



EECS 230 Deep Learning

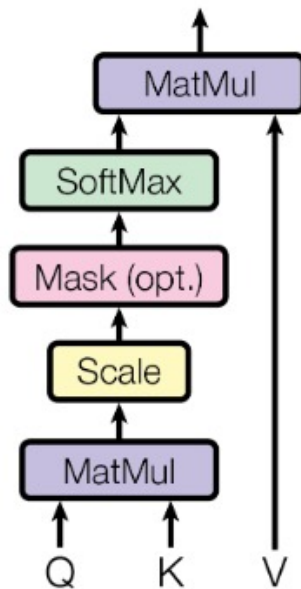
Lecture 13: Training LLM

Recap: Attention Operation

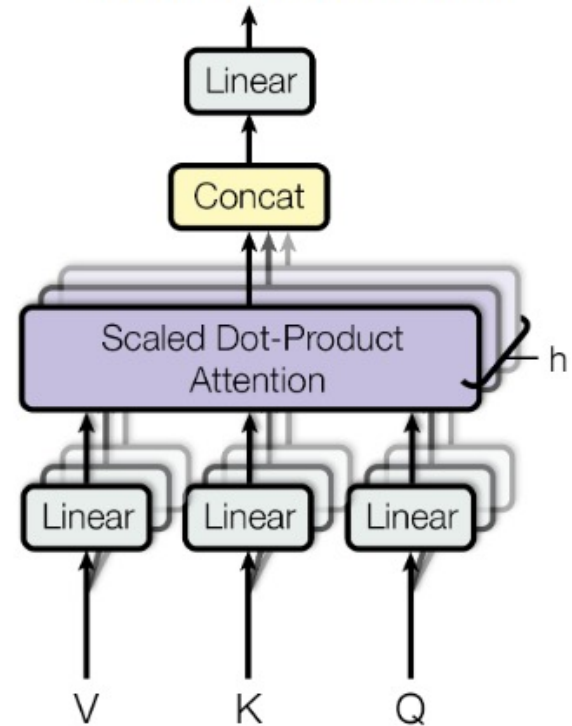
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Recap: Multi-head attention

Scaled Dot-Product Attention

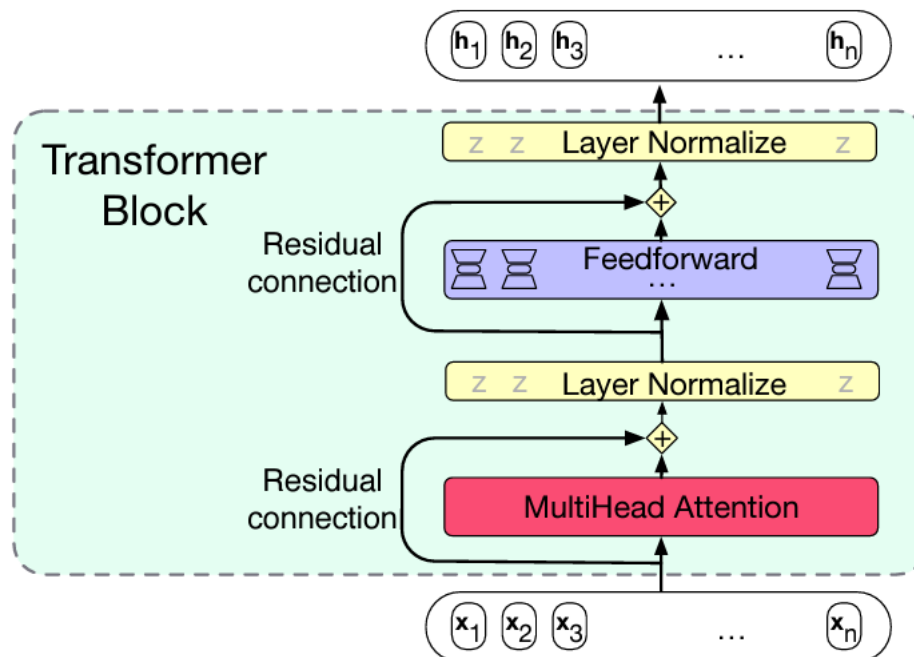


Multi-Head Attention

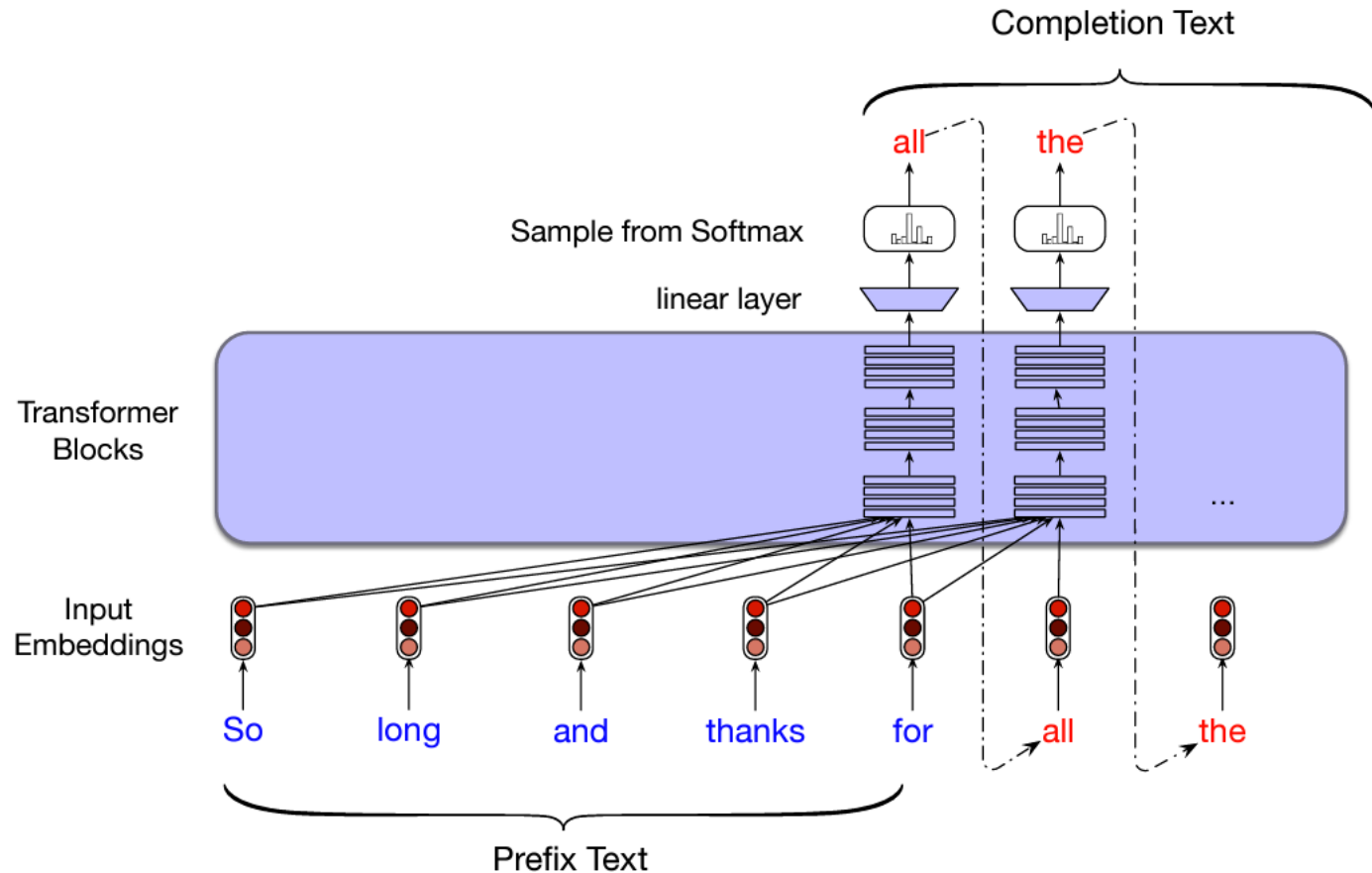


Recap: From Attention to Transformer Block

- ❑ A transformer block has
 - ❑ **Self Attention**: information exchange between tokens
 - ❑ **Feed forward network**: Information transform within tokens
 - ❑ E.g. a multi-layer perceptron with 1 hidden layer
 - ❑ **Normalization** (Layer normalization)
 - ❑ **Residual connection**



Recap:Transformer-based Large Language Model

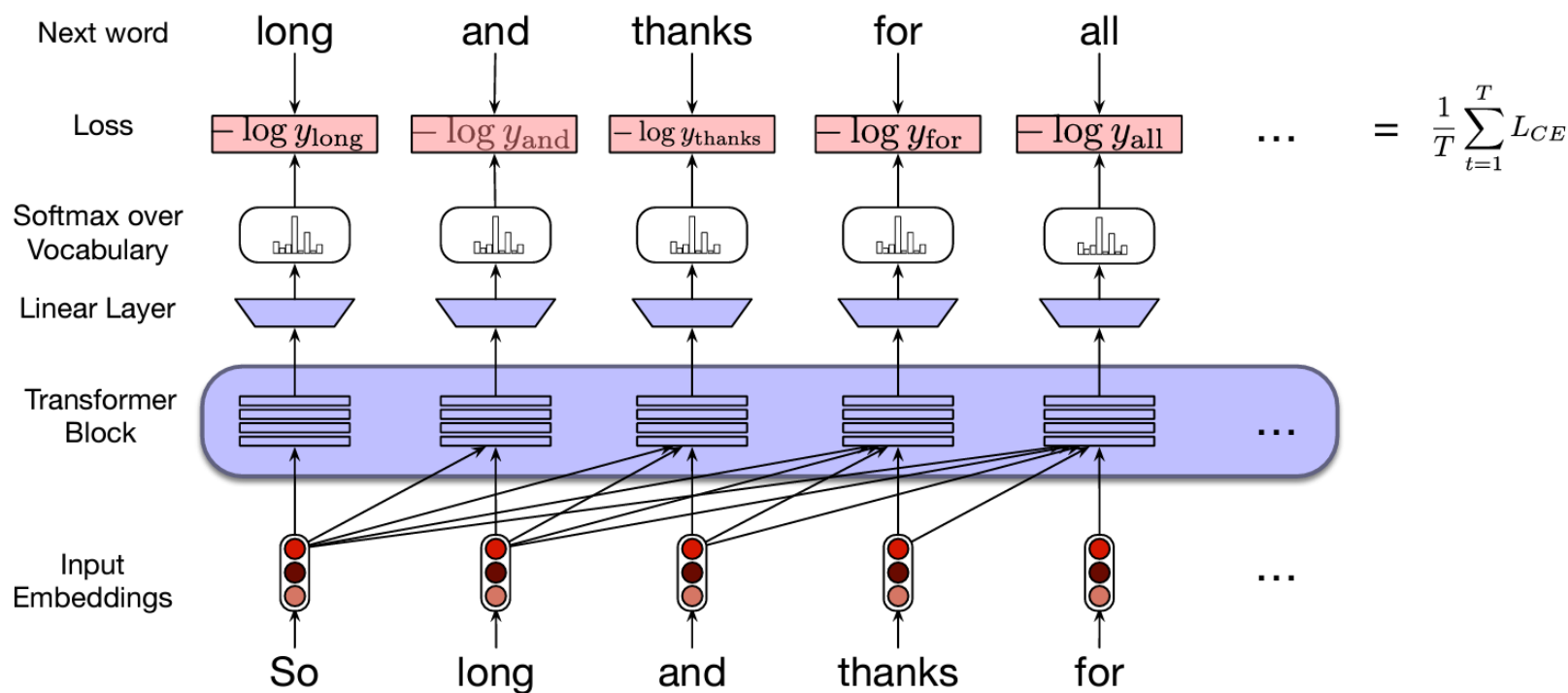


Outline

- ❑ Training LLM
- ❑ LLM beyond language modeling
- ❑ Fine-tuning LLM

Self-supervised training algorithm

- Take a corpus of text as training material
- Minimize cross-entropy loss for next-word prediction



Training corpora for large language models

❑ commoncrawl.org

- ❑ Over 250 billion pages spanning 17 years.
- ❑ Free and open corpus since 2007.
- ❑ Cited in over 10,000 research papers.
- ❑ 3–5 billion new pages added each month.

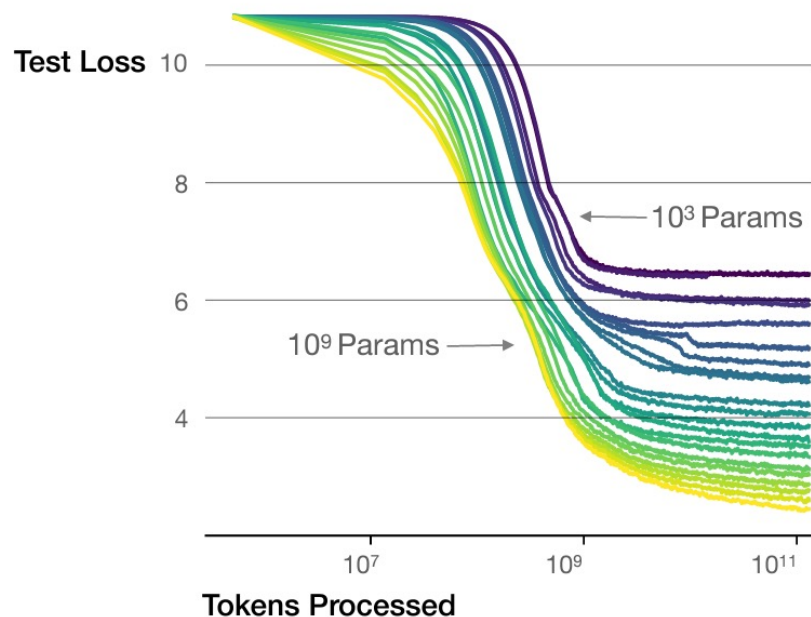
❑ Training data for GPT-3

- ❑ web (429 billion tokens)
- ❑ some text from books (67 billion tokens)
- ❑ Wikipedia (3 billion tokens).

Scaling laws

- Cross-entropy loss scales as a power-law with
 - model size
 - dataset size
 - amount of compute

Larger models require **fewer samples** to reach the same performance



$$L(N) = \left(\frac{N_c}{N} \right)^{\alpha_N}$$
$$L(D) = \left(\frac{D_c}{D} \right)^{\alpha_D}$$
$$L(C) = \left(\frac{C_c}{C} \right)^{\alpha_C}$$

Conditional generation using LLM

- ❑ Many practical NLP tasks can be cast as word prediction
- ❑ For example, sentiment classification
 - ❑ The sentiment of the sentence “I like Jackie Chan” is:
- ❑ Question answering
 - ❑ Q: Who wrote the book “The Origin of Species”? A:
- ❑ Summarization
 - ❑ Append token tl;dr

How large are LLMs?

Model	Creators	Year of release	Training Data (# tokens)	Model Size (# parameters)
GPT-2	OpenAI	2019	~10 billion (40Gb)	1.5 billion
GPT-3 (cf. ChatGPT)	OpenAI	2020	300 billion	175 billion
PaLM	Google	2022	780 billion	540 billion
Chinchilla	DeepMind	2022	1.4 trillion	70 billion
LaMDA (cf. Bard)	Google	2022	1.56 trillion	137 billion
LLaMA	Meta	2023	1.4 trillion	65 billion
LLaMA-2	Meta	2023	2 trillion	70 billion
GPT-4	OpenAI	2023	?	?



Fine-tuning LLM

Few-shot Learning with LLMs

❑ Suppose you have...

❑ a dataset $D = \{(x_i, y_i)\}, i=1\dots N$, and N is rather small (i.e. few-shot setting)

❑ a very large (billions of parameters) pre-trained language model

❑ There are two ways to “learn”

Option A: Supervised fine-tuning

❑ Definition: fine-tune the LLM on the training data using...

- ❑ a standard supervised objective
- ❑ backpropagation to compute gradients
- ❑ your favorite optimizer (e.g. Adam)

❑ Pros:

- ❑ fits into the standard ML recipe
- ❑ still works if N is large

❑ Cons:

- ❑ backpropagation requires $\sim 3x$ the memory and computation time as the forward computation
- ❑ you might not have access to the model weights at all

Option B: In-context learning

❑ Definition:

- ❑ 1. feed training examples to the LLM as a prompt
- ❑ 2. allow the LLM to infer patterns in the training examples during inference (i.e. decoding)
- ❑ 3. take the output of the LLM following the prompt as its prediction

❑ Cons:

- ❑ the prompt may be very long and Transformer LMs require $O(N^2)$ time/space where N = length of context

❑ Pros:

- ❑ no backpropagation required and only one pass through the training data
- ❑ does not require model weights, only API access

Fine-Tuning vs. In-Context Learning

- ❑ Why would we ever bother with fine-tuning if it's so inefficient?
- ❑ Even for very large LMs, fine-tuning often beats in-context learning

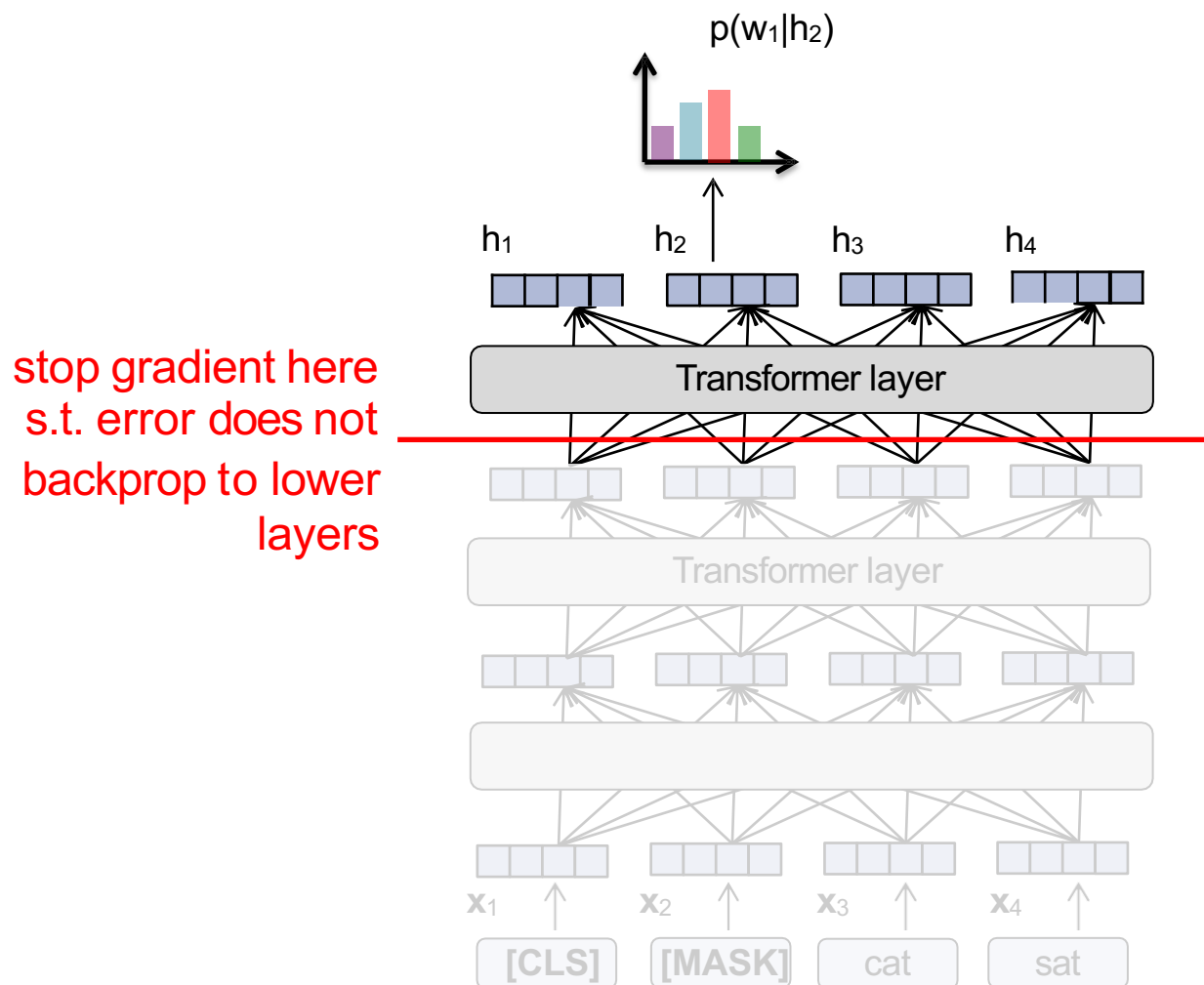
Method	MNLI-m (Val. Acc./%)	RTE (Val. Acc./%)
GPT-3 Few-Shot	40.6	69.0
GPT-3 Fine-Tuned	89.5	85.4

Parameter Efficient Fine-Tuning

- ❑ Goal: perform fine-tuning of fewer parameters, but achieve performance on a downstream task that is comparable to finetuning of all parameters
- ❑ Various approaches:
 - ❑ Subset: Pick a subset of the parameters and fine-tune only those
 - ❑ Adapters: add additional layers that have few parameters and tune only the parameters of those layers, keeping all others fixed
 - ❑ LoRA: learn a small delta for each of the parameter matrices with the delta chosen to be low rank

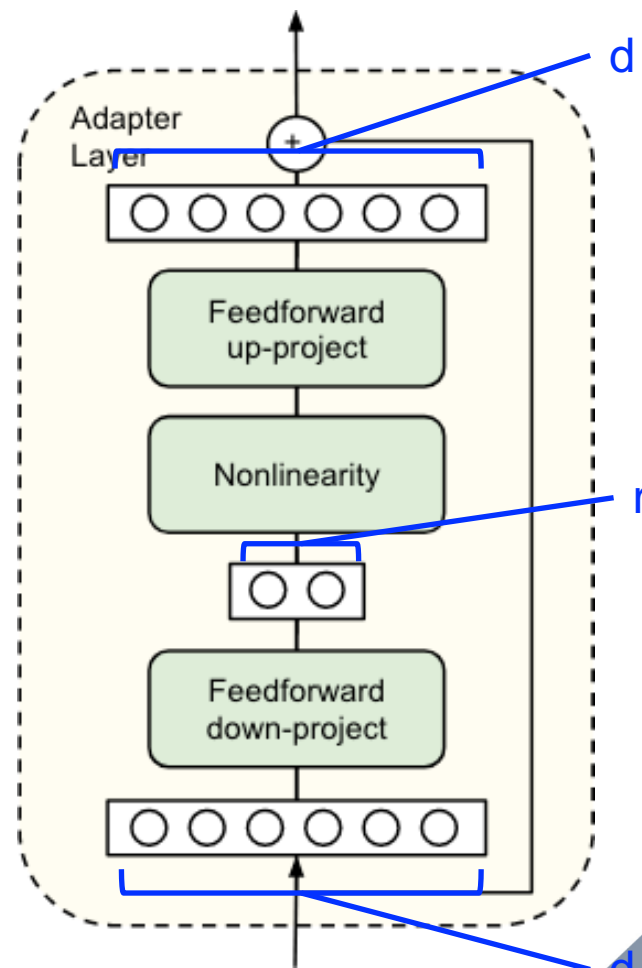
Fine-Tuning the Top Layers Only

- keep all parameters fixed except for the top K layers



Adapter Module

- ❑ An adapter layer is simply a feed forward neural network with one hidden layer, and a residual connection
- ❑ For input dimension, d , the adapter layer also has output dimension d , but bottlenecks to a lower dimension r in the middle



Adapter for transformer

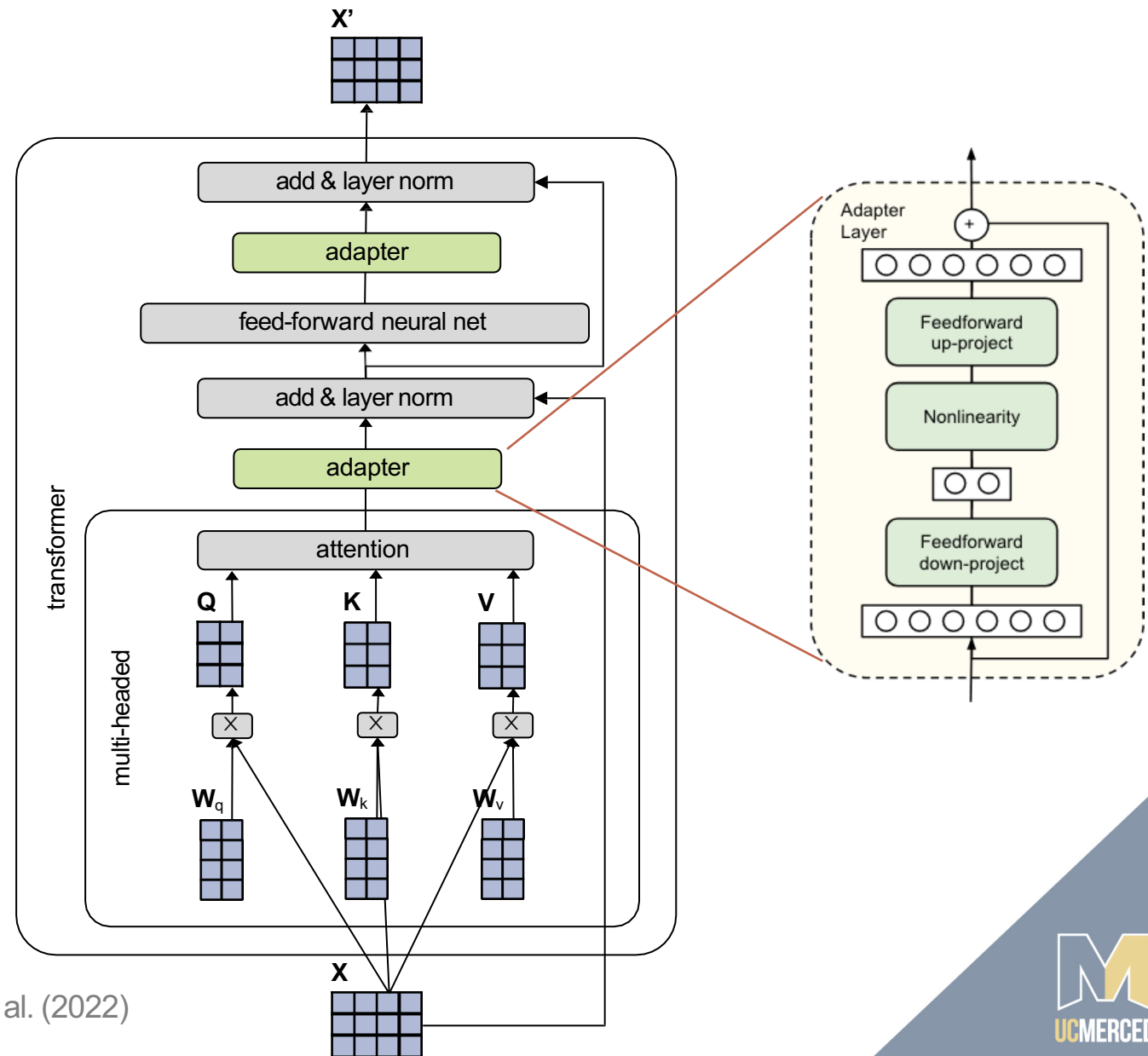
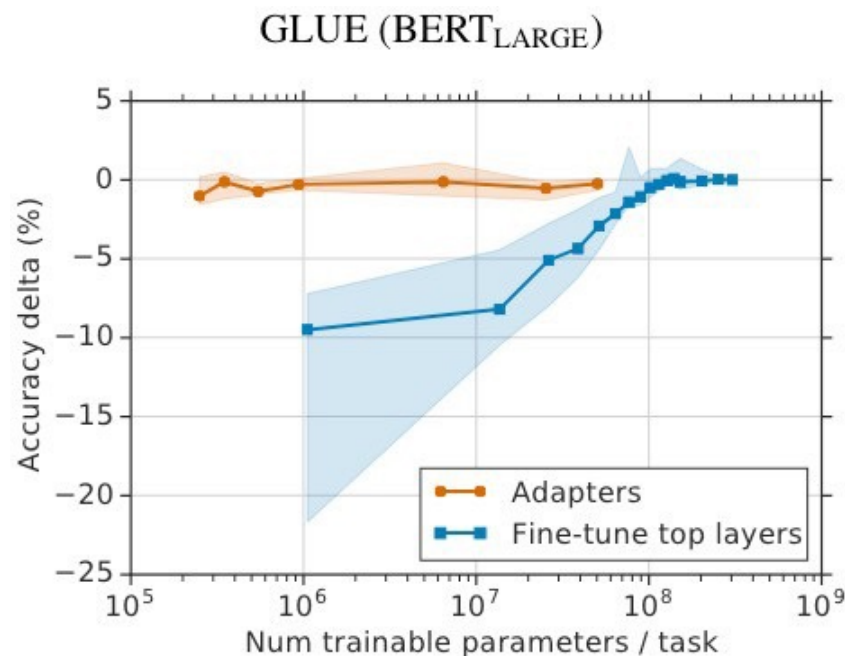


Figure inspired by He et al. (2022)

Adapter Results

- **Pretrained Model:** BERT_{LARGE}
- **Baseline Method:** fine-tune only the top K layers of BERT_{LARGE}
- Adapters achieve nearly the performance (i.e. 0% delta) of full fine-tuning but with substantially fewer parameters
- Sometimes adapters even outperform full fine-tuning



Low-Rank Adaptation (LORA)

- ❑ Adapters and related methods introduce inference latency at test time that is non-trivial

Batch Size	32	16	1
Sequence Length	512	256	128
$ \Theta $	0.5M	11M	11M
Fine-Tune/LoRA	1449.4 \pm 0.8	338.0 \pm 0.6	19.8 \pm 2.7
Adapter ^L	1482.0 \pm 1.0 (+2.2%)	354.8 \pm 0.5 (+5.0%)	23.9 \pm 2.1 (+20.7%)
Adapter ^H	1492.2 \pm 1.0 (+3.0%)	366.3 \pm 0.5 (+8.4%)	25.8 \pm 2.2 (+30.3%)

Table 1: Inference latency of a single forward pass in GPT-2 medium measured in milliseconds, averaged over 100 trials. We use an NVIDIA Quadro RTX8000. “ $|\Theta|$ ” denotes the number of trainable parameters in adapter layers. Adapter^L and Adapter^H are two variants of adapter tuning, which we describe in [Section 5.1](#). The inference latency introduced by adapter layers can be significant in an online, short-sequence-length scenario. See the full study in [Appendix E](#).

LoRA

Key Idea

- Keep the original pretrained parameters \mathbf{W}_0 fixed during fine-tuning
- Learn an additive modification to those parameters $\Delta\mathbf{W}$
- Define $\Delta\mathbf{W}$ via a low rank decomposition:

$$\Delta\mathbf{W} = \mathbf{B}\mathbf{A}$$

where $\mathbf{B}\mathbf{A}$ has rank r , which is **much less** than the input dimension k or the output dimension d

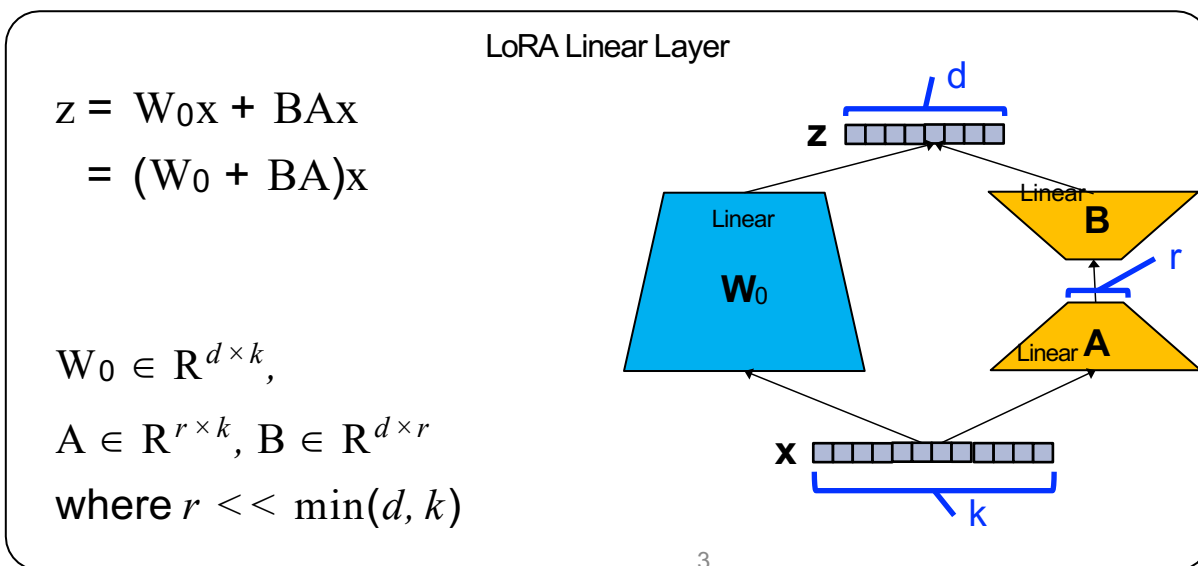
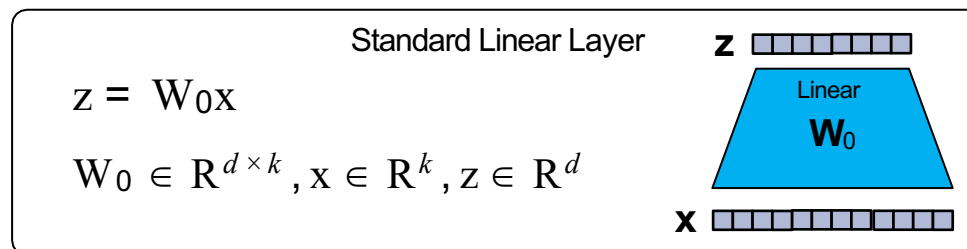


Figure inspired by He et al. (2022)

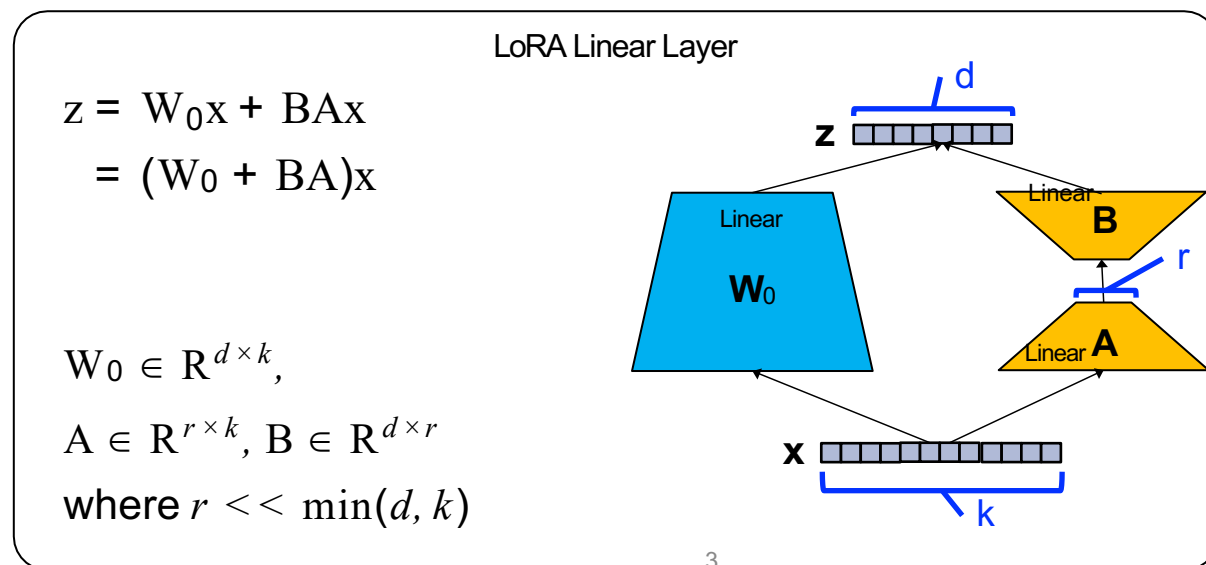
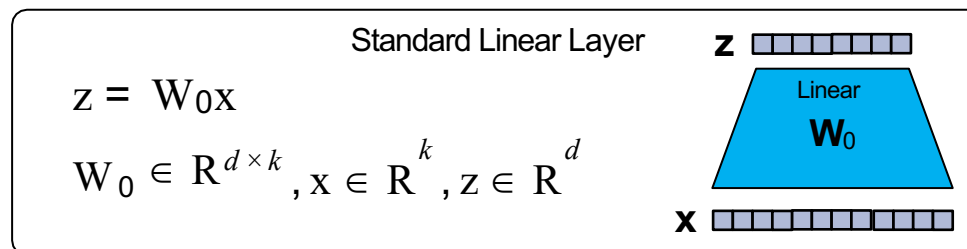
LoRA

Initialization

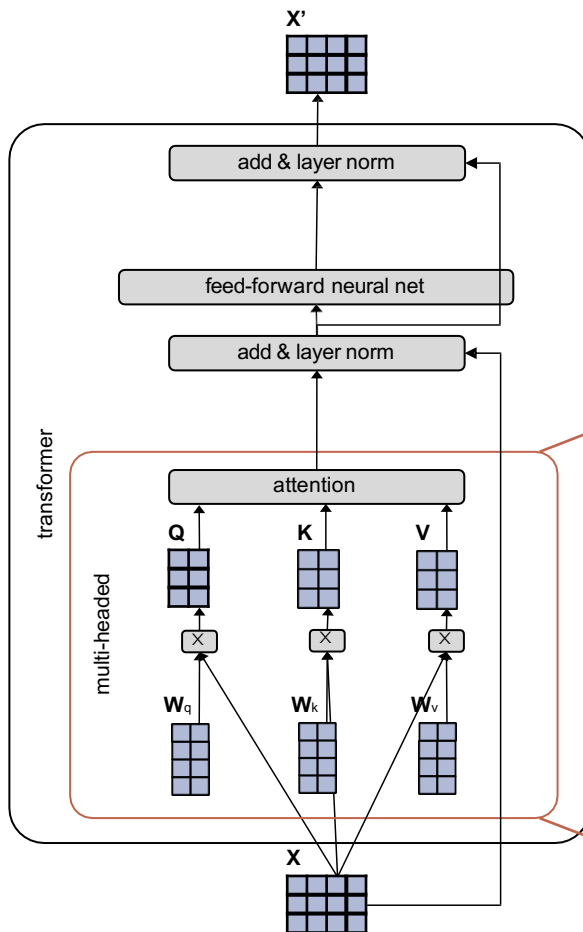
- We initialize the trainable parameters:
 $A_{ij} \sim N(0, \sigma^2), \forall i, j$
 $B = 0$
- This ensures that, at the start of fine tuning, the parameters have their pretrained values:

$$\Delta W = BA = 0$$

$$W_0 + BA = W_0$$



Transformer Layer



multi-headed attention

$$X^l = \text{concat}(X^{l(1)}, \dots, X^{l(h)})$$

$$X^{l(i)} = \text{softmax} \left(\frac{Q^{(i)} (K^{(i)})^T}{\sqrt{d_k}} \right) + M^l V^{(i)}$$

$$Q^{(i)} = X W_q^{(i)}$$

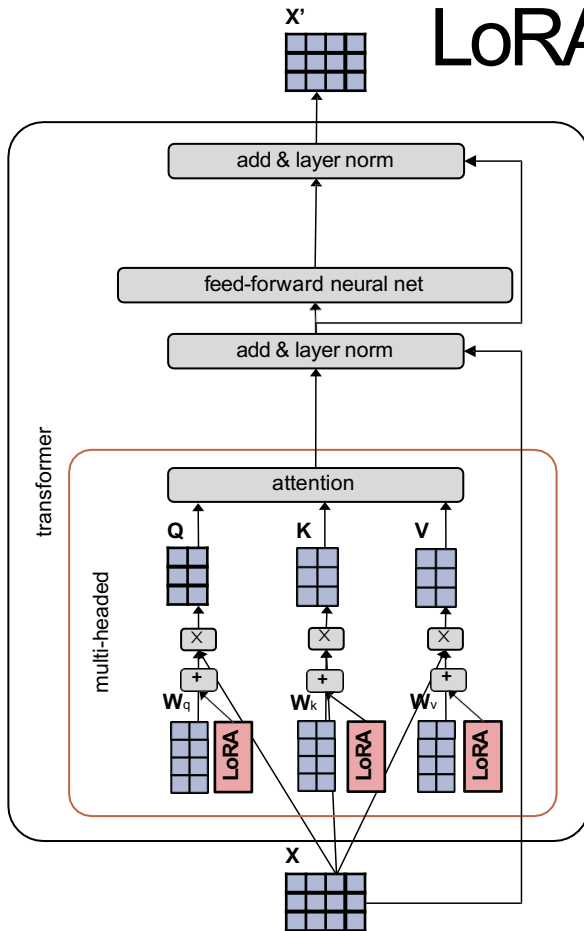
$$K^{(i)} = X W_k^{(i)}$$

$$V^{(i)} = X W_v^{(i)}$$

$$X = [x_1, \dots, x_N]^T$$

LoRA for Transformer

- LoRA linear layers could replace every linear layer in the Transformer layer
- But the original paper only applies LoRA to the attention weights



$$z = W_0 x + B A x$$

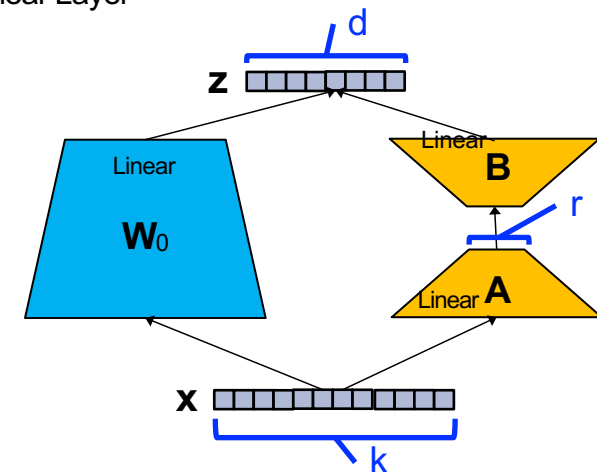
$$= (W_0 + B A) x$$

$$W_0 \in \mathbb{R}^{d \times k},$$

$$A \in \mathbb{R}^{r \times k}, B \in \mathbb{R}^{d \times r}$$

where $r \ll \min(d, k)$

LoRA Linear Layer



LoRA for Transformer

- LoRA linear layers could replace every linear layer in the Transformer layer
- But the original paper only applies LoRA to the attention weights

$$Q = \text{LoRALinear}(X; W_q, A_q, B_q)$$

$$K = \text{LoRALinear}(X; W_k, A_k, B_k)$$

$$V = \text{LoRALinear}(X; W_v, A_v, B_v)$$

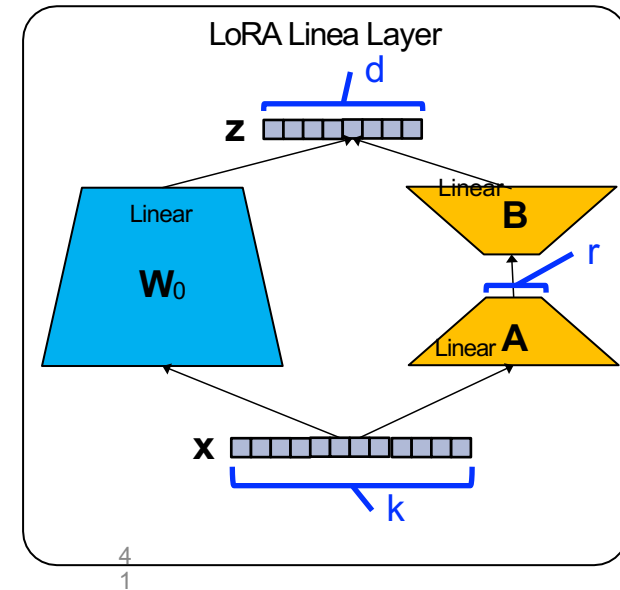
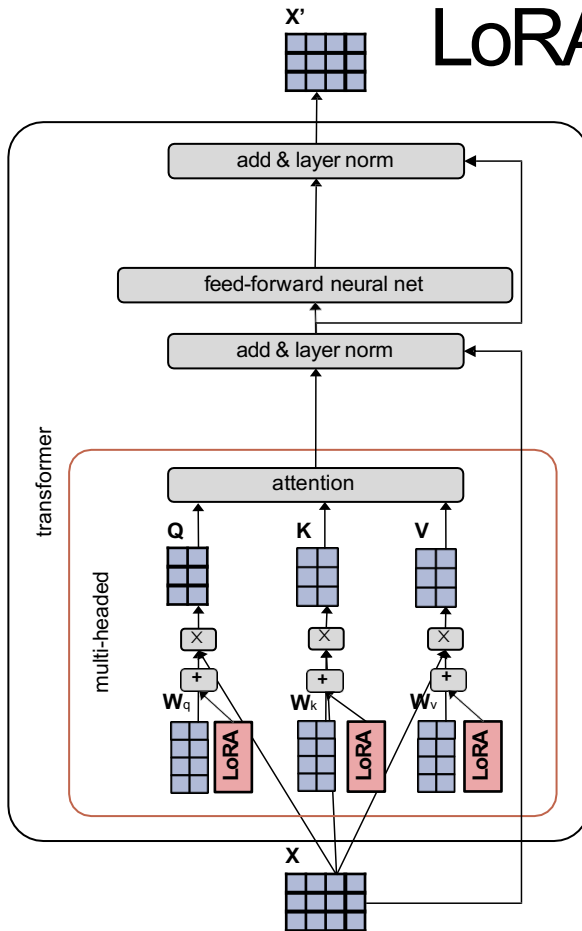
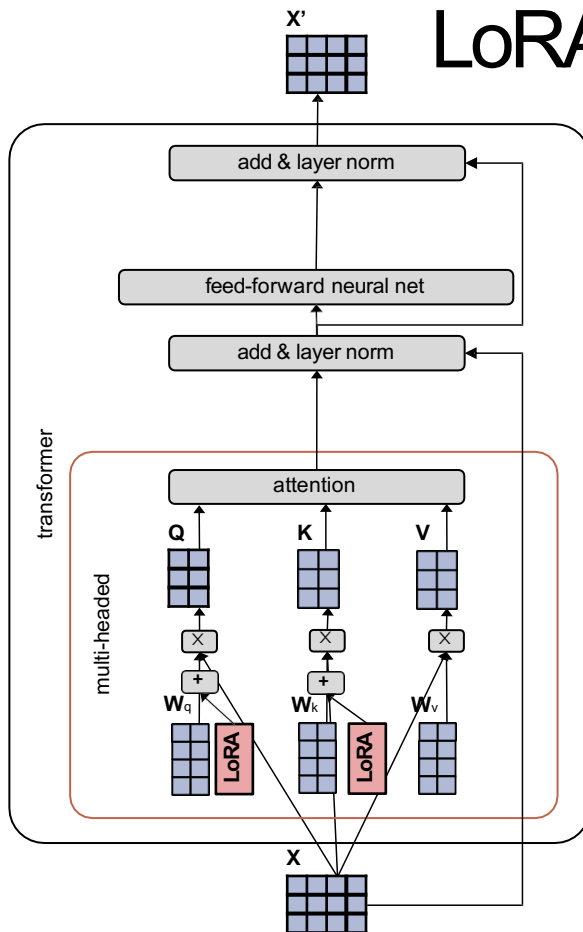


Figure inspired by He et al. (2022)

LoRA for Transformer



- LoRA linear layers could replace every linear layer in the Transformer layer
- But the original paper only applies LoRA to the attention weights
- Empirically, for GPT-3, they also find that it is most efficient to include LoRA **only on the query and key** linear layers:

Weight Type Rank r	# of Trainable Parameters = 18M						
	W_q 8	W_k 8	W_v 8	W_o 8	W_q, W_k 4	W_q, W_v 4	W_q, W_k, W_v, W_o 2
WikiSQL ($\pm 0.5\%$)	70.4	70.0	73.0	73.2	71.4	73.7	73.7
MultiNLI ($\pm 0.1\%$)	91.0	90.8	91.0	91.3	91.3	91.3	91.7

Table 5: Validation accuracy on WikiSQL and MultiNLI after applying LoRA to different types of attention weights in GPT-3, given the same number of trainable parameters. Adapting both W_q and W_v gives the best performance overall. We find the standard deviation across random seeds to be consistent for a given dataset, which we report in the first column.

Summary

- ❑ Training LLM

- ❑ Parameter-efficient Fine-tuning

 - ❑ Subset

 - ❑ Adaptor

 - ❑ LoRA