

Revisiting Parameter-Efficient Tuning: Are We Really There Yet?

Anonymous ACL submission

Abstract

Parameter-efficient tuning (PETuning) methods have been deemed by many as the new paradigm for using pretrained language models (PLMs). By tuning just a fraction amount of parameters comparing to full model finetuning, PETuning methods claim to have achieved performance on par with or even better than finetuning. In this work, we take a step back and re-examine these PETuning methods by conducting the first comprehensive investigation into the training and evaluation of PETuning methods. We found the problematic validation and testing practice in current studies, when accompanied by the instability nature of PETuning methods, has led to unreliable conclusions. When being compared under a truly fair evaluation protocol, PETuning cannot yield consistently competitive performance while finetuning remains to be the best-performing method in medium- and high-resource settings. We delve deeper into the cause of the instability and observed that the number of trainable parameters and training iterations are two main factors: reducing trainable parameters and prolonging training iterations may lead to higher stability in PETuning methods.¹

1 Introduction

Pretrained Language Models (PLMs) such as BERT (Devlin et al., 2019) and RoBERTa (Liu et al., 2019) have orchestrated tremendous progress in NLP in the past few years, achieving state-of-the-art on a large variety of benchmarks such as GLUE (Wang et al., 2018) and SuperGLUE (Wang et al., 2019). Most successful applications of PLMs follow the pretraining-and-finetuning transfer learning paradigm (Devlin et al., 2019), where PLMs are used as backbone to be combined with additional parameters and finetuned on downstream tasks in an end-to-end manner. Whilst being simple and effective, such paradigm requires task-specific tuning

¹Our code is available at <https://github.com/AnonymousARRsubmission/PETuning>.

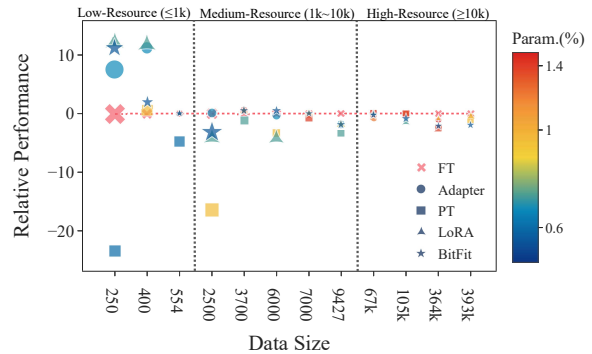


Figure 1: The relative performance difference of PETuning methods, i.e., Adapter, prefix tuning (PT), LoRA, BitFit, comparing with the full finetuning (FT) over different training data size of 12 tasks from GLUE and SuperGLUE. The tasks and their split into the three resource bands are illustrated in Appx. §B.1. The size of each point denotes the standard deviation and the colours of PETuning methods denote the percentage of trainable parameters over different tasks compared to full finetuning. The key takeaway message is that PETuning methods outperform finetuning only in the low-resource tasks but remain on par or behind in medium and high-resource settings.

of the full model that consists of hundreds of millions (Devlin et al., 2019; Liu et al., 2019), or even billions (Radford et al., 2019; Brown et al., 2020; Raffel et al., 2020) of parameters for each task, which is time-consuming and resource-intensive.

To avoid full model finetuning, there has been a surge of studies on **Parameter-Efficient Tuning** (PETuning) methods, which aim to tune the PLMs by adjusting lightweight trainable parameters while keeping most pretrained parameters frozen. Various ways have been used in these PETuning methods to introduce the lightweight trainable parameters. Adapter (Houlsby et al., 2019; Pfeiffer et al., 2020) is one of these that injects a small portion of *model-level* parameters within each transformer (Vaswani et al., 2017) layer of the pretrained language model. Prompt-tuning (Qin and Eisner,

2021; Liu et al., 2021c; Lester et al., 2021) is another one that introduces trainable continuous embeddings into the original sequences of input word embeddings to augment the PLMs on the *feature level*. Diff-pruning (Guo et al., 2021) learns and updates additional sparse diff-vector for all pretrained parameters, and LoRA (Hu et al., 2022) learns low-rank matrices to approximate the updated matrices, both of which update the PLMs on the *parameter level*. Moreover, BitFit (Zaken et al., 2021) *partially tunes* the bias terms of PLMs, without even introducing any new parameters. More details for these methods can be seen in Appx. §A.

Given the various exciting progresses of PETuning methods that all seem to demonstrate their competitive performance with higher training efficiency, the idea that PETuning could be a new general paradigm in place of full finetuning for transfer learning in NLP becomes never more tempting (Liu et al., 2021a). We, however, argue that current evidences are insufficient to support the complete overthrow of full finetuning. First, we point out that the current evaluation strategy, i.e., the development set is used for both early stopping and reporting results, used in a number of studies for PETuning (Lester et al., 2021; Vu et al., 2021; Liu et al., 2021b; Pfeiffer et al., 2021) does not provide fair model comparisons. This essentially causes data leakage that results in misleading conclusions (§2). Second, statistical significance is rarely reported when comparing PETuning methods. This is an especially crucial issue as we show that the finetuning and PETuning processes are inherently unstable due to various randomness, such as weight initialization and training data order (§3.3).

To fairly compare these tuning strategies, this study conducts a comprehensive re-examination on the effectiveness of PETuning methods. **Our main contributions** are: **1)** We conduct controlled experiments (§2) and reveal the fundamental flaw of the current evaluation scheme (i.e., its failure to assess generalisation) and how that leads to misinterpretations of the progress in the field. **2)** We offer a more reliable practice for model selection that is not prone to overfitting. **3)** We revisit the performance of PETuning in comparison with finetuning across tasks with various, and have reached very different conclusions on different data scales. **4)** We conduct the first comprehensive study to investigate the stability of off-the-shelf PETuning methods and identify the main contributing factors.

Key Findings: **1)** Finetuning cannot be fully replaced so far, since there is no PETuning method that can consistently outperform finetuning across all tasks and settings. We conclude that PETuning may be more suitable for low-resource tasks, but struggle on medium-resource tasks and fall behind finetuning across the board on high-resource tasks (see Figure 1). **2)** All the PETuning methods unanimously show instability across different random seeds similar to finetuning (Dodge et al., 2020), where the randomness comes from both weight initialisation and training data order. **3)** We found prompt-tuning lags far behind finetuning, which is a very different conclusion from previous studies. We show that prompt-tuning is highly unstable and cannot robustly and consistently re-produce its reported competitive performance (usually reported as a single run or the optimal run across multiple episodes (Lester et al., 2021; Liu et al., 2021b)) in our fair evaluation setup. **4)** Within each PETuning method, reducing the size of trainable parameters is likely to yield better stability (but not necessary to yield better or poorer performance). **5)** The stability of PETuning methods is substantially proportional to the scale of training data, and we further highlight the most crucial factor behind is the number of training iterations.

For the rest of the paper, we begin with the analysis on why the current evaluation protocol can be flawed (§2), and follow with a rigorous re-examination with a fairer protocol to benchmark the performance and stability of PETuning (§3).

2 The Unreliability of Misused Early Stopping

GLUE² and SuperGLUE³ have become the *de facto* benchmarks for verifying model effectiveness in Natural Language Understanding. For the sake of validity and fairness of the evaluation, the labels of test set in these benchmarks are not released. Instead, web portals are provided for submitting and evaluating the prediction results. Due to the limited number of allowed evaluation submissions to these benchmarks, a large number of works have followed a common practice that the model performance is only reported and compared based on the development sets rather than the real test sets, where the development set is treated as the “test set” (Lester et al., 2021; Vu et al., 2021; Liu et al.,

²<https://gluebenchmark.com>.

³<https://super.gluebenchmark.com>.

2021b; Pfeiffer et al., 2021).

While this practice is a convenient approximation of model performance as it allows quickly-obtained results from large-scaled experiments, there has been a serious problem in this setting: a single set is often used for both validating and testing the model. Therefore, the reported results under such setting might come from overly-optimistic checkpoints since early stopping is applied on the same set. We argue that such ‘trick’ practice breaches the standard train/development/test paradigm and compromises fair and rigorous comparison, leading to unreliable conclusions and misunderstandings of the examined models.

To verify our concern, in this section, we scrutinise this trick practice by comparing it with a rigorous evaluation protocol with strictly separated sets for validating and testing respectively, and thus provide comprehensive analyses to reveal the effects of the misused early stopping strategy.

Evaluation Setup. We adopt Roberta_{base} (Liu et al., 2019) as our base model, and experiment on the RTE dataset, which is a textual entailment dataset included in both GLUE and SuperGLUE. We divide the original development set of the RTE dataset by a 50%/50% split⁴ (denoted by *dev.1* and *dev.2* respectively), and compare the performance over finetuning and four PETuning methods, i.e., Adapter, prefix tuning (PT), LoRA, BitFit⁵. In particular, we use the *dev.2* set as the test set, and use the *dev.1* set or the *dev.2* set as the development set for model selection, respectively (denoted by RTE_{1-2} or RTE_{2-2}). We set the number of epochs to 50 and early stop when validation scores do not improve for 10 consecutive epochs following Mao et al. (2021). Both *evaluation loss* and *accuracy* are used as the stopping metrics⁶.

Results and Analyses. From Table 1, we can see that using a single set as both the development and test sets (i.e. RTE_{2-2} , which is the “trick” setting) can substantially boost the performances of PETuning models, comparing with using two sepa-

⁴A normal way to create new dev. set is to separate part of training set while using original dev. set as test set, as what we do in §3.1. However, such newly created dev. set may be of a different distribution from the new test set, where the comparison between using new dev. set and test set for early stopping might be unfair.

⁵We choose prefix tuning as the representative of prompt-tuning. See Appx. §B.2 for more details of the compared models.

⁶See Appx. §B.3 for the full hyperparameters settings.

	<i>Evaluation loss</i>		<i>Accuracy</i>	
	RTE_{1-2}	RTE_{2-2}	RTE_{1-2}	RTE_{2-2}
FT	78.89 ± 1.36	78.89 ± 1.36	79.28 ± 1.9	79.62 ± 2.22
Adapter	75.1 ± 1.60	76.3 ± 4.26	76.55 ± 3.57	78.42 ± 3.7
PT	57.55 ± 2.71	66.19 ± 8.51	57.84 ± 4.85	67.19 ± 11.37
LoRA	75.22 ± 2.77	75.94 ± 3.39	75.11 ± 3.3	77.7 ± 4.57
BitFit	70.79 ± 10.38	71.3 ± 10.19	66.76 ± 12.98	68.2 ± 13.72

Table 1: Mean and standard deviation results with different dev./test splits for RTE task across 20 runs. *Evaluation loss* and *accuracy* are the stopping metrics. Bold denotes the highest mean value for corresponding method with specific stopping metric.

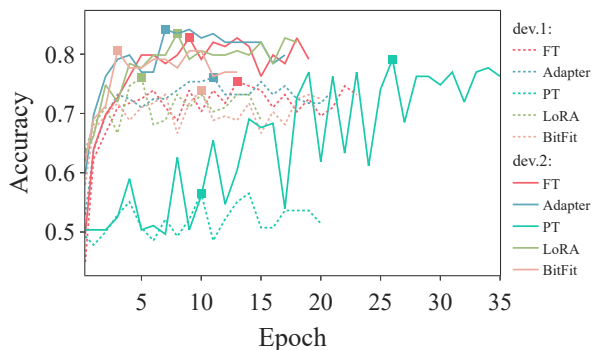


Figure 2: Comparing early stopped points selected by RTE_{1-2} and RTE_{2-2} , i.e., checkpoints with the best accuracy scores from *dev.1* and *dev.2* over training epochs. The markers denote the epochs selected by early stopping. Comparing the two checkpoint results on dev.2 (i.e. test performance), the RTE_{2-2} (same set for test and dev) checkpoint usually shows higher performance than the checkpoint selected in RTE_{1-2} by a large gap.

rate ones (i.e., RTE_{1-2}). Particularly, prefix tuning (PT) gains $\sim 10\%$ improvements on both the evaluation loss and accuracy stopping metrics. However, such performance boost does not mean genuine improvement in terms of better generalisation.

To demonstrate this in a more intuitive way, we plot the evaluation performance on development sets (i.e. *dev.1* and *dev.2* respectively) over training steps in Figure 2⁷. For each model, its early stopped epochs over the two sets are drastically different, suggesting that there is significant behavioural difference of the models across sets and best checkpoint selected on one set does not necessarily generalise well on the other set. In fact, the ability of models to mitigate such gap (e.g., from the best-performing checkpoints on *dev.1* to the best-performing ones on unseen *dev.2*) precisely denotes corresponding ability of generalisation,

⁷The best-performing runs of RTE_{1-2} and RTE_{2-2} are used for this visualisation.

which is essentially the criteria to measure the models’ effectiveness (Raschka, 2018). However, the evaluation scheme RTE_{2-2} , i.e., the misused early stopping, reuses the test set multiple times during training stage, which is tantamount to *leaking* the test information to erase this gap, resulting in unreliable evaluation.

This observation motivates us to re-examine these PETuning methods with a fairer evaluation.

3 Experiments with Fair Evaluation

In this section, we use a fairer evaluation protocol that strictly separates development and test sets. Based on this protocol, we conduct extensive experiments to investigate the effectiveness of PETuning methods (concluded in Figure 1). First, we experiment over a wide range of tasks under various levels of resource abundance to fully compare the performance of PETuning with finetuning (§3.2). Further, we provide in-depth analyses for the instability of PETuning methods, investigating the possible causes and provide practical suggestions of using PETuning methods (§3.3).

3.1 Experimental Setup

Data Setup. We conduct experiments on 12 datasets from GLUE and SuperGLUE, which are divided into three levels according to their sizes: (1) *low-resource* (< 1k data points), including CB (de Marneffe et al., 2019), COPA (Roemmele et al., 2011), and WSC (Levesque et al., 2012); (2) *medium-resource* (1k ~10k data points), including RTE (Wang et al., 2018), MRPC (Dolan and Brockett, 2005), WiC (Pilehvar and Camacho-Collados, 2019), STS-B (Cer et al., 2017), and BoolQ (Clark et al., 2019); (3) *high-resource* (> 10k data points), including SST-2 (Wang et al., 2018), MNLI (Williams et al., 2018), QNLI (Wang et al., 2018), and QQP⁸.

Since using a single set for both early stopping and testing could result in unreliable results (§2), we use separate development and test sets for all our experiments. Specifically, the original training set of each dataset is split into new train set and development set by a 90%/10% proportion, and the original development set is used as the test set.⁹

⁸<https://quoradata.quora.com/First-Quora-Dataset-Release-Question-Pairs>

⁹Ideally, the standard train-dev-test splits of GLUE and SuperGLUE should be used. However, due to the amount of experiments and evaluations need to be done in our ultra-large-scale investigation, we create our own splits instead of

Evaluation Setup. We again experiment with the Roberta_{base} model (Liu et al., 2019) on our 12 datasets. All experimental results are reported across 20 runs for low- and medium-resource tasks, and 10 runs for high-resource tasks with different random seeds, respectively. We train for 50 epochs and early stop the training when evaluation loss do not decrease for 10 consecutive epochs.¹⁰

3.2 Analysis of Performance

From the average performance for all tasks in Table 2, we can observe that most of the PETuning methods (i.e., Adapter, LoRA, and BitFit) indeed have some performance gains when compared with finetuning. It is known that PETuning methods have far better tuning efficiency, with significantly less tuning parameters (< 2% of full model parameters), comparing with full finetuning (Mao et al., 2021). However, it remains questionable to conclude that PETuning methods are more advantageous as the overall comparison may neglect important divergences in the wide range of tasks with different scales of training data. To provide a finer-grained view for the comparison between finetuning and PETuning, we group the results of the 12 tasks in Table 2 into low-, medium-, and high-resource tasks. Whilst most PETuning methods outperform finetuning on low-resource settings, the best PETuning is merely comparable to finetuning in medium-resource tasks and lags behind finetuning in high-resource tasks. We summarise the trend in Table 3, and provide more detailed analyses in the following.

Adapter & LoRA & BitFit only perform better on low-resource tasks. From Table 2, we observe that Adapter, LoRA, and BitFit obtain outstanding performance on the low-resource tasks and significantly outperform finetuning by large margins¹¹ (especially LoRA obtains ~8% performance gains on average). However, the trend changes when training data size gets larger. For the medium-resource tasks, only Adapter and BitFit can maintain a comparable performance with finetuning. LoRA and prefix tuning lags behind substantially. For the high-resource setting, finetuning performs consistently better than all PETuning methods¹². In

submitting models to the learderboards.

¹⁰See Appx. §B.3 for the full hyperparameters settings.

¹¹Similar observation for Adapter was previously reported in He et al. (2021). We extend it to more PETuning methods.

¹²These findings are also observed on the same task with different number of training instances. See Appx. §C.1 for

Dataset↓, Model→	FT	Adapter	PT	LoRA	BitFit
<i>Low-Resource</i>					
CB	70.00 \pm 13.32	77.49 \pm 13.20	46.55 \downarrow \pm 5.74	82.05 \uparrow \pm 9.62	81.12 \uparrow \pm 8.94
COPA	54.70 \pm 3.36	65.90 \uparrow \pm 5.42	55.35 \pm 5.07	66.4 \uparrow \pm 9.05	56.65 \pm 3.72
WSC	63.46 \pm 0.0	63.46 \pm 0.0	58.7 \downarrow \pm 4.69	63.46 \pm 0.0	63.46 \pm 0.0
Avg. (Low)	62.72 \pm 4.58	68.95 \uparrow \pm 4.83	53.53 \downarrow \pm 3.22	70.64 \uparrow \pm 4.32	67.08 \uparrow \pm 3.57
<i>Medium-Resource</i>					
RTE	73.77 \pm 3.17	73.88 \pm 1.88	57.36 \downarrow \pm 8.01	69.69 \downarrow \pm 7.89	70.67 \pm 10.77
MRPC	90.54 \pm 1.05	91.06 \pm 0.63	89.35 \downarrow \pm 1.31	91.03 \pm 0.95	91.06 \pm 0.71
WiC	65.47 \pm 2.04	65.12 \pm 1.88	62.12 \downarrow \pm 1.32	61.29 \downarrow \pm 6.7	66.0 \pm 1.41
STS-B	90.42 \pm 0.26	90.23 \downarrow \pm 0.1	89.64 \pm 0.39	90.47 \pm 0.11	90.44 \pm 0.15
BoolQ	78.75 \pm 0.72	76.93 \pm 0.92	75.44 \pm 0.47	76.92 \pm 1.33	76.9 \pm 0.84
Avg. (Medium)	79.79 \pm 0.99	79.44 \pm 0.74	74.78 \downarrow \pm 1.62	77.88 \downarrow \pm 2.02	79.01 \pm 2.22
<i>High-Resource</i>					
SST-2	94.15 \pm 0.0	93.34 \downarrow \pm 0.31	94.15 \pm 0.0	94.15 \pm 0.0	93.92 \pm 0.07
QNLI	92.40 \pm 0.12	92.31 \pm 0.09	92.31 \pm 0.27	91.00 \pm 0.69	91.60 \pm 1.01
QQP	91.38 \pm 0.06	90.28 \pm 0.0	88.90 \downarrow \pm 0.32	90.45 \downarrow \pm 0.17	89.28 \downarrow \pm 0.0
MNLI	87.42 \pm 0.20	86.88 \downarrow \pm 0.17	86.30 \downarrow \pm 0.08	86.96 \pm 0.24	85.50 \downarrow \pm 0.32
Avg. (High)	91.34 \pm 0.09	90.70 \downarrow \pm 0.12	90.42 \downarrow \pm 0.14	90.64 \downarrow \pm 0.21	90.08 \downarrow \pm 0.29
Avg. (All)	79.37	80.57	74.68	80.32	79.72

Table 2: Mean and standard deviation results for each of the 12 tasks. We report the f1 score for CB and MRPC, Pearson correlation for STS-B, and accuracy for other tasks (matched accuracy for MNLI). Higher is better for all metrics. One-tailed t-test is used for the comparison between PETuning and finetuning. One PETuning method that outperforms (\uparrow) or falls behind (\downarrow) finetuning when accepting the corresponding alternative hypothesis, where p-value < 0.05 that means they have significant differences.

	Low	Medium	High
Adapter	\nearrow	\rightarrow	\searrow
PT	\searrow	\searrow	\searrow
LoRA	\nearrow	\searrow	\searrow
BitFit	\nearrow	\rightarrow	\searrow

Table 3: Performance comparison between PETuning and finetuning on low-, medium-, and high-resource settings, respectively. Arrows indicate whether corresponding PETuning method significantly outperforms finetuning (\nearrow), falls behind (\searrow), or their results across multiple runs without significant differences (\rightarrow).

particular, among the PETuning methods, Adapter obtains the highest scores on high-resource tasks. This result suggests that low-resource would be the main factor that PETuning methods could outperform the full finetuning.

Prefix tuning consistently underperforms finetuning. According to Table 2 and Table 3, finetuning beats prefix tuning by large margins on most tasks across multiple runs, contradicting to what has been reported in Liu et al. (2021b). One possible reason is that prefix tuning is highly unstable more details.

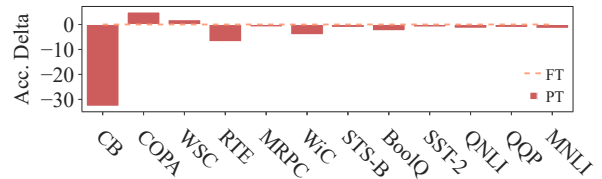


Figure 3: Relative performance differences of prefix tuning (PT) over full finetuning (FT) on the upper bounds of multi-run results. PT achieves close upper bounds compared with FT on most of the 12 tasks.

to train and thus may have exploited the misused early stopping more than other PETuning methods (see Figure 2). Besides early stopping protocol, previous works on prefix tuning only report their result of a single run (Liu et al., 2021b; Lester et al., 2021; Vu et al., 2021), which might lead to biased conclusion. In Figure 3 we further plot the upper bounds of these runs, and we indeed observe that the optimal run from prefix tuning achieves competitive performance compared with finetuning on many tasks. However, the results in Table 2 verify that this competitiveness would plummet across different runs by varying the random seeds. Such instability of prefix tuning leads to its poor average performance in our experiments. We further discuss this in §3.3.

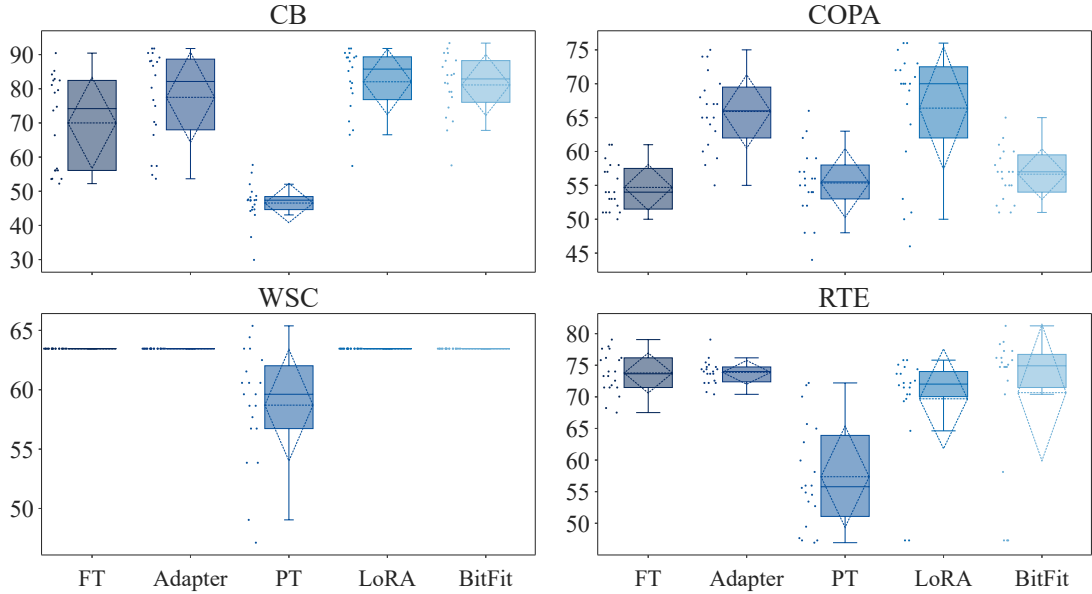


Figure 4: The experimental results over 20 different random seeds across CB, COPA, WSC, and RTE datasets, where finetuning and PETuning methods show large instability. The dashed rhombuses denote the mean (horizontal dashed line) and standard deviation (vertical distance).

	WI	DO	Global
FT	55.40 \pm 4.55	55.35 \pm 3.32	54.70 \pm 3.36
Adapter	67.15 \pm 5.40	66.35 \pm 7.36	65.90 \pm 5.42
PT	55.00 \pm 5.13	54.75 \pm 4.97	55.35 \pm 5.07
LoRA	63.60 \pm 7.93	64.60 \pm 8.56	66.40 \pm 9.05
BitFit	58.40 \pm 2.29	56.00 \pm 4.00	56.65 \pm 3.72

Table 4: Performance over 20 runs on COPA task, controlled by global random seeds, weight initialization (WI) random seeds, and data order (DO) random seeds, respectively. (Visualised in Figure 12 in the Appendix.)

Finetuning cannot be fully replaced. To summarise, PETuning has exceptional performance in resource-poor scenarios and usually outperform the more expensive full-model finetuning. However, when dataset size increases, finetuning regains dominance in medium- and high-resource setups. This indicates that finetuning cannot be fully replaced so far. We also delved deeper into understanding why finetuning lags behind PETuning on low-resource settings. Our investigation points to the different fitting capabilities of finetuning and PETuning. Specifically, finetuning is more prone to overfitting on low-resource tasks¹³.

3.3 Analysis of Stability

By revisiting the results in Table 2, we can observe that both finetuning and all PETuning methods ex-

hibit large standard deviations on several tasks, i.e., CB, COPA, WSC, and RTE. To further understand this phenomenon, in Figure 4, we visualise the performance distribution of 20 runs of finetuning and PETuning methods on these tasks. Surprisingly, large fluctuations are seen on all four tasks across all methods, where the margins between the lower and upper bounds could reach over 30%. While Dodge et al. (2020); Mosbach et al. (2021) have previously identified such variation exists for finetuning, our experiments further validate that such instability also occurs in all PETuning methods and could be even more prominent in certain tasks.

This level of instability severely hampers the application of PETuning and there is a pressing need to understand the underlying cause. However, to the best of our knowledge, no previous studies have systematically discussed the instability issue in PETuning methods. In this section, we provide the first comprehensive investigation on this matter. While instability is measured as the performance differences between random seeds, we further disentangle two randomness sources (weight initialisation and training data order) to better describe model instability. We then investigate two factors that might affect model instability: (1) trainable parameter size; and (2) training data size and training iterations. Through controlled experiments, we find that model instability is reflected by both changing data order and changing weight initialisa-

¹³See Appx. §C.2 for more details and analyses.

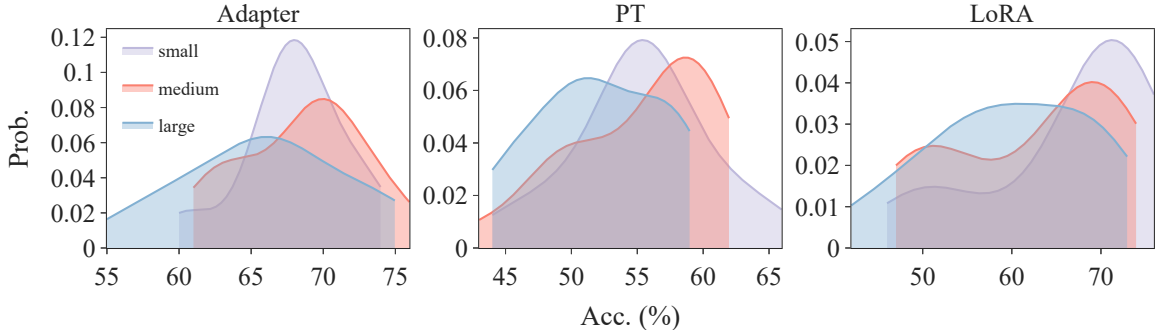


Figure 5: Performance probability density curves of Adapter, prefix tuning (PT), and LoRA over small, medium, and large parameter scales on COPA task across 20 runs. (See the numerical results and analyses in Appx. §C.3.)

379 tion. Reducing model size and increasing training
 380 iteration seems to have positive impact on model
 381 stability. We discuss all these points in detail in the
 382 followings.

383 **Weight initialisation and data order work to-**
 384 **gether.** Instability is measured from performance
 385 changes due to randomness introduced by random
 386 seeds. Two key things impacted by random seeds
 387 are (a) the initialisation of trainable weights (in-
 388 cluding extra parameters of PETuning methods and
 389 the classification head), and (b) the order of train-
 390 ing data fed to the model. To disentangle these two
 391 factors, following the setting in Dodge et al. (2020),
 392 we use two separate random seeds to control weight
 393 initialisation and training data order respectively,
 394 comparing with using one global random seed to
 395 control these two factors simultaneously.

396 Table 4 demonstrates that each of the two factors
 397 could individually lead to large standard deviations,
 398 which means the instability of PETuning methods
 399 are sensitive to either training data order, or weight
 400 initialisation, or both. This observation indicates
 401 that the sources of instability for PETuning can be
 402 multifaceted – isolating and enhancing stability via
 403 controlling individual factor can be challenging.¹⁴

404 **Models with fewer trainable parameters are**
 405 **more stable.** To investigate the impact of model
 406 size on model stability, we define three sizes, *small*,
 407 *medium*, and *large*, for each PETuning method.
 408 The three sizes correspond to the reduction factor
 409 of {64, 16, 2} for Adapter¹⁵, the prompt length of
 410 {32, 64, 128} for prefix tuning, and the rank of {8,

¹⁴Prior works mainly focused on obtaining better prior (e.g., prompt/weight initialisation) to improve model performance/stability but did not touch upon the multifaceted nature of instability (Pfeiffer et al., 2021; Lester et al., 2021; Vu et al., 2021).

¹⁵The smaller the reduction factor, the more parameters the model has.

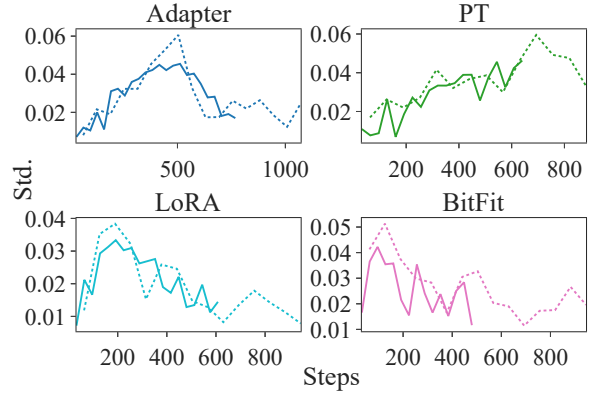


Figure 6: Standard deviations of data size in {1k (solid line), 2k (dashed line)} over training steps on WiC task across 20 runs. (See that on BoolQ task in Figure 13 in the Appendix.)

411 16, 32} for LoRA. We conduct a set of controlled
 412 experiments on the COPA task where PETuning
 413 methods exhibit high instability. We perform 20
 414 runs for each setting and use kernel density esti-
 415 mation (KDE) (Chen, 2017) to estimate the probability
 416 density curves of the multi-run results.

417 As shown in Figure 5, for all PETuning methods,
 418 we consistently observe that the probability density
 419 curves would be progressively flatter (having lower
 420 peak) as the number of parameters increase from
 421 small to large. This suggests that more trainable
 422 parameters for PETuning leads to a wider range of
 423 performance distribution, resulting in higher insta-
 424 bility. That said, when it comes to model perfor-
 425 mance, the best-performing model usually is not
 426 the smallest one. We conjecture that models with
 427 fewer trainable parameters converge quickly to the
 428 rough global minima but could be underfitting the
 429 real data manifold.

430 **Data size does not affect instability directly, but**
 431 **training iterations do.** Figure 1 and Table 2

432 suggest that PETuning methods almost always
433 have larger standard deviations on lower-resource
434 tasks¹⁶. To investigate if training data size directly
435 affects the stability of PETuning, inspired by Mos-
436 bach et al. (2021), we compare models that are
437 trained with randomly sampled 1k and 2k training
438 instances from WiC training set and validated with
439 an another separately sampled 1k development set.

440 In Figure 6, we observe that the solid and dashed
441 lines of each PETuning method are substantially
442 intertwined, which means the standard deviations
443 (instability) of PETuning methods trained by 1k
444 or 2k samples would not have significant differ-
445 ences with the same number of steps. Instead, the
446 culprit, leading to the discrepancy of instability
447 across different training data sizes, is essentially
448 the number of training iterations. As shown in Fig-
449 ure 6, the standard deviations of PETuning methods
450 have an initial ascent stage where models are fit-
451 ting to the training data and thus having fluctuating
452 performance. After the ascent stage, the standard
453 deviations substantially decrease as the number of
454 training iterations (steps) get larger. With number
455 of epochs being fixed, the total number of itera-
456 tions on small datasets is small and the standard
457 deviation has yet to decrease, causing the higher in-
458 stability in lower-resource tasks. In particular, due
459 to the weaker fitting capabilities (Ding et al., 2022),
460 prefix tuning (PT) has a longer ascent stage, which
461 might need more training iterations to obtain more
462 stable performance. That said, prolong the training
463 on small datasets do not necessarily enhance model
464 performance, and the best checkpoint may still be
465 only appearing when the standard deviation is high.

466 4 Related Work

467 **Instability of finetuning PLMs.** While our study
468 is, to the best of our knowledge, the first to system-
469 atically investigate PETuning instability, prior stud-
470 ies have looked into the instability of finetuning
471 PLMs. Dodge et al. (2020) illustrated the inherent
472 instability of finetuning by controlling the random
473 seeds and provided a new early stopping strategy
474 to improve instability. Lee et al. (2020) proposed a
475 new regularisation method by mixing two models
476 based on dropout to prevent catastrophic forgetting
477 and to improve instability. More recently, Mosbach
478 et al. (2021) revisited the hypotheses of finetuning
479 instability proposed by previous studies and found
480 that optimisation difficulties can lead to vanishing

481 gradients, which further causes finetuning instabil-
482 ity. Zhang et al. (2021) also reveal that optimisation
483 significantly affects the instabilities in few-sample
484 fine-tuning.

485 **Analysis of PETuning.** As PETuning methods
486 have become a prominent research direction, a
487 great number of studies are dedicated to analyse the
488 characteristics of these methods. He et al. (2021)
489 investigate the effectiveness of Adapter across dif-
490 ferent scales and Han et al. (2021) provide a robust
491 strategy for training Adapter. Recently, He et al.
492 (2022) and Mao et al. (2021) propose a unified view
493 to connect various PETuning methods. However,
494 there has not been reliable validation and compar-
495 ison for off-the-shelf PETuning methods in terms
496 of stability and effectiveness, and this is where our
497 paper bridges the gap.

498 5 Conclusion

499 This work conducted a rigorous re-examination
500 on the current parameter-efficient tuning (PETun-
501 ing) methods. We demonstrated that performing
502 early stopping and evaluation on the same dataset
503 (a common practice used in many past studies)
504 could lead to unreliable conclusions. This issue is
505 more pronounced when accompanied by the insta-
506 bility nature of PETuning, leading to the inflated
507 results and overly optimistic estimates of PETuning
508 approaches. We re-evaluate these PETuning meth-
509 ods on the performance and stability aspects on a
510 rigorous evaluation protocol that strictly separates
511 validation and test sets. By the fine-grained compar-
512 ison between PETuning and finetuning on the per-
513 formance, we found that PETuning methods are not
514 consistently competitive with finetuning, namely
515 prefix tuning performing poorly across tasks and
516 most PETuning methods perform worse than fine-
517 tuning on higher-resource settings. By system-
518 atically investigating the instability of PETuning
519 methods, we found that models’ instability is sensi-
520 tive to both weight initialisation and training data
521 order. We identify two major factors behind such
522 instability: 1) models with fewer parameters are
523 more stable within each PETuning method (but not
524 across models); 2) more training iterations can usu-
525 ally reduce instability. Our overall re-examination
526 conclude that finetuning still cannot be fully re-
527 placed by PETuning so far, and there are many key
528 challenges for PETuning in terms of both perfor-
529 mance and instability, which need to be addressed
530 in future work.

¹⁶This is further confirmed in Appx. §C.1.

531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587

References

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language Models are Few-Shot Learners](#). In *Advances in Neural Information Processing Systems*, pages 1877–1901.

Daniel Cer, Mona Diab, Eneko Agirre, Iñigo Lopez-Gazpio, and Lucia Specia. 2017. [SemEval-2017 Task 1: Semantic Textual Similarity Multilingual and Crosslingual Focused Evaluation](#). In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 1–14.

Yen-Chi Chen. 2017. [A tutorial on kernel density estimation and recent advances](#). *Biostatistics & Epidemiology*, (1):161–187.

Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. 2019. [BoolQ: Exploring the Surprising Difficulty of Natural Yes/No Questions](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2924–2936.

Marie-Catherine de Marneffe, Mandy Simons, and Judith Tonhauser. 2019. [The CommitmentBank: Investigating projection in naturally occurring discourse](#). *Proceedings of Sinn und Bedeutung*, (2):107–124.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4171–4186.

Ning Ding, Yujia Qin, Guang Yang, Fuchao Wei, Zonghan Yang, Yusheng Su, Shengding Hu, Yulin Chen, Chi-Min Chan, Weize Chen, et al. 2022. [Delta Tuning: A Comprehensive Study of Parameter Efficient Methods for Pre-trained Language Models](#). *arXiv preprint arXiv:2203.06904*.

Jesse Dodge, Gabriel Ilharco, Roy Schwartz, Ali Farhadi, Hannaneh Hajishirzi, and Noah Smith. 2020. [Fine-Tuning Pretrained Language Models: Weight Initializations, Data Orders, and Early Stopping](#). *arXiv preprint arXiv:2002.06305*.

William B. Dolan and Chris Brockett. 2005. [Automatically Constructing a Corpus of Sentential Paraphrases](#). In *Proceedings of the Third International Workshop on Paraphrasing*.

Yuxian Gu, Xu Han, Zhiyuan Liu, and Minlie Huang. 2021. [PPT: Pre-trained prompt tuning for few-shot learning](#). *arXiv preprint arXiv:2109.04332*. 588
589
590

Demi Guo, Alexander Rush, and Yoon Kim. 2021. [Parameter-Efficient Transfer Learning with Diff Pruning](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing*, pages 4884–4896. 591
592
593
594
595
596

Wenjuan Han, Bo Pang, and Ying Nian Wu. 2021. [Robust Transfer Learning with Pretrained Language Models through Adapters](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing*, pages 854–861. 597
598
599
600
601
602
603

Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. 2022. [Towards a Unified View of Parameter-Efficient Transfer Learning](#). In *International Conference on Learning Representations*. 604
605
606
607
608

Ruidan He, Linlin Liu, Hai Ye, Qingyu Tan, Bosheng Ding, Liying Cheng, Jiawei Low, Lidong Bing, and Luo Si. 2021. [On the Effectiveness of Adapter-based Tuning for Pretrained Language Model Adaptation](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing*, pages 2208–2222. 609
610
611
612
613
614
615
616

Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. [Parameter-Efficient Transfer Learning for NLP](#). In *Proceedings of the 36th International Conference on Machine Learning*, pages 2790–2799. 617
618
619
620
621
622

Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. [LoRA: Low-Rank Adaptation of Large Language Models](#). In *International Conference on Learning Representations*. 623
624
625
626
627

Cheolhyoung Lee, Kyunghyun Cho, and Wanmo Kang. 2020. [Mixout: Effective Regularization to Finetune Large-scale Pretrained Language Models](#). In *8th International Conference on Learning Representations*. 628
629
630
631

Jaejun Lee, Raphael Tang, and Jimmy Lin. 2019. [What Would Elsa Do? Freezing Layers During Transformer Fine-Tuning](#). *arXiv preprint arXiv:1911.03090*. 632
633
634
635

Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. [The Power of Scale for Parameter-Efficient Prompt Tuning](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3045–3059. 636
637
638
639
640

Hector J. Levesque, Ernest Davis, and Leora Morgenstern. 2012. [The Winograd Schema Challenge](#). In

643			
644		<i>Proceedings of the Thirteenth International Conference on Principles of Knowledge Representation and Reasoning</i> , page 552–561.	
645			
646	Xiang Lisa Li and Percy Liang. 2021. Prefix-Tuning: Optimizing Continuous Prompts for Generation . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing</i> , pages 4582–4597.		
647			
648			
649			
650			
651			
652	Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2021a. Pre-train, Prompt, and Predict: A Systematic Survey of Prompting Methods in Natural Language Processing . <i>arXiv preprint arXiv:2107.13586</i> .		
653			
654			
655			
656			
657	Xiao Liu, Kaixuan Ji, Yicheng Fu, Zhengxiao Du, Zhilin Yang, and Jie Tang. 2021b. P-Tuning v2: Prompt Tuning Can Be Comparable to Fine-tuning Universally Across Scales and Tasks . <i>arXiv preprint arXiv:2110.07602</i> .		
658			
659			
660			
661			
662	Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. 2021c. GPT Understands, Too . <i>arXiv preprint arXiv:2103.10385</i> .		
663			
664			
665	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A Robustly Optimized Bert Pretraining Approach . <i>arXiv preprint arXiv:1907.11692</i> .		
666			
667			
668			
669			
670	Yuning Mao, Lambert Mathias, Rui Hou, Amjad Almahairi, Hao Ma, Jiawei Han, Wen-tau Yih, and Madian Khabsa. 2021. Unipelt: A Unified Framework for Parameter-Efficient Language Model Tuning . <i>arXiv preprint arXiv:2110.07577</i> .		
671			
672			
673			
674			
675	Zaiqiao Meng, Fangyu Liu, Thomas Clark, Ehsan Shareghi, and Nigel Collier. 2021. Mixture-of-Partitions: Infusing Large Biomedical Knowledge Graphs into BERT . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 4672–4681.		
676			
677			
678			
679			
680			
681	Marius Mosbach, Maksym Andriushchenko, and Dietrich Klakow. 2021. On the Stability of Fine-tuning BERT: Misconceptions, Explanations, and Strong Baselines . In <i>9th International Conference on Learning Representations</i> .		
682			
683			
684			
685			
686	Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé, Kyunghyun Cho, and Iryna Gurevych. 2021. AdapterFusion: Non-Destructive Task Composition for Transfer Learning . In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 487–503.		
687			
688			
689			
690			
691			
692	Jonas Pfeiffer, Andreas Rücklé, Clifton Poth, Aishwarya Kamath, Ivan Vulić, Sebastian Ruder, Kyunghyun Cho, and Iryna Gurevych. 2020. AdapterHub: A Framework for Adapting Transformers . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 46–54.		
693			
694			
695			
696			
697			
698			
		Mohammad Taher Pilehvar and Jose Camacho-Collados. 2019. WiC: the Word-in-Context Dataset for Evaluating Context-Sensitive Meaning Representations . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 1267–1273.	699
			700
			701
			702
			703
			704
			705
		Guanghui Qin and Jason Eisner. 2021. Learning How to Ask: Querying LMs with Mixtures of Soft Prompts . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 5203–5212.	706
			707
			708
			709
			710
			711
		Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language Models are Unsupervised Multitask Learners . <i>OpenAI blog</i> , (8):9.	712
			713
			714
			715
		Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer . <i>J. Mach. Learn. Res.</i> , pages 140:1–140:67.	716
			717
			718
			719
			720
			721
		Sebastian Raschka. 2018. Model Evaluation, Model Selection, and Algorithm Selection in Machine Learning . <i>arXiv preprint arXiv:1811.12808</i> .	722
			723
			724
		Melissa Roemmele, Cosmin Adrian Bejan, and Andrew S Gordon. 2011. Choice of Plausible Alternatives: An Evaluation of Commonsense Causal Reasoning . In <i>AAAI spring symposium: logical formalizations of commonsense reasoning</i> , pages 90–95.	725
			726
			727
			728
			729
		Timo Schick and Hinrich Schütze. 2021. It’s Not Just Size That Matters: Small Language Models Are Also Few-Shot Learners . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 2339–2352.	730
			731
			732
			733
			734
			735
		Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need . In <i>Advances in Neural Information Processing Systems</i> .	736
			737
			738
			739
			740
		Tu Vu, Brian Lester, Noah Constant, Rami Al-Rfou, and Daniel Cer. 2021. Spot: Better Frozen Model Adaptation through Soft Prompt Transfer . <i>arXiv preprint arXiv:2110.07904</i> .	741
			742
			743
			744
		Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2019. SuperGLUE: A Stickier Benchmark for General-Purpose Language Understanding Systems . In <i>Advances in Neural Information Processing Systems</i> , volume 32. Curran Associates, Inc.	745
			746
			747
			748
			749
			750
			751
		Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. GLUE: A Multi-Task Benchmark and Analysis Platform for	752
			753
			754

755 [Natural Language Understanding](#). In *Proceedings of*
756 *the 2018 EMNLP Workshop BlackboxNLP: Analyz-*
757 *ing and Interpreting Neural Networks for NLP*, pages
758 353–355.

759 Adina Williams, Nikita Nangia, and Samuel Bowman.
760 2018. [A Broad-Coverage Challenge Corpus for Sen-](#)
761 [tence Understanding through Inference](#). In *Proceed-*
762 *ings of the 2018 Conference of the North Ameri-*
763 *can Chapter of the Association for Computational*
764 *Linguistics: Human Language Technologies*, pages
765 1112–1122.

766 Elad Ben Zaken, Shauli Ravfogel, and Yoav Gold-
767 berg. 2021. [Bitfit: Simple parameter-efficient](#)
768 [fine-tuning for transformer-based masked language-](#)
769 [models](#). *arXiv preprint arXiv:2106.10199*.

770 Tianyi Zhang, Felix Wu, Arzoo Katiyar, Kilian Q Wein-
771 berger, and Yoav Artzi. 2021. [Revisiting Few-sample](#)
772 [BERT Fine-tuning](#). In *International Conference on*
773 *Learning Representations*.

774 Zexuan Zhong, Dan Friedman, and Danqi Chen. 2021.
775 [Factual Probing Is \[MASK\]: Learning vs. Learning](#)
776 [to Recall](#). In *Proceedings of the 2021 Conference*
777 *of the North American Chapter of the Association*
778 *for Computational Linguistics: Human Language*
779 *Technologies*, pages 5017–5033.

Appendix

A PETuning Methods

PETuning methods are unique in keeping (most) pretrained parameters of PLMs frozen and finetuning only light-weight additional parameters or a fraction of the PLM’s parameters for downstream tasks.

To achieve efficient tuning of PLMs, existing PETuning methods are generally designed by two different manners: (1) training additional parameters on different levels of PLMs, including model-level (Appx. §A.1), feature-level (Appx. §A.2), and the parameter-level (Appx. §A.3), or (2) tuning partial parameters of the base model (Appx. §A.4). Figure 7 shows the difference of these PETuning methods.

A.1 Model-Level

Adapter-Tuning. *Adapters* (Houlsby et al., 2019; Pfeiffer et al., 2020, 2021; Meng et al., 2021) are a type of PETuning approaches that insert small newly initialised parameter modules on the model-level (i.e., each transformer layer) of PLMs. In particular, these adapter modules are normally moulded by a two-layer feed-forward neural network with a bottleneck: (1) a down-projection with $\mathbf{W}_{\text{down}} \in \mathbb{R}^{d \times r}$ to project the input \mathbf{h}_i to a lower-dimensional space specified by bottleneck dimension r ; (2) an up-projection with $\mathbf{W}_{\text{up}} \in \mathbb{R}^{r \times d}$ to project back to the input size. Mathematically, the adapter can be defined as:

$$\mathbf{h}_a = \mathbf{W}_{\text{up}}^\top f\left(\mathbf{W}_{\text{down}}^\top \mathbf{h}_i\right), \quad (1)$$

where \mathbf{h}_a is the output and $f(\cdot)$ is the activation function. During the finetuning, the model only updates the parameters of the adapter modules while keeps the underlying pretrained model fixed.

A.2 Feature-Level

Prompt-Tuning. Prompt-Tuning (Lester et al., 2021) is another type of PETuning approaches that introduce additional tunable parameters on the feature-level. Specifically, prompt-tuning introduces additional tunable prefix (or suffix) vectors, namely prompts (Zhong et al., 2021; Schick and Schütze, 2021), to extend the input text features (or the input of each transformer layer (Li and Liang, 2021; Liu et al., 2021b)), and tunes only the prompts. Besides its simplicity and lightness, prompt-tuning could achieve on par performance,

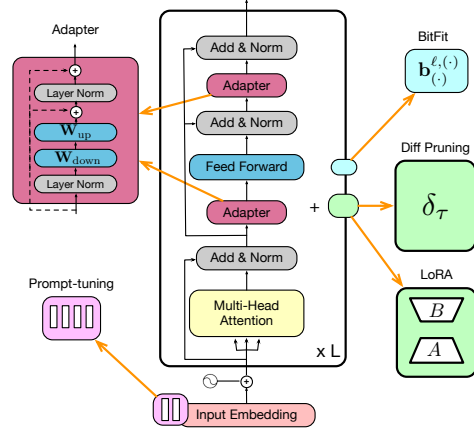


Figure 7: Different PETuning methods by adjusting trainable parameter on model level (Adapter), feature level (Prompt-tuning), parameter level (Diff Pruning and LoRA), and partial-tuning level (BitFit).

particularly in billions-size PLMs, and even better performance, comparing with the full finetuning (Liu et al., 2021b).

A.3 Parameter-Level

Diff-Pruning. Diff-pruning (Guo et al., 2021) works on all parameters of PLMs, which aims to learn additional trainable sparse parameters for the entire PLMs. Specifically, for the pretrained parameters Θ , diff-pruning reparameterizes the task-specific model parameters Θ_τ as:

$$\Theta_\tau = \Theta + \delta_\tau, \quad (2)$$

where δ_τ denotes the trainable diff vector, which is regularised to be sparse.

LoRA. LoRA (Hu et al., 2022) focuses on the updating procedure of the language model parameters. For a pretrained weight matrix $\mathbf{W} \in \mathbb{R}^{d \times k}$, LoRA uses trainable low-rank matrices to approximate the updates ($\Delta\mathbf{W}$) by:

$$\mathbf{W} + \Delta\mathbf{W} = \mathbf{W} + \mathbf{B}\mathbf{A}, \quad (3)$$

where $\mathbf{B} \in \mathbb{R}^{d \times r}$, $\mathbf{A} \in \mathbb{R}^{r \times k}$, and the rank $r \ll \min(d, k)$.

A.4 Partial Finetuning

BitFit. *Partial finetuning* aims to tune a fraction of PLMs parameters without introducing any additional ones. For example, Lee et al. (2019) only tunes the top layers, however, which usually performs much worse than full finetuning. With the principle of efficiency and effectiveness, BitFit (Zaken et al., 2021) turns to tune the bias terms of PLMs to obtain competitive performance.

B General Experimental Setup

In this section, we illustrate the general task, compared methods, and hyperparameter settings for the following experiments. Apart from that, in §2 and §3, we will additionally illustrate their specific data and evaluation setups, respectively.

B.1 Task Setup

In order to extensively compare the performance and stability of PETuning methods with the full-model finetuning, we select a full set of 12 tasks across low-, medium- and high-resource scales of GLUE and SuperGLUE, including natural language inference (CB, RTE, MNLI, QNLI), question answering (COPA, BoolQ), paraphrasing (MRPC, QQP), sentiment analysis (SST-2), sentence similarity (STS-B), word sense disambiguation (WiC), and coreference resolution (WSC) tasks. According to the dataset sizes, we divide these tasks into three levels:

- **Low-Resource:** the tasks with training data size smaller than 1k, including CB, COPA, and WSC.
- **Medium-Resource:** the tasks with training data size between 1k and 10k, including RTE, MRPC, WiC, STS-B, and BoolQ.
- **High-Resource:** the tasks with training data size larger than 10k, including SST-2, QNLI, QQP, and MNLI.

B.2 Compared Methods

We have chosen four representative PETuning methods: **Adapter**, **Prompt-tuning**, **LoRA**, and **BitFit**, across all levels. For Adapter, we use the Pfeiffer architecture (Pfeiffer et al., 2020) since it has reported better performance than others. For Prompt-tuning, due to the poor performance of standard prompt-tuning (Lester et al., 2021) on small PLMs, e.g., base versions of Bert and Roberta, we adopt the settings of **prefix tuning** (Li and Liang, 2021) to add continuous prompts for each transformer layer of PLMs. For LoRA & BitFit, we take the architectures from their origin papers (Hu et al., 2022; Zaken et al., 2021).

B.3 Hyperparameter Setup

We adopt Roberta_{base} as the base model released by Huggingface¹⁷. The grid search is used to se-

¹⁷<https://github.com/huggingface/transformers>

lect the learning rate from {1e-6, 1e-5, 5e-5, 1e-4, 5e-4, 1e-3, 5e-3, 1e-2} and batch size from {16, 32}. We search the reduction factor from {2, 16, 64} following (Pfeiffer et al., 2021) for Adapter, the prompt length from {8, 16, 32, 64} for prefix tuning, and the scaling factor α and rank from {8, 16} for LoRA following its origin paper. There are many studies focusing on achieving better initialization by post pretraining for PETuning methods such as Adapter (Pfeiffer et al., 2021) and prompt (Vu et al., 2021; Gu et al., 2021), however, to be a fair comparison, the extra parameters of all PETuning methods are initialized randomly.

We set the number of epochs to 50 and adopt the early stopping strategy with the patience of 10 worse-performing epochs on our new development set following (Mao et al., 2021). In particular, for §2, to fully investigate the effects of early stopping on the task RTE, we use both *evaluation loss* and *accuracy* as the stopping metrics; for §3, due to the variety of evaluation metrics for the tasks, we use the *evaluation loss* as the common stopping metric.

C Additional Experiments and Analyses

C.1 The Same Task with Different Training Data Sizes

To make the above conclusions in §3.2 more convincing, we conduct fine-grained experiments following (He et al., 2021). Specifically, we separately sample 500, 5k and 50k training instances from the original training data as representatives of low-, medium- and high-resource settings, in addition to draw another 1k samples as development set for each task. We report experimental results for WiC, STS-B, BoolQ, SST-2, QNLI, QQP, and MNLI, which have more than 6k training samples, and following the settings illustrated in §3.1.

Confirming our conclusions, in Table 5, we obtain fully consistent findings with §3.2 and Table 3, that prefix tuning consistently falls behind finetuning on various-resources tasks; Adapter&LoRA&BitFit significantly outperforms finetuning on low-resource tasks; Adapter&BitFit keep competitive with finetuning and LoRA lags behind; and all PETuning methods falls behind on high-resource tasks.

In addition, we plot the mean and std. values with different data scales on the same task in Figure 8, to confirm the std. is substantially proportional to training data size.

Model↓, Dataset→	WiC	STS-B	BoolQ	SST-2	QNLI	QQP	MNLI	Avg.
500								
FT	56.12 \pm 2.13	83.81 \pm 1.34	62.17 \pm 0.0	87.89 \pm 1.42	77.54 \pm 6.5	74.21 \pm 3.11	58.61 \pm 6.21	71.47 \pm 1.44
Adapter	58.36 \uparrow \pm 3.98	85.74 \uparrow \pm 0.64	61.56 \pm 1.51	89.12 \uparrow \pm 0.84	79.85 \uparrow \pm 1.55	75.91 \uparrow \pm 1.01	60.92 \uparrow \pm 2.86	73.07 \uparrow \pm 0.84
PT	56.25 \pm 1.07	73.98 \downarrow \pm 3.79	57.55 \downarrow \pm 4.43	82.01 \downarrow \pm 2.75	70.77 \downarrow \pm 5.16	66.37 \downarrow \pm 1.42	36.35 \downarrow \pm 1.51	63.32 \downarrow \pm 1.28
LoRA	58.64 \uparrow \pm 1.16	85.06 \uparrow \pm 0.83	60.48 \pm 4.74	89.1 \uparrow \pm 0.78	80.86 \uparrow \pm 0.7	72.43 \pm 2.9	63.1 \uparrow \pm 2.57	72.81 \uparrow \pm 0.84
BitFit	57.92 \uparrow \pm 2.4	84.39 \uparrow \pm 1.98	62.16 \pm 0.04	88.38 \pm 8.19	78.38 \pm 2.74	73.13 \pm 1.63	61.47 \uparrow \pm 2.08	72.26 \uparrow \pm 1.23
5k								
FT	66.98 \pm 1.26	90.76 \pm 0.03	73.81 \pm 1.11	92.66 \pm 0.92	87.12 \pm 0.37	83.72 \pm 0.17	78.18 \pm 0.82	81.89 \pm 0.44
Adapter	65.15 \pm 2.85	90.24 \pm 0.05	73.51 \pm 1.63	92.66 \pm 0.25	86.69 \pm 0.49	82.97 \pm 0.32	78.0 \pm 1.17	81.32 \pm 0.71
PT	66.25 \pm 1.29	89.14 \pm 0.66	62.32 \downarrow \pm 0.14	91.7 \pm 0.92	85.83 \downarrow \pm 1.32	80.11 \downarrow \pm 1.52	77.78 \pm 0.42	79.02 \downarrow \pm 0.83
LoRA	62.46 \downarrow \pm 1.38	90.57 \pm 0.15	71.9 \downarrow \pm 0.41	92.35 \pm 0.7	87.49 \pm 0.39	83.03 \pm 0.36	77.67 \pm 0.26	80.78 \downarrow \pm 0.45
BitFit	67.61 \pm 0.49	90.37 \pm 0.19	75.08 \uparrow \pm 0.56	92.35 \pm 0.56	86.47 \pm 0.28	82.9 \pm 0.25	78.87 \pm 0.1	81.95 \pm 0.15
50k								
FT	-	-	-	93.46 \pm 0.18	90.07 \pm 0.22	88.36 \pm 0.18	84.71 \pm 0.41	89.15 \pm 0.27
Adapter	-	-	-	93.02 \pm 0.25	89.03 \downarrow \pm 0.17	86.67 \downarrow \pm 0.26	84.03 \pm 0.56	88.19 \downarrow \pm 0.18
PT	-	-	-	93.32 \pm 0.47	88.23 \downarrow \pm 0.59	85.21 \downarrow \pm 0.79	82.96 \downarrow \pm 0.20	87.43 \downarrow \pm 0.41
LoRA	-	-	-	93.35 \pm 0.27	89.49 \pm 0.13	87.20 \downarrow \pm 0.33	83.26 \downarrow \pm 0.11	88.33 \downarrow \pm 0.20
BitFit	-	-	-	92.99 \pm 0.38	89.00 \downarrow \pm 0.09	87.51 \downarrow \pm 0.14	83.25 \downarrow \pm 0.08	88.19 \downarrow \pm 0.12

Table 5: Mean and standard deviation results for the 7 tasks by 500, 5k, and 50k samples of training data sets across 20 runs.

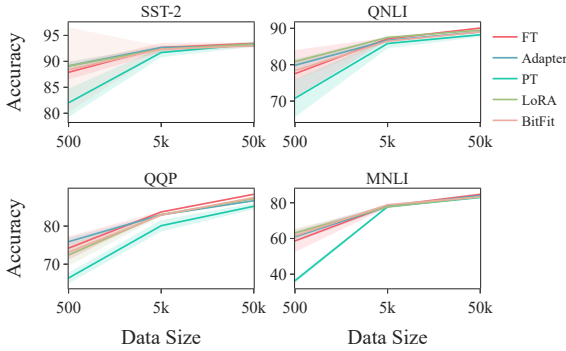


Figure 8: Performance over data scale in 500, 5k, 50k on SST-2, QNLI, QQP, and MNLI. The shaded regions are the standard deviations.

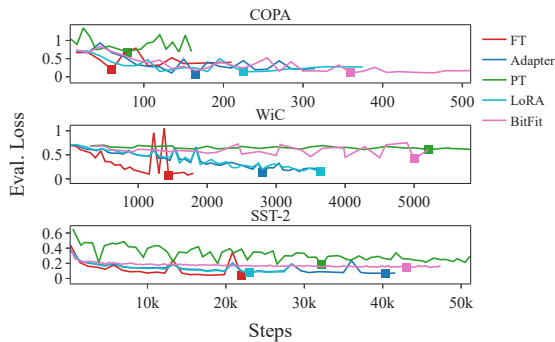


Figure 9: Evaluation loss over training steps on COPA, WiC, and SST-2.

	Small	Medium	Large
Adapter	68.0 \pm 3.42	68.45 \pm 4.17	65.9 \pm 5.42
PT	55.35 \pm 4.71	55.35 \pm 5.07	52.1 \pm 5.24
LoRA	66.4 \pm 9.05	62.4 \pm 8.99	59.8 \pm 9.24

Table 6: Performance over 20 runs on COPA task, controlled by global random seeds, weight initialization (WI) random seeds, and data order (DO) random seeds, respectively.

C.2 Finetuning is More Prone to Overfit

To investigate the reasons behind the performance differences of finetuning and PETuning under different training resources, we plot the evaluation loss over training steps for COPA, WiC, and SST-2, as the representatives of low-, medium-, and high-resource tasks, respectively in Figure 9. We observe that finetuning always converges faster than PETuning, especially on low-resource task COPA, where the training steps are less than 100. One possible explanation for the aforementioned differences is that finetuning may be more prone to overfit on low-resource settings with fewer training iterations, resulting in poorer performance.

C.3 High Stability on Fewer Trainable Parameters.

The probability density curves (Figure 5 and Figure 10) have statistically confirmed PETuning methods tend to exhibit higher stability with fewer train-

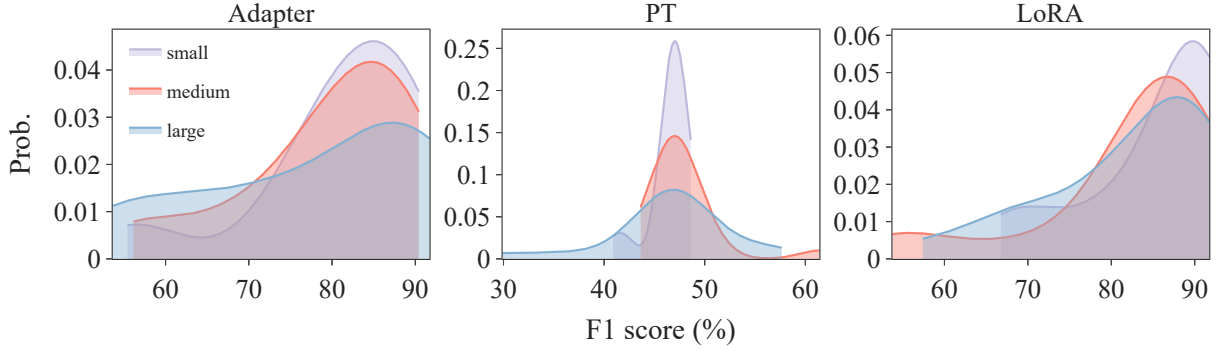


Figure 10: Performance probability density curves of Adapter, prefix tuning (PT), and LoRA over small, medium, and large parameter scales on CB task across 20 runs.

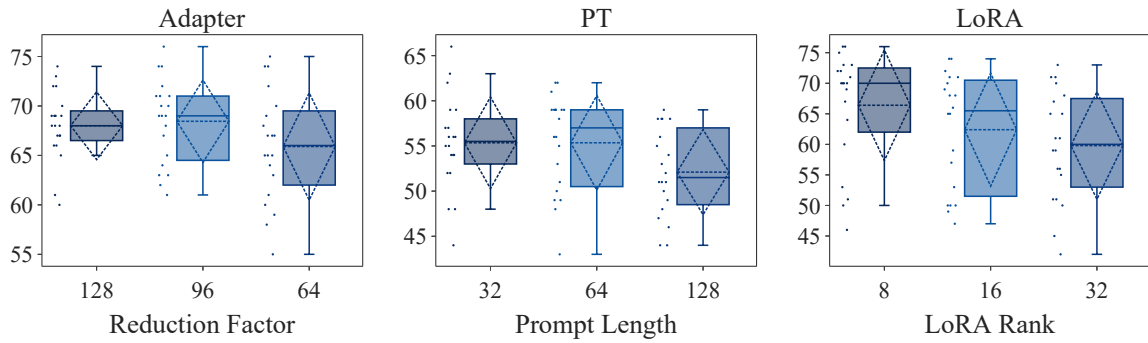


Figure 11: Performance over Adapter, prefix tuning (PT), and LoRA over small, medium, and large parameter scales on COPA task.

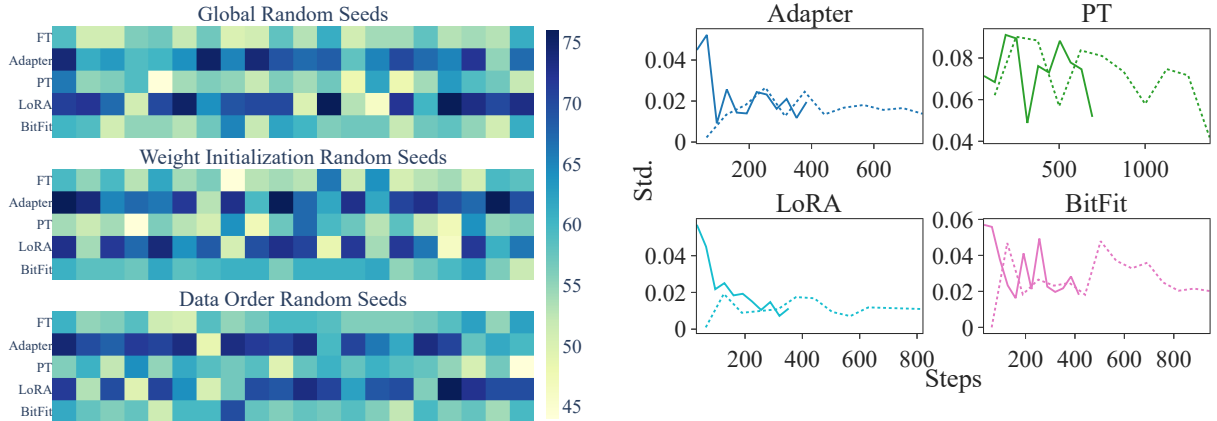


Figure 12: Performance over 20 runs on RTE, controlled by global random seeds, weight initialization (WI) random seeds, and data order (DO) random seeds, respectively.

Figure 13: Standard deviations of data size in {1k (solid line), 2k (dashed line)} over training steps on BoolQ task across 20 runs.

able parameters. In Table 6, we also directly list the numerical results of Adapter, PT, and LoRA over small, medium, and large parameter scales across 20 runs. While the results substantially support our conclusion that Adapter and PT achieve lowest standard deviations on the small parameter scale, except LoRA obtain slightly lower std. on the

medium one. To gain more understanding about the multi-run results, we visualise them in Figure 11. Confirming our conclusion, we can observe that PETuning methods indeed show a trend towards clustering points and smaller boxes on small parameter scale, which means probably higher stability. However, there are also likely to generalise outliers on small parameter scale as shown in Figure 11, especially under our limited 20 runs, which

977
978
979
980
981
982
983
984
985

986 could lead to the increasing variance. This special
987 case might result in the inconsistent phenomenon
988 of LoRA with other PETuning methods, nonethe-
989 less, the results and phenomena in [Table 6](#) and [Fig-
990 ure 11](#) generally further support the conclusion that
991 PETuning are likely to have high stability on fewer
992 trainable parameters.