

Lecture 13: Fine Tuning Alternatives

Andrei Arsene Simion

Columbia University

April 26, 2023

Outline

- ▶ Adapters
- ▶ Prefix-Tuning
- ▶ LORA

Outline

- ▶ **Adapters**
- ▶ Prefix-Tuning
- ▶ LORA

Fine Tuning

- ▶ Given a transformer, what options do we have to use it for a specific task?
- ▶ Let's focus on Classification and BERT, where we add a MLP on top of the CLS token
- ▶ Option1: We can optimize just the "head" MLP parameters
 - ▶ Using BERT as a feature extractor
 - ▶ This is relatively cheap
 - ▶ We do not, however, get SOTA on tasks
- ▶ Option 2: We can optimize ALL of the parameters (BERT + MLP)
 - ▶ This is fine-tuning
 - ▶ We can get SOTA
 - ▶ But, the big problem is that this is quite expensive
- ▶ Is there something in-between?
- ▶ We saw one example: ULM-Fit

Adapter

- ▶ The idea of an adapter is to insert a few critical layers in a transformer and allow gradients only on parameters from these layers
 - ▶ If you pick the architecture right, you can save time on fine tuning
 - ▶ Generally, this will increase the inference cost
 - ▶ Big idea: insert a small MLP, tune the LayerNorm parameters
 - ▶ Let d_{model} be the embedding dimensions, and $d_{adapter}$ be the MLP
 - ▶ Total gradient parameters per transformer block:
 - ▶ MLP: $2d_{model}d_{adapter} + d_{model} + d_{adapter}$
 - ▶ LayerNorm: $4d_{model}$
 - ▶ Use $d_{adapter} \ll d_{model}$

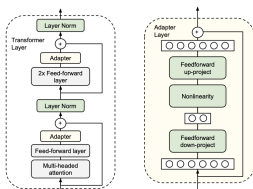


Figure 2. Architecture of the adapter module and its integration with the Transformer. **Left:** We add the adapter module twice to each Transformer layer: after the projection following multi-headed attention and after the two feed-forward layers. **Right:** The adapter consists of a bottleneck which contains few parameters relative to the attention and feedforward layers in the original model. The adapter also contains a skip-connection. During adapter tuning, the green layers are trained on the downstream data, this includes the adapter, the layer normalization parameters, and the final classification layer (not shown in the figure).

Adapter vs Fine Tuning

- ▶ Adapter-based tuning requires training two orders of magnitude fewer parameters to fine-tuning, while attaining similar performance

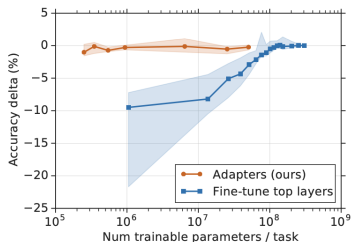


Figure 1. Trade-off between accuracy and number of trained task-specific parameters, for adapter tuning and fine-tuning. The y-axis is normalized by the performance of full fine-tuning, details in Section 3. The curves show the 20th, 50th, and 80th performance percentiles across nine tasks from the GLUE benchmark. Adapter-based tuning attains a similar performance to full fine-tuning with two orders of magnitude fewer trained parameters.

Adapter $d_{adapter}$ depends on task

- GLUE: Different $d_{adapter}$ might be optimal per task

| | Total num params | Trained params / task | CoLA | SST | MRPC | STS-B | QQP | MNLI _m | MNLI _{mm} | QNLI | RTE | Total |
|-----------------------|---------------------|--------------------------|------|------|------|-------|------|-------------------|--------------------|------|------|-------|
| BERT _{LARGE} | 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 |

Table 1. Results on GLUE test sets scored using the GLUE evaluation server. MRPC and QQP are evaluated using F1 score. STS-B is evaluated using Spearman's correlation coefficient. CoLA is evaluated using Matthew's Correlation. The other tasks are evaluated using accuracy. Adapter tuning achieves comparable overall score (80.0) to full fine-tuning (80.4) using 1.3× parameters in total, compared to 9×. Fixing the adapter size to 64 leads to a slightly decreased overall score of 79.6 and slightly smaller model.

Adapter on other tasks

- ▶ BERT Baselines are stronger for Classification than other baselines
- ▶ Variable fine-tuning: fine-tune only the top k layers, and freeze the remainder
- ▶ Each task as a different setting
- ▶ Fine-tuning might not be optimal; it can be overkill

| Dataset | No BERT baseline | BERT _{BASE} Fine-tune | BERT _{BASE} Variable FT | BERT _{BASE} Adapters |
|------------------------------------|---------------------|-----------------------------------|-------------------------------------|----------------------------------|
| 20 newsgroups | 91.1 | 92.8 ± 0.1 | 92.8 ± 0.1 | 91.7 ± 0.2 |
| Crowdfower airline | 84.5 | 83.6 ± 0.3 | 84.0 ± 0.1 | 84.5 ± 0.2 |
| Crowdfower corporate messaging | 91.9 | 92.5 ± 0.5 | 92.4 ± 0.6 | 92.9 ± 0.3 |
| Crowdfower disasters | 84.9 | 85.3 ± 0.4 | 85.3 ± 0.4 | 84.1 ± 0.2 |
| Crowdfower economic news relevance | 81.1 | 82.1 ± 0.0 | 78.9 ± 2.8 | 82.5 ± 0.3 |
| Crowdfower emotion | 36.3 | 38.4 ± 0.1 | 37.6 ± 0.2 | 38.7 ± 0.1 |
| Crowdfower global warming | 82.7 | 84.2 ± 0.4 | 81.9 ± 0.2 | 82.7 ± 0.3 |
| Crowdfower political audience | 81.0 | 80.9 ± 0.3 | 80.7 ± 0.8 | 79.0 ± 0.5 |
| Crowdfower political bias | 76.8 | 75.2 ± 0.9 | 76.5 ± 0.4 | 75.9 ± 0.3 |
| Crowdfower political message | 43.8 | 38.9 ± 0.6 | 44.9 ± 0.6 | 44.1 ± 0.2 |
| Crowdfower primary emotions | 33.5 | 36.9 ± 1.6 | 38.2 ± 1.0 | 33.9 ± 1.4 |
| Crowdfower progressive opinion | 70.6 | 71.6 ± 0.5 | 75.9 ± 1.3 | 71.7 ± 1.1 |
| Crowdfower progressive stance | 54.3 | 63.8 ± 1.0 | 61.5 ± 1.3 | 60.6 ± 1.4 |
| Crowdfower US economic performance | 75.6 | 75.3 ± 0.1 | 76.5 ± 0.4 | 77.3 ± 0.1 |
| Customer complaint database | 54.5 | 55.9 ± 0.1 | 56.4 ± 0.1 | 55.4 ± 0.1 |
| News aggregator dataset | 95.2 | 96.3 ± 0.0 | 96.5 ± 0.0 | 96.2 ± 0.0 |
| SMS spam collection | 98.5 | 99.3 ± 0.2 | 99.3 ± 0.2 | 95.1 ± 2.2 |
| Average | 72.7 | 73.7 | 74.0 | 73.3 |
| Total number of params | — | 17× | 9.9× | 1.19× |
| Trained params/task | — | 100% | 52.9% | 1.14% |

Table 2. Test accuracy for additional classification tasks. In these experiments we transfer from the BERT_{BASE} model. For each task and algorithm, the model with the best validation set accuracy is chosen. We report the mean test accuracy and s.e.m. across runs with different random seeds.

Adapter results: other cases

- ▶ Fine-tune just the top k BERT layers or just tune LayerNorm?
- ▶ LayerNorm tuning alone is suboptimal

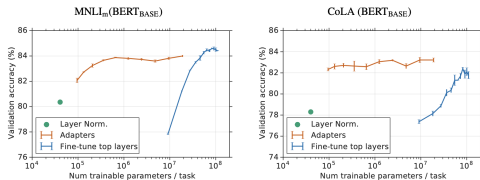


Figure 4. Validation set accuracy versus number of trained parameters for three methods: (i) Adapter tuning with an adapter sizes 2^n for $n = 0 \dots 9$ (orange). (ii) Fine-tuning the top k layers for $k = 1 \dots 12$ (blue). (iii) Tuning the layer normalization parameters only (green). Error bars indicate ± 1 s.e.m. across three random seeds.

Adapter Results

► Aggregate results for Adapters vs fine tuning

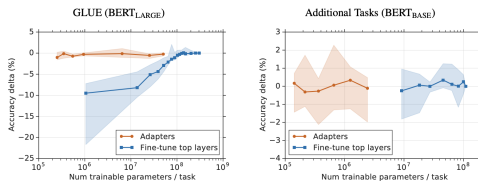


Figure 3. Accuracy versus the number of trained parameters, aggregated across tasks. We compare adapters of different sizes (orange) with fine-tuning the top n layers, for varying n (blue). The lines and shaded areas indicate the 20th, 50th, and 80th percentiles across tasks. For each task and algorithm, the best model is selected for each point along the curve. For GLUE, the validation set accuracy is reported. For the additional tasks, we report the test-set accuracies. To remove the intra-task variance in scores, we normalize the scores for each model and task by subtracting the performance of full fine-tuning on the corresponding task.

Adapter Results

- ▶ Adapters also work on SQUAD based tasks

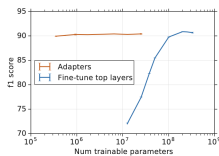


Figure 5. Validation accuracy versus the number of trained parameters for SQuAD v1.1. Error bars indicate the s.e.m. across three seeds, using the best hyperparameters.

Adapter Robustness

- ▶ Ablating one layer's Adaptor makes no difference
- ▶ Adapters on the lower layers have a smaller impact than the higher-layers - **those are the crucial ones**
 - ▶ This indicates that adapters perform well because they automatically prioritize higher layers
 - ▶ One intuition is that the lower layers extract lower-level features that are shared among tasks, while the higher layers build features that are unique to different task
- ▶ The performance of adapters is robust for standard deviations below 10^{-2}
 - ▶ However, when the initialization is too large, performance degrades

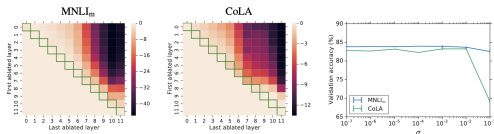


Figure 6. **Left, Center:** Ablation of trained adapters from continuous layer spans. The heatmap shows the relative decrease in validation accuracy to the fully trained adapted model. The y and x axes indicate the first and last layers ablated (inclusive), respectively. The diagonal cells, highlighted in green, indicate ablation of a single layer's adapters. The cell in the top-right indicates ablation of all adapters. Cells in the lower triangle are meaningless, and are set to 0%, the best possible relative performance. **Right:** Performance of BERT_{BASE} using adapters with different initial weight magnitudes. The x-axis is the standard deviation of the initialization distribution.

Outline

- ▶ Adapters
- ▶ Prefix-Tuning
- ▶ LORA

In-Context learning

- ▶ Recall in-context learning from GPT3
 - ▶ 0% of parameters need to be tuned
 - ▶ But, you can't exploit large training sets
 - ▶ Manually written prompts may be suboptimal
- ▶ Contrast this with fine-tuning, which requires 100 % of all parameters and is expensive

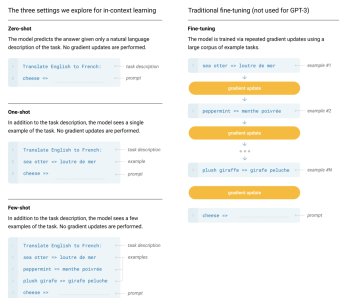


Figure 2.1: Zero-shot, one-shot and few-shot, contrasted with traditional fine-tuning. The panels above show four methods for performing a task with a language model – fine-tuning is the traditional method, whereas zero-, one-, and few-shot, which we study in this work, require the model to perform the task with only forward passes at test time. We typically present the model with a few dozen examples in the few shot setting. Exact phrasings for all task descriptions, examples and prompts can be found in Appendix G.

Prefix-Tuning

- ▶ Can we combine in-context learning and fine-tuning?
 - ▶ It needs to be lightweight
 - ▶ It needs to allow us to exploit a large training set via the trainable prefix
 - ▶ In-context learning only works with large models like GPT3, we'd like something which works with smaller models ...

Fine Tuning vs Adapters vs Prefix-Tuning

- ▶ Tuning top k layers ~ 20 % of parameters
 - ▶ Performance drop with respect to fine tuning
- ▶ Adapters ~ 3 -4 % of parameters
 - ▶ Moderately light weight
 - ▶ Comparable performance to fine-tuning
- ▶ Prefix-tuning ~ 0.1 % of parameters
 - ▶ Very lightweight
 - ▶ Comparable performance to fine-tuning

Prefix-Tuning for different Tasks

- ▶ Prefix-tuning adds a prefix to each transformer block (at each level)
- ▶ For each task, you have a different prefix you optimize over
- ▶ The Language Model's parameters are fixed

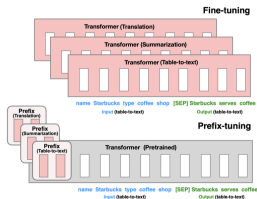


Figure 1: Fine-tuning (top) updates all Transformer parameters (the red Transformer box) and requires storing a full model copy for each task. We propose prefix-tuning (bottom), which freezes the Transformer parameters and only optimizes the prefix (the red prefix blocks). Consequently, we only need to store the prefix for each task, making prefix-tuning modular and space-efficient. Note that each vertical block denote transformer activations at one time step.

Prefix Tuning Architectures

- ▶ Prefix tuning inserts task-specific parameters for different tasks / architectures
- ▶ For example, $z = (PREFIX; x; y)$
- ▶ Put another way, prefix basically overrides the Self Attention Mechanism and this prefix as at every level
- ▶ So, we have h_i

$$h_i = \begin{cases} P_{\theta}(i, :), & \text{if } i \in P_{idx}, \\ LM_{\phi}(z_i, h_{<i}), & \text{otherwise} \end{cases}$$

- ▶ Because of optimizer difficulty, P_{θ} is a MLP

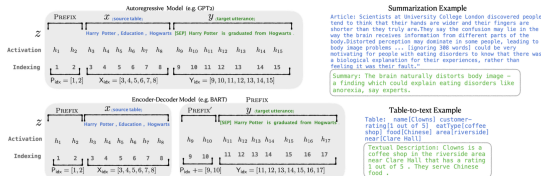


Figure 2: An annotated example of prefix-tuning using an autoregressive LM (top) and an encoder-decoder model (bottom). The prefix activations $\forall i \in P_{idx}, h_i$ are drawn from a trainable matrix P_{θ} . The remaining activations are computed by the Transformer.

Prefix-Tuning: Low Data Settings

- ▶ In low data settings, prefix tuning does better than fine tuning
- ▶ Here, they sample small datasets and explore how well their model does
- ▶ Both methods tend to undergenerate (missing table contents) in low data regimes, prefix-tuning tends to be more faithful than fine-tuning
- ▶ For example, fine-tuning (100, 200) falsely claims a low customer rating while the true rating is average, whereas prefix-tuning (100, 200) generates a description that is faithful to the table.

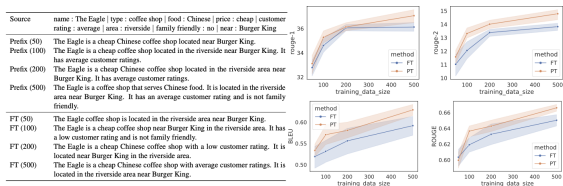


Figure 3: (Left) qualitative examples in lowdata settings. (Right) prefix-tuning (orange) outperforms fine-tuning (blue) in low-data regimes in addition to requiring many fewer parameters. The top two plots correspond to summarization, measured by ROUGE-1 and ROUGE-2. The bottom two plots correspond to table-to-text, measured by BLEU and ROUGE-L. The x-axis is the training size and the y-axis is the evaluation metric (higher is better).

Prefix-Tuning: Summarization (BART)

- ▶ For summarization tasks, prefix tuning does better than fine tuning
- ▶ ROUGE \sim Recall-BLEU

| | R-1 \uparrow | R-2 \uparrow | R-L \uparrow |
|-------------------------------|----------------|----------------|----------------|
| FINE-TUNE(Lewis et al., 2020) | 45.14 | 22.27 | 37.25 |
| PREFIX(2%) | 43.80 | 20.93 | 36.05 |
| PREFIX(0.1%) | 42.92 | 20.03 | 35.05 |

Table 2: Metrics for summarization on XSUM. Prefix-tuning slightly underperforms fine-tuning.

| | news-to-sports | | | within-news | | |
|-----------|----------------|----------------|----------------|----------------|----------------|----------------|
| | R-1 \uparrow | R-2 \uparrow | R-L \uparrow | R-1 \uparrow | R-2 \uparrow | R-L \uparrow |
| FINE-TUNE | 38.15 | 15.51 | 30.26 | 39.20 | 16.35 | 31.15 |
| PREFIX | 39.23 | 16.74 | 31.51 | 39.41 | 16.87 | 31.47 |

Table 3: Extrapolation performance on XSUM. Prefix-tuning outperforms fine-tuning on both news-to-sports and within-news splits.

Prefix Tuning: Prefix Length

- ▶ Longer prefix length means better performance, but the sweet spot depends on the task

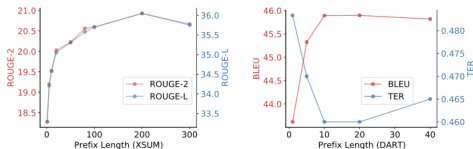


Figure 4: Prefix length vs. performance on summarization (left) and table-to-text (right). Performance increases as the prefix length increases up to a threshold (200 for summarization and 10 for table-to-text) and then a slight performance drop occurs. Each plot reports two metrics (on two vertical axes).

Prefix Tuning: Overall results

- Overall, prefix tuning was superior to fine tuning / adapters

| | E2E | | | | | WebNLG | | | | | | | | | DART | | | | | |
|-------------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| | BLEU | NIST | MET | R-L | CIDEr | BLEU | | | MET | | | TER ↓ | | | BLEU | MET | TER ↓ | Mover | BERT | BLEURT |
| | | | | | | S | U | A | S | U | A | S | U | A | | | | | | |
| GPT-2 _{MEDIUM} | | | | | | | | | | | | | | | | | | | | |
| FINE-TUNE | 68.2 | 8.62 | 46.2 | 71.0 | 2.47 | 64.2 | 27.7 | 46.5 | 0.45 | 0.30 | 0.38 | 0.33 | 0.76 | 0.53 | 46.2 | 0.39 | 0.46 | 0.50 | 0.94 | 0.39 |
| FT-TOP2 | 68.1 | 8.59 | 46.0 | 70.8 | 2.41 | 53.6 | 18.9 | 36.0 | 0.38 | 0.23 | 0.31 | 0.49 | 0.99 | 0.72 | 41.0 | 0.34 | 0.56 | 0.43 | 0.93 | 0.21 |
| ADAPTER(3%) | 68.9 | 8.71 | 46.1 | 71.3 | 2.47 | 60.4 | 48.3 | 54.9 | 0.43 | 0.38 | 0.41 | 0.35 | 0.45 | 0.39 | 45.2 | 0.38 | 0.46 | 0.50 | 0.94 | 0.39 |
| ADAPTER(0.1%) | 66.3 | 8.41 | 45.0 | 69.8 | 2.40 | 54.5 | 45.1 | 50.2 | 0.39 | 0.36 | 0.38 | 0.40 | 0.46 | 0.43 | 42.4 | 0.36 | 0.48 | 0.47 | 0.94 | 0.33 |
| PREFIX(0.1%) | 69.7 | 8.81 | 46.1 | 71.4 | 2.49 | 62.9 | 45.6 | 55.1 | 0.44 | 0.38 | 0.41 | 0.35 | 0.49 | 0.41 | 46.4 | 0.38 | 0.46 | 0.50 | 0.94 | 0.39 |
| GPT-2 _{LARGE} | | | | | | | | | | | | | | | | | | | | |
| FINE-TUNE | 68.5 | 8.78 | 46.0 | 69.9 | 2.45 | 65.3 | 43.1 | 55.5 | 0.46 | 0.38 | 0.42 | 0.33 | 0.53 | 0.42 | 47.0 | 0.39 | 0.46 | 0.51 | 0.94 | 0.40 |
| Prefix | 70.3 | 8.85 | 46.2 | 71.7 | 2.47 | 63.4 | 47.7 | 56.3 | 0.45 | 0.39 | 0.42 | 0.34 | 0.48 | 0.40 | 46.7 | 0.39 | 0.45 | 0.51 | 0.94 | 0.40 |
| SOTA | 68.6 | 8.70 | 45.3 | 70.8 | 2.37 | 63.9 | 52.8 | 57.1 | 0.46 | 0.41 | 0.44 | - | - | - | - | - | - | - | - | - |

Table 1: Metrics (higher is better, except for TER) for table-to-text generation on E2E (left), WebNLG (middle) and DART (right). With only 0.1% parameters, Prefix-tuning outperforms other lightweight baselines and achieves a comparable performance with fine-tuning. The best score is boldfaced for both GPT-2_{MEDIUM} and GPT-2_{LARGE}.

Prefix Tuning: Variants

- ▶ Prefix: $[PREFIX; x, y]$
- ▶ Infix: $[x; INFIX; y]$
- ▶ Embedding only: only have the tuned weights on the first layer, not every layer
 - ▶ Embedding only tuning is the worse then either Infix or Prefix tuning
 - ▶ Infix tuning is better than embedding only, but it is inferior to prefix tuning
 - ▶ Is this because the prefix makes x have the prefix in it's representation also, so this representation is "stronger" before it gets to y and gets stronger still when it gets to y ?

| | E2E | | | | |
|------------------------------------|------|------|------|-------|-------|
| | BLEU | NIST | MET | ROUGE | CIDEr |
| PREFIX | 69.7 | 8.81 | 46.1 | 71.4 | 2.49 |
| Embedding-only: EMB-{PrefixLength} | | | | | |
| EMB-1 | 48.1 | 3.33 | 32.1 | 60.2 | 1.10 |
| EMB-10 | 62.2 | 6.70 | 38.6 | 66.4 | 1.75 |
| EMB-20 | 61.9 | 7.11 | 39.3 | 65.6 | 1.85 |
| Infix-tuning: INFIX-{PrefixLength} | | | | | |
| INFIX-1 | 67.9 | 8.63 | 45.8 | 69.4 | 2.42 |
| INFIX-10 | 67.2 | 8.48 | 45.8 | 69.9 | 2.40 |
| INFIX-20 | 66.7 | 8.47 | 45.8 | 70.0 | 2.42 |

Table 4: Intrinsic evaluation of Embedding-only (§7.2) and Infixing (§7.3). Both Embedding-only ablation and Infix-tuning underperforms full prefix-tuning.

Outline

- ▶ Adapters
- ▶ Prefix-Tuning
- ▶ LORA

LORA

- ▶ Both adapters and prefix-tuning have some issues
 - ▶ Adapters introduce high inference latency
 - ▶ Prefix-tuning is difficult to optimize, and performance changes non-monotonically in trainable parameters
 - ▶ Also, reserving a part of the sequence length for adaptation necessarily reduces the sequence length available to process a downstream task; might this lead to degradation in performance?

LORA

- ▶ LORA: For each matrix W in a transformer, for example a key matrix W_k , set $W_k = W_k + BA$ where B, A have low rank $r \ll \min\{d_k, d_{model}\}$ and $B \in \mathbb{R}^{d_k \times r}$, $A \in \mathbb{R}^{r \times d_{model}}$
- ▶ LORA: Low-Rank Adaptation
- ▶ Initialize B to all zeros and A to Gaussian noise so $BA = 0$ at the start but gradient updates are not trivially "stuck"

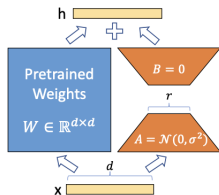


Figure 1: Our reparametrization. We only train A and B .

Benefits of LORA

- ▶ Why would LORA be a good idea?
 - ▶ Generalized full fine-tuning - if we select r to the right choice and have a LORA set B, A for each matrix, we basically have fine tuning
 - ▶ Adapters converge to an MLP
 - ▶ Prefix models cannot handle longer sequences
 - ▶ Benefit of LORA: **no additional layers, so no additional inference latency**
 - ▶ Another benefit: at test time, we can use $W = W + BA$, just a new matrix
 - ▶ Yet another: we can swap in different BA for different tasks, so it's easy to "adapt" the model
- ▶ For their paper, the authors focus on just the Self Attention parameters like W_k or W_q ; they freeze all the other parameters

LoRA vs Adapters

- For online short sequences, Adapters can be slow because they effectively make layers deeper

| 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 B](#).

LORA on ROBERTA (NLU)

- For Text Classification and NLU (understanding) tasks, LORA was generally superior

| 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 \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 _{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 \pm 2.2 | 90.3 \pm 1 | 85.4 |
| RoB _{base} (LoRA) | 0.3M | 87.5 \pm 3 | 95.1\pm2 | 89.7 \pm 7 | 63.4 \pm 1.2 | 93.3\pm3 | 90.8 \pm 1 | 86.6\pm7 | 91.5\pm2 | 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\pm2 | 96.2 \pm 5 | 90.9\pm1.2 | 68.2\pm1.9 | 94.9\pm3 | 91.6 \pm 1 | 87.4\pm2.5 | 92.6\pm2 | 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\pm2 | 91.9\pm1 | 83.8 \pm 2.9 | 92.1 \pm 7 | 88.4 |
| RoB _{large} (Adpt ^P)† | 0.8M | 90.5\pm3 | 96.6\pm2 | 89.7 \pm 1.2 | 67.8 \pm 2.5 | 94.8\pm3 | 91.7 \pm 2 | 80.1 \pm 2.9 | 91.9 \pm 4 | 87.9 |
| RoB _{large} (Adpt ^H)† | 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 _{large} (Adpt ^H)† | 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 _{large} (LoRA)† | 0.8M | 90.6\pm2 | 96.2 \pm 5 | 90.2\pm1.0 | 68.2 \pm 1.9 | 94.8\pm3 | 91.6 \pm 2 | 85.2\pm1.1 | 92.3\pm5 | 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\pm2 | 96.9 \pm 2 | 92.6\pm6 | 72.4\pm1.1 | 96.0\pm1 | 92.9\pm1 | 94.9\pm4 | 93.0\pm2 | 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.

LORA on GPT2 (NLG)

- For GPT2, LORA is generally superior so it works well for text generation

| 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 ^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.

LoRA on various number of parameters

- ▶ The authors observe a significant performance drop when we use more than 256 special tokens for prefix-embedding tuning or more than 32 special tokens for prefix-layer tuning

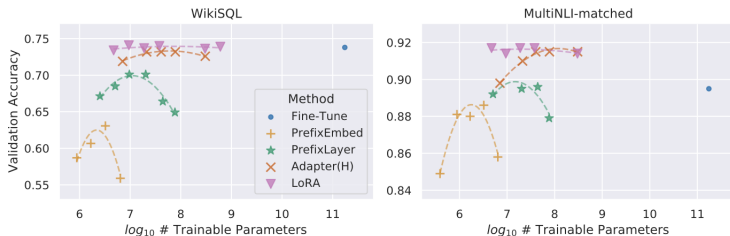


Figure 2: GPT-3 175B validation accuracy vs. number of trainable parameters of several adaptation methods on WikiSQL and MNLI-matched. LoRA exhibits better scalability and task performance. See [Section F.2](#) for more details on the plotted data points.

LORA and GPT3

► LORA scales well on larger models

| Model&Method | # Trainable Parameters | WikiSQL | MNLI-m | SAMSum |
|-------------------------------|------------------------|-------------|-------------|-----------------------|
| | | Acc. (%) | Acc. (%) | 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 $\pm 0.5\%$, MNLI-m around $\pm 0.1\%$, and SAMSum around $\pm 0.2/\pm 0.2/\pm 0.1$ for the three metrics.

LoRA: Which matrices?

- ▶ Given a parameter budget constraint, which subset of weight matrices in a pre-trained Transformer should we adapt to maximize downstream performance?
- ▶ Fix the number of trainable parameters: who should get LoRA matrices B , A ?
- ▶ From their study: it seems like the ideal matrices to apply this method to is a combination of W_k and W_v , not just one matrix alone

| | # 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 | 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.

LORA: Which r to use?

- Is the “optimal” adaptation matrix BA really rank- deficient?
If so, what is a good rank to use in practice?
- Even a very small rank does well: it is more of a choice of which matrices to apply the method to

| | Weight Type | $r = 1$ | $r = 2$ | $r = 4$ | $r = 8$ | $r = 64$ |
|--------------------------|----------------------|---------|---------|---------|---------|----------|
| WikiSQL($\pm 0.5\%$) | W_q | 68.8 | 69.6 | 70.5 | 70.4 | 70.0 |
| | W_q, W_v | 73.4 | 73.3 | 73.7 | 73.8 | 73.5 |
| | W_q, W_k, W_v, W_o | 74.1 | 73.7 | 74.0 | 74.0 | 73.9 |
| MultiNLI ($\pm 0.1\%$) | W_q | 90.7 | 90.9 | 91.1 | 90.7 | 90.7 |
| | W_q, W_v | 91.3 | 91.4 | 91.3 | 91.6 | 91.4 |
| | W_q, W_k, W_v, W_o | 91.2 | 91.7 | 91.7 | 91.5 | 91.4 |

Table 6: Validation accuracy on WikiSQL and MultiNLI with different rank r . To our surprise, a rank as small as one suffices for adapting both W_q and W_v on these datasets while training W_q alone needs a larger r . We conduct a similar experiment on GPT-2 in [Section H.2](#).

LORA: $BA = \Delta W$ vs W

- ▶ What is the connection between BA and W ? Does BA highly correlate with W ? How large is BA compared to W ?
- ▶ Core idea: use the Frobenius norm of a matrix
$$\|M\| = \sqrt{\sum_i \sigma_i^2}$$
- ▶ SVD: For any M , we have $M = U\Sigma V$
- ▶ We can use U , V to project a matrix into another matrix's space
- ▶ We can use U , V to understand "correlation" between matrices

LORA: $BA = \Delta W$ vs W

- ▶ First, BA has a stronger correlation with W compared to a random matrix, indicating that BA amplifies some features that are already in W
- ▶ Second, instead of repeating the top singular directions of W , BA only amplifies directions that are not emphasized in W
- ▶ Third, the amplification factor is rather huge: $21.5 \sim 6.91/0.32$ for $r = 4$
- ▶ Hypothesis: **LORA matrix potentially amplifies the important features for specific downstream tasks that were learned but not emphasized in the general pre-training model**

| | $r = 4$ | | | $r = 64$ | | |
|-----------------------------|---------------------------|-------|--------|---------------------------|-------|--------|
| | ΔW_q | W_q | Random | ΔW_q | W_q | Random |
| $\ U^\top W_q V^\top\ _F =$ | 0.32 | 21.67 | 0.02 | 1.90 | 37.71 | 0.33 |
| $\ W_q\ _F = 61.95$ | $\ \Delta W_q\ _F = 6.91$ | | | $\ \Delta W_q\ _F = 3.57$ | | |

Table 7: The Frobenius norm of $U^\top W_q V^\top$ where U and V are the left/right top r singular vector directions of either (1) ΔW_q , (2) W_q , or (3) a random matrix. The weight matrices are taken from the 48th layer of GPT-3.

References

- ▶ LORA
- ▶ Adapters
- ▶ Prefix Tuning