# Paper Sharing

## LoRA：Low-Rank Adaptation of Large Language Models
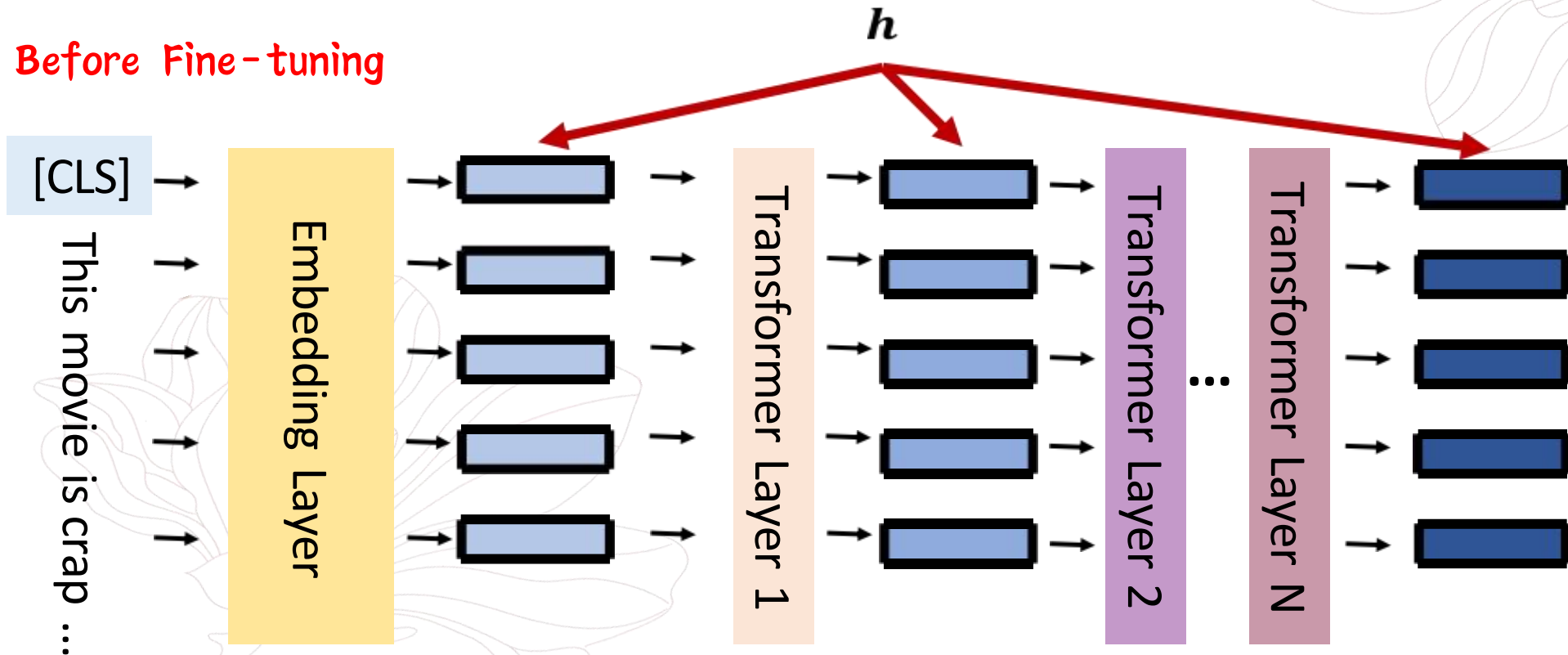
Lecturer：Yuxin Wu

2024.1.18

# Fine-tuning

- ## What is standard fine-tuning really doing?
- Modify the hidden representations ( $h$ ) of the PLM such that it can perform well on downstream task

$h$

Before Fine-tuning

[CLS] This movie is crap ...

Embedding Layer

Transformer Layer 1

Transformer Layer 2

...

Transformer Layer N

# Fine-tuning

- ## What is standard fine-tuning really doing?
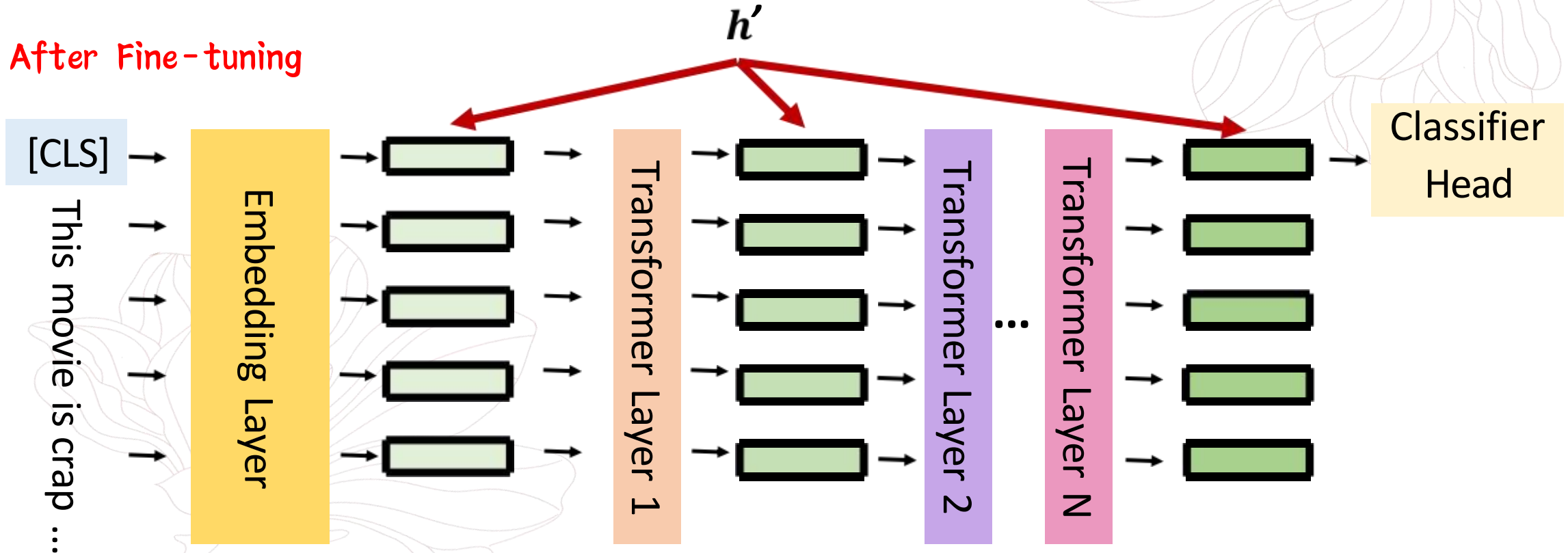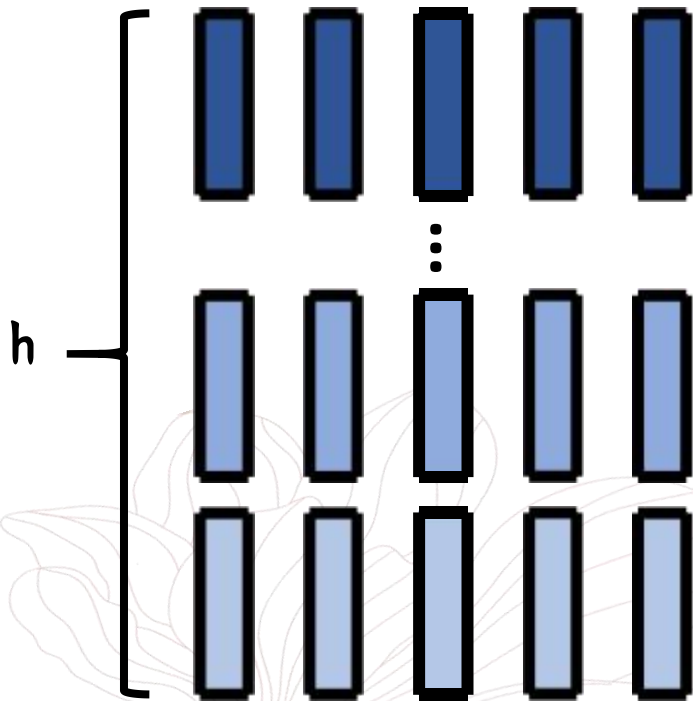- Modify the hidden representations ( $h$ ) of the PLM such that it can perform well on downstream task

# Parameter-Efficient Fine-tuning

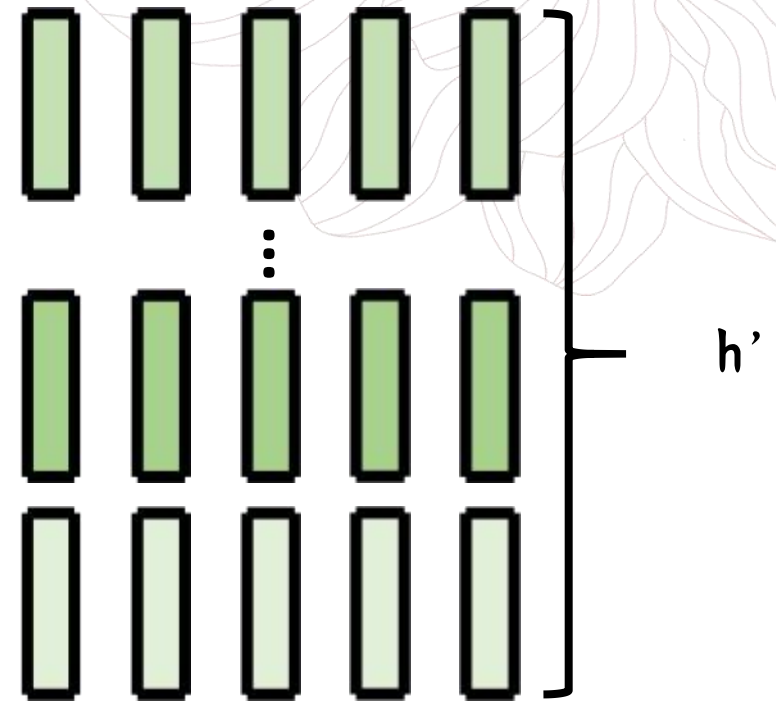- Fine-tuning = modifying the hidden representation based on a PLM



Before Fine-tuning

After Fine-tuning

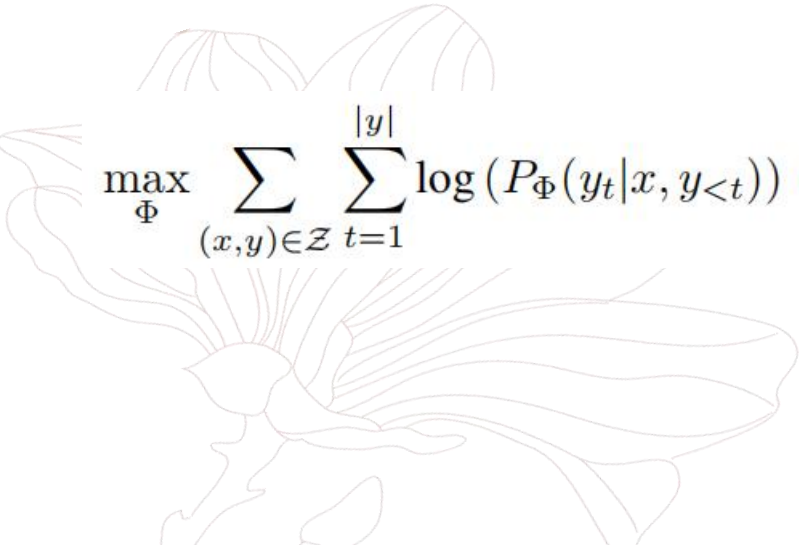h: hidden representation calculated by the original PLM

h': hidden representation calculated by the pre-trained model

# Parameter-Efficient Fine-tuning

- However

Finetune LLM is quite expensive
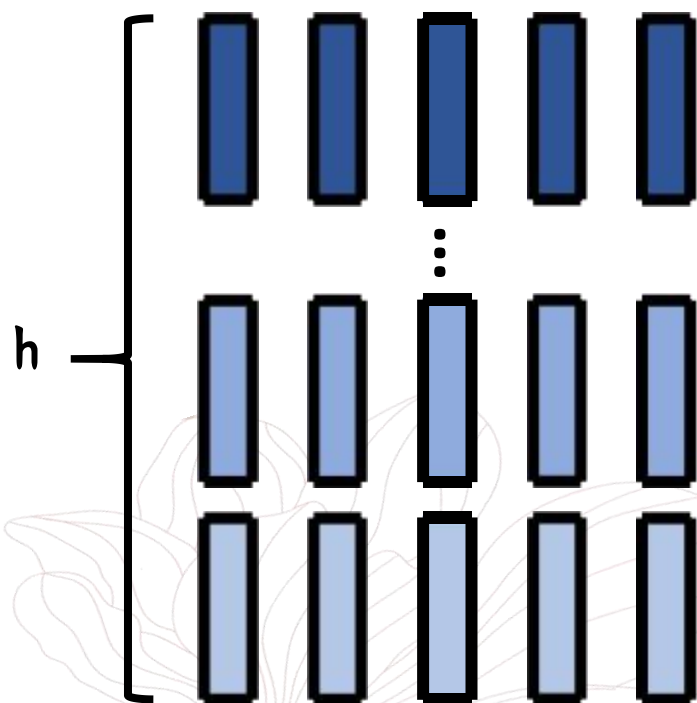
A more parameter-efficient approach is needed

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

$$\implies$$

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

# Parameter-Efficient Fine-tuning

- Fine-tuning = modifying the hidden representation based on a PLM

Before Fine-tuning

After Fine-tuning



h

h'

$= h + \Delta h$

h: hidden representation calculated by the original PLM

h': hidden representation calculated by the pre-trained model

# Parameter-Efficient Fine-tuning: Adapter

- Use special submodules to modify hidden representations



Before Fine-tuning

After Fine-tuning

Adapter

h

h'

$= h + \Delta h$

h: hidden representation calculated by the original PLM

h' : hidden representation calculated by the pre-trained model

# Parameter-Efficient Fine-tuning: Adapter

- Adapters: small trainable submodules inserted in transformers
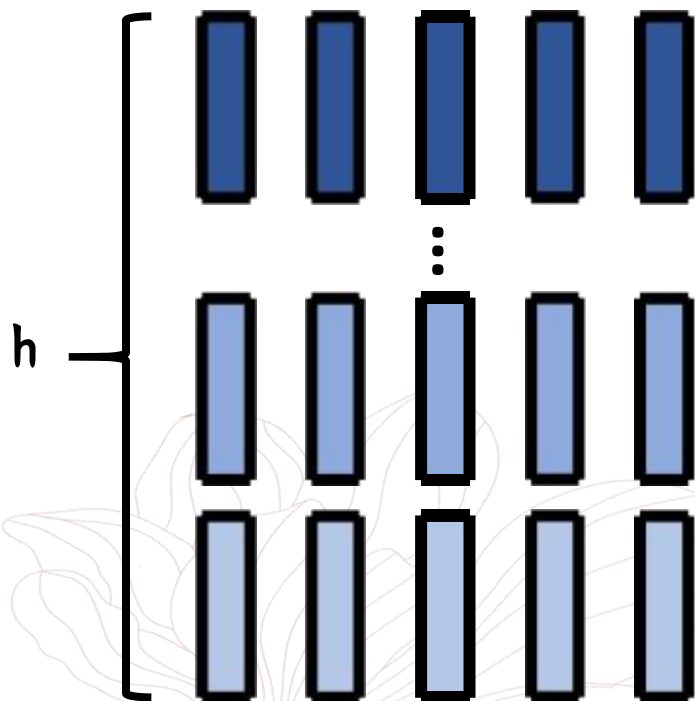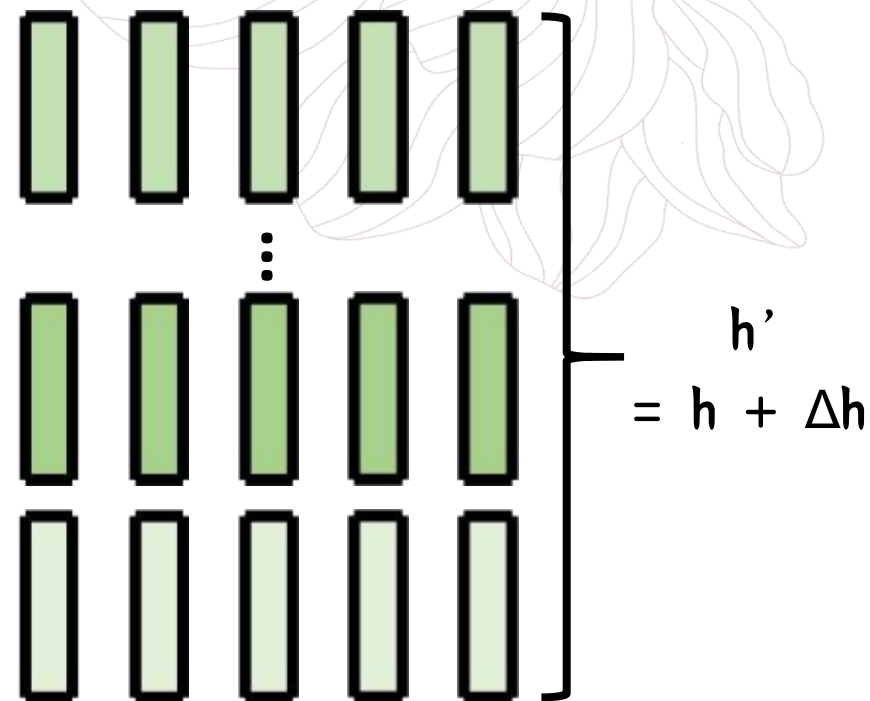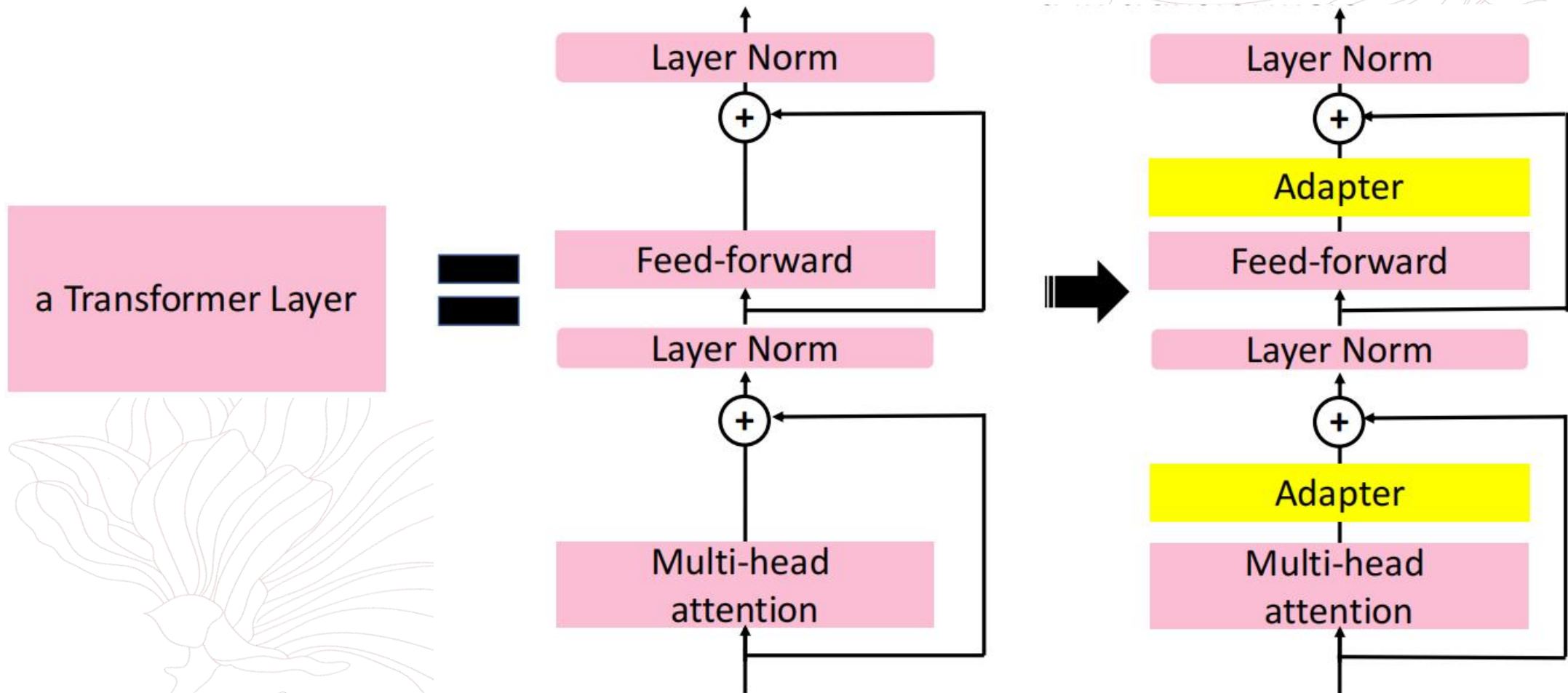
# Parameter-Efficient Fine-tuning: Adapter

- Adapters: Introduce Inference Latency

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

Table 1: Infernece 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 B.

# Parameter-Efficient Fine-tuning: Prefix

- Prefix

    - difficult to optimize

- reserving a part of the sequence length for adaptation necessarily reduces the sequence length available to process a downstream task

# Parameter-Efficient Fine-tuning: LoRA

- Use special submodules to modify hidden representations



Before Fine-tuning

After Fine-tuning

LoRA

h

h'

$= h + \Delta h$

h: hidden representation
calculated by the original PLM

h' : hidden representation
calculated by the pre-trained model

# Parameter-Efficient Fine-tuning: LoRA

- LoRA: Low-Rank Adaptation of Large Language Models



transformer layer

# Parameter-Efficient Fine-tuning: LoRA

- LoRA: Low-Rank Adaptation of Large Language Models

# Parameter-Efficient Fine-tuning: LoRA

- LoRA

# Parameter-Efficient Fine-tuning: LoRA

- ## LoRA

  - <span style="color:red">Low-Rank Adaptation</span> of Large Language Models
  - Motivation: Downstream fine-tunings have low intrinsic dimension
  - Weight after fine-tuning = $w_o$ (pre-trained weight) $+\Delta w$ (updates to the weight)
  - Hypothesis: The updates to the weight ($\Delta w$) also gives a low intrinsic rank
  - Fine-tuned weight = $w_o + \Delta w = w_o + BA$, rank $r \ll \min(d_{FFW}, d_{model})$

# Parameter-Efficient Fine-tuning: LoRA

- LoRA

LoRA possesses several key advantages.

- A pre-trained model can be shared and used to build many small LoRA modules for different tasks. We can freeze the shared model and efficiently switch tasks by replacing the matrices $A$ and $B$ in Figure 1, reducing the storage requirement and task-switching overhead significantly.

- LoRA makes training more efficient and lowers the hardware barrier to entry by up to 3 times when using adaptive optimizers since we do not need to calculate the gradients or maintain the optimizer states for most parameters. Instead, we only optimize the injected, much smaller low-rank matrices.

- Our simple linear design allows us to merge the trainable matrices with the frozen weights when deployed, *introducing no inference latency* compared to a fully fine-tuned model, by construction.

- LoRA is orthogonal to many prior methods and can be combined with many of them, such as prefix-tuning. We provide an example in Appendix E.

# Experiments

- LoRA

| 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$_{\pm.0}$ | 94.2$_{\pm.1}$ | 88.5$_{\pm1.1}$ | 60.8$_{\pm.4}$ | 93.1$_{\pm.1}$ | 90.2$_{\pm.0}$ | 71.5$_{\pm2.7}$ | 89.7$_{\pm.3}$ | 84.4 |
| RoB$_{base}$ (Adpt$^D$)* | 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$_{\pm2.2}$ | 90.3$_{\pm.1}$ | 85.4 |
| RoB$_{base}$ (LoRA) | 0.3M | 87.5$_{\pm.3}$ | 95.1$_{\pm.2}$ | 89.7$_{\pm.7}$ | 63.4$_{\pm1.2}$ | 93.3$_{\pm.3}$ | 90.8$_{\pm.1}$ | 86.6$_{\pm.7}$ | 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$_{\pm.2}$ | 96.2$_{\pm.5}$ | 90.9$_{\pm1.2}$ | 68.2$_{\pm1.9}$ | 94.9$_{\pm.3}$ | 91.6$_{\pm.1}$ | 87.4$_{\pm2.5}$ | 92.6$_{\pm.2}$ | 89.0 |
| RoB$_{large}$ (Adpt$^P$)† | 3.0M | 90.2$_{\pm.3}$ | 96.1$_{\pm.3}$ | 90.2$_{\pm.7}$ | 68.3$_{\pm1.0}$ | 94.8$_{\pm.2}$ | 91.9$_{\pm.1}$ | 83.8$_{\pm2.9}$ | 92.1$_{\pm.7}$ | 88.4 |
| RoB$_{large}$ (Adpt$^P$)† | 0.8M | 90.5$_{\pm.3}$ | 96.6$_{\pm.2}$ | 89.7$_{\pm1.2}$ | 67.8$_{\pm2.5}$ | 94.8$_{\pm.3}$ | 91.7$_{\pm.2}$ | 80.1$_{\pm2.9}$ | 91.9$_{\pm.4}$ | 87.9 |
| RoB$_{large}$ (Adpt$^H$)† | 6.0M | 89.9$_{\pm.5}$ | 96.2$_{\pm.3}$ | 88.7$_{\pm2.9}$ | 66.5$_{\pm4.4}$ | 94.7$_{\pm.2}$ | 92.1$_{\pm.1}$ | 83.4$_{\pm1.1}$ | 91.0$_{\pm1.7}$ | 87.8 |
| RoB$_{large}$ (Adpt$^H$)† | 0.8M | 90.3$_{\pm.3}$ | 96.3$_{\pm.5}$ | 87.7$_{\pm1.7}$ | 66.3$_{\pm2.0}$ | 94.7$_{\pm.2}$ | 91.5$_{\pm.1}$ | 72.9$_{\pm2.9}$ | 91.5$_{\pm.5}$ | 86.4 |
| RoB$_{large}$ (LoRA)† | 0.8M | 90.6$_{\pm.2}$ | 96.2$_{\pm.5}$ | 90.2$_{\pm1.0}$ | 68.2$_{\pm1.9}$ | 94.8$_{\pm.3}$ | 91.6$_{\pm.2}$ | 85.2$_{\pm1.1}$ | 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$_{\pm.2}$ | 92.6$_{\pm.6}$ | 72.4$_{\pm1.1}$ | 96.0$_{\pm.1}$ | 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.

| Model & Method | # Trainable Parameters | 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$^L$)* | 0.37M | 66.3 | 8.41 | 45.0 | 69.8 | 2.40 |
| GPT-2 M (Adapter$^L$)* | 11.09M | 68.9 | 8.71 | 46.1 | 71.3 | 2.47 |
| GPT-2 M (Adapter$^H$) | 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$^{Top2}$)* | 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 | 70.4$_{\pm.1}$ | 8.85$_{\pm.02}$ | 46.8$_{\pm.2}$ | 71.8$_{\pm.1}$ | 2.53$_{\pm.02}$ |
| GPT-2 L (FT)* | 774.03M | 68.5 | 8.78 | 46.0 | 69.9 | 2.45 |
| GPT-2 L (Adapter$^L$) | 0.88M | 69.1$_{\pm.1}$ | 8.68$_{\pm.03}$ | 46.3$_{\pm.0}$ | 71.4$_{\pm.2}$ | 2.49$_{\pm.0}$ |
| GPT-2 L (Adapter$^L$) | 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 | 70.4$_{\pm.1}$ | 8.89$_{\pm.02}$ | 46.8$_{\pm.2}$ | 72.0$_{\pm.2}$ | 2.47$_{\pm.02}$ |

Table 3: GPT-2 medium (M) and large (L) with different adaptation methods on the E2E NLG Challenge. For all metrics, higher is better. LoRA outperforms several baselines with comparable or fewer trainable parameters. Confidence intervals are shown for experiments we ran. * indicates numbers published in prior works.

| Model&Method | # Trainable Parameters | WikiSQL Acc. (%) | MNLI-m Acc. (%) | SAMSum R1/R2/RL |
|---|---|---|---|---|
| GPT-3 (FT) | 175,255.8M | 73.8 | 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$^H$) | 7.1M | 71.9 | 89.8 | 53.0/28.9/44.8 |
| GPT-3 (Adapter$^H$) | 40.1M | 73.2 | 91.5 | 53.2/29.0/45.1 |
| GPT-3 (LoRA) | 4.7M | 73.4 | 91.7 | 53.8/29.8/45.9 |
| GPT-3 (LoRA) | 37.7M | 74.0 | 91.6 | 53.4/29.2/45.1 |

Table 4: Performance of different adaptation methods on GPT-3 175B. We report the logical form validation accuracy on WikiSQL, validation accuracy on MultiNLI-matched, and Rouge-1/2/L on SAMSum. LoRA performs better than prior approaches, including full fine-tuning. The results on WikiSQL have a fluctuation around ±0.5%, MNLI-m around ±0.1%, and SAMSum around ±0.2/±0.2/±0.1 for the three metrics.

# Thank you for listening

主讲人：吴雨欣

2024.1.18