| Method | Dataset | MNLI | SST-2 | MRPC | CoLA | QNLI | QQP | RTE | STS-B |
|---|---|---|---|---|---|---|---|---|---|
| | Optimizer | | | | AdamW | | | | |
| | Warmup Ratio | | | | 0.06 | | | | |
| | LR Schedule | | | | Linear | | | | |
| RoBERTa base LoRA | Batch Size | 16 | 16 | 16 | 32 | 32 | 16 | 32 | 16 |
| | # Epochs | 30 | 60 | 30 | 80 | 25 | 25 | 80 | 40 |
| | Learning Rate | 5E-04 | 5E-04 | 4E-04 | 4E-04 | 4E-04 | 5E-04 | 5E-04 | 4E-04 |
| | LoRA Config. | | | | $r_q = r_v = 8$ | | | | |
| | LoRA $\alpha$ | | | | 8 | | | | |
| | Max Seq. Len. | | | | 512 | | | | |
| RoBERTa large LoRA | Batch Size | 4 | 4 | 4 | 4 | 4 | 4 | 8 | 8 |
| | # Epochs | 10 | 10 | 20 | 20 | 10 | 20 | 20 | 30 |
| | Learning Rate | 3E-04 | 4E-04 | 3E-04 | 2E-04 | 2E-04 | 3E-04 | 4E-04 | 2E-04 |
| | LoRA Config. | | | | $r_q = r_v = 8$ | | | | |
| | LoRA $\alpha$ | | | | 16 | | | | |
| | Max Seq. Len. | 128 | 128 | 512 | 128 | 512 | 512 | 512 | 512 |
| RoBERTa large LoRA† | Batch Size | | | | 4 | | | | |
| | # Epochs | 10 | 10 | 20 | 20 | 10 | 20 | 20 | 10 |
| | Learning Rate | 3E-04 | 4E-04 | 3E-04 | 2E-04 | 2E-04 | 3E-04 | 4E-04 | 2E-04 |
| | LoRA Config. | | | | $r_q = r_v = 8$ | | | | |
| | LoRA $\alpha$ | | | | 16 | | | | |
| | Max Seq. Len. | | | | 128 | | | | |
| RoBERTa large Adpt$^P$ (3M)† | Batch Size | | | | 32 | | | | |
| | # Epochs | 10 | 20 | 20 | 20 | 10 | 20 | 20 | 20 |
| | Learning Rate | 3E-05 | 3E-05 | 3E-04 | 3E-04 | 3E-04 | 3E-04 | 3E-04 | 3E-04 |
| | Bottleneck $r$ | | | | 64 | | | | |
| | Max Seq. Len. | | | | 128 | | | | |
| RoBERTa large Adpt$^P$ (0.8M)† | Batch Size | | | | 32 | | | | |
| | # Epochs | 5 | 20 | 20 | 20 | 10 | 20 | 20 | 20 |
| | Learning Rate | 3E-04 | 3E-04 | 3E-04 | 3E-04 | 3E-04 | 3E-04 | 3E-04 | 3E-04 |
| | Bottleneck $r$ | | | | 16 | | | | |
| | Max Seq. Len. | | | | 128 | | | | |
| RoBERTa large Adpt$^H$ (6M)† | Batch Size | | | | 32 | | | | |
| | # Epochs | 10 | 5 | 10 | 10 | 5 | 20 | 20 | 10 |
| | Learning Rate | 3E-05 | 3E-04 | 3E-04 | 3E-04 | 3E-04 | 3E-04 | 3E-04 | 3E-04 |
| | Bottleneck $r$ | | | | 64 | | | | |
| | Max Seq. Len. | | | | 128 | | | | |
| RoBERTa large Adpt$^H$ (0.8M)† | Batch Size | | | | 32 | | | | |
| | # Epochs | 10 | 5 | 10 | 10 | 5 | 20 | 20 | 10 |
| | Learning Rate | 3E-04 | 3E-04 | 3E-04 | 3E-04 | 3E-04 | 3E-04 | 3E-04 | 3E-04 |
| | Bottleneck $r$ | | | | 8 | | | | |
| | Max Seq. Len. | | | | 128 | | | | |

Table 9: The hyperparameters we used for RoBERTa on the GLUE benchmark.

### D.3 GPT-2

We train all of our GPT-2 models using AdamW (Loshchilov & Hutter, 2017) with a linear learning rate schedule for 5 epochs. We use the batch size, learning rate, and beam search beam size described in Li & Liang (2021). Accordingly, we also tune the above hyperparameters for LoRA. We report the mean over 3 random seeds; the result for each run is taken from the best epoch. The hyperparameters used for LoRA in GPT-2 are listed in Table 11. For those used for other baselines, see Li & Liang (2021).

### D.4 GPT-3

For all GPT-3 experiments, we train using AdamW (Loshchilov & Hutter, 2017) for 2 epochs with a batch size of 128 samples and a weight decay factor of 0.1. We use a sequence length of 384 for

| Method | Dataset | MNLI | SST-2 | MRPC | CoLA | QNLI | QQP | RTE | STS-B |
|---|---|---|---|---|---|---|---|---|---|
| | Optimizer | | | | AdamW | | | | |
| | Warmup Ratio | | | | 0.1 | | | | |
| | LR Schedule | | | | Linear | | | | |
| DeBERTa XXL LoRA | Batch Size | 8 | 8 | 32 | 4 | 6 | 8 | 4 | 4 |
| | # Epochs | 5 | 16 | 30 | 10 | 8 | 11 | 11 | 10 |
| | Learning Rate | 1E-04 | 6E-05 | 2E-04 | 1E-04 | 1E-04 | 1E-04 | 2E-04 | 2E-04 |
| | Weight Decay | 0 | 0.01 | 0.01 | 0 | 0.01 | 0.01 | 0.01 | 0.1 |
| | CLS Dropout | 0.15 | 0 | 0 | 0.1 | 0.1 | 0.2 | 0.2 | 0.2 |
| | LoRA Config. | | | | $r_q = r_v = 8$ | | | | |
| | LoRA $\alpha$ | | | | 8 | | | | |
| | Max Seq. Len. | 256 | 128 | 128 | 64 | 512 | 320 | 320 | 128 |

Table 10: The hyperparameters for DeBERTa XXL on tasks included in the GLUE benchmark.

| Dataset | E2E | WebNLG | DART |
|---|---|---|---|
| | | Training | |
| Optimizer | | AdamW | |
| Weight Decay | 0.01 | 0.01 | 0.0 |
| Dropout Prob | 0.1 | 0.1 | 0.0 |
| Batch Size | | 8 | |
| # Epoch | | 5 | |
| Warmup Steps | | 500 | |
| Learning Rate Schedule | | Linear | |
| Label Smooth | 0.1 | 0.1 | 0.0 |
| Learning Rate | | 0.0002 | |
| Adaptation | | $r_q = r_v = 4$ | |
| LoRA $\alpha$ | | 32 | |
| | | Inference | |
| Beam Size | | 10 | |
| Length Penalty | 0.9 | 0.8 | 0.8 |
| no repeat ngram size | | 4 | |

Table 11: The hyperparameters for GPT-2 LoRA on E2E, WebNLG and DART.

WikiSQL (Zhong et al., 2017), 768 for MNLI (Williams et al., 2018), and 2048 for SAMSum (Gliwa et al., 2019). We tune learning rate for all method-dataset combinations. See Section D.4 for more details on the hyperparameters used. For prefix-embedding tuning, we find the optimal $l_p$ and $l_i$ to be 256 and 8, respectively, totalling $3.2M$ trainable parameters. We use $l_p = 8$ and $l_i = 8$ for prefix-layer tuning with $20.2M$ trainable parameters to obtain the overall best performance. We present two parameter budgets for LoRA: 4.7M ($r_q = r_v = 1$ or $r_v = 2$) and 37.7M ($r_q = r_v = 8$ or $r_q = r_k = r_v = r_o = 2$). We report the best validation performance from each run. The training hyperparameters used in our GPT-3 experiments are listed in Table 12.

# E  COMBINING LORA WITH PREFIX TUNING

LoRA can be naturally combined with existing prefix-based approaches. In this section, we evaluate two combinations of LoRA and variants of prefix-tuning on WikiSQL and MNLI.

**LoRA+PrefixEmbed (LoRA+PE)** combines LoRA with prefix-embedding tuning, where we insert $l_p + l_i$ special tokens whose embeddings are treated as trainable parameters. For more on prefix-embedding tuning, see Section 5.1.

**LoRA+PrefixLayer (LoRA+PL)** combines LoRA with prefix-layer tuning. We also insert $l_p + l_i$ special tokens; however, instead of letting the hidden representations of these tokens evolve natu-

| Hyperparameters | Fine-Tune | PreEmbed | PreLayer | BitFit | Adapter[H] | LoRA |
|---|---|---|---|---|---|---|
| Optimizer | | | | AdamW | | |
| Batch Size | | | | 128 | | |
| # Epoch | | | | 2 | | |
| Warmup Tokens | | | | 250,000 | | |
| LR Schedule | | | | Linear | | |
| Learning Rate | 5.00E-06 | 5.00E-04 | 1.00E-04 | 1.6E-03 | 1.00E-04 | 2.00E-04 |

Table 12: The training hyperparameters used for different GPT-3 adaption methods. We use the same hyperparameters for all datasets after tuning learning rate.

rally, we replace them after every Transformer block with an input agnostic vector. Thus, both the embeddings and subsequent Transformer block activations are treated as trainable parameters. For more on prefix-layer tuning, see Section 5.1.

In Table 15, we show the evaluation results of LoRA+PE and LoRA+PL on WikiSQL and MultiNLI. First of all, LoRA+PE significantly outperforms both LoRA and prefix-embedding tuning on WikiSQL, which indicates that LoRA is somewhat orthogonal to prefix-embedding tuning. On MultiNLI, the combination of LoRA+PE doesn't perform better than LoRA, possibly because LoRA on its own already achieves performance comparable to the human baseline. Secondly, we notice that LoRA+PL performs slightly worse than LoRA even with more trainable parameters. We attribute this to the fact that prefix-layer tuning is very sensitive to the choice of learning rate and thus makes the optimization of LoRA weights more difficult in LoRA+PL.

# F   ADDITIONAL EMPIRICAL EXPERIMENTS

## F.1   ADDITIONAL EXPERIMENTS ON GPT-2

We also repeat our experiment on DART (Nan et al., 2020) and WebNLG (Gardent et al., 2017) following the setup of Li & Liang (2021). The result is shown in Table 13. Similar to our result on E2E NLG Challenge, reported in Section 5, LoRA performs better than or at least on-par with prefix-based approaches given the same number of trainable parameters.

| Method | # Trainable Parameters | DART | | |
|---|---|---|---|---|
| | | BLEU↑ | MET↑ | TER↓ |
| | | GPT-2 Medium | | |
| Fine-Tune | 354M | 46.2 | **0.39** | **0.46** |
| Adapter[L] | 0.37M | 42.4 | 0.36 | 0.48 |
| Adapter[L] | 11M | 45.2 | 0.38 | **0.46** |
| FT[Top2] | 24M | 41.0 | 0.34 | 0.56 |
| PrefLayer | 0.35M | 46.4 | 0.38 | **0.46** |
| LoRA | 0.35M | **47.1**$_{\pm.2}$ | **0.39** | **0.46** |
| | | GPT-2 Large | | |
| Fine-Tune | 774M | 47.0 | **0.39** | 0.46 |
| Adapter[L] | 0.88M | 45.7$_{\pm.1}$ | 0.38 | 0.46 |
| Adapter[L] | 23M | 47.1$_{\pm.1}$ | **0.39** | **0.45** |
| PrefLayer | 0.77M | 46.7 | 0.38 | **0.45** |
| LoRA | 0.77M | **47.5**$_{\pm.1}$ | **0.39** | **0.45** |

Table 13: GPT-2 with different adaptation methods on DART. The variances of MET and TER are less than 0.01 for all adaption approaches.