

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 <sup>L</sup>	$1482.0 \pm 1.0$ (+2.2%)	$354.8 \pm 0.5$ (+5.0%)	$23.9 \pm 2.1$ (+20.7%)
Adapter <sup>H</sup>	$1492.2 \pm 1.0$ (+3.0%)	$366.3 \pm 0.5$ (+8.4%)	$25.8 \pm 2.2$ (+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<sup>L</sup> and Adapter<sup>H</sup> 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.

## 4 OUR METHOD

We describe the simple design of LoRA and its practical benefits. The principles outlined here apply to any dense layers in deep learning models, though we only focus on certain weights in Transformer language models in our experiments as the motivating use case.

### 4.1 LOW-RANK-PARAMETRIZED UPDATE MATRICES

A neural network contains many dense layers which perform matrix multiplication. The weight matrices in these layers typically have full-rank. When adapting to a specific task, Aghajanyan et al. (2020) shows that the pre-trained language models have a low “intrinsic dimension” and can still learn efficiently despite a random projection to a smaller subspace. Inspired by this, we hypothesize the updates to the weights also have a low “intrinsic rank” during adaptation. For a pre-trained weight matrix  $W_0 \in \mathbb{R}^{d \times k}$ , we constrain its update by representing the latter with a low-rank decomposition  $W_0 + \Delta W = W_0 + BA$ , where  $B \in \mathbb{R}^{d \times r}$ ,  $A \in \mathbb{R}^{r \times k}$ , and the rank  $r \ll \min(d, k)$ . During training,  $W_0$  is frozen and does not receive gradient updates, while  $A$  and  $B$  contain trainable parameters. Note both  $W_0$  and  $\Delta W = BA$  are multiplied with the same input, and their respective output vectors are summed coordinate-wise. For  $h = W_0x$ , our modified forward pass yields:

$$h = W_0x + \Delta Wx = W_0x + BAx \quad (3)$$

We illustrate our reparametrization in Figure 1. We use a random Gaussian initialization for  $A$  and zero for  $B$ , so  $\Delta W = BA$  is zero at the beginning of training. We then scale  $\Delta Wx$  by  $\frac{\alpha}{r}$ , where  $\alpha$  is a constant in  $r$ . When optimizing with Adam, tuning  $\alpha$  is roughly the same as tuning the learning rate if we scale the initialization appropriately. As a result, we simply set  $\alpha$  to the first  $r$  we try and do not tune it. This scaling helps to reduce the need to retune hyperparameters when we vary  $r$  (Yang & Hu, 2021).

**A Generalization of Full Fine-tuning.** A more general form of fine-tuning allows the training of a subset of the pre-trained parameters. LoRA takes a step further and does not require the accumulated gradient update to weight matrices to have full-rank during adaptation. This means that when applying LoRA to all weight matrices and training all biases<sup>2</sup>, we roughly recover the expressiveness of full fine-tuning by setting the LoRA rank  $r$  to the rank of the pre-trained weight matrices. In other words, as we increase the number of trainable parameters<sup>3</sup>, training LoRA roughly converges to training the original model, while adapter-based methods converges to an MLP and prefix-based methods to a model that cannot take long input sequences.

**No Additional Inference Latency.** When deployed in production, we can explicitly compute and store  $W = W_0 + BA$  and perform inference as usual. Note that both  $W_0$  and  $BA$  are in  $\mathbb{R}^{d \times k}$ . When we need to switch to another downstream task, we can recover  $W_0$  by subtracting  $BA$  and then adding a different  $B'A'$ , a quick operation with very little memory overhead. Critically, this

<sup>2</sup>They represent a negligible number of parameters compared to weights.

<sup>3</sup>An inevitability when adapting to hard tasks.

---

guarantees that we do not introduce any additional latency during inference compared to a fine-tuned model by construction.

#### 4.2 APPLYING LORA TO TRANSFORMER

In principle, we can apply LoRA to any subset of weight matrices in a neural network to reduce the number of trainable parameters. In the Transformer architecture, there are four weight matrices in the self-attention module ( $W_q, W_k, W_v, W_o$ ) and two in the MLP module. We treat  $W_q$  (or  $W_k, W_v$ ) as a single matrix of dimension  $d_{model} \times d_{model}$ , even though the output dimension is usually sliced into attention heads. We limit our study to **only adapting the attention weights** for downstream tasks and freeze the MLP modules (so they are not trained in downstream tasks) both for simplicity and parameter-efficiency. We further study the effect on adapting different types of attention weight matrices in a Transformer in Section 7.1. We leave the empirical investigation of adapting the MLP layers, LayerNorm layers, and biases to a future work.

**Practical Benefits and Limitations.** The most significant benefit comes from the reduction in memory and storage usage. For a large Transformer trained with Adam, we reduce that VRAM usage by up to 2/3 if  $r \ll d_{model}$  as we do not need to store the optimizer states for the frozen parameters. On GPT-3 175B, we reduce the VRAM consumption during training from 1.2TB to 350GB. With  $r = 4$  and only the query and value projection matrices being adapted, the checkpoint size is reduced by roughly  $10,000 \times$  (from 350GB to 35MB)<sup>4</sup>. This allows us to train with significantly fewer GPUs and avoid I/O bottlenecks. Another benefit is that we can switch between tasks while deployed at a much lower cost by only swapping the LoRA weights as opposed to all the parameters. This allows for the creation of many customized models that can be swapped in and out on the fly on machines that store the pre-trained weights in VRAM. We also observe a 25% speedup during training on GPT-3 175B compared to full fine-tuning<sup>5</sup> as we do not need to calculate the gradient for the vast majority of the parameters.

LoRA also has its limitations. For example, it is not straightforward to batch inputs to different tasks with different  $A$  and  $B$  in a single forward pass, if one chooses to absorb  $A$  and  $B$  into  $W$  to eliminate additional inference latency. Though it is possible to not merge the weights and dynamically choose the LoRA modules to use for samples in a batch for scenarios where latency is not critical.

### 5 EMPIRICAL EXPERIMENTS

We evaluate the downstream task performance of LoRA on RoBERTa (Liu et al., 2019), DeBERTa (He et al., 2021), and GPT-2 (Radford et al., b), before scaling up to GPT-3 175B (Brown et al., 2020). Our experiments cover a wide range of tasks, from natural language understanding (NLU) to generation (NLG). Specifically, we evaluate on the GLUE (Wang et al., 2019) benchmark for RoBERTa and DeBERTa. We follow the setup of Li & Liang (2021) on GPT-2 for a direct comparison and add WikiSQL (Zhong et al., 2017) (NL to SQL queries) and SAMSum (Gliwa et al., 2019) (conversation summarization) for large-scale experiments on GPT-3. See Appendix C for more details on the datasets we use. We use NVIDIA Tesla V100 for all experiments.

#### 5.1 BASELINES

To compare with other baselines broadly, we replicate the setups used by prior work and reuse their reported numbers whenever possible. This, however, means that some baselines might only appear in certain experiments.

**Fine-Tuning (FT)** is a common approach for adaptation. During fine-tuning, the model is initialized to the pre-trained weights and biases, and all model parameters undergo gradient updates. A simple variant is to update only some layers while freezing others. We include one such baseline reported in prior work (Li & Liang, 2021) on GPT-2, which adapts just the last two layers (**FT<sup>Top2</sup>**).

---

<sup>4</sup>We still need the 350GB model during deployment; however, storing 100 adapted models only requires  $350\text{GB} + 35\text{MB} * 100 \approx 354\text{GB}$  as opposed to  $100 * 350\text{GB} \approx 35\text{TB}$ .

<sup>5</sup>For GPT-3 175B, the training throughput for full fine-tuning is 32.5 tokens/s per V100 GPU; with the same number of weight shards for model parallelism, the throughput is 43.1 tokens/s per V100 GPU for LoRA.

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 <sub>±.0</sub>	94.2 <sub>±.1</sub>	88.5 <sub>±1.1</sub>	60.8 <sub>±.4</sub>	93.1 <sub>±.1</sub>	90.2 <sub>±.0</sub>	71.5 <sub>±2.7</sub>	89.7 <sub>±.3</sub>	84.4
RoB <sub>base</sub> (Adpt <sup>D</sup> )*	0.9M	87.3 <sub>±.1</sub>	94.7 <sub>±.3</sub>	88.4 <sub>±.1</sub>	62.6 <sub>±.9</sub>	93.0 <sub>±.2</sub>	90.6 <sub>±.0</sub>	75.9 <sub>±2.2</sub>	90.3 <sub>±.1</sub>	85.4
RoB <sub>base</sub> (LoRA)	0.3M	87.5 <sub>±.3</sub>	<b>95.1</b> <sub>±.2</sub>	89.7 <sub>±.7</sub>	63.4 <sub>±1.2</sub>	<b>93.3</b> <sub>±.3</sub>	90.8 <sub>±.1</sub>	<b>86.6</b> <sub>±.7</sub>	<b>91.5</b> <sub>±.2</sub>	<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</b> <sub>±.2</sub>	96.2 <sub>±.5</sub>	<b>90.9</b> <sub>±1.2</sub>	<b>68.2</b> <sub>±1.9</sub>	<b>94.9</b> <sub>±.3</sub>	91.6 <sub>±.1</sub>	<b>87.4</b> <sub>±2.5</sub>	<b>92.6</b> <sub>±.2</sub>	<b>89.0</b>
RoB <sub>large</sub> (Adpt <sup>P</sup> )†	3.0M	90.2 <sub>±.3</sub>	96.1 <sub>±.3</sub>	90.2 <sub>±.7</sub>	<b>68.3</b> <sub>±1.0</sub>	<b>94.8</b> <sub>±.2</sub>	<b>91.9</b> <sub>±.1</sub>	83.8 <sub>±2.9</sub>	92.1 <sub>±.7</sub>	88.4
RoB <sub>large</sub> (Adpt <sup>P</sup> )†	0.8M	<b>90.5</b> <sub>±.3</sub>	<b>96.6</b> <sub>±.2</sub>	89.7 <sub>±1.2</sub>	67.8 <sub>±2.5</sub>	<b>94.8</b> <sub>±.3</sub>	91.7 <sub>±.2</sub>	80.1 <sub>±2.9</sub>	91.9 <sub>±.4</sub>	87.9
RoB <sub>large</sub> (Adpt <sup>H</sup> )†	6.0M	89.9 <sub>±.5</sub>	96.2 <sub>±.3</sub>	88.7 <sub>±2.9</sub>	66.5 <sub>±4.4</sub>	94.7 <sub>±.2</sub>	92.1 <sub>±.1</sub>	83.4 <sub>±1.1</sub>	91.0 <sub>±1.7</sub>	87.8
RoB <sub>large</sub> (Adpt <sup>H</sup> )†	0.8M	90.3 <sub>±.3</sub>	96.3 <sub>±.5</sub>	87.7 <sub>±1.7</sub>	66.3 <sub>±2.0</sub>	94.7 <sub>±.2</sub>	91.5 <sub>±.1</sub>	72.9 <sub>±2.9</sub>	91.5 <sub>±.5</sub>	86.4
RoB <sub>large</sub> (LoRA)†	0.8M	<b>90.6</b> <sub>±.2</sub>	96.2 <sub>±.5</sub>	<b>90.2</b> <sub>±1.0</sub>	68.2 <sub>±1.9</sub>	<b>94.8</b> <sub>±.3</sub>	91.6 <sub>±.2</sub>	<b>85.2</b> <sub>±1.1</sub>	<b>92.3</b> <sub>±.5</sub>	<b>88.6</b>
DeBERT <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
DeBERT <sub>XXL</sub> (LoRA)	4.7M	<b>91.9</b> <sub>±.2</sub>	96.9 <sub>±.2</sub>	<b>92.6</b> <sub>±.6</sub>	<b>72.4</b> <sub>±1.1</sub>	<b>96.0</b> <sub>±.1</sub>	<b>92.9</b> <sub>±.1</sub>	<b>94.9</b> <sub>±.4</sub>	<b>93.0</b> <sub>±.2</sub>	<b>91.3</b>

Table 2: RoBERTa<sub>base</sub>, RoBERTa<sub>large</sub>, and DeBERTa<sub>XXL</sub> 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.

**Bias-only or BitFit** is a baseline where we only train the bias vectors while freezing everything else. Contemporarily, this baseline has also been studied by BitFit (Zaken et al., 2021).

**Prefix-embedding tuning (PreEmbed)** inserts special tokens among the input tokens. These special tokens have trainable word embeddings and are generally not in the model’s vocabulary. Where to place such tokens can have an impact on performance. We focus on “prefixing”, which prepends such tokens to the prompt, and “infixing”, which appends to the prompt; both are discussed in Li & Liang (2021). We use  $l_p$  (resp.  $l_i$ ) denote the number of prefix (resp. infix) tokens. The number of trainable parameters is  $|\Theta| = d_{model} \times (l_p + l_i)$ .

**Prefix-layer tuning (PreLayer)** is an extension to prefix-embedding tuning. Instead of just learning the word embeddings (or equivalently, the activations after the embedding layer) for some special tokens, we learn the activations after every Transformer layer. The activations computed from previous layers are simply replaced by trainable ones. The resulting number of trainable parameters is  $|\Theta| = L \times d_{model} \times (l_p + l_i)$ , where  $L$  is the number of Transformer layers.

**Adapter tuning** as proposed in Houlsby et al. (2019) inserts adapter layers between the self-attention module (and the MLP module) and the subsequent residual connection. There are two fully connected layers with biases in an adapter layer with a nonlinearity in between. We call this original design **Adapter<sup>H</sup>**. Recently, Lin et al. (2020) proposed a more efficient design with the adapter layer applied only after the MLP module and after a LayerNorm. We call it **Adapter<sup>L</sup>**. This is very similar to another design proposed in Pfeiffer et al. (2021), which we call **Adapter<sup>P</sup>**. We also include another baseline called AdapterDrop (Rücklé et al., 2020) which drops some adapter layers for greater efficiency (**Adapter<sup>D</sup>**). We cite numbers from prior works whenever possible to maximize the number of baselines we compare with; they are in rows with an asterisk (\*) in the first column. In all cases, we have  $|\Theta| = \hat{L}_{Adpt} \times (2 \times d_{model} \times r + r + d_{model}) + 2 \times \hat{L}_{LN} \times d_{model}$  where  $\hat{L}_{Adpt}$  is the number of adapter layers and  $\hat{L}_{LN}$  the number of trainable LayerNorms (e.g., in Adapter<sup>L</sup>).

**LoRA** adds trainable pairs of rank decomposition matrices in parallel to existing weight matrices. As mentioned in Section 4.2, we only apply LoRA to  $W_q$  and  $W_v$  in most experiments for simplicity. The number of trainable parameters is determined by the rank  $r$  and the shape of the original weights:  $|\Theta| = 2 \times \hat{L}_{LoRA} \times d_{model} \times r$ , where  $\hat{L}_{LoRA}$  is the number of weight matrices we apply LoRA to.