



# LoRA: Low-Rank Adaptation of Large Language Models

Edward Hu\*    Yelong Shen\*    Phillip Wallis    Zeyuan Allen-Zhu  
Yuanzhi Li    Shean Wang    Lu Wang    Weizhu Chen  
Microsoft Corporation

ICLR 2022

HUMANE Lab 최종현

2025.03.06 랩 세미나

# Background

---

- Language models are getting larger and larger
- Needs fine-tuning to adapt to different domain / purpose
- Fine-tuning LLMs are expensive in terms of memory and compute (e.g., GPT-3)
- Pre-trained models have a low intrinsic dimensionality – most of the parameters aren't needed for effective adaptation (Li et al., 2018)

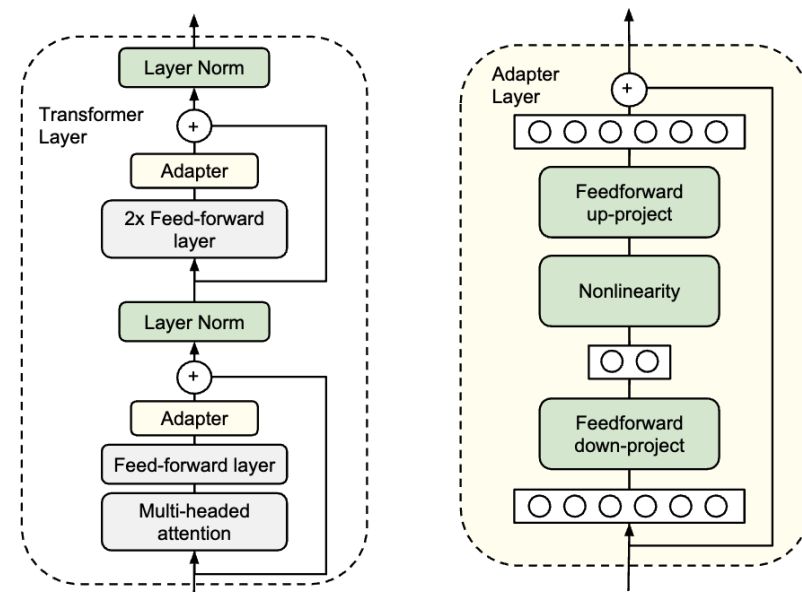
# Previous Methods

---

- Full Fine-Tuning (FT)
  - Updates all parameters of the pre-trained model
  - Best possible performance for a given task
  - No inference latency
  - Requires lots of VRAM and computation power
  - Storage overhead
  - Not scalable for multiple tasks – training and storing different fine-tuned models is infeasible

# Previous Methods

- Adapter Layers
  - Small trainable modules inserted into each Transformer layer
  - Original model weights are frozen – only adapter parameters are updated
  - Bottleneck structure
    - Down-projection: reduces dimensionality
    - Non-linearity: ReLU, GeLU
    - Up-projection: restores original size
  - Storage cost per task
  - Inference latency – each adapter introduces extra computations



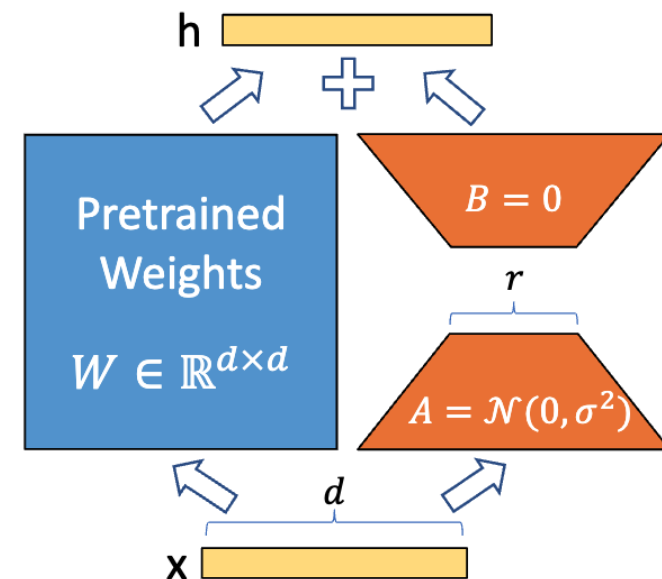
# Previous Methods

---

- BitFit (Bias-Only Fine-Tuning)
  - Only update bias terms in the Transformer layers are updated
  - Lightweight alternative to full fine-tuning
  - Weaker performance
- Prefix-tuning
  - Learn a set of trainable “prefix” tokens
  - Keeps the original model unchanged
  - Hard to optimize

# What is LoRA?

- Low-Rank Adaptation
- Freezes the pre-trained model
- Injects trainable low-rank matrices (A and B) into Transformer layers
- Without LoRA:  $h = W_0x$
- With LoRA:  $h = W_0x + \Delta Wx = W_0x + BAx$



# How LoRA works

---

- Pre-trained model's weight:  $W_0 \in \mathbb{R}^{d \times k}$
- Accumulated gradient update (Core idea)
  - $\Delta W = BA$
  - $B \in \mathbb{R}^{d \times r}$  (initialized with 0)
  - $A \in \mathbb{R}^{r \times k}$  (initialized with random Gaussian)
  - $r \ll \min(d, k)$  –  $r$  is decided with experiments
- $h = W_0x + \Delta Wx = W_0x + BAx$
- Only update  $\Delta W$ , thus only needing to learn  $d \times r + r \times k$  instead of  $d \times k$

# How LoRA works

---

- Scaling factor -  $\frac{\alpha}{r}$ 
  - $r$  – same as previous  $r$  (used in matrices  $A$  and  $B$ )
  - $\alpha$  – same as  $r$  above (but can be tuned separately from  $r$ )
- $h = W_0x + \frac{\alpha}{r}\Delta Wx = W_0x + \frac{\alpha}{r}BAx$



# How LoRA works

---

- $h = W_0x + \frac{\alpha}{r}\Delta Wx = W_0x + \frac{\alpha}{r}BAx$
- $h = W_qx + \frac{\alpha}{r}\Delta Wx = W_qx + \frac{\alpha}{r}BAx$
- $h = W_kx + \frac{\alpha}{r}\Delta Wx = W_kx + \frac{\alpha}{r}BAx$
- $h = W_vx + \frac{\alpha}{r}\Delta Wx = W_vx + \frac{\alpha}{r}BAx$
- $h = W_Ox + \frac{\alpha}{r}\Delta Wx = W_Ox + \frac{\alpha}{r}BAx$

# Choosing $r$

- ' $r$ ' is selective with experiments
- Putting all parameters in  $\Delta W_q, \Delta W_k$  results in lower performance
- Adapting both yields the best result
- Preferable to adapt more weight matrices than a single type of weights with larger rank

	# of Trainable Parameters = 18M						
Weight Type Rank $r$	$W_q$ 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	<b>73.7</b>	<b>73.7</b>
MultiNLI ( $\pm 0.1\%$ )	91.0	90.8	91.0	91.3	91.3	91.3	<b>91.7</b>

# Training

---

- Standard Objective

$$\max_{\Phi} \sum_{(x,y) \in \mathcal{Z}} \sum_{t=1}^{|y|} \log (P_{\Phi}(y_t|x, y_{<t}))$$

- Model updates all its parameters
- Parameters need to be stored and maintained for each task
- If the pre-trained model is large (e.g., GPT-3), it's computationally expensive

- LoRA Objective

$$\max_{\Theta} \sum_{(x,y) \in \mathcal{Z}} \sum_{t=1}^{|y|} \log (p_{\Phi_0 + \Delta\Phi(\Theta)}(y_t|x, y_{<t}))$$

- Pre-trained weights are kept frozen
- Only low-rank matrices – A and B (collectively denoted as  $\Theta$ )
- Same loss, fewer parameters
- During backpropagation, gradients are computed with respect to A and B (since  $W_0$  is frozen)

# Result

Model & Method	# Trainable Parameters	E2E NLG Challenge				
		BLEU	NIST	MET	ROUGE-L	CIDEr
GPT-2 M (FT)*	354.92M	68.2	8.62	46.2	71.0	2.47
GPT-2 M (Adapter <sup>L</sup> )*	0.37M	66.3	8.41	45.0	69.8	2.40
GPT-2 M (Adapter <sup>L</sup> )*	11.09M	68.9	8.71	46.1	71.3	2.47
GPT-2 M (Adapter <sup>H</sup> )	11.09M	67.3 $\pm$ .6	8.50 $\pm$ .07	46.0 $\pm$ .2	70.7 $\pm$ .2	2.44 $\pm$ .01
GPT-2 M (FT <sup>Top2</sup> )*	25.19M	68.1	8.59	46.0	70.8	2.41
GPT-2 M (PreLayer)*	0.35M	69.7	8.81	46.1	71.4	2.49
GPT-2 M (LoRA)	0.35M	<b>70.4<math>\pm</math>.1</b>	<b>8.85<math>\pm</math>.02</b>	<b>46.8<math>\pm</math>.2</b>	<b>71.8<math>\pm</math>.1</b>	<b>2.53<math>\pm</math>.02</b>
GPT-2 L (FT)*	774.03M	68.5	8.78	46.0	69.9	2.45
GPT-2 L (Adapter <sup>L</sup> )	0.88M	69.1 $\pm$ .1	8.68 $\pm$ .03	46.3 $\pm$ .0	71.4 $\pm$ .2	<b>2.49<math>\pm</math>.0</b>
GPT-2 L (Adapter <sup>L</sup> )	23.00M	68.9 $\pm$ .3	8.70 $\pm$ .04	46.1 $\pm$ .1	71.3 $\pm$ .2	2.45 $\pm$ .02
GPT-2 L (PreLayer)*	0.77M	70.3	8.85	46.2	71.7	2.47
GPT-2 L (LoRA)	0.77M	<b>70.4<math>\pm</math>.1</b>	<b>8.89<math>\pm</math>.02</b>	<b>46.8<math>\pm</math>.2</b>	<b>72.0<math>\pm</math>.2</b>	2.47 $\pm$ .02

# Result

Model&Method	# Trainable Parameters	WikiSQL	MNLI-m	SAMSum
		Acc. (%)	Acc. (%)	R1/R2/RL
GPT-3 (FT)	175,255.8M	<b>73.8</b>	89.5	52.0/28.0/44.5
GPT-3 (BitFit)	14.2M	71.3	91.0	51.3/27.4/43.5
GPT-3 (PreEmbed)	3.2M	63.1	88.6	48.3/24.2/40.5
GPT-3 (PreLayer)	20.2M	70.1	89.5	50.8/27.3/43.5
GPT-3 (Adapter <sup>H</sup> )	7.1M	71.9	89.8	53.0/28.9/44.8
GPT-3 (Adapter <sup>H</sup> )	40.1M	73.2	<b>91.5</b>	53.2/29.0/45.1
GPT-3 (LoRA)	4.7M	73.4	<b>91.7</b>	<b>53.8/29.8/45.9</b>
GPT-3 (LoRA)	37.7M	<b>74.0</b>	<b>91.6</b>	53.4/29.2/45.1

# Advantages

---

- Compute and memory efficient
- Reduce VRAM usage by up to 2/3 (on GPT-3, 1.2TB to 350GB)
- Checkpoint size reduced from 350GB to 35MB
- Can train with significantly fewer GPUs and avoid I/O bottlenecks
- Switch between tasks while deployed at a much lower cost by only swapping the LoRA weights

# Summary

---

- LoRA is an efficient way to train a pre-trained model
- Freezes pre-trained model and injects trainable matrices
- Can reduce memory footprint and computation
- Achieves similar or higher performance compared to fine-tuning

# Q&A