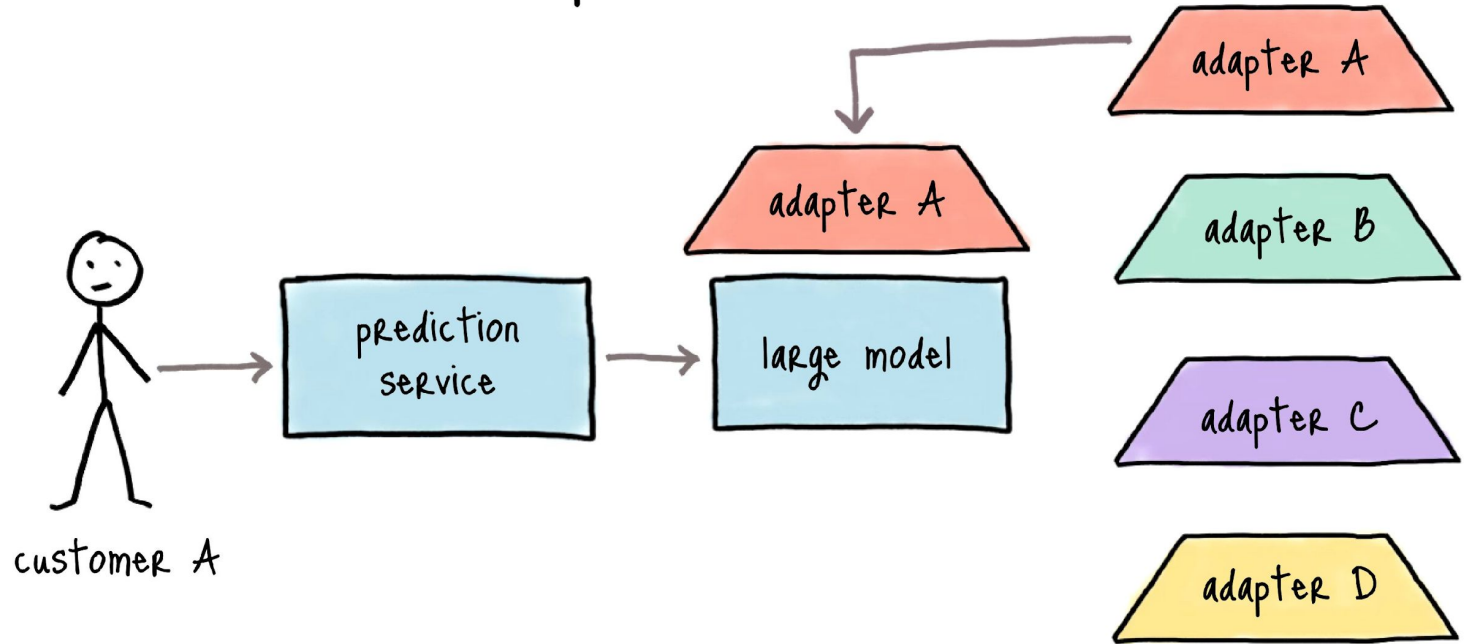


# LORA: LOW-RANK ADAPTATION OF LARGE LANGUAGE MODELS

Edward J. Hu and Yelong Shen and Phillip Wallis and Zeyuan Allen-Zhu and Yuanzhi Li and Shean Wang  
and Lu Wang and Weizhu Chen

Presented by,  
Anagha M B, Gokul P V, and Chandu Dadi

# personalized models



# Challenges with Existing Adaptation Methods

## **Adapter Layers:**

Small modules inserted into transformer layers. Their purpose is to modify the model's behavior for specific tasks without fine-tuning the entire network.

**Advantages:** Parameter-efficient; supports multi-task learning

## **Problems:**

- **Inference Latency:** Additional depth increases sequential computations.
- **Scalability Issues:** Requires GPU synchronization (e.g., AllReduce) in large-scale deployments.

## Prompt Tuning:

Modifies **input tokens** with task-specific prefixes.

### Example:

Input sequence: "Translate English to French: What is your name?"

With prompt tuning: "[TASK\_SPECIFIC] Translate English to French: What is your name?"

**Advantages:** Keeps model weights frozen; lightweight.

### Problems:

- **Optimization Difficulty:** Training prompt embeddings is unstable.
- **Sequence Length Reduction:** Prefix tokens reduce space for task input.

# Eckart-Young Theorem

Given a matrix  $A \in \mathbb{R}^{m \times n}$  and its singular value decomposition:

$$A = U \Sigma V^T$$

where  $\Sigma = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_r)$ , with  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$  (singular values), the best rank- $k$  approximation of  $A$ , in the Frobenius norm, is:

$$A_k = U_k \Sigma_k V_k^T$$

where  $\Sigma_k$  is the diagonal matrix containing the top  $k$  singular values, and  $U_k$ ,  $V_k$  are the corresponding singular vectors.

# Key Ideas of LoRA

- Instead of updating all model parameters, LoRA injects trainable low-rank matrices into the pretrained weights.
- Pretrained weights remain frozen, reducing computational overhead.
- Parameter updates are represented as the product of two low-rank matrices.
- Gradients scaled by  $\alpha/r$

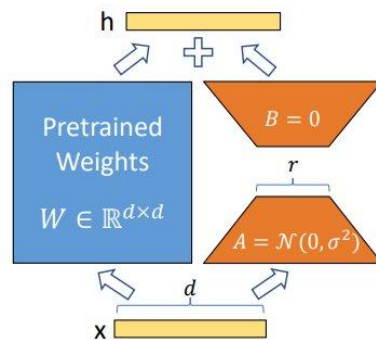


Figure 1: Our reparametrization. We only train  $A$  and  $B$ .

$$h = W_0 x + \Delta W x = W_0 x + B A x$$

# LoRA Memory Efficiency

- The low-rank structure significantly reduces memory overhead and computational complexity during training.

Storage:  $d \times k \rightarrow r(d+k)$

- Example: For  $d=512$ ,  $k=512$ ,  $r=8$
  - Original: 262,144 parameters
  - LoRA: 8,192 parameters (96.9% reduction)
- Maintains expressiveness of full-rank updates with far fewer trainable parameters.



[https://colab.research.google.com/drive/1DzArDRxgcPiu32fwWd3BgREy9SKRcOJ\\_](https://colab.research.google.com/drive/1DzArDRxgcPiu32fwWd3BgREy9SKRcOJ_)



<https://colab.research.google.com/drive/142Wy4B5cVUWnzS2JaSt9jPsyVcYTtB3O>



# Applications of LoRA

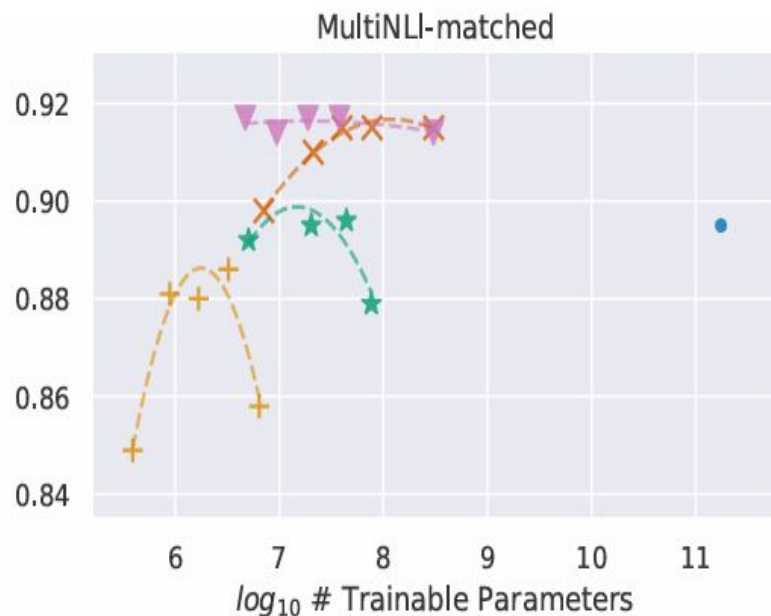
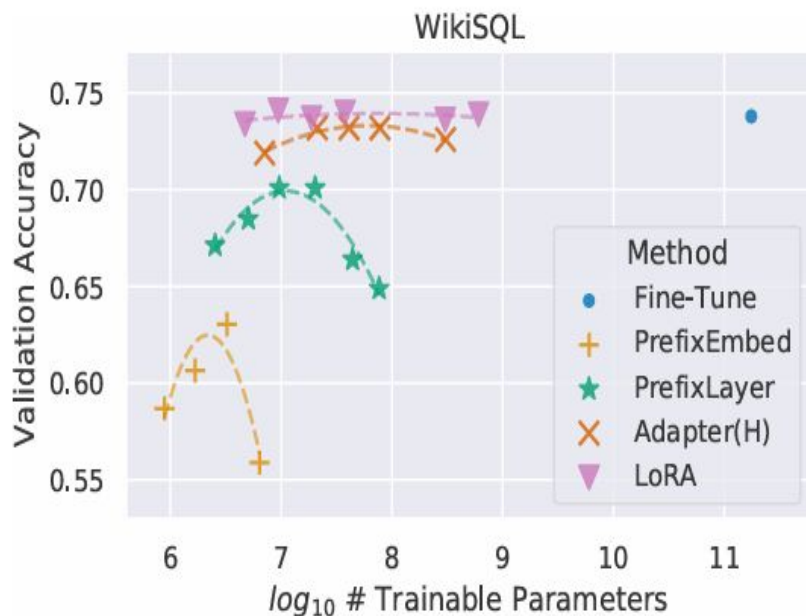
1. Natural Language Processing (NLP)
  - a. **Text Classification:** Fine-tuning LLMs like RoBERTa, BERT, or GPT for sentiment analysis, spam detection, or topic classification with minimal parameter updates.
  - b. **Question Answering (QA):** Tuning models like RoBERTa or GPT for domain-specific Q&A applications, such as customer support or medical FAQs.
2. Personalized AI and Chatbots
3. Generative AI

# EMPIRICAL EXPERIMENTS

Model & Method	# Trainable Parameters	MNLI	SST-2	MRPC	CoLA	QNLI	QQP	RTE	STS-B	Avg.
RoB <sub>base</sub> (FT)*	125.0M	<b>87.6</b>	94.8	90.2	<b>63.6</b>	92.8	<b>91.9</b>	78.7	91.2	86.4
RoB <sub>base</sub> (BitFit)*	0.1M	84.7	93.7	<b>92.7</b>	62.0	91.8	84.0	81.5	90.8	85.2
RoB <sub>base</sub> (Adpt <sup>D</sup> )*	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 <sub>base</sub> (Adpt <sup>D</sup> )*	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 <sub>base</sub> (LoRA)	0.3M	87.5 $\pm$ .3	<b>95.1<math>\pm</math>.2</b>	89.7 $\pm$ .7	63.4 $\pm$ 1.2	<b>93.3<math>\pm</math>.3</b>	90.8 $\pm$ .1	<b>86.6<math>\pm</math>.7</b>	<b>91.5<math>\pm</math>.2</b>	<b>87.2</b>
RoB <sub>large</sub> (FT)*	355.0M	90.2	<b>96.4</b>	<b>90.9</b>	68.0	94.7	<b>92.2</b>	86.6	92.4	88.9
RoB <sub>large</sub> (LoRA)	0.8M	<b>90.6<math>\pm</math>.2</b>	96.2 $\pm$ .5	<b>90.9<math>\pm</math>1.2</b>	<b>68.2<math>\pm</math>1.9</b>	<b>94.9<math>\pm</math>.3</b>	91.6 $\pm$ .1	<b>87.4<math>\pm</math>2.5</b>	<b>92.6<math>\pm</math>.2</b>	<b>89.0</b>
RoB <sub>large</sub> (Adpt <sup>P</sup> )†	3.0M	90.2 $\pm$ .3	96.1 $\pm$ .3	90.2 $\pm$ .7	<b>68.3<math>\pm</math>1.0</b>	<b>94.8<math>\pm</math>.2</b>	<b>91.9<math>\pm</math>.1</b>	83.8 $\pm$ 2.9	92.1 $\pm$ .7	88.4
RoB <sub>large</sub> (Adpt <sup>P</sup> )†	0.8M	<b>90.5<math>\pm</math>.3</b>	<b>96.6<math>\pm</math>.2</b>	89.7 $\pm$ 1.2	67.8 $\pm$ 2.5	<b>94.8<math>\pm</math>.3</b>	91.7 $\pm$ .2	80.1 $\pm$ 2.9	91.9 $\pm$ .4	87.9
RoB <sub>large</sub> (Adpt <sup>H</sup> )†	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 <sub>large</sub> (Adpt <sup>H</sup> )†	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 <sub>large</sub> (LoRA)†	0.8M	<b>90.6<math>\pm</math>.2</b>	96.2 $\pm$ .5	<b>90.2<math>\pm</math>1.0</b>	68.2 $\pm$ 1.9	<b>94.8<math>\pm</math>.3</b>	91.6 $\pm$ .2	<b>85.2<math>\pm</math>1.1</b>	<b>92.3<math>\pm</math>.5</b>	<b>88.6</b>
DeB <sub>XXL</sub> (FT)*	1500.0M	91.8	<b>97.2</b>	92.0	72.0	<b>96.0</b>	92.7	93.9	92.9	91.1
DeB <sub>XXL</sub> (LoRA)	4.7M	<b>91.9<math>\pm</math>.2</b>	96.9 $\pm$ .2	<b>92.6<math>\pm</math>.6</b>	<b>72.4<math>\pm</math>1.1</b>	<b>96.0<math>\pm</math>.1</b>	<b>92.9<math>\pm</math>.1</b>	<b>94.9<math>\pm</math>.4</b>	<b>93.0<math>\pm</math>.2</b>	<b>91.3</b>

# EMPIRICAL EXPERIMENTS

GPT-3 175 B



# CONCLUSION

Fine-tuning enormous language models is expensive in terms of the hardware required and the storage/switching cost for hosting independent instances for different tasks.

LoRA, an efficient adaptation strategy that neither introduces inference latency nor reduces input sequence length while retaining high model quality.

# References

- <https://arxiv.org/pdf/2106.09685>
- <https://magazine.sebastianraschka.com/p/practical-tips-for-finetuning-llms>
- <https://towardsdatascience.com/understanding-lora-low-rank-adaptation-for-finetuning-large-models-936bce1a07c6>
- <https://github.com/huggingface/smol-course>

Thank you!