

# FPT: Improving Prompt Tuning Efficiency via Progressive Training

Anonymous ACL submission

## Abstract

Recently, prompt tuning (PT) has gained increasing attention as a parameter-efficient way of tuning pre-trained language models (PLMs). Despite extensively reducing the number of tunable parameters and achieving satisfying performance, PT is training-inefficient due to its slow convergence. To improve PT’s training efficiency, we first make some novel observations about the prompt transferability of “partial PLMs”, which are defined by compressing a PLM in depth or width. We observe that the soft prompts learned by different partial PLMs of various sizes are similar in the parameter space, implying that these soft prompts could potentially be transferred among partial PLMs. Inspired by these observations, we propose Fast Prompt Tuning (FPT), which starts by conducting PT using a small-scale partial PLM, then progressively expands its depth and width until the full-model size. After each expansion, we recycle the previously learned soft prompts as initialization for the enlarged partial PLM and then proceed PT. We demonstrate the feasibility of FPT on 5 tasks and show that FPT could save over 30% training computations while achieving comparable performance.

## 1 Introduction

The emergence of pre-trained language models (PLMs) has broken the glass ceiling for various NLP tasks (Min et al., 2021). Versatile semantic and syntactic knowledge acquired during pre-training could be leveraged when PLMs are adapted towards a specific downstream task to boost performance. The de facto strategy for such an adaptation is full-parameter fine-tuning, which is computationally expensive and profligate since it requires tuning and storing all the parameters in the PLM for each downstream task.

To remedy this, several parameter-efficient tuning algorithms are proposed in place of the vanilla fine-tuning (Houlsby et al., 2019; Li and Liang,

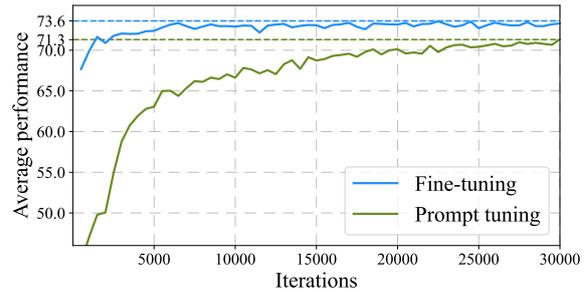


Figure 1: Average performance growth of  $T5_{\text{LARGE}}$  on 5 investigated tasks in this paper comparing fine-tuning and PT. The convergence speed of PT is much slower than fine-tuning in terms of training steps.

2021; Hu et al., 2022; Zaken et al., 2021), among which prompt tuning (PT) (Lester et al., 2021) has gained increasing attention recently. PT prepends a few *virtual tokens* to the input text, these tokens are tuned during training while all the other PLM parameters remain frozen. Despite its simple form, PT has been demonstrated to achieve remarkable performance in various NLP tasks. Especially when the scale of the PLM becomes extremely huge, PT could achieve comparable performance to fine-tuning (Lester et al., 2021). Despite extensively reducing the number of tunable parameters and achieving satisfying performance, PT is criticized to be training-inefficient due to the slow convergence as illustrated in Figure 1, and such incompetence would limit the practical application of PT. Hence in this paper, we explore how to improve PT’s training efficiency.

Our motivation is based on novel observations about the prompt transferability among “partial PLMs”. Here a partial PLM is defined by compressing a PLM in depth or width, which is implemented by dropping several layers or masking part of the connections in the feed-forward network (FFN) in each Transformer (Vaswani et al., 2017) layer. We observe that the soft prompts of the same task learned by different partial PLMs of various

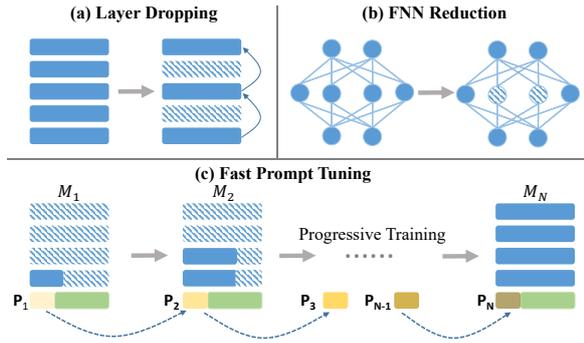


Figure 2: The framework of Fast Prompt Tuning (FPT). The top part (a,b) shows two methods to construct a partial PLM. The bottom part (c) shows FPT’s training process, we conduct PT on a partial PLM, progressively expand its size and transfer the trained prompts.

sizes tend to be close in the parameter space, implying that these soft prompts could potentially be transferred among different partial PLMs.

Inspired by the above observations, we propose Fast Prompt Tuning (FPT), which starts by conducting PT using a small-scale partial PLM to obtain the corresponding soft prompts. After that, we progressively expand the partial PLM’s depth and width until the full-model size by rehabilitating the dropped layers and masked neurons. After each expansion, we recycle the previously learned soft prompts as initialization for the enlarged PLM and then proceed PT. Since the partial PLM requires fewer computations for each step, keeping the total training steps unchanged, we could reduce the overall computations consumed, and in the meantime, achieve comparable PT performance. In experiments, we demonstrate the feasibility of FPT on 5 NLP tasks. The experimental results show that FPT could save around 30% training computations and achieve satisfying downstream performance.

## 2 Prompt Tuning on a Partial PLM

### 2.1 Prompt Tuning

For a given input sequence  $\mathcal{X} = \{x_1, x_2, \dots, x_n\}$  and its target label  $\mathcal{Y}$ , PT first converts  $\mathcal{X}$  into a matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$ , where  $d$  is the hidden size. After that, PT prepends  $l$  tunable soft prompt tokens  $\mathbf{P} \in \mathbb{R}^{l \times d}$  before  $\mathbf{X}$ , creating a new input matrix  $[\mathbf{P}; \mathbf{X}] \in \mathbb{R}^{(l+n) \times d}$ , which is then processed by the PLM. The training objective is to maximize  $\mathcal{P}(\mathcal{Y} | [\mathbf{P}; \mathbf{X}])$ , where only  $\mathbf{P}$  is optimized during training and the parameters of PLM are frozen. Although PT is applied to the entire PLM by default, in this section, we investigate how the performance

would become if we conduct PT on a partial PLM, i.e., only part of the parameters in the PLM participate in the computation.

### 2.2 Partial PLM Construction

Using partial parameters in a PLM is typically applied to reduce the **inference computation** for fine-tuning, such as early exit (Teerapittayanon et al., 2016; Xin et al., 2020) and model pruning (Chen et al., 2020; Sun et al., 2020), which assume that the features produced by a part of a PLM may already suffice to classify some input examples. In this paper, we investigate its application in reducing the **training computation** of PT, and propose to construct partial PLMs by shrinking the original PLM in both depth and width, as illustrated in Figure 2 (a, b). Details are listed in appendix B.

**Layer Dropping.** Based on previous findings (Clark et al., 2019; Jawahar et al., 2019) that adjacent layers in PLMs generally have similar functionalities, we hypothesize that removing part of these layers may not significantly hurt the overall performance, and we propose to drop a PLM’s layers uniformly to construct a partial PLM consisting of fewer layers than the original PLM. After that, we directly build connections among the remaining layers keeping the original order, which is found empirically to work well although such connections do not exist during pre-training.

**FFN Reduction.** Jaszczur et al. (2021) indicate that only part of the neurons in the FFN layers will be activated for a given input. Such a sparse activation phenomenon inspires us to reduce the computation in FFN by shrinking the width of the FFN layer. Specifically, the FFN layer consists of two fully connected networks with a nonlinear activation function  $\sigma$ , and it processes an input representation  $\mathbf{x} \in \mathbb{R}^d$  as:  $\text{FFN}(\mathbf{x}) = \sigma(\mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2$ , where  $\mathbf{W}_1 \in \mathbb{R}^{d \times d'}$  and  $\mathbf{W}_2 \in \mathbb{R}^{d' \times d}$  are the weight matrices,  $\mathbf{b}_1 \in \mathbb{R}^{d'}$  and  $\mathbf{b}_2 \in \mathbb{R}^d$  are the bias terms. We abandon a portion of  $\mathbf{W}_1 / \mathbf{W}_2$ ’s columns / rows (i.e., reducing  $d'$ ) by masking the neurons that are seldom activated. In practice, before training, we feed a few downstream examples prepended by randomly initialized soft prompts into the full-size PLM and record the neuron activation of each dimension of  $d'$ .

**Compound Reduction.** Since the above methods are compatible with each other, we try to combine them to form a partial PLM smaller than the

	Enc./Dec. Layer	FFN Dimension	MNLI (Acc)	QQP (Acc)	SQUAD2.0 (EM)	ReCoRD (EM)	XSUM (ROUGE-L)	Avg.	$\Delta$	FLOPs
<b>PT</b>	24 / 24	2,816	<b>86.07</b>	<b>87.26</b>	<b>76.09</b>	<b>81.46</b>	<b>26.65</b>	<b>71.51</b>	-	100%
<b>LD</b>	6 / 6	2,816	60.34	78.29	48.14	24.75	17.40	45.78	-25.73	30%
	12 / 12	2,816	63.90	80.64	52.87	39.09	19.69	51.24	-20.27	54%
	18 / 18	2,816	81.41	86.05	63.97	59.87	22.51	62.76	-8.75	77%
<b>FR</b>	24 / 24	704	78.18	85.19	66.68	62.90	22.46	63.08	-8.43	58%
	24 / 24	1,408	82.62	86.45	72.65	74.59	24.61	68.19	-3.32	72%
	24 / 24	2,112	<u>84.93</u>	<u>86.77</u>	<u>74.73</u>	<u>79.52</u>	<u>26.12</u>	70.41	-1.10	86%
<b>CR</b>	6 / 6	704	62.53	78.38	48.62	23.99	16.49	46.00	-25.51	20%
	12 / 12	1,408	64.09	78.91	50.90	29.50	18.88	48.45	-23.06	40%
	18 / 18	2,112	80.63	86.32	63.42	58.97	22.18	62.30	-9.21	66%

Table 1: Average results for partial PLM PT on  $T5_{LARGE}$  with layer dropping (**LD**), FFN Reduction (**FR**), and compound reduction (**CR**).  $\Delta$  denotes the performance degeneration compared with vanilla **PT** of each setting.

original PLM in both depth and width.

### 2.3 Observations

To explore PT’s performance on a partial PLM, we conduct experiments on  $T5_{LARGE}$  (Raffel et al., 2020). We choose 5 representative NLP datasets in English, covering the tasks of natural language inference (MNLI (Williams et al., 2018)), paraphrase (QQP (link)), reading comprehension (SQUAD2.0 (Rajpurkar et al., 2018) and RECoRD (Zhang et al., 2018)), and summarization (XSUM (Narayan et al., 2018)). For both layer dropping and FFN reduction, we evaluate the performance when we reduce the number of Transformer layers or FFN intermediate dimension to  $\{\frac{1}{4}, \frac{1}{2}, \frac{3}{4}\}$ . We train all models using the same steps and the details are described in appendix B.

**Overall Performance.** The overall results are shown in Table 1. We observe that for each method, **despite abandoning a large portion of parameters, a partial PLM reserves most of the PT performance of the full-size PLM.** As expected, the performance becomes better when more parameters are retained. In addition, we find that the performance drop is less sensitive to FFN reduction than layer dropping. Specifically, there is only 1.10% performance drop on average when 25% neurons are masked. These results indicate that the resulting partial PLM still retains most of the functionalities of the original PLM.

**Prompt Embedding Visualization.** Taking a step further, we visualize the learned prompt embeddings of different partial PLMs using t-SNE (Van der Maaten and Hinton, 2008) in Figure 3, and describe the details in appendix C. We observe that **for the same task, the soft prompts obtained by different partial PLMs tend to form**

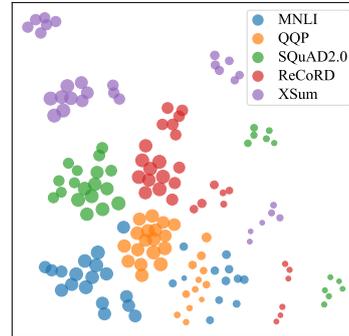


Figure 3: Visualization of 5 investigated tasks’ soft prompts of different partial PLMs. A marker with a larger size means the performance of the corresponding soft prompts on the partial PLM is better.

**a compact cluster in the parameter space.** This phenomenon implies that the soft prompts corresponding to the same task (1) have a great potential of transferring among different partial PLMs, and (2) could serve as a better initialization that leads to faster convergence. Apart from the visualization, we further report the cosine similarity of the learned prompts in appendix D to verify the above phenomenon from another aspect.

## 3 Fast Prompt Tuning

In this section, we propose Fast Prompt Tuning (FPT), which aims at accelerating PT via *progressive training* (Gong et al., 2019). Progressive training is typically leveraged for improving pre-training efficiency, instead, we focus on its application in PLM’s downstream adaptation.

### 3.1 Methodology

Formally speaking, as visualized in Figure 2 (c), we split the original PT training process into  $N$  stages. We start with a small-size partial PLM  $\mathcal{M}_1$  and

Method		MNLI (Acc)	QQP (Acc)	SQUAD2.0 (EM)	RECORD (EM)	XSUM (ROUGE-L)	Avg.	FLOPs	Improve $\uparrow$
T5 <sub>LARGE</sub>	PT	86.07	<b>87.26</b>	76.09	<b>81.46</b>	<b>26.65</b>	<b>71.51</b>	100%	-
	FPT <sub>LD</sub>	85.72	86.51	75.89	80.23	26.27	70.92	72%	0.11
	FPT <sub>FR</sub>	<b>86.49</b>	87.11	<b>76.26</b>	81.07	26.55	71.50	83%	<b>0.49</b>
	FPT <sub>CR</sub>	85.13	86.40	75.63	81.00	26.21	70.87	<b>65%</b>	0.38
T5 <sub>XL</sub>	PT	89.00	88.20	81.08	<b>88.48</b>	<b>30.53</b>	75.46	100%	-
	FPT <sub>LD</sub>	88.99	88.09	<b>82.18</b>	88.06	30.52	<b>75.57</b>	86%	<b>0.78</b>
	FPT <sub>FR</sub>	88.84	<b>88.21</b>	81.74	88.46	30.52	75.55	84%	0.76
	FPT <sub>CR</sub>	<b>89.18</b>	87.34	80.88	87.82	30.43	75.13	<b>74%</b>	0.48

Table 2: Performance of the vanilla PT and three variants of our method. FPT<sub>LD</sub>, FPT<sub>FR</sub>, and FPT<sub>CR</sub> refer to constructing partial PLMs by layer dropping, FFN reduction, and compound reduction. The ‘‘FLOPs’’ column shows the average relative FLOPs consumed on 5 tasks compared with the full-size PLM. The column ‘‘Improve $\uparrow$ ’’ denotes the performance improvement of each FPT<sub>\*</sub> method over PT when PT uses the same FLOPs as FPT<sub>\*</sub>.

then progressively rehabilitate its depth and width until the full-size model  $\mathcal{M}_N$ , creating a series of partial PLMs  $\{\mathcal{M}_i\}_{i=1}^{N-1}$  with growing sizes. The architectures of the partial PLMs are constructed using the same method in § 2.2.

During each training stage  $i$ , we conduct PT on a partial PLM  $\mathcal{M}_i$  and obtain the learned soft prompts  $\mathbf{P}_i$ . Based on the observation that  $\mathcal{M}_i$  retains a large portion of functionalities of the full-size PLM  $\mathcal{M}_N$ , we conjecture that  $\mathcal{M}_i$  could serve as a perfect substitute for  $\mathcal{M}_N$  and learn how to deal with the downstream task. In addition, considering that the soft prompts learned by different partial PLMs are close in the parameter space, we could transfer the knowledge learned by  $\mathcal{M}_i$  to  $\mathcal{M}_{i+1}$  through recycling  $\mathbf{P}_i$ . Specifically, after each model expansion, we directly use  $\mathbf{P}_i$  as initialization for training  $\mathcal{M}_{i+1}$  in the next stage. Since for each partial PLM, fewer parameters participate in both the forward and backward process, the computations could be reduced. Keeping the total number of training steps the same, FPT could accelerate training compared with the vanilla PT.

### 3.2 Experiments and Analyses

We follow most of the experimental settings in § 2 and also describe the training details in appendix B. We report FLOPs for the vanilla PT and FPT to compare training efficiency. We evaluate both T5<sub>LARGE</sub> and T5<sub>XL</sub> (a larger T5 model) on each task and train for 30k and 15k steps, respectively. We test FPT’s performance when we progressively expand the model’s depth, width, and both of them. Unless otherwise specified, for most of FPT’s methods, we split the training process into 4 stages. Each of the first three stages takes 20% steps, while the last stage takes 40% steps.

**Results.** We list the results in Table 2, from which we observe that (1) on average, all three variants of FPT achieve comparable performance with PT and utilize fewer computations (e.g., FPT<sub>CR</sub> saves around 30% FLOPs). On several tasks (e.g., MNLI and SQUAD2.0), FPT even exceeds PT’s performance; (2) combining both layer dropping and FFN reduction (i.e., FPT<sub>CR</sub>) is more training-efficient. However, we also observe that saving more computations generally leads to poorer performance. Among all three variants of FPT, FPT<sub>FR</sub> strikes the best balance between the performance and training efficiency; (3) moreover, we compare both PT and FPT’s performance when PT consumes the same computations as each variant of FPT. As reflected in the column ‘‘Improve $\uparrow$ ’’, controlling the training computations the same, our FPT outperforms PT, and the improvement is more significant for T5<sub>XL</sub> than T5<sub>LARGE</sub>, showing that FPT has a great potential to apply to large-scale PLMs. We also verify the effectiveness of our partial model construction designs in appendix E, and show in appendix F that the performance of FPT is not sensitive to the duration of each training stage. We leave the explorations on other tasks and the effect of training budgets as future work.

## 4 Conclusion

In this work, towards improving PT’s training efficiency, we first make several insightful observations by conducting PT on partial PLMs, and then propose FPT based on the observations. The results on 5 datasets demonstrate the feasibility of FPT in saving the training computations. Being the first attempt towards accelerating PT, we encourage future work to design more sophisticated algorithms to further improve PT’s training efficiency.

280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
  
296  
297  
298  
299  
  
300  
301  
302  
303  
304  
305  
306  
  
307  
308  
309  
310  
311  
312  
313  
  
314  
315  
316  
317  
318  
319  
320  
321  
322  
  
323  
324  
325  
326  
327  
328  
329  
  
330  
331  
332  
333  
334  
335  
336  
337

## References

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared 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 M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher 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 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.

Cheng Chen, Yichun Yin, Lifeng Shang, Xin Jiang, Yujia Qin, Fengyu Wang, Zhi Wang, Xiao Chen, Zhiyuan Liu, and Qun Liu. 2021. [bert2bert: Towards reusable pretrained language models](#).

Tianlong Chen, Jonathan Frankle, Shiyu Chang, Sijia Liu, Yang Zhang, Zhangyang Wang, and Michael Carbin. 2020. [The lottery ticket hypothesis for pre-trained BERT networks](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.

Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. 2019. [What does BERT look at? an analysis of BERT’s attention](#). In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 276–286, Florence, Italy. Association for Computational Linguistics.

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, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Linyuan Gong, Di He, Zhuohan Li, Tao Qin, Liwei Wang, and Tie-Yan Liu. 2019. [Efficient training of BERT by progressively stacking](#). In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pages 2337–2346. PMLR.

Xiaotao Gu, Liyuan Liu, Hongkun Yu, Jing Li, Chen Chen, and Jiawei Han. 2021. [On the transformer growth for progressive BERT training](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5174–5180, Online. Association for Computational Linguistics.

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*. 338–342

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, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pages 2790–2799. PMLR. 343–351

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*. 352–358

Sebastian Jaszczur, Aakanksha Chowdhery, Afroz Mohiuddin, Lukasz Kaiser, Wojciech Gajewski, Henryk Michalewski, and Jonni Kanerva. 2021. [Sparse is enough in scaling transformers](#). In *Advances in Neural Information Processing Systems*. 359–361

Ganesh Jawahar, Benoît Sagot, and Djamel Seddah. 2019. [What does BERT learn about the structure of language?](#) In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3651–3657, Florence, Italy. Association for Computational Linguistics. 362–367

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, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics. 368–374

Xiang Lisa Li and Percy Liang. 2021. [Prefix-tuning: Optimizing continuous prompts for generation](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4582–4597, Online. Association for Computational Linguistics. 375–381

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](#). 382–387

Rabeeh Karimi Mahabadi, James Henderson, and Sebastian Ruder. 2021. [Compacter: Efficient low-rank hypercomplex adapter layers](#). In *Advances in Neural Information Processing Systems*. 388–391

Bonan Min, Hayley Ross, Elior Sulem, Amir Pouran Ben Veyseh, Thien Huu Nguyen, Oscar Sainz,

394	Eneko Agirre, Ilana Heinz, and Dan Roth. 2021. <a href="#">Recent advances in natural language processing via large pre-trained language models: A survey</a> . <i>ArXiv preprint</i> , abs/2111.01243.	453
395		454
396		455
397		456
398	Shashi Narayan, Shay B. Cohen, and Mirella Lapata. 2018. <a href="#">Don't give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization</a> . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 1797–1807, Brussels, Belgium. Association for Computational Linguistics.	457
399		458
400		459
401		460
402		461
403		462
404		463
405	Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Z. Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. <a href="#">Pytorch: An imperative style, high-performance deep learning library</a> . In <i>Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada</i> , pages 8024–8035.	464
406		465
407		466
408		467
409		468
410		469
411		470
412		471
413		472
414		473
415		474
416		475
417		476
418	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. <a href="#">Exploring the limits of transfer learning with a unified text-to-text transformer</a> . <i>Journal of Machine Learning Research</i> , 21(140):1–67.	477
419		478
420		479
421		480
422		481
423		482
424	Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. <a href="#">Know what you don't know: Unanswerable questions for SQuAD</a> . In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 784–789, Melbourne, Australia. Association for Computational Linguistics.	483
425		484
426		485
427		486
428		487
429		488
430		489
431	Andreas Rücklé, Gregor Geigle, Max Glockner, Tilman Beck, Jonas Pfeiffer, Nils Reimers, and Iryna Gurevych. 2021. <a href="#">AdapterDrop: On the efficiency of adapters in transformers</a> . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 7930–7946, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	490
432		491
433		492
434		493
435		494
436		495
437		496
438		497
439	Noam Shazeer and Mitchell Stern. 2018. <a href="#">Adafactor: Adaptive learning rates with sublinear memory cost</a> . In <i>Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholm, Sweden, July 10-15, 2018</i> , volume 80 of <i>Proceedings of Machine Learning Research</i> , pages 4603–4611. PMLR.	498
440		499
441		500
442		501
443		502
444		503
445		504
446	Zhiqing Sun, Hongkun Yu, Xiaodan Song, Renjie Liu, Yiming Yang, and Denny Zhou. 2020. <a href="#">MobileBERT: a compact task-agnostic BERT for resource-limited devices</a> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 2158–2170, Online. Association for Computational Linguistics.	505
447		506
448		507
449		508
450		509
451		510
452		
	Surat Teerapittayanon, Bradley McDanel, and Hsiang-Tsung Kung. 2016. <a href="#">Branchynet: Fast inference via early exiting from deep neural networks</a> . In <i>2016 23rd International Conference on Pattern Recognition (ICPR)</i> , pages 2464–2469. IEEE.	
	Laurens Van der Maaten and Geoffrey Hinton. 2008. <a href="#">Visualizing data using t-sne</a> . <i>Journal of machine learning research</i> , 9(11).	
	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. <a href="#">Attention is all you need</a> . In <i>Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA</i> , pages 5998–6008.	
	Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. <a href="#">A broad-coverage challenge corpus for sentence understanding through inference</a> . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 1112–1122, New Orleans, Louisiana. Association for Computational Linguistics.	
	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. <a href="#">Transformers: State-of-the-art natural language processing</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 38–45, Online. Association for Computational Linguistics.	
	Ji Xin, Raphael Tang, Jaejun Lee, Yaoliang Yu, and Jimmy Lin. 2020. <a href="#">DeeBERT: Dynamic early exiting for accelerating BERT inference</a> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 2246–2251, Online. Association for Computational Linguistics.	
	Zhilin Yang, Zihang Dai, Yiming Yang, Jaime G. Carbonell, Ruslan Salakhutdinov, and Quoc V. Le. 2019. <a href="#">Xlnet: Generalized autoregressive pretraining for language understanding</a> . In <i>Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada</i> , pages 5754–5764.	
	Elad Ben Zaken, Shauli Ravfogel, and Yoav Goldberg. 2021. <a href="#">Bitfit: Simple parameter-efficient fine-tuning for transformer-based masked language-models</a> . <i>ArXiv preprint</i> , abs/2106.10199.	
	Sheng Zhang, Xiaodong Liu, Jingjing Liu, Jianfeng Gao, Kevin Duh, and Benjamin Van Durme. 2018. <a href="#">ReCoRD: Bridging the gap between human and machine commonsense reading comprehension</a> .	

## 511 Appendices

### 512 A Related Work

513 **Prompt Tuning.** PLMs have achieved excellent  
514 performance on many NLP tasks relying on their  
515 powerful natural language understanding and gener-  
516 ation capabilities (Devlin et al., 2019; Liu et al.,  
517 2019; Yang et al., 2019). However, with the emer-  
518 gence of large-scale PLMs like T5 (Raffel et al.,  
519 2020) and GPT-3 (Brown et al., 2020), tuning all  
520 the parameters of a PLM (i.e., full-parameter fine-  
521 tuning), which requires huge storage and mem-  
522 ory costs, is not flexible for real-world applica-  
523 tions on massive downstream tasks. Therefore,  
524 parameter-efficient tuning methods (Houlsby et al.,  
525 2019; Mahabadi et al., 2021; Hu et al., 2022; Zaken  
526 et al., 2021; He et al., 2022; Rücklé et al., 2021)  
527 attract more and more attention, among which  
528 prompt tuning (PT) (Lester et al., 2021) is a sim-  
529 ple and effective one. By prepending a few train-  
530 able embeddings before the input sequence, PT can  
531 achieve comparable performance to full-parameter  
532 fine-tuning. With the size of PLM getting larger,  
533 the performance of PT gets closer to vanilla fine-  
534 tuning (Lester et al., 2021), showing great potential  
535 to utilize extremely large PLMs in future. How-  
536 ever, due to the slow convergence shown in Fig-  
537 ure 1, PT’s training efficiency becomes a serious  
538 drawback and may limit its practical application.

539 **Progressive Training.** Considering that pre-  
540 training usually requires tremendous computational  
541 resources, researchers propose *progressive train-*  
542 *ing* to improve the training efficiency (Gong et al.,  
543 2019). Progressive training starts training using a  
544 shallow model, and gradually grows the depth of  
545 the model along the training process by replicat-  
546 ing existing layers (parameter recycling). In this  
547 way, the pre-training efficiency can be improved  
548 a lot. To further improve training efficiency, later  
549 works propose to progressively grow PLMs in both  
550 depth and width (Gu et al., 2021), and design better  
551 initialization methods to inherit the functionality  
552 of existing models (Chen et al., 2021). Instead of  
553 leveraging progressive training during the process  
554 of pre-training, we apply it to PLM’s downstream  
555 adaptation, with a focus on PT. Furthermore, con-  
556 ventional progressive training duplicates existing  
557 parameters to grow a PLM’s size until the full-  
558 model’s size. Instead, we have already obtained  
559 a full-size PLM, and propose to construct partial  
560 models with growing sizes by dropping / masking

existing parameters.

### 561 B Implementation Details

562 Our implementation is based on PyTorch (Paszke  
563 et al., 2019) and transformers (Wolf et al., 2020).  
564 The experiments are conducted with 8 NVIDIA  
565 32GB V100 GPUs, and each experiment requires  
566 fewer than 10 hours to finish.

567 **Partial PLM Construction.** As mentioned in  
568 § 2.2, we design three methods to construct par-  
569 tial PLMs. Specifically, for **layer dropping**, we  
570 select layers uniformly. For example, to select 3  
571 layers out of a 24-layer PLM, we will select layer  
572 {1, 12, 24} to construct the partial PLM. For **FFN**  
573 **reduction**, to estimate the activation of each neu-  
574 ron (dimension) in FFN layer  $l$ , we first randomly  
575 sample 1,000 examples to form a small dataset  $\mathcal{D}$ .  
576 We prepend each example  $\mathcal{X}$  (without the label) in  
577  $\mathcal{D}$  with randomly initialized soft prompts and feed  
578 it into the full-size PLM to obtain the input repre-  
579 sentation  $\mathbf{x}^l$  of FFN layer. After that, we obtain the  
580 activation score of each neuron using the follow-  
581 ing equation  $S = \sum_{\mathcal{X} \in \mathcal{D}} \sum_{i=1}^{|\mathcal{X}|} |\sigma(\mathbf{x}_i^l \mathbf{W}_1^l + \mathbf{b}_1^l)|$ ,  
582 where  $\mathbf{W}_1^l, \mathbf{b}_1^l$  are the parameters in FFN layer  $l$ ,  
583 and  $|\mathcal{X}|$  denotes the sequence length. The neurons  
584 (dimensions) with smaller activation score (i.e., sel-  
585 dom activated) will be masked. Note that the T5  
586 model is composed of both an encoder and a de-  
587 coder, due to the difference of the input length and  
588 output length on various tasks, the computation  
589 overload of the encoder and decoder may vary a  
590 lot. Therefore, for the tasks (MNLI and QQP) that  
591 have a lighter computation overload on the decoder  
592 (i.e., small output length), shrinking the decoder  
593 model size has little impact on saving the computa-  
594 tional costs, hence we retain the whole decoder un-  
595 der this setting; for other three tasks (SQUAD2.0,  
596 RECORD and XSUM), the output length on de-  
597 coder is much longer and we conduct partial PLM  
598 construction on both the encoder and decoder. We  
599 calculate FLOPs for each experiment using the *pt-*  
600 *flops* tool <sup>1</sup>, and report the average FLOPs of 5  
601 tasks in Table 1 and Table 2.

602 **Partial PLM Prompt Tuning.** We use T5<sub>LARGE</sub>  
603 for our experiments of partial PLM PT. Following  
604 Lester et al. (2021), we leverage the LM-adapted  
605 version of T5 checkpoints, which are additionally  
606 trained for 100k steps. The adapted version of T5  
607

<sup>1</sup><https://github.com/sovrasov/flops-counter.pytorch>

Partial PLMs	Layer Dropping		FFN Reduction		Compound Reduction		Training Steps	
	$L_{enc}/L_{dec}$	$d'$	$L_{enc}/L_{dec}$	$d'$	$L_{enc}/L_{dec}$	$d'$		
T5 <sub>LARGE</sub>	$\mathcal{M}_1$	6 / 6	2,816	24 / 24	704	6 / 6	704	6,000
	$\mathcal{M}_2$	12 / 12	2,816	24 / 24	1,408	12 / 12	1,408	6,000
	$\mathcal{M}_3$	18 / 18	2,816	24 / 24	2,112	18 / 18	2,112	6,000
	$\mathcal{M}_4$	24 / 24	2,816	24 / 24	2,816	24 / 24	2,816	12,000
T5 <sub>XL</sub>	$\mathcal{M}_1$	18 / 18	5,120	24 / 24	1,280	18 / 18	1,280	4,000
	$\mathcal{M}_2$	18 / 18	5,120	24 / 24	2,560	18 / 18	2,560	4,000
	$\mathcal{M}_3$	18 / 18	5,120	24 / 24	3,840	18 / 18	3,840	4,000
	$\mathcal{M}_4$	24 / 24	5,120	24 / 24	5,120	24 / 24	5,120	8,000

Table 3: Architecture details of the partial PLMs on the three construction methods for both T5<sub>LARGE</sub> and T5<sub>XL</sub>.  $L_{enc}$  and  $L_{dec}$  denote the number of layers of the encoder and decoder of the partial model  $\mathcal{M}_i$ , respectively. We also list the training steps for each stage in the last column.

has been demonstrated to achieve stable and better PT performance. For the implementation of PT, we set the prompt length to 20 and randomly initialize the soft prompts. The optimizer is chosen as Adafactor (Shazeer and Stern, 2018) and the learning rate is set to 0.3. We choose a batch size of 32 and the greedy decoding to generate the predictions. The training steps are set to 30k to ensure that PT will not get stuck in a local optimum. We run all the experiments 3 times with different random seeds and report the average results.

**Fast Prompt Tuning.** For the implementations of FPT, we train T5<sub>LARGE</sub> / T5<sub>XL</sub> with a total step of 30k / 15k. The number of training steps of T5<sub>XL</sub> is chosen smaller than T5<sub>LARGE</sub> since we find empirically that T5<sub>XL</sub> converges faster than T5<sub>LARGE</sub>. As mentioned in § 3.2, unless otherwise specified, we split the whole training process into 4 stages. Each of the first three stages takes 20% of the training steps, while the last stage (full-model PT) takes 40% training stages. Except for layer dropping on T5<sub>XL</sub>, we find that a partial PLM, with fewer than 12 layers in either the encoder or decoder, achieves poor PT performance. Therefore, we only use two training stages where the first stage takes 60% training steps and the second stage takes 40% training steps. More detailed settings about the partial model construction are shown in Table 3. The experiments with T5<sub>LARGE</sub> are run three times with different random seeds and the average results are reported while experiments with T5<sub>XL</sub> are conducted once due to their huge computation consumption.

## C Prompt Embedding Visualization

In Figure 3, we visualize the soft prompts of different partial PLMs and tasks in Table 1. The embed-

ding used for visualization is derived by averaging soft prompt along the token length dimension. As described in § 2.3, we run each experiment three times with different random seeds to get stable results. Therefore, we plot 30 points (3 runs  $\times$  (9 partial PLM + 1 full-size PLM)) for each task in Figure 3. And the size of the marker in the figure denotes the performance of the soft prompts on corresponding partial PLMs. Larger size indicates better performance. We can observe that soft prompts with better performance will be easier to form a compact cluster.

## D Prompt Embedding Similarity

To further gain insights on the transferability of the soft prompts learned by T5<sub>LARGE</sub>'s different partial PLMs defined in Table 3, in addition to the visualization conducted in § 2.3, we calculate the average cosine similarity of the soft prompts corresponding to different tasks in Table 4. Specifically, for different partial PLMs  $\mathcal{M}_1, \mathcal{M}_2, \dots, \mathcal{M}_{N-1}$  and the full-size PLM  $\mathcal{M}_N$ , we conduct PT with each model  $\mathcal{M}_i$  on the task  $\mathcal{T}_j$  and obtain the corresponding soft prompts  $\mathbf{P}_i^j \in \mathbb{R}^{l \times d}$ . Then we average  $\mathbf{P}_i^j$  along the token length dimension, and get a vector  $\bar{\mathbf{P}}_i^j \in \mathbb{R}^d$ . After that, we calculate  $\mathcal{S}(\mathcal{T}_j^P, \mathcal{T}_k^F)$  (average cosine similarity between (1) task  $j$ 's partial PLMs' prompts and (2) task  $k$ 's full PLM's prompts) using the following metric:

$$\mathcal{S}(\mathcal{T}_j^P, \mathcal{T}_k^F) = \frac{1}{N-1} \sum_{i=1}^{N-1} \frac{\bar{\mathbf{P}}_i^j \cdot \bar{\mathbf{P}}_N^k}{\|\bar{\mathbf{P}}_i^j\| \|\bar{\mathbf{P}}_N^k\|} \quad (1)$$

From the results in Table 4, we observe that the highest similarity is achieved when  $j = k$ , showing that the prompts of the partial PLMs are closer to the same task's prompts of the full-size model. This phenomenon is aligned with the observation

$\mathcal{T}^P \backslash \mathcal{T}^F$	MNLI	QQP	SQUAD2.0	ReCoRD	XSUM
MNLI	<b>0.249</b>	0.131	0.175	0.109	0.139
QQP	0.145	<b>0.177</b>	0.135	0.103	0.126
SQUAD2.0	0.202	0.143	<b>0.286</b>	0.190	0.202
ReCoRD	0.167	0.119	0.219	<b>0.224</b>	0.195
XSUM	0.164	0.128	0.237	0.188	<b>0.301</b>

Table 4: Average prompt similarity ( $\mathcal{S}(\mathcal{T}_j^P, \mathcal{T}_k^F)$ ) among different tasks. The highest score in each row is **highlighted**.

Selection Method	Performance	Relative FLOPs
<i>Layer Dropping</i>		
Uniform	<b>70.92</b>	72%
Last	69.91	
<i>FFN Reduction</i>		
Activation	<b>71.50</b>	84%
Random	66.80	

Table 5: Average performance on 5 investigated tasks using different strategies of layer dropping and FFN reduction on T5<sub>LARGE</sub>.

in Figure 3, implying that on the same task, the soft prompts learned by partial PLMs could be potentially transferred to the full-size PLM.

## E Effect of Partial Model Construction Designs for FPT

We construct a partial PLM by dropping a few layers or masking some neurons. As mentioned in § 2.2, for layer dropping, we retain the layers uniformly; for FFN reduction, we mask the neurons that are less likely to be activated. How to select the retained parameters is essential to the performance of FPT. To demonstrate this, in Table 5, we experiment FPT with another strategy for layer dropping and FFN reduction, respectively.

For layer dropping, we compare our strategy of dropping layers uniformly (denoted as *Uniform*) with dropping the last few layers (denoted as *Last*). For both methods, we retain the same number of layers. For example, in order to select 3 layers from a 24-layer PLM, the *Uniform* strategy will retain the layer {1, 12, 24}, and the *Last* strategy will retain the layer {1, 2, 3}. From Table 5, we can derive that the *Uniform* strategy is slightly better than the *Last* strategy. We hypothesize the reason is that the overall functionalities of a PLM are uniformly distributed among different layers, and adjacent layers tend to have similar functionalities. Therefore, retaining layers uniformly tends to re-

serve more functionalities than only retaining the first few layers.

For FFN reduction, we compare our strategy of masking neurons based on the activation score (denoted as *Activation*) with randomly masking neurons (denoted as *Random*). For the *Activation* strategy, we feed 1000 samples prepended by randomly initialized soft prompts into the PLM, and then record the activation score of neurons along each dimension. The results in Table 5 show that the *Activation* strategy significantly outperforms the *Random* strategy, demonstrating the effectiveness of our method. Randomly masking neurons may abandon those highly activated (most informative) ones, which hinders PT’s convergence. We also find empirically the activation score of each neuron in FFN layer may vary a lot across different tasks, which means different neurons may respond differently to the input. This phenomenon also implies that there may exist some “functional partitions” in the FFN layers of PLMs.

## F Effect of Duration for Each Training Stage

To show the effects of the duration of each training stage, following Gong et al. (2019), we conduct experiments on MNLI using T5<sub>LARGE</sub> with three proposed variants of FPT, and evaluate the effects of training duration for the last two stages.

Specifically, for layer dropping of FPT, we conduct PT on the 18-layer partial PLM for 15k steps, and save the learned soft prompts every 3k steps to get  $15/3 = 5$  sets of soft prompts. Then using each of these 5 soft prompts as the initialization, we conduct PT with the full-size PLM for 3k steps. We report the validation performance and compare FPT with the vanilla PT. For FFN reduction and compound reduction of FPT, we conduct similar experiments except that we start from a partial PLM using different construction methods.

The results are shown in Figure 4, from which

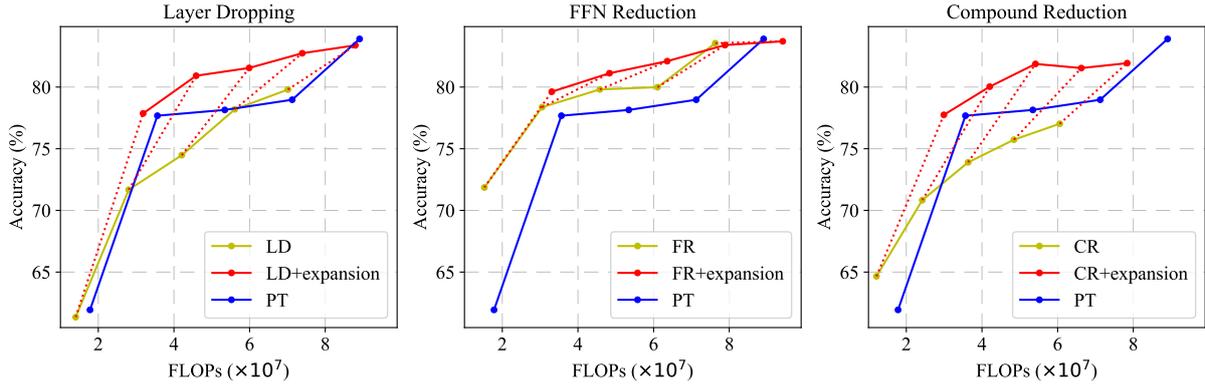


Figure 4: The validation performance on MNLI with different training duration for the last two stages. We conduct this ablation study for each of the three variants of FPT. We compare our FPT with different expanding time (red line) with the vanilla PT (blue line) and PT without model expansion (yellow line). Each red dot is connected with a yellow dot by a dashed line, indicating it is initialized by the yellow dot and optimized by conducting PT on full-size PLM.

746 we can see that expanding the partial PLM's size  
747 and then conducting PT (i.e., the red line) performs  
748 better than only conducting PT on the partial PLM  
749 (i.e., the yellow line). In addition, comparing our  
750 FPT (i.e., the red line) with the vanilla PT (i.e., the  
751 blue line), there is a specific threshold  $s'$  of training  
752 steps, if we expand the partial PLM before  $s'$ ,  
753 the training efficiency can be improved compared  
754 with the vanilla PT; however, after  $s'$ , expanding  
755 the partial PLM and continuing PT on it does not  
756 bring consistent improvement over the vanilla PT.  
757 In general, expanding the partial PLM between 3k  
758 steps and 12k steps works well for all three variants  
759 of FPT, indicating that within a reasonably broad  
760 range, the performance improvement of FPT is not  
761 sensitive to the duration of each training stage. We  
762 aim to explore how to decide the optimal duration  
763 for each training stage in future to make our FPT  
764 more practical.