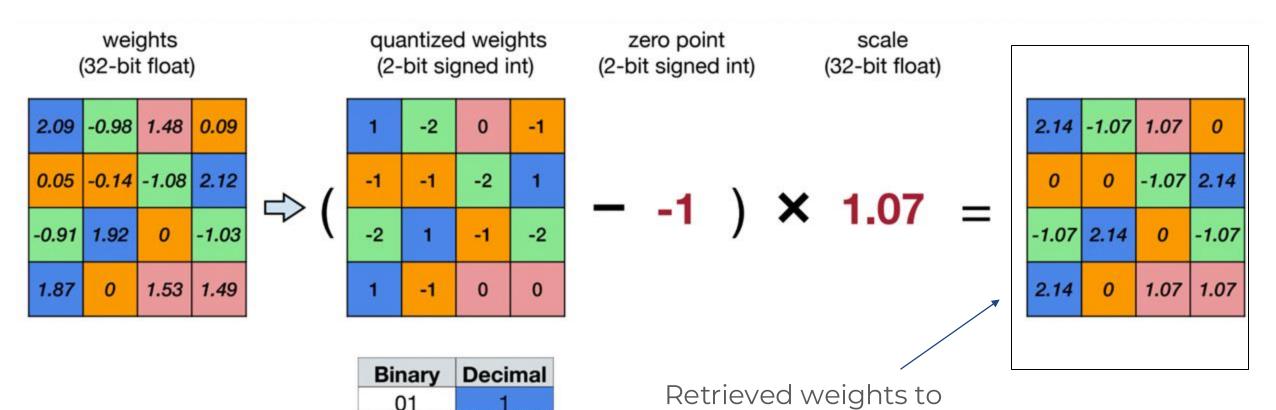
0

-2

Virtually all modern methods use this (as opposed to KMeans Quantization)

be used at inference

time



Should we use KMeans/Linear Quantization on all weights at the same time?

Should we use KMeans/Linear Quantization on all weights at the same time?

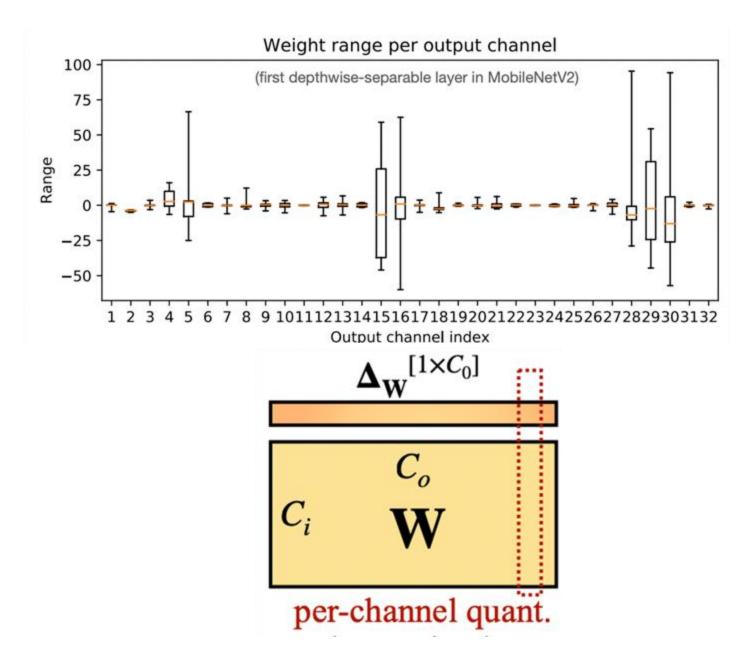
No. We should group them.

The size of the grouping we choose is denoted as the granularity.



Weight Granularity

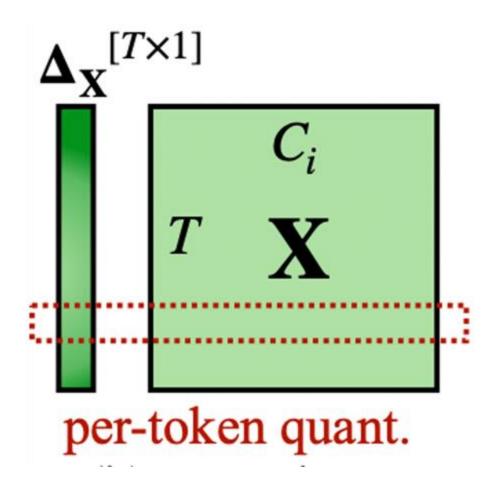
- Weight matrices will often have different variances along each output channel
- High variance in weights means that applying linear quantization will result in large performance degradation
- □ To fix this, we can perform linear quantization along each channel of the weight tensor separately



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



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

When do we perform this quantization?



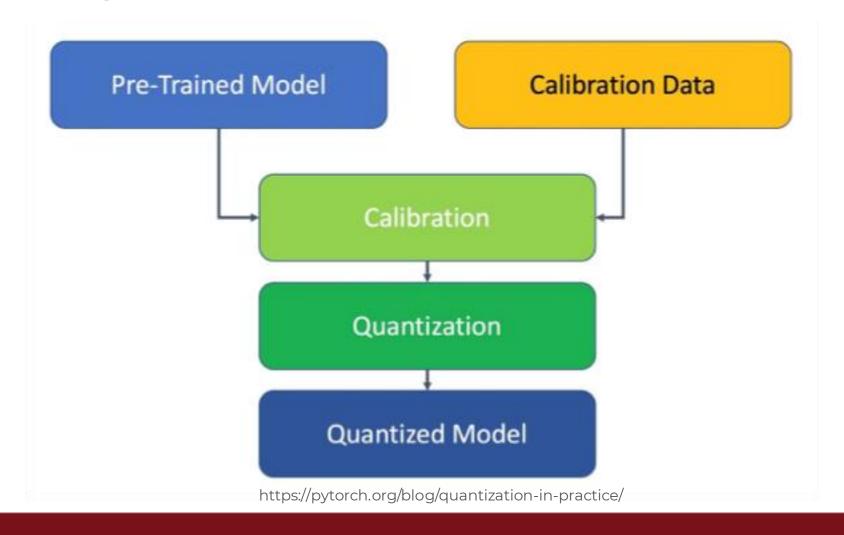
When do we perform this quantization?

Typically, this is done after first completely training a model (pretraining → SFT → post-training)



Post Training Quantization (PTQ)

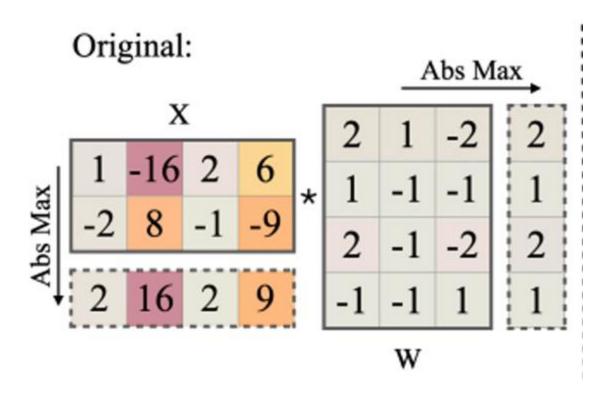
Quantize after training



What are some of the modern methods for performing quantization?



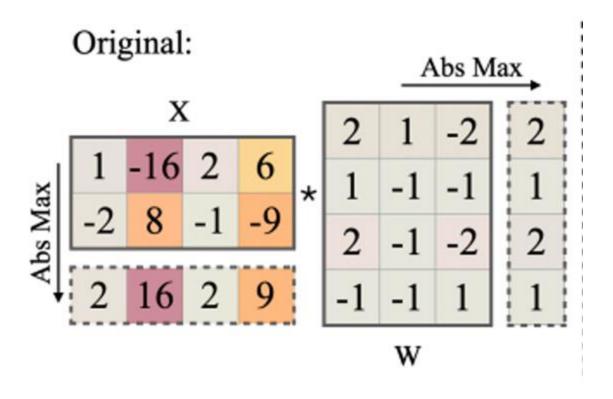
AWQ (W4A16)



Observation: ?

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

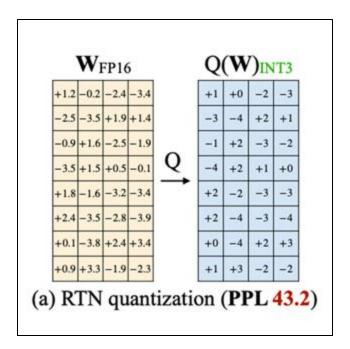
AWQ (W4A16)

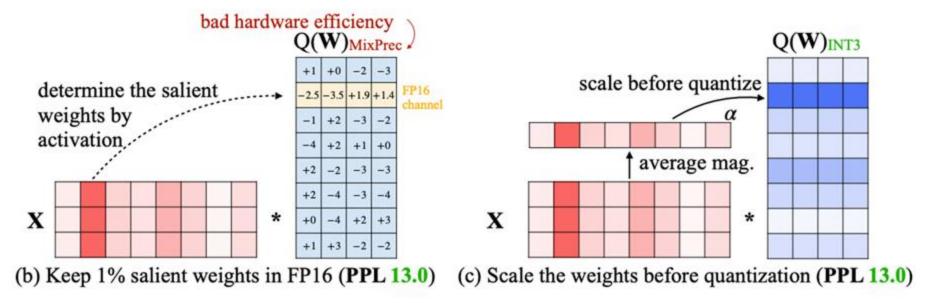


Observation: High variance channels are fixed in activations in LLM FFN layers-weights have relatively little difference in variance

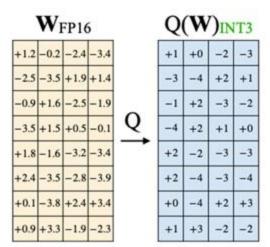
SmoothQuant: Accurate and Efficient Post-Training Quantization for Large Language Models[Xiao et. al., ICML 2023]

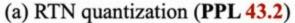
Normal quantization on LLMs performs poorly due to outliers in the model's hidden state

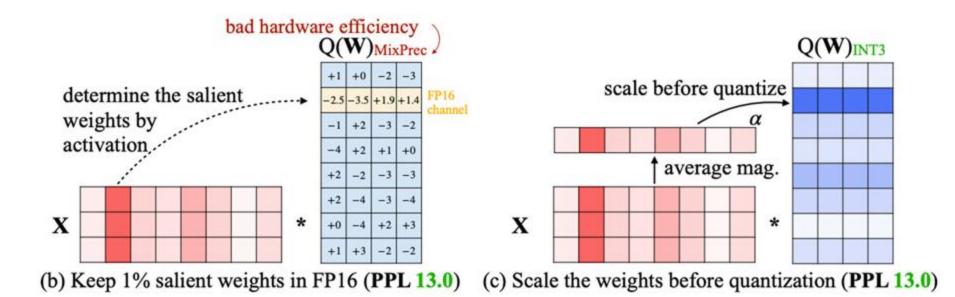




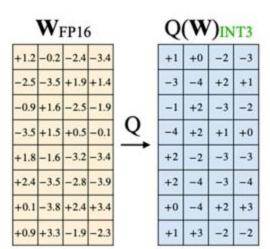
AWQ (c below) scales weights to handle the channels separately

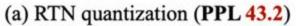


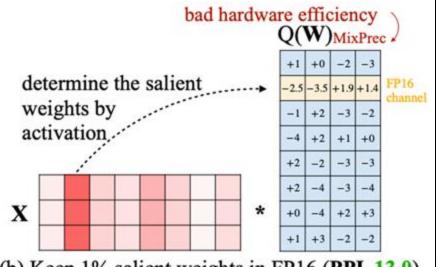




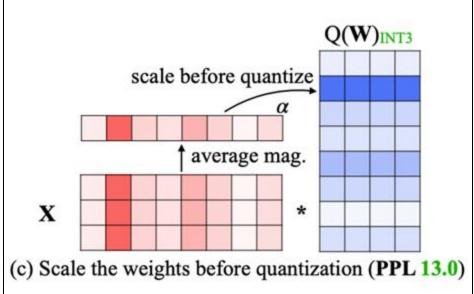
Uses a calibration batch of data to determine what the outlier channels are along with a perchannel scaling factor for these channels



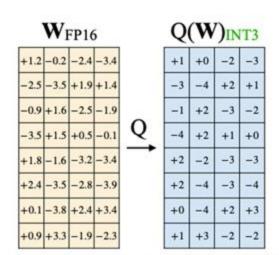




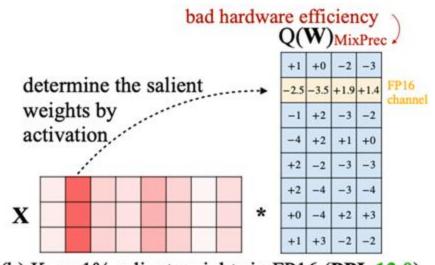
(b) Keep 1% salient weights in FP16 (PPL 13.0)



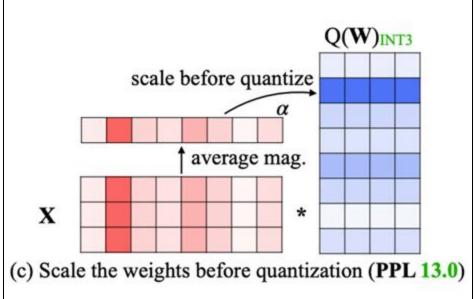
AWQ is weight-only → multiplication is still performed in BF16/FP16.



(a) RTN quantization (PPL 43.2)



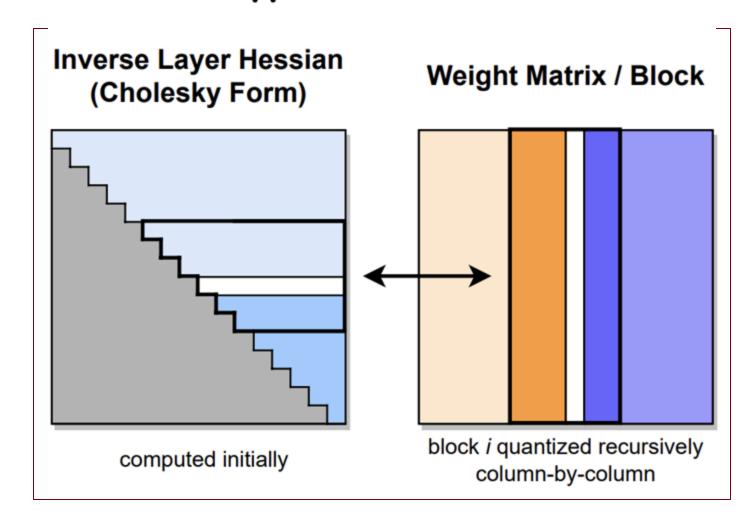
(b) Keep 1% salient weights in FP16 (PPL 13.0)



GPTQ

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

- Quantizes a block of each weight matrix at a time.
- Uses a loss term + the inverse hessian for the layer to update the nonquantized weights
- Takes a long time



GPTQ: Accurate Post-Training Quantization for Generative Pre-Trained Transformers [Frantar et al., arXiv 2022]

GGUF Quantization

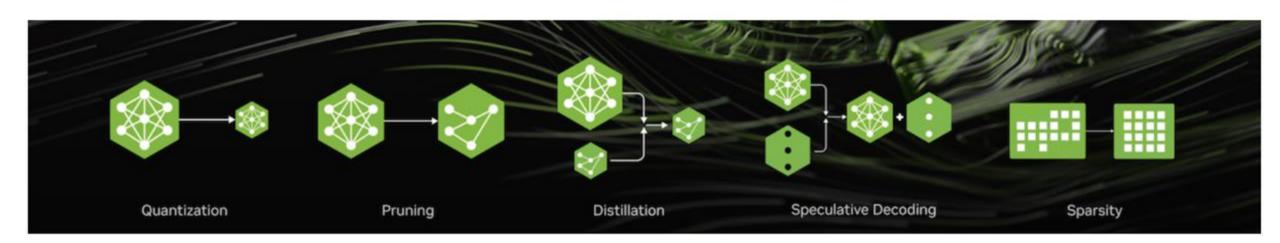
- No model retraining as is done with GPTQ
- A model is quantized by taking blocks of the input matrix, then performing linear quantization on each block separately
- Uses several different methods based on level of quantization
- Will come with names such as IQ2_XS-IQ4_XS, Q2_K_S-Q5_K_S
- You can use <u>unsloth</u> to create models
 with this type
 GPTQ: Accurate Post-Training Quantization for Generative Pre-Trained Transformers [Frantar et al., arXiv 2022]



,

FP8 Quantization (Nvidia)

- Offered by the <u>NVIDIA TensorRT Model Optimizer</u> library
- Performs a separate per-channel post-training quantization
- Supported by hardware
- This library also performs a number of other functions for improving model efficiency



Modern Quantization (Open Source Models)

Oftentimes, open source LLM releases from model providers will release versions of models that use several of the above quantization techniques

```
Qwen/Qwen3-235B-A22B-Thinking-2507-FP8

    Text Generation - .: 235B - Updated Jul 30 - ± 17.9k - ♥ 64

Qwen/Qwen3-235B-A22B-Thinking-2507

▼ Text Generation • ... 235B • Updated Aug 17 • ± 47.3k • • • ○ 377

Qwen/Qwen3-235B-A22B-Instruct-2507-FP8

    Text Generation → ...i 235B → Updated Sep 17 → ± 97.2k → ♥ 134

0wen/Owen3-235B-A22B-Instruct-2507

▼ Text Generation • ... 235B • Updated Sep 17 • ± 108k • • ○ 714

Owen/Owen3-30B-A3B-Thinking-2507-FP8

    Text Generation - .:: 31B - Updated Jul 30 - ± 20.2k - ♥ 45

  Qwen/Qwen3-30B-A3B-Thinking-2507

▼ Text Generation - ...: 31B - Updated Aug 17 - ± 267k - * - ♥ - ♥ 315

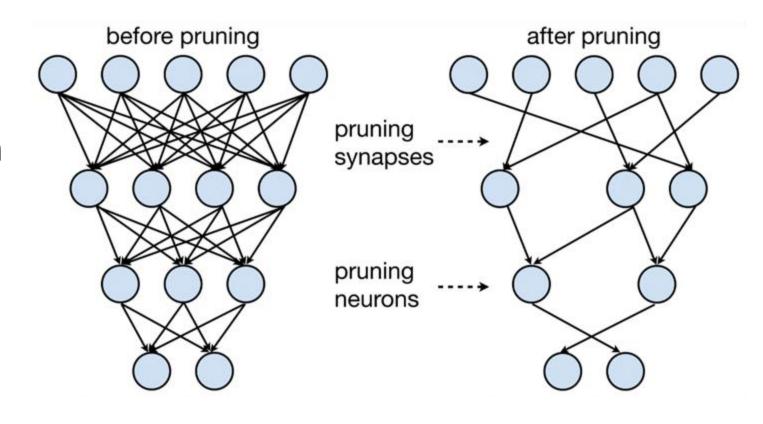
Qwen/Qwen3-30B-A3B-Instruct-2507-FP8
Text Generation • .:: 31B • Updated Sep 17 • ± 486k • ♥ 92
```

Efficient LLMs

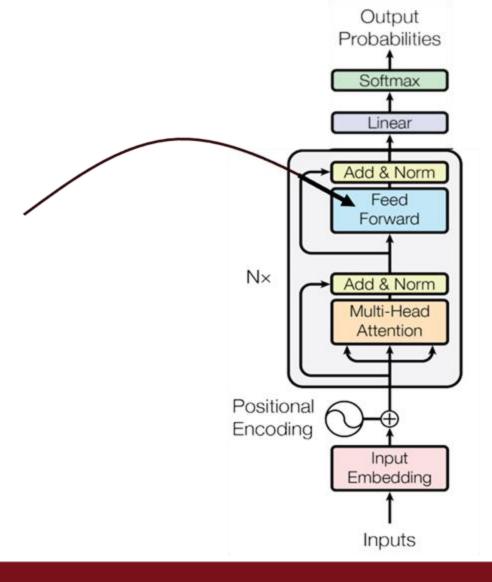
- Quantization
- Sparsity
- ☐ Long Context
- ☐ Serving & Systems
- Distillation
- ☐ Parameter Efficient Fine-Tuning
- Alternative Paradigms

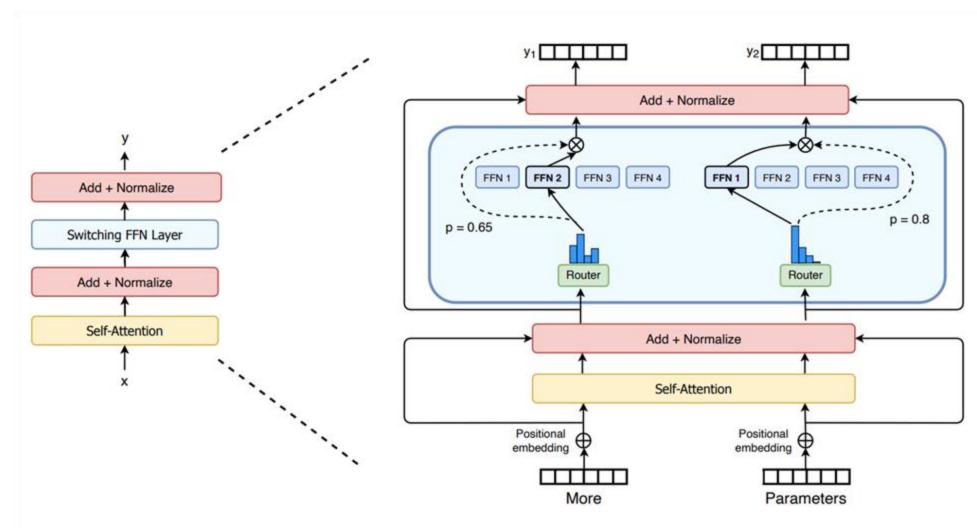
Sparsity

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

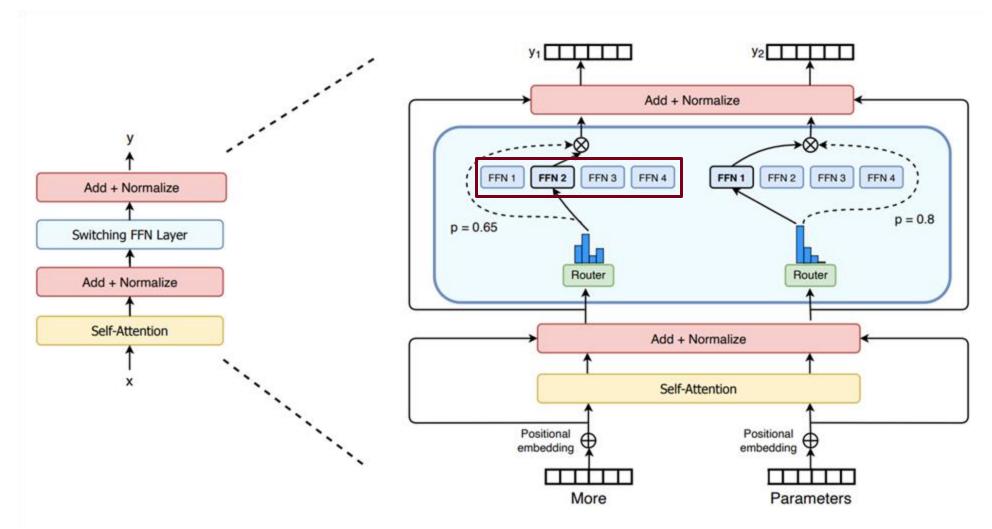


Replace FFN layers in traditional transformers with a switching FFN layer (more generally called an MoE layer)

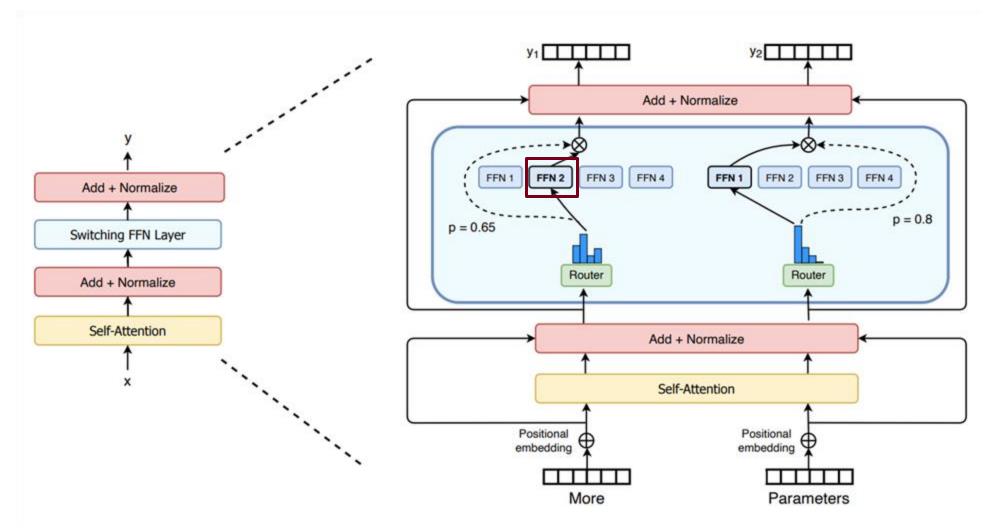




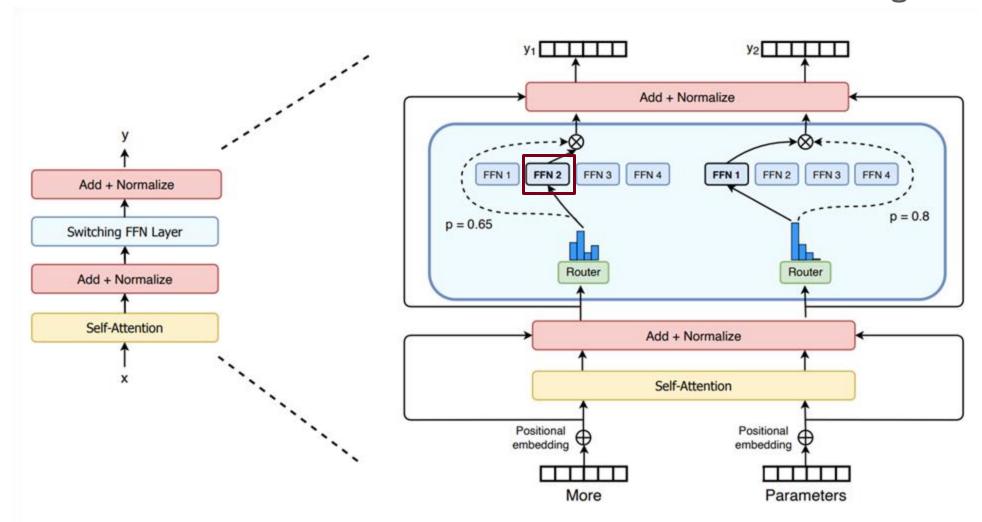
Four FFN layers



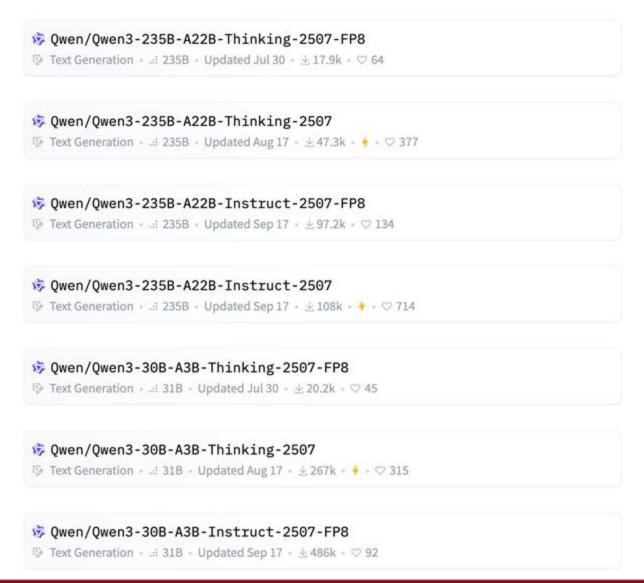
Only one is used per token



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

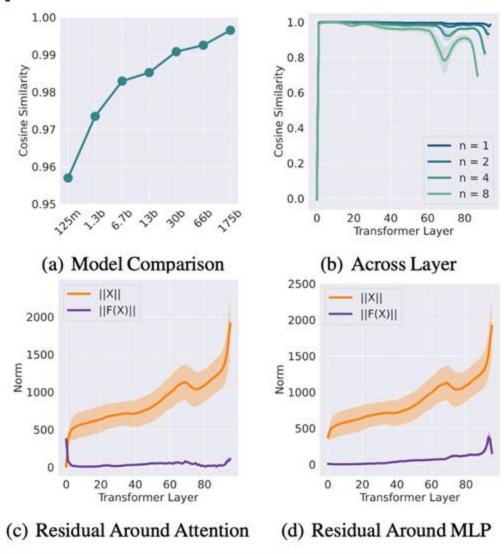


Most of the best-performing open source models produced today are mixture of experts models



Deja Vu: Contextual Sparsity

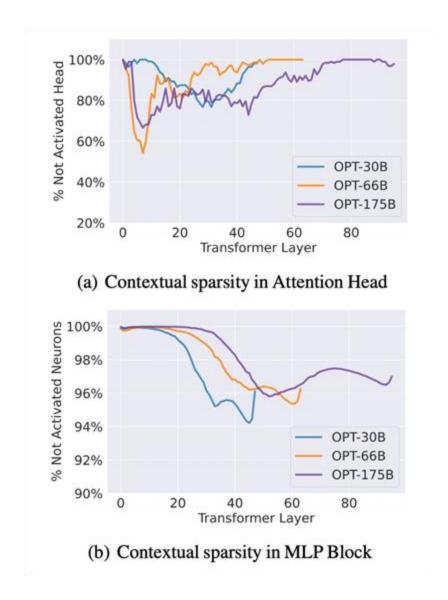
Observation 1: Model activations change very little between consecutive layers of a network



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

Deja Vu: Contextual Sparsity

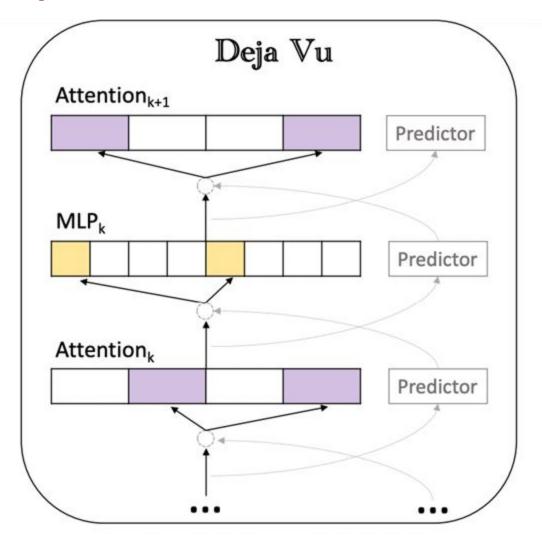
Observation 2: Most attention heads and most neurons are not used



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

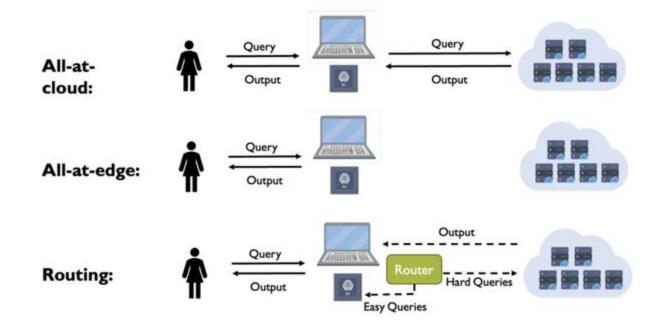


Deja Vu: Contextual Sparsity for Efficient LLMs at Inference Time [Liu et al., 2023]



Model Routing

Route easier queries to smaller models, harder queries to larger models

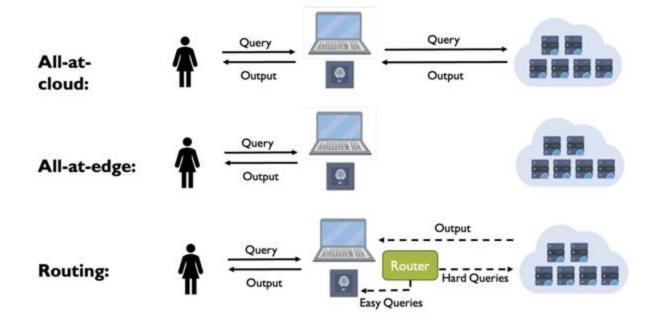


Hybrid LLM: Cost-Efficient and Quality-Aware Query Routing [Ding et al., 2024]

Model Routing

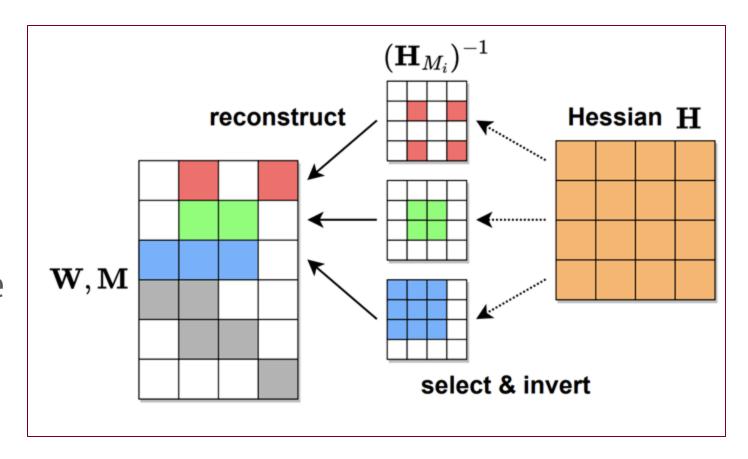
Valuable for companies like OpenAl/Anthropic when they serve multiple models of different capabilities & sizes

Route easier queries to smaller models, harder queries to larger models



SparseGPT

- Similar to GPTQ, uses a reconstruction loss term to prune entries from each weight matrix
- This can be taken advantage of with either 2:4, 4:8 sparsity on Nvidia GPUs



Efficient LLMs

- Quantization
- Sparsity
- Long Context
- ☐ Serving & Systems
- Distillation
- ☐ Parameter Efficient Fine-Tuning
- Alternative Paradigms

What is a central issue with having longer context models?



Attention(Q, K, V) = softmax
$$\left(\frac{QK^{\top}}{\sqrt{d}}\right)V$$

What is a central issue with having longer context models?

The storage costs are linear in generation length and quadratic in computation



Attention(Q, K, V) = softmax
$$\left(\frac{QK^{\top}}{\sqrt{d}}\right)V$$

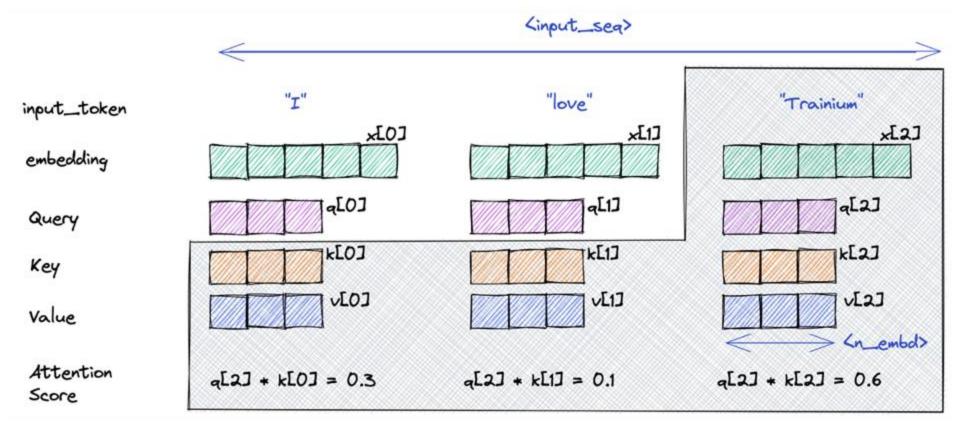
What is a central issue with having longer context models?

The **storage** costs are linear in generation length and quadratic in **computation**



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-Ilm-inference-with-neuron.html

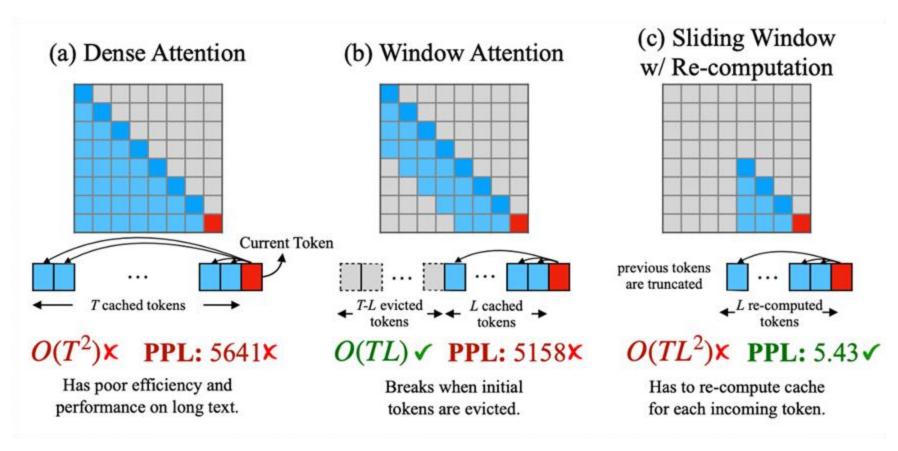
The KV-Cache

In total, we must store

Batch_size * seq_len * num_heads * num_layers * emb_dim * 2

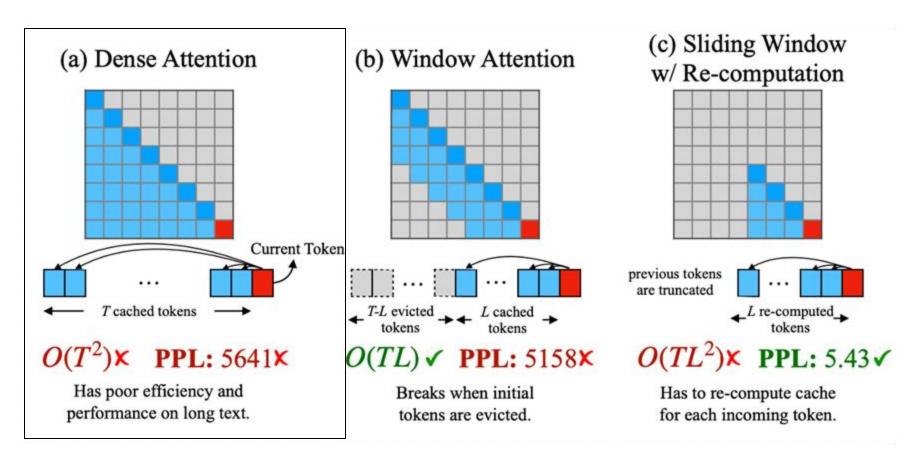
separate values in the kv cache

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

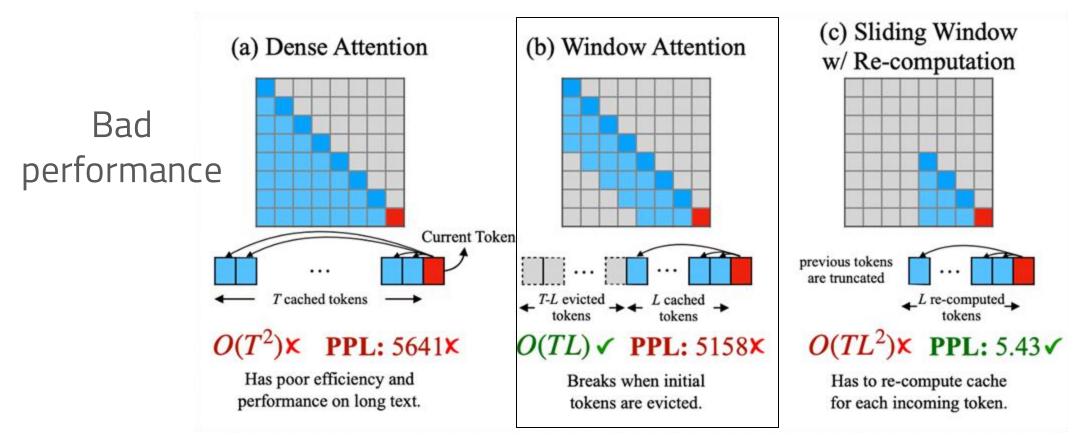


How can we extend models to have much longer context length at minimal cost?

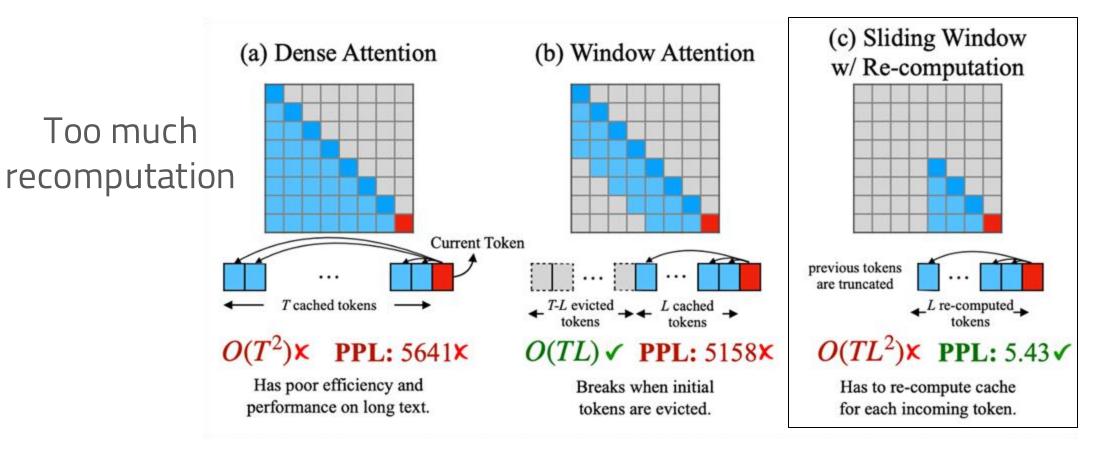
Too much storage



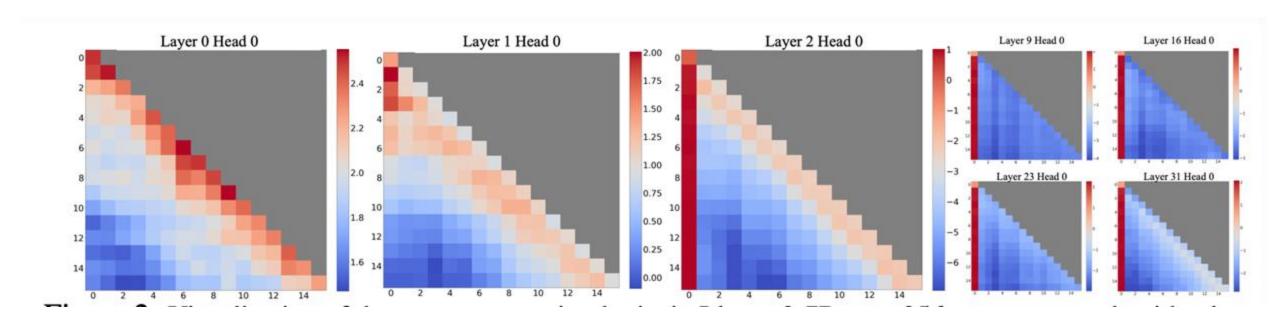
How can we extend models to have much longer context length at minimal cost?

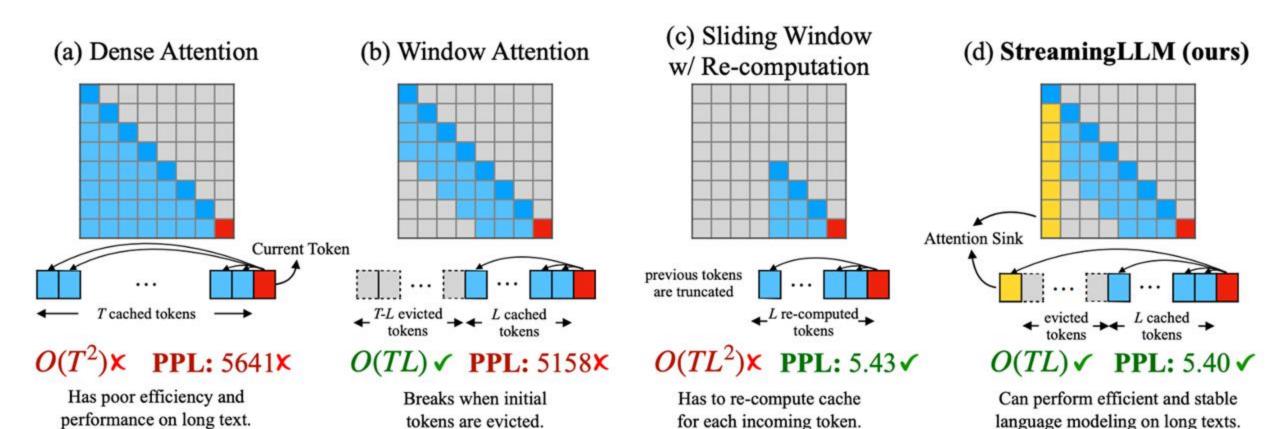


How can we extend models to have much longer context length at minimal cost?



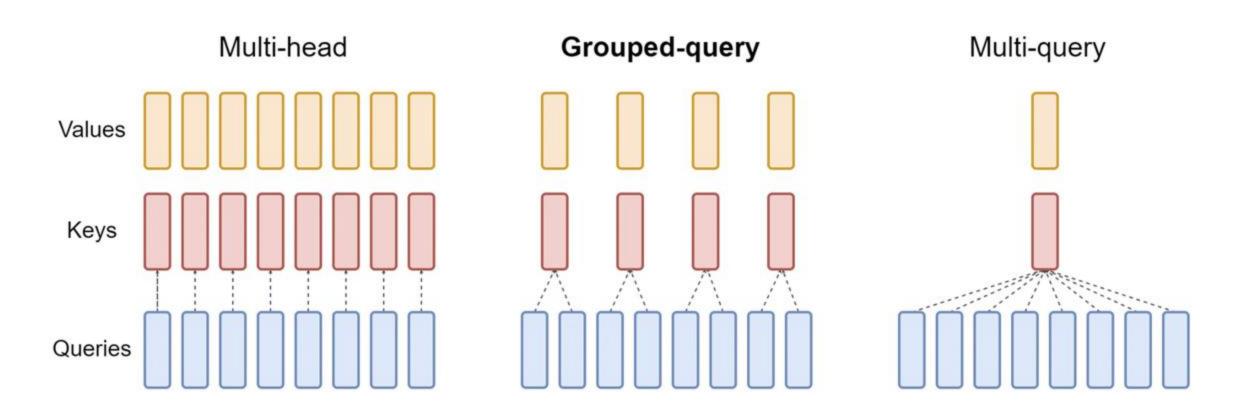
Observation: Most attention is either placed on the first token or to tokens that the model has recently seen.







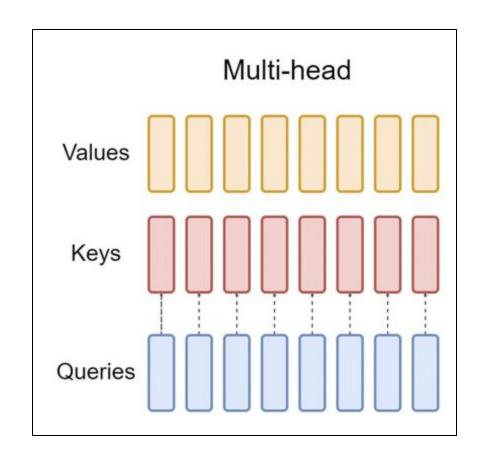
MHA/GQA/MQA

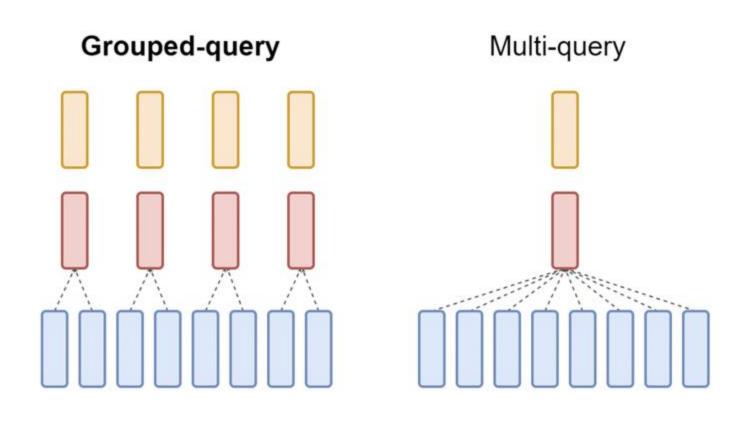


GQA: Training Generalized Multi-Query Transformer Models from Multi-Head Checkpoints

MHA/GQA/MQA

Each attention head calculates separate keys and values for each token

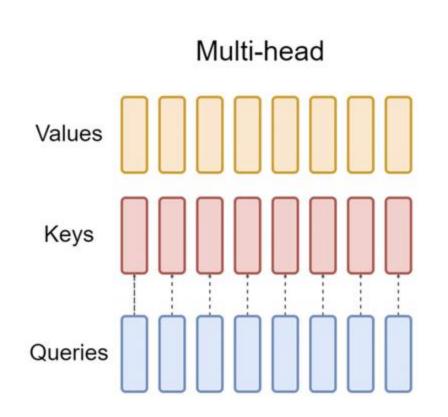


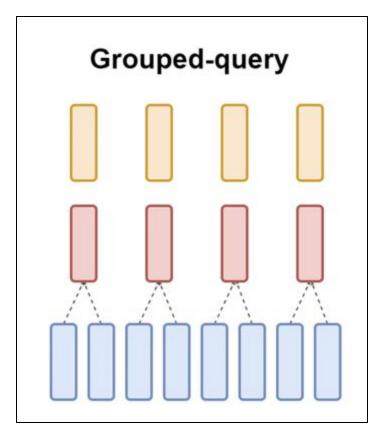


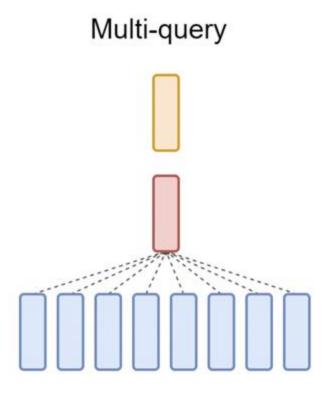
GQA: Training Generalized Multi-Query Transformer Models from Multi-Head Checkpoints

MHA/GQA/MQA

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



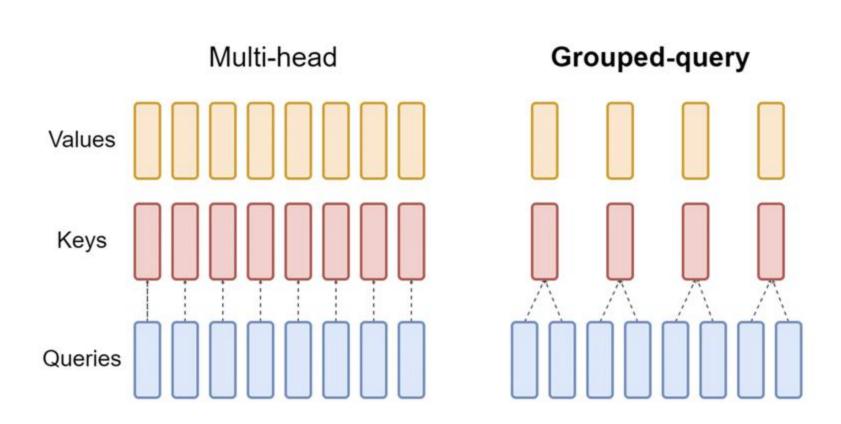


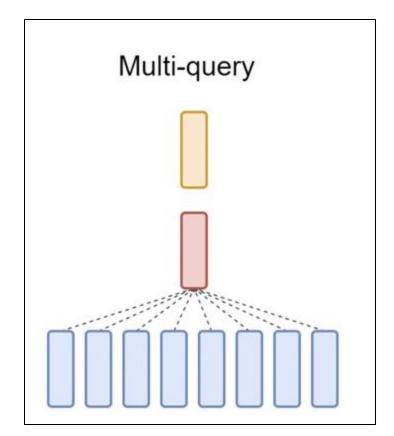


GQA: Training Generalized Multi-Query Transformer Models from Multi-Head Checkpoints

MHA/GQA/MQA

Attention heads share the same keys and values for each token

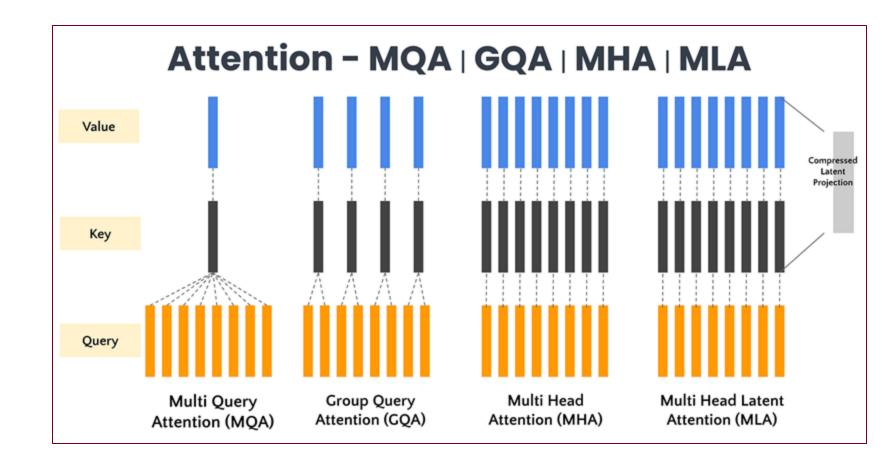




GQA: Training Generalized Multi-Query Transformer Models from Multi-Head Checkpoints

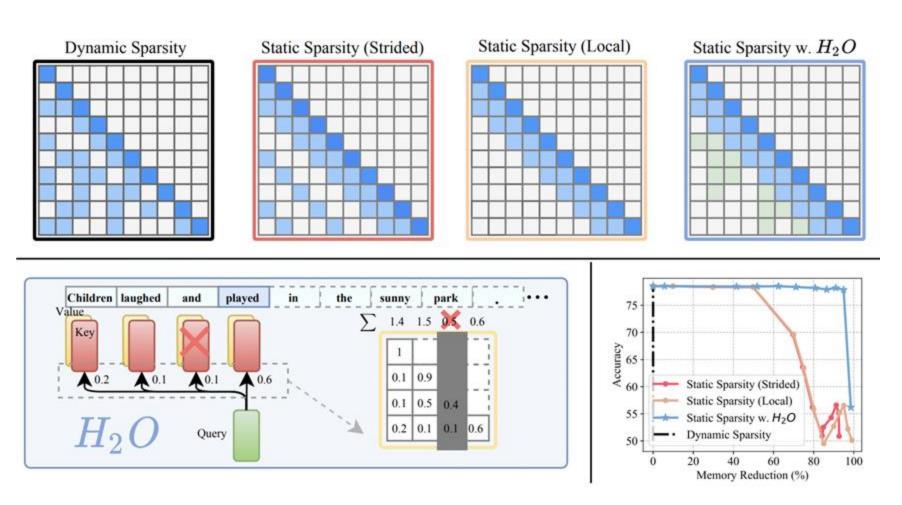
MHA Latent Query Attention

Some modern models
(such as DeepSeek)
train with what is
called latent query
attention → Keys +
Values are compressed
to a lower dimension to
reduce storage costs



H₂0: Heavy Hitter Oracle

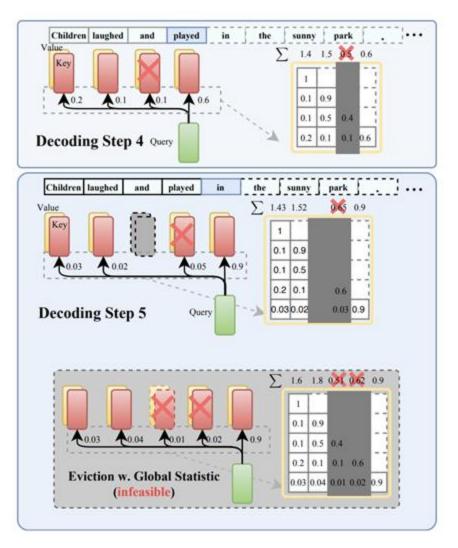
Evict all but k-highest cumulative attention tokens from cache



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

H₂0: 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]

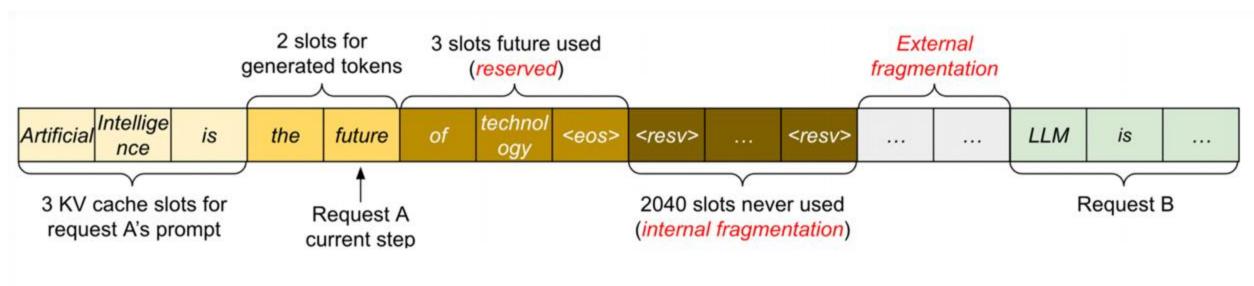
Efficient LLMs

- Quantization
- Sparsity
- ☐ Long Context
- Serving & Systems
- Distillation
- ☐ Parameter Efficient Fine-Tuning
- Alternative Paradigms

PagedAttention

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

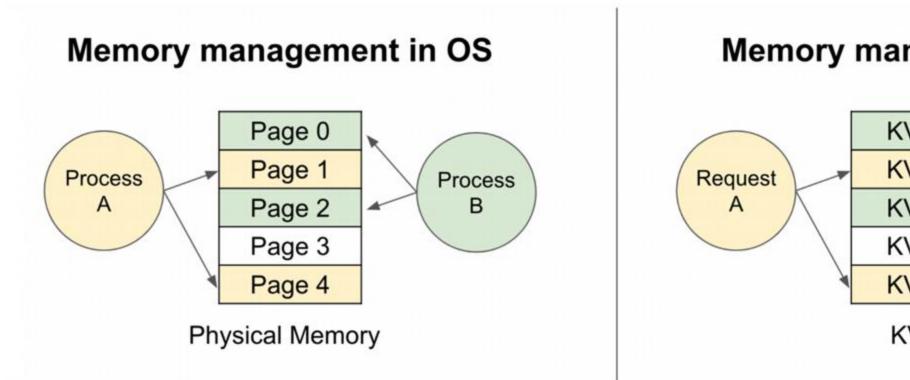
-Originally, most systems just assign fixed sized blocks of memory to each incoming request. How to improve?

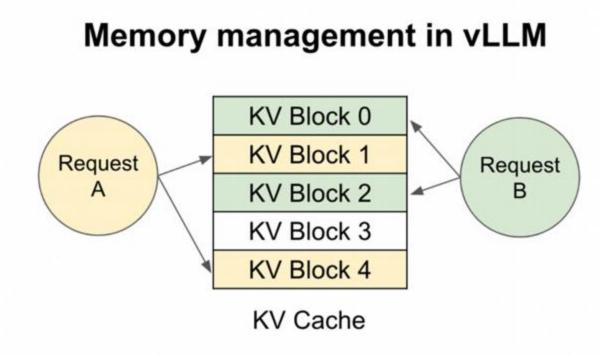


Efficient Memory Management for Large Language Model Serving with PagedAttention (Kwon et al., 2023)

PagedAttention

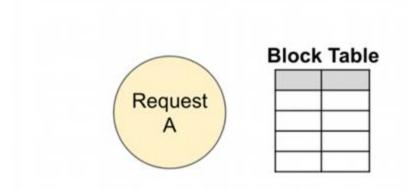
Let's adopt a similar approach to that found in virtual memory!





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

PagedAttention



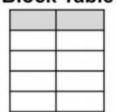
Logical KV blocks

Alan	Turing	is	а		
computer	scientist	and	mathema tician		
renowned					

Physical KV blocks

computer	scientist	and	mathem atician		
Artificial	Intellige nce	is	the		
renowned					
future	of	technolog y			
Alan	Turing	is	а		

Block Table





Logical KV blocks

Artificial	Intelligence	is	the
future	of	technology	

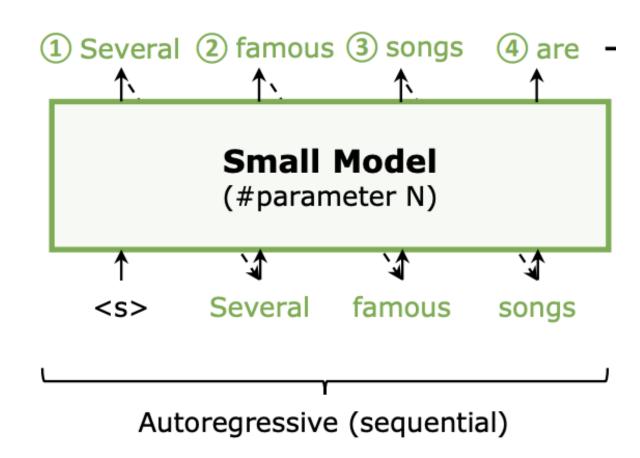
Efficient Memory Management for Large Language Model Serving with PagedAttention (Kwon et al., 2023)

Examples of Modern Inference Systems

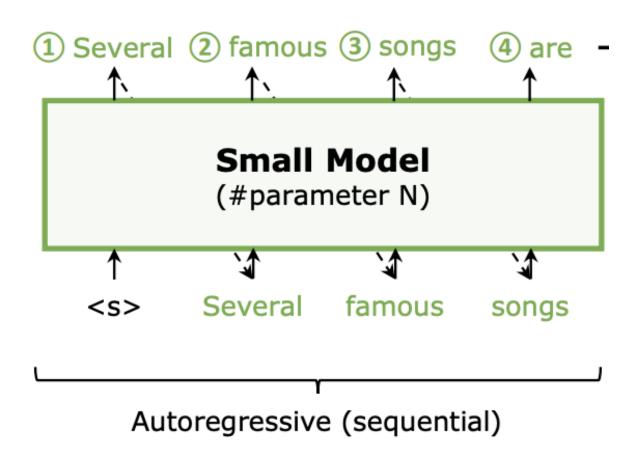
vLLM, sglang, and TensorRT-LLM are all examples of inference engines which take advantage of this approach. If you are serving models in any setting, you should use one of the above (rather than directly instantiating some huggingface or pytorch model)



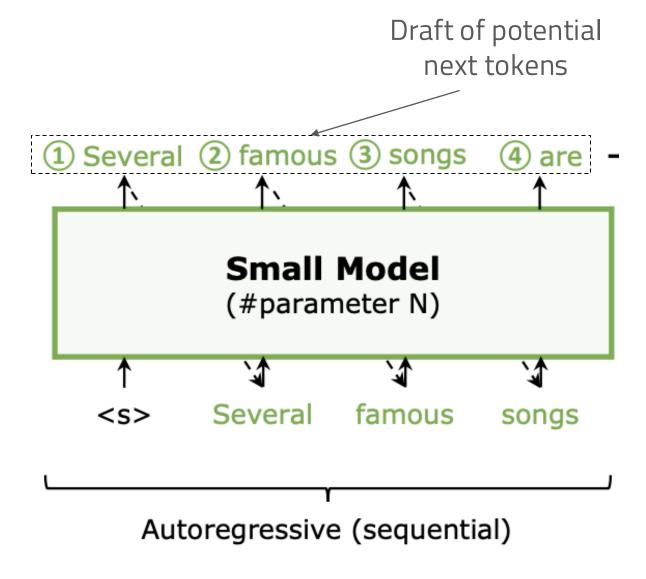
- In standard autoregressive decoding, we are only using each parameter one time when the batch size is 1
- This means standard decoding has a *low* arithmetic intensity and is memory bound
- We have a bunch more compute we could be getting for free given how massively parallel GPUs are



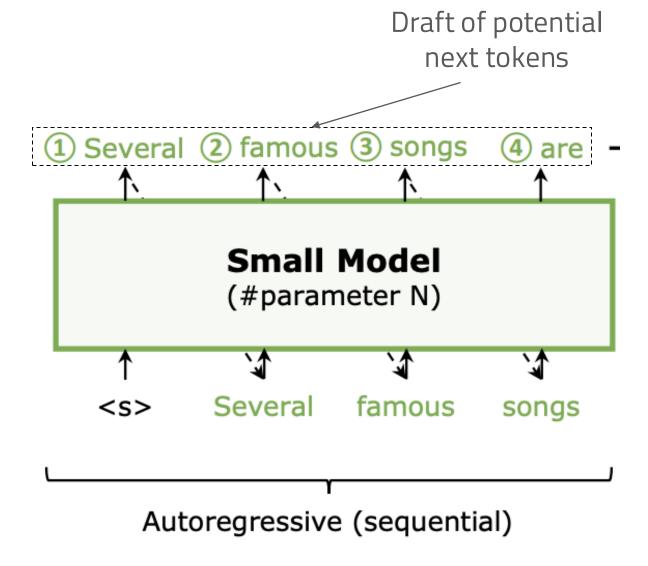
 Speculative decoding resolves this by 'speculating' multiple tokens into the future with a smaller, cheaper model



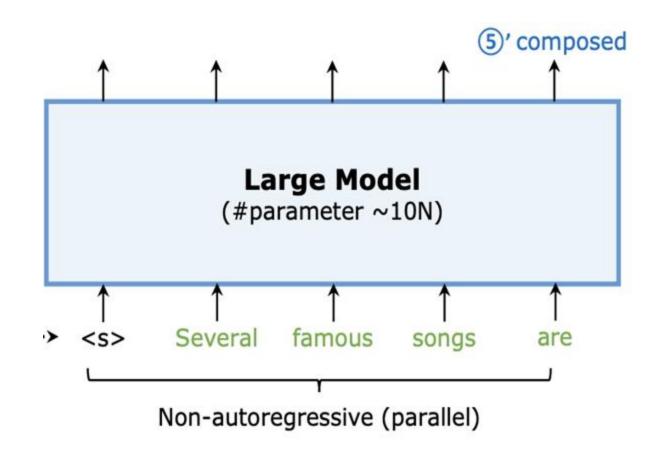
 Speculative decoding resolves this by 'speculating' multiple tokens into the future with a smaller, cheaper model



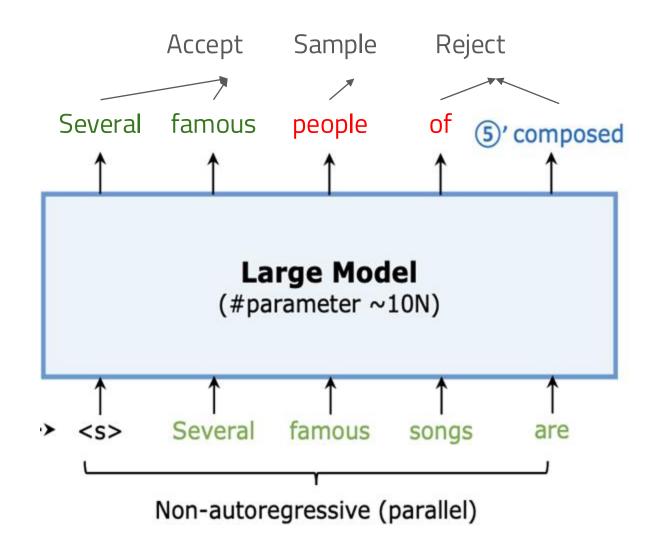
- Speculative decoding resolves this by 'speculating' multiple tokens into the future with a smaller, cheaper model
- We can now send this set of tokens on to a much larger model to verify the sequence



- 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

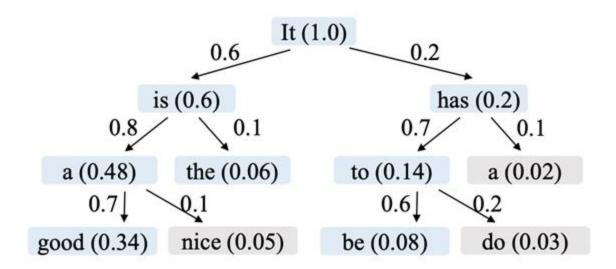


- Because the sequence will be run in parallel the arithmetic intensity will be proportional to the number of draft tokens
- We run each token through and see if the output of the large model matches that of the smaller, draft model
- We accept the matching tokens



Advanced Approaches

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



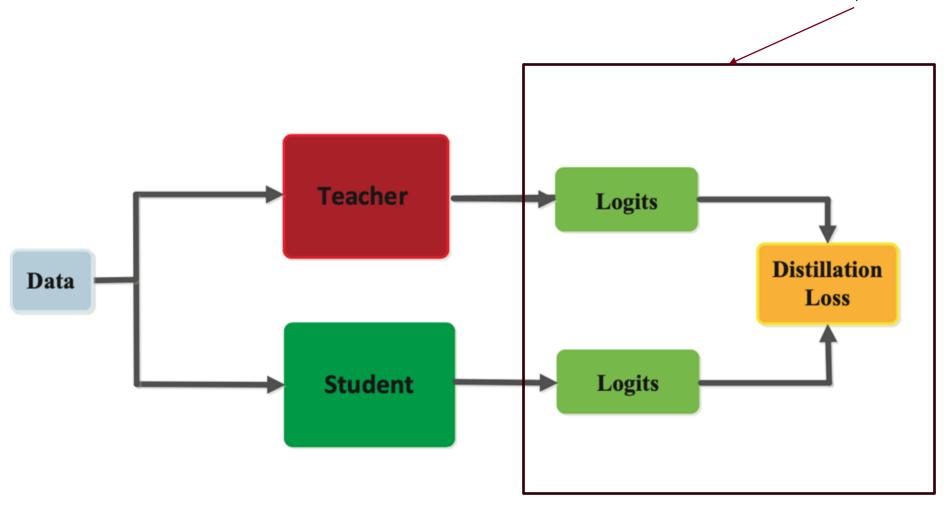


Efficient LLMs

- Quantization
- Sparsity
- ☐ Long Context
- ☐ Serving & Systems
- Distillation
- ☐ Parameter Efficient Fine-Tuning
- Alternative Paradigms

Logit Distillation

Don't use labels, use the modeling distribution of a better 'Teacher' model to train a smaller, 'Student' model



Logit Distillation

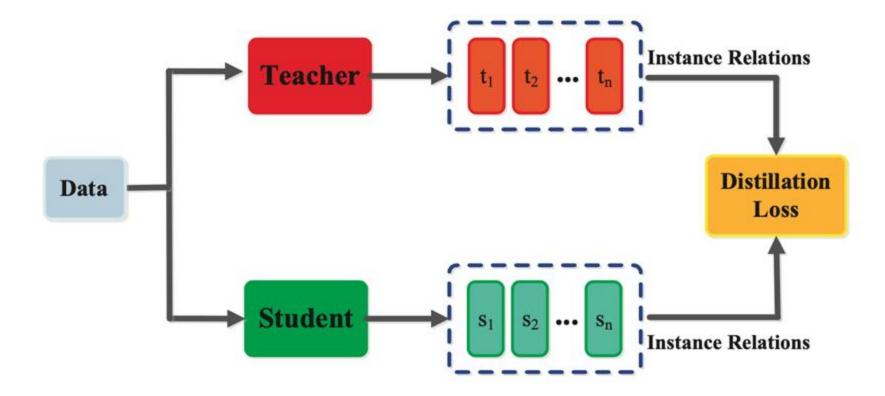




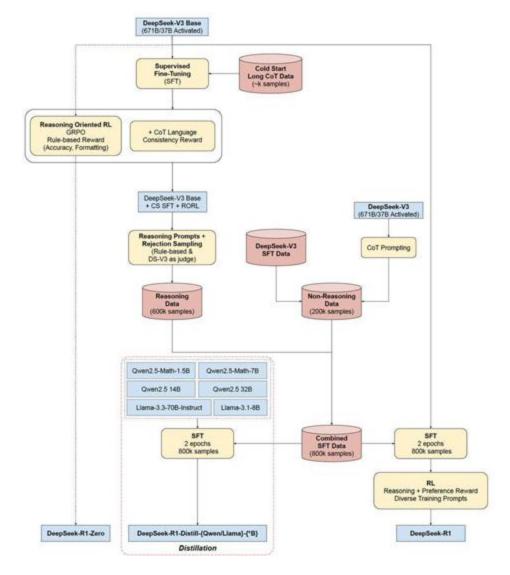
Both Gemma3 & Llama4 use logit distillation during pretraining

Layer Wise Distillation (LWD)

This can be expanded to trying to match the hidden states of the student and teacher model as well

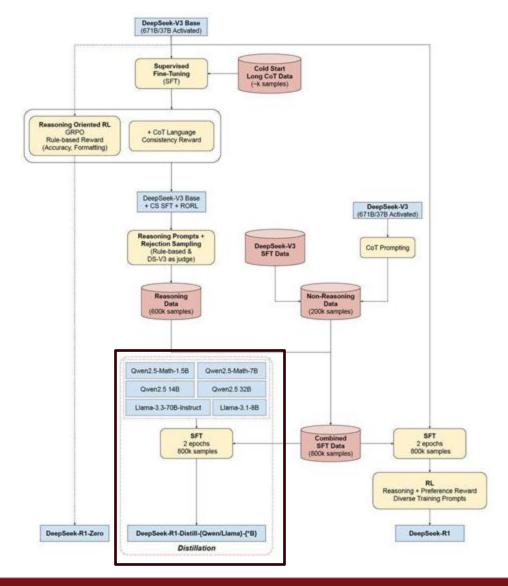


Output Tokens as Distillation (Synthetic Data)



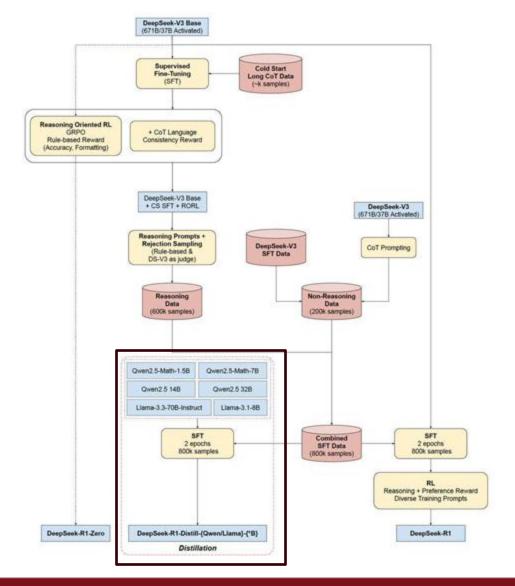
DeepSeekR1 outputs tokens for a given set of problems. In this case 'distillation', just means performing Supervised fine tuning on the tokens create by the larger V3 Models

Output Tokens as Distillation (Synthetic Data)



DeepSeekR1 outputs tokens for a given set of problems. In this case 'distillation', just means performing Supervised fine tuning on the tokens create by the larger V3 Models

Output Tokens as Distillation (Synthetic Data)



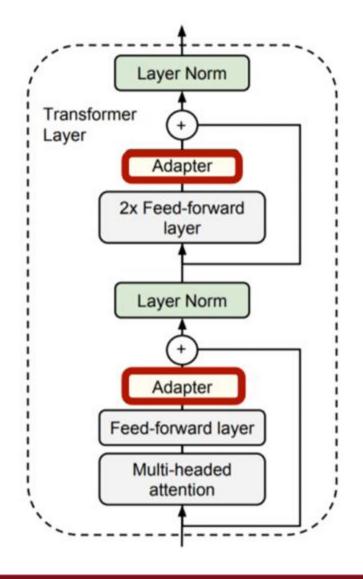
DeepSeekR1 outputs tokens for a given set of problems. In this case 'distillation', just means performing Supervised fine tuning on the tokens create by the larger V3 Models

Most small-model training today employs similar methods as these

Efficient LLMs

- Quantization
- Sparsity
- ☐ Long Context
- ☐ Serving & Systems
- Distillation
- □ Parameter Efficient Fine-Tuning
- Alternative Paradigms

Adapter



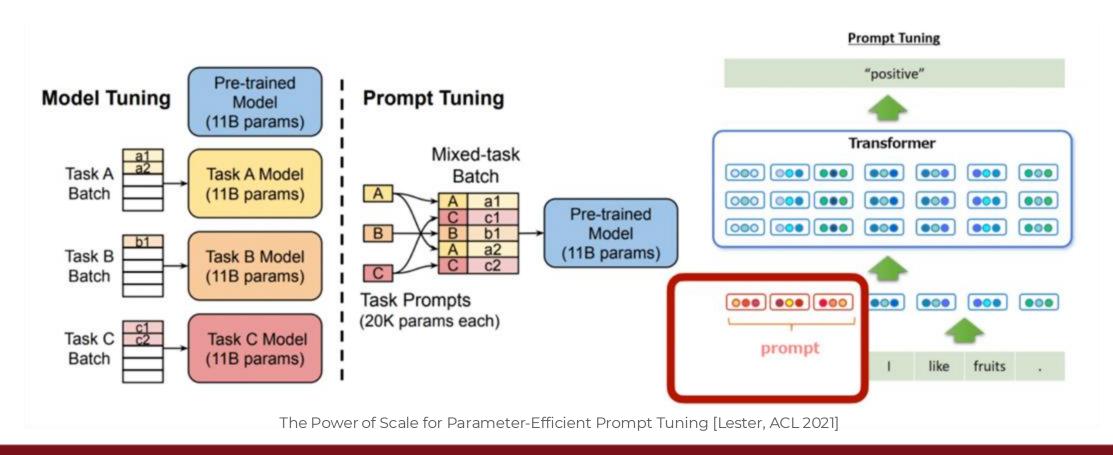
Add trainable layers after each feedforward layer

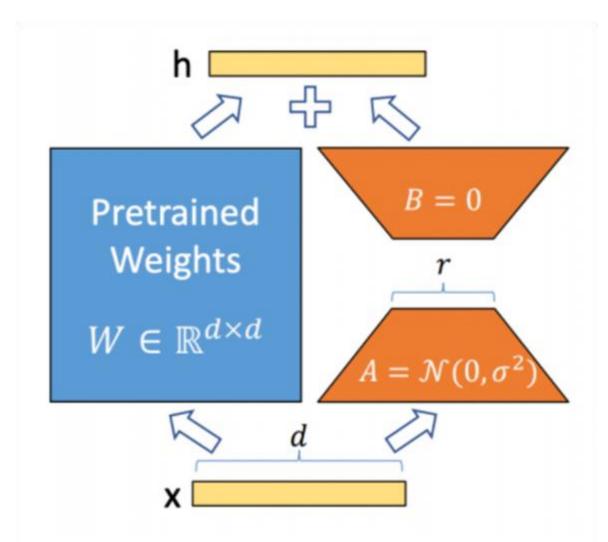
	Total num params	Trained params / task	CoLA	SST	MRPC	STS-B	QQP	$MNLI_{m}$	$MNLI_{mm}$	QNLI	RTE	Total
BERTLARGE	9.0×	100%	60.5	94.9	89.3	87.6	72.1	86.7	85.9	91.1	70.1	80.4
Adapters (8-256)	1.3×	3.6%	59.5	94.0	89.5	86.9	71.8	84.9	85.1	90.7	71.5	80.0
Adapters (64)	1.2×	2.1%	56.9	94.2	89.6	87.3	71.8	85.3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

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

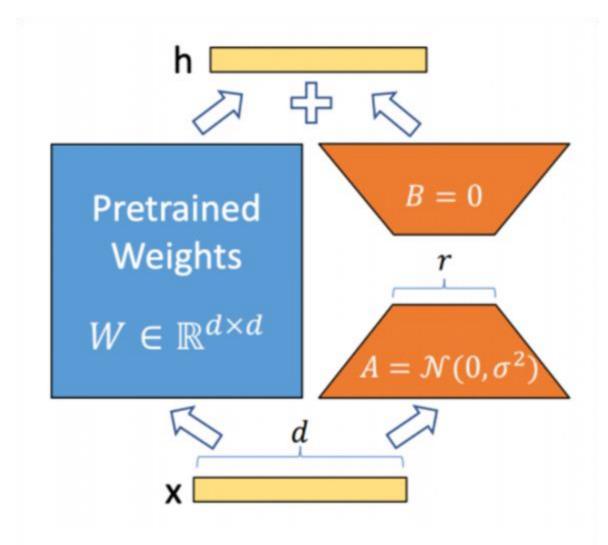
Prompt Tuning (Soft Prompting)

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

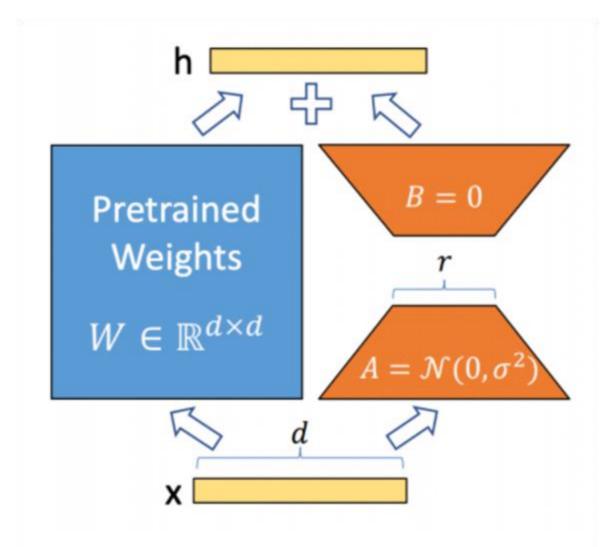




- ☐ Hypothesizes that fine-tuning results in only low rank updates
- ☐ Thus, we may approximate the updates themselves as low-rank and train on this low-rank approximation directly



Normal Finetuning:



Normal Finetuning:

LoRA Finetuning:

$$h = Wx + BAx$$

Update B,A Leave W unchanged

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±.0	$94.2_{\pm .1}$	$88.5_{\pm 1.1}$	$60.8 \scriptstyle{\pm.4}$	$93.1_{\pm .1}$	$90.2_{\pm .0}$	$71.5_{\pm 2.7}$	$89.7_{\pm .3}$	84.4
$RoB_{base} (Adpt^{D})*$	0.9M	87.3±.1	94.7±.3	$88.4_{\pm .1}$	62.6±.9	93.0±.2	$90.6_{\pm .0}$	$75.9_{\pm 2.2}$	$90.3_{\pm .1}$	85.4
RoB _{base} (LoRA)	0.3M	87.5 _{±.3}	$\textbf{95.1}_{\pm .2}$	$89.7_{\pm.7}$	$63.4_{\pm 1.2}$	$\textbf{93.3}_{\pm.3}$	$90.8 \scriptstyle{\pm .1}$	$\textbf{86.6}_{\pm.7}$	$\textbf{91.5}_{\pm .2}$	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±.2	$96.2 \scriptstyle{\pm .5}$	$\textbf{90.9}_{\pm 1.2}$	$\textbf{68.2}\scriptstyle{\pm 1.9}$	$\textbf{94.9}_{\pm.3}$	$91.6 \scriptstyle{\pm .1}$	$\textbf{87.4} \scriptstyle{\pm 2.5}$	$\textbf{92.6}_{\pm .2}$	89.0
RoB _{large} (Adpt ^P)†	3.0M	90.2±.3	96.1±.3	90.2±.7	68.3 _{±1.0}	94.8±.2	91.9 _{±.1}	83.8 _{±2.9}	92.1 _{±.7}	88.4
RoB _{large} (Adpt ^P)†	0.8M	90.5±.3	$96.6_{\pm .2}$	$89.7_{\pm 1.2}$	$67.8_{\pm 2.5}$	94.8±.3	$91.7_{\pm .2}$	$80.1_{\pm 2.9}$	$91.9_{\pm .4}$	87.9
RoB _{large} (Adpt ^H)†	6.0M	89.9±.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±.3	96.3±.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±.2	$96.2 \scriptstyle{\pm .5}$	$\textbf{90.2}_{\pm 1.0}$	$68.2{\scriptstyle\pm1.9}$	$\textbf{94.8}_{\pm.3}$	$91.6 \scriptstyle{\pm .2}$	$\textbf{85.2}{\scriptstyle\pm1.1}$	$\textbf{92.3}_{\pm.5}$	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	$91.9_{\pm .2}$	$96.9_{\pm.2}$	$\textbf{92.6}_{\pm.6}$	$72.4_{\pm 1.1}$	$\textbf{96.0}_{\pm.1}$	$\textbf{92.9}_{\pm.1}$	$94.9_{\pm .4}$	$93.0_{\pm .2}$	91.3

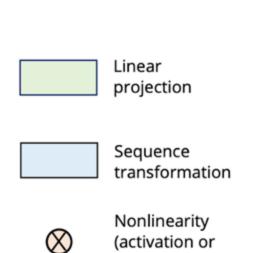
Table 2: RoBERTa_{base}, RoBERTa_{large}, and DeBERTa_{XXL} with different adaptation methods on the GLUE benchmark. We report the overall (matched and mismatched) accuracy for MNLI, Matthew's correlation for CoLA, Pearson correlation for STS-B, and accuracy for other tasks. Higher is better for all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

Efficient LLMs

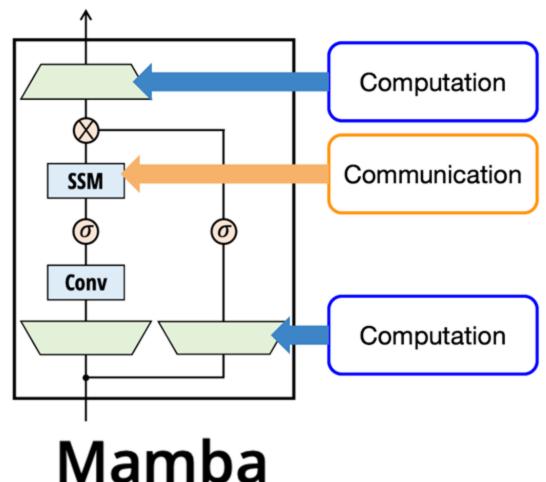
- Quantization
- Sparsity
- ☐ Long Context
- ☐ Serving & Systems
- Distillation
- ☐ Parameter Efficient Fine-Tuning
- Alternative Paradigms

State Space Model (SSM)

Linearity Strikes back



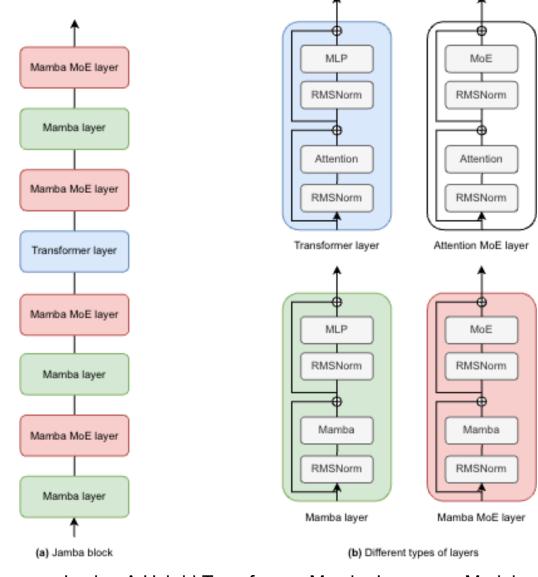
multiplication)



Mamba: Linear-Time Sequence Modeling with Selective State Spaces (Gu et. Al)

SSM

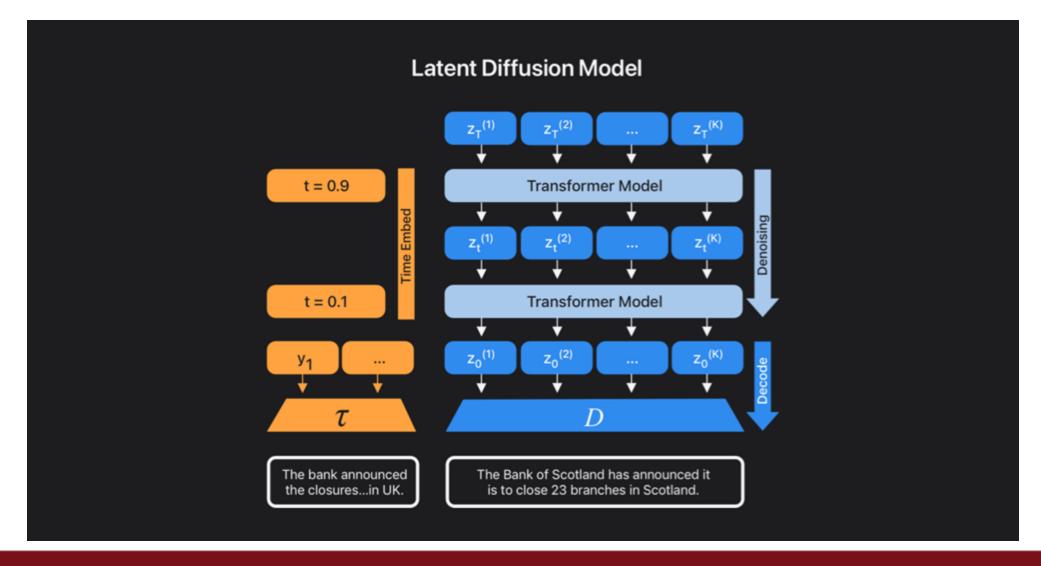
- SSMs are growing in popularity
- Linear complexity makes them excellent candidates for very long tasks
- If they can begin outperforming transformers in practical settings, these could displace the transformer architecture



Jamba: A Hybrid Transformer-Mamba Language Model

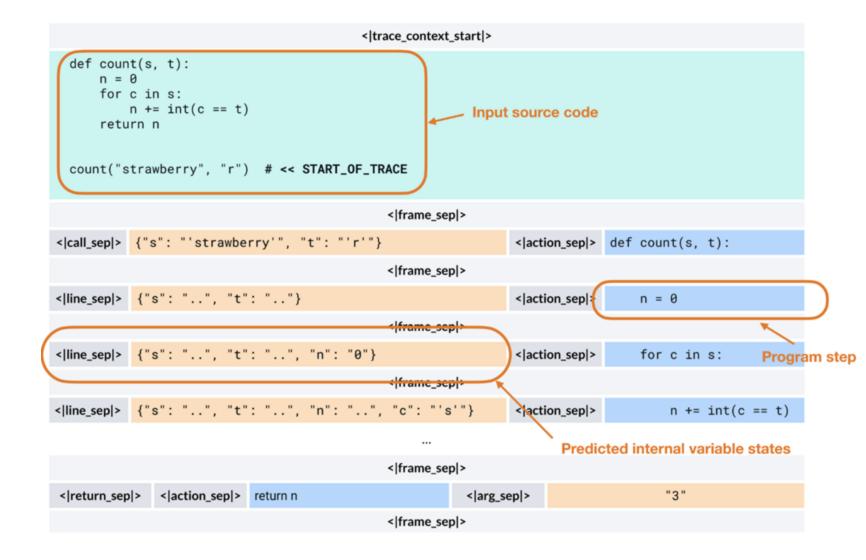
Text Diffusion

Predict multiple tokens as the same time using a diffusion model



World Models

Use some dense environment signals to predict what the consequences of each token prediction is



Recursive Transformers

- ☐ Use transformers recursively to iteratively refine model outputs
- ☐ This approach is still mostly useful in toy settings at the moment

