

Method	WebNLG								
	BLEU↑			MET↑			TER↓		
	U	S	A	U	S	A	U	S	A
GPT-2 Medium									
Fine-Tune (354M)	27.7	<b>64.2</b>	46.5	.30	<b>.45</b>	.38	.76	<b>.33</b>	.53
Adapter <sup>L</sup> (0.37M)	45.1	54.5	50.2	.36	.39	.38	.46	.40	.43
Adapter <sup>L</sup> (11M)	<b>48.3</b>	60.4	54.9	<b>.38</b>	.43	<b>.41</b>	<b>.45</b>	.35	<b>.39</b>
FT <sup>Top2</sup> (24M)	18.9	53.6	36.0	.23	.38	.31	.99	.49	.72
Prefix (0.35M)	45.6	62.9	55.1	<b>.38</b>	.44	<b>.41</b>	.49	.35	.40
LoRA (0.35M)	46.7 <sub>±.4</sub>	62.1 <sub>±.2</sub>	<b>55.3<sub>±.2</sub></b>	<b>.38</b>	.44	<b>.41</b>	.46	<b>.33</b>	<b>.39</b>
GPT-2 Large									
Fine-Tune (774M)	43.1	65.3	55.5	.38	<b>.46</b>	.42	.53	.33	.42
Adapter <sup>L</sup> (0.88M)	<b>49.8<sub>±.0</sub></b>	61.1 <sub>±.0</sub>	56.0 <sub>±.0</sub>	.38	.43	.41	<b>.44</b>	.35	.39
Adapter <sup>L</sup> (23M)	49.2 <sub>±.1</sub>	64.7 <sub>±.2</sub>	<b>57.7<sub>±.1</sub></b>	<b>.39</b>	<b>.46</b>	<b>.43</b>	.46	.33	.39
Prefix (0.77M)	47.7	63.4	56.3	<b>.39</b>	.45	.42	.48	.34	.40
LoRA (0.77M)	48.4 <sub>±.3</sub>	64.0 <sub>±.3</sub>	57.0 <sub>±.1</sub>	<b>.39</b>	.45	.42	.45	<b>.32</b>	<b>.38</b>

Table 14: GPT-2 with different adaptation methods on WebNLG. The variances of MET and TER are less than 0.01 for all the experiments we ran. “U” indicates unseen categories, “S” indicates seen categories, and “A” indicates all categories in the test set of WebNLG.

## F.2 ADDITIONAL EXPERIMENTS ON GPT-3

We present additional runs on GPT-3 with different adaptation methods in Table 15. The focus is on identifying the trade-off between performance and the number of trainable parameters.

## F.3 LOW-DATA REGIME

To evaluate the performance of different adaptation approaches in the low-data regime, we randomly sample 100, 1k and 10k training examples from the full training set of MNLI to form the low-data MNLI-*n* tasks. In Table 16, we show the performance of different adaptation approaches on MNLI-*n*. To our surprise, PrefixEmbed and PrefixLayer performs very poorly on MNLI-100 dataset, with PrefixEmbed performing only slightly better than random chance (37.6% vs. 33.3%). PrefixLayer performs better than PrefixEmbed but is still significantly worse than Fine-Tune or LoRA on MNLI-100. The gap between prefix-based approaches and LoRA/Fine-tuning becomes smaller as we increase the number of training examples, which might suggest that prefix-based approaches are not suitable for low-data tasks in GPT-3. LoRA achieves better performance than fine-tuning on both MNLI-100 and MNLI-Full, and comparable results on MNLI-1k and MNLI-10K considering the ( $\pm 0.3$ ) variance due to random seeds.

The training hyperparameters of different adaptation approaches on MNLI-*n* are reported in Table 17. We use a smaller learning rate for PrefixLayer on the MNLI-100 set, as the training loss does not decrease with a larger learning rate.

## G MEASURING SIMILARITY BETWEEN SUBSPACES

In this paper we use the measure  $\phi(A, B, i, j) = \psi(U_A^i, U_B^j) = \frac{\|U_A^{i\top} U_B^j\|_F^2}{\min\{i, j\}}$  to measure the subspace similarity between two column orthonormal matrices  $U_A^i \in \mathbb{R}^{d \times i}$  and  $U_B^j \in \mathbb{R}^{d \times j}$ , obtained by taking columns of the left singular matrices of  $A$  and  $B$ . We point out that this similarity is simply a reverse of the standard Projection Metric that measures distance between subspaces Ham & Lee (2008).

Method	Hyperparameters	# Trainable Parameters	WikiSQL	MNLI-m
Fine-Tune	-	175B	73.8	89.5
PrefixEmbed	$l_p = 32, l_i = 8$	0.4 M	55.9	84.9
	$l_p = 64, l_i = 8$	0.9 M	58.7	88.1
	$l_p = 128, l_i = 8$	1.7 M	60.6	88.0
	$l_p = 256, l_i = 8$	3.2 M	63.1	88.6
	$l_p = 512, l_i = 8$	6.4 M	55.9	85.8
PrefixLayer	$l_p = 2, l_i = 2$	5.1 M	68.5	89.2
	$l_p = 8, l_i = 0$	10.1 M	69.8	88.2
	$l_p = 8, l_i = 8$	20.2 M	70.1	89.5
	$l_p = 32, l_i = 4$	44.1 M	66.4	89.6
	$l_p = 64, l_i = 0$	76.1 M	64.9	87.9
Adapter <sup>H</sup>	$r = 1$	7.1 M	71.9	89.8
	$r = 4$	21.2 M	73.2	91.0
	$r = 8$	40.1 M	73.2	91.5
	$r = 16$	77.9 M	73.2	91.5
	$r = 64$	304.4 M	72.6	91.5
LoRA	$r_v = 2$	4.7 M	73.4	<b>91.7</b>
	$r_q = r_v = 1$	4.7 M	73.4	91.3
	$r_q = r_v = 2$	9.4 M	73.3	91.4
	$r_q = r_k = r_v = r_o = 1$	9.4 M	74.1	91.2
	$r_q = r_v = 4$	18.8 M	73.7	91.3
	$r_q = r_k = r_v = r_o = 2$	18.8 M	73.7	<b>91.7</b>
	$r_q = r_v = 8$	37.7 M	73.8	<b>91.6</b>
	$r_q = r_k = r_v = r_o = 4$	37.7 M	74.0	<b>91.7</b>
	$r_q = r_v = 64$	301.9 M	73.6	91.4
	$r_q = r_k = r_v = r_o = 64$	603.8 M	73.9	91.4
LoRA+PE	$r_q = r_v = 8, l_p = 8, l_i = 4$	37.8 M	75.0	91.4
	$r_q = r_v = 32, l_p = 8, l_i = 4$	151.1 M	<b>75.9</b>	91.1
	$r_q = r_v = 64, l_p = 8, l_i = 4$	302.1 M	<b>76.2</b>	91.3
LoRA+PL	$r_q = r_v = 8, l_p = 8, l_i = 4$	52.8 M	72.9	90.2

Table 15: Hyperparameter analysis of different adaptation approaches on WikiSQL and MNLI. Both prefix-embedding tuning (PrefixEmbed) and prefix-layer tuning (PrefixLayer) perform worse as we increase the number of trainable parameters, while LoRA’s performance stabilizes. Performance is measured in validation accuracy.

Method	MNLI(m)-100	MNLI(m)-1k	MNLI(m)-10k	MNLI(m)-392K
GPT-3 (Fine-Tune)	60.2	<b>85.8</b>	88.9	89.5
GPT-3 (PrefixEmbed)	37.6	75.2	79.5	88.6
GPT-3 (PrefixLayer)	48.3	82.5	85.9	89.6
GPT-3 (LoRA)	<b>63.8</b>	85.6	<b>89.2</b>	<b>91.7</b>

Table 16: Validation accuracy of different methods on subsets of MNLI using GPT-3 175B. MNLI- $n$  describes a subset with  $n$  training examples. We evaluate with the full validation set. LoRA performs exhibits favorable sample-efficiency compared to other methods, including fine-tuning.

To be concrete, let the singular values of  $U_A^{i\top} U_B^j$  to be  $\sigma_1, \sigma_2, \dots, \sigma_p$  where  $p = \min\{i, j\}$ . We know that the Projection Metric Ham & Lee (2008) is defined as:

$$d(U_A^i, U_B^j) = \sqrt{p - \sum_{i=1}^p \sigma_i^2} \in [0, \sqrt{p}]$$

Hyperparameters	Adaptation	MNLI-100	MNLI-1k	MNLI-10K	MNLI-392K
Optimizer	-			AdamW	
Warmup Tokens	-			250,000	
LR Schedule	-			Linear	
Batch Size	-	20	20	100	128
# Epoch	-	40	40	4	2
Learning Rate	FineTune			5.00E-6	
	PrefixEmbed	2.00E-04	2.00E-04	4.00E-04	5.00E-04
	PrefixLayer	5.00E-05	5.00E-05	5.00E-05	1.00E-04
	LoRA			2.00E-4	
Adaptation-Specific	PrefixEmbed $l_p$	16	32	64	256
	PrefixEmbed $l_i$			8	
	PrefixTune			$l_p = l_i = 8$	
	LoRA			$r_q = r_v = 8$	

Table 17: The hyperparameters used for different GPT-3 adaptation methods on MNLI(m)-n.

where our similarity is defined as:

$$\phi(A, B, i, j) = \psi(U_A^i, U_B^j) = \frac{\sum_{i=1}^p \sigma_i^2}{p} = \frac{1}{p} \left( 1 - d(U_A^i, U_B^j)^2 \right)$$

This similarity satisfies that if  $U_A^i$  and  $U_B^j$  share the same column span, then  $\phi(A, B, i, j) = 1$ . If they are completely orthogonal, then  $\phi(A, B, i, j) = 0$ . Otherwise,  $\phi(A, B, i, j) \in (0, 1)$ .

## H ADDITIONAL EXPERIMENTS ON LOW-RANK MATRICES

We present additional results from our investigation into the low-rank update matrices.

### H.1 CORRELATION BETWEEN LORA MODULES

See Figure 6 and Figure 7 for how the results presented in Figure 3 and Figure 4 generalize to other layers.

### H.2 EFFECT OF $r$ ON GPT-2

We repeat our experiment on the effect of  $r$  (Section 7.2) in GPT-2. Using the E2E NLG Challenge dataset as an example, we report the validation loss and test metrics achieved by different choices of  $r$  after training for 26,000 steps. We present our result in Table 18. The optimal rank for GPT-2 Medium is between 4 and 16 depending on the metric used, which is similar to that for GPT-3 175B. Note that the relationship between model size and the optimal rank for adaptation is still an open question.

### H.3 CORRELATION BETWEEN $W$ AND $\Delta W$

See Figure 8 for the normalized subspace similarity between  $W$  and  $\Delta W$  with varying  $r$ .

Note again that  $\Delta W$  does not contain the top singular directions of  $W$ , since the similarity between the top 4 directions in  $\Delta W$  and the top-10% of those in  $W$  barely exceeds 0.2. This gives evidence that  $\Delta W$  contains those “task-specific” directions that are otherwise *not* emphasized in  $W$ .

An interesting next question to answer, is how “strong” do we need to amplify those task-specific directions, in order for the model adaptation to work well?