

Global Prompt Cell: A Portable Control Module for Effective Prompt Tuning

Chi Liu, Haochun Wang, Nuwa Xi, Sendong Zhao, and Bing Qin

Harbin Institute of Technology
{cliu, hcwang, nwxi, sdzhao, bqin}@ir.hit.edu.cn

Abstract. As a novel approach to tuning pre-trained models, prompt tuning involves freezing the parameters in downstream tasks while inserting trainable embeddings into inputs in the first layer. However, previous methods have mainly focused on the initialization of prompt embeddings. The strategy of training and utilizing prompt embeddings in a reasonable way has become a limiting factor in the effectiveness of prompt tuning. To address this issue, we introduce the Global Prompt Cell (GPC), a portable control module for prompt tuning that selectively preserves prompt information across all encoder layers. Our experimental results demonstrate a 5.8% improvement on SuperGLUE datasets compared to vanilla prompt tuning.

1 Introduction

Prompt-based methods can be classified into two categories: discrete prompt tuning [15,1] and continuous prompt tuning [8]. Discrete prompt tuning transforms the task into a “fill-in-the-blank” format and then utilizes a pre-trained language model to predict the answer, which operates similarly to the masked language model (MLM) [2]. In subsequent research, continuous prompt tuning [8] introduced soft prompts (i.e., prompt embeddings) to replace manual templates, which consist of special tokens with adjustable embeddings. We refer to continuous prompt tuning as “prompt tuning” for simplicity.

Vanilla prompt tuning, which concatenates prompt embeddings with input tokens in the first layer and updates only the parameters of the prompt embeddings during the training phase, has several limitations [8,21]. Firstly, since the effectiveness of prompt embeddings is highly related to the length, it is necessary to use prompt embeddings with hundreds of tokens in length to achieve better downstream task performance, as suggested by [8] and [21]. However, long prompt embeddings also reduce the possible length of input text. In addition, prompt learning requires a longer time to converge compared to full fine-tuning [22], and its effectiveness still has significant room for improvement.

These drawbacks are due to the traditional approach of inserting prompt embeddings into the input layer and concatenating them with token embeddings for model input in prompt learning. However, prompt embeddings have significant differences compared to token embeddings. Firstly, prompt embeddings have not

undergone pre-training and require more optimization steps compared to token embeddings. Secondly, prompt embeddings do not have semantic information but serve as task-specific vectors to guide the model for downstream tasks. Finally, in pre-trained models with multiple layers, prompt learning freezes the parameters during training, updating only the bottom-level prompt embeddings, which can cause long-distance backpropagation to result in vanishing gradients and slow convergence [18]. These reasons indicate the need to design better training and utilization methods for prompt embeddings, instead of using the same training method as token embeddings.

The RNNs contain a hidden unit that can preserve important information during sequence iterations and guide the model to output results [19]. RNNs can be unfolded into a long sequence, and through the hidden unit, the model can selectively integrate information from different times from the beginning to the end. The input at the initial time can guide the judgment at the final time, and this approach is beneficial in alleviating the problem of vanishing gradients, reducing forgetting, and speeding up convergence.

Prompt embeddings suffer from information loss when passing through each layer of transformers. Although residual connection modules exist within each layer, these modules operate on the entire sequence, including the token embeddings, rather than the prompt embeddings alone. Therefore, due to the reasons mentioned earlier, we need to design a dynamic information fusion mechanism specifically for the part corresponding to the prompt embeddings. Inspired by RNNs, we propose the Global Prompt Cell (GPC) to address the aforementioned issues.

The Global Prompt Cell (GPC) consists of two units: the remembering unit and the forgetting unit. Remembering unit should be applied to the prompt embeddings before passing through the transformer layer because information loss occurs after passing through the layers. Therefore, we use remembering unit to selectively remember certain information. On the other hand, forgetting unit should be applied to the prompt embeddings after passing through the transformer layers. The forgetting unit can selectively forget some information in the latest prompt embeddings to fuse with the previous prompt embeddings.

By utilizing the same remembering unit and forgetting unit for every model layer, our approach gathers information regarding prompt embeddings from all layers, enabling it to guide the model to achieve better downstream task performance. As a result, prompt embeddings are no longer simply a vector concatenated with token embeddings, but also a control module that aids PTMs in making better decisions. Moreover, since GPC only acts on the prompt embeddings updates between layers, it can be considered a plug-in module and requires only a small number of additional parameters. Lastly, we eliminate the verbalizer in vanilla prompt tuning and use classification heads for downstream tasks, reducing the difficulty of selecting optimal verbalizers for various tasks.

To summarize, our contributions are three-fold:

1. We propose a new training and utilization method for prompt embeddings called the Global Prompt Cell (GPC). To our knowledge, GPC is the first

method that specifically aims to improve the training and utilization of prompt embeddings, reduce forgetting during training, and ultimately enhance downstream task performance. Experiments prove its effectiveness and show its validity in architecture with ablation study.

2. Our method simplifies prompt tuning by discarding the verbalizer, which further reduces the time and computing consumption to select the optimal verbalizer, and achieves even better results.
3. Our method can be viewed as an easy-to-implement plug-in module with only a few additional parameters, which makes GPC both model-agnostic and task-agnostic.

2 Related Work

2.1 Pre-trained Language Model

Recently substantial works have shown that pre-trained models (PTMs) can learn universal language representations through pre-trained tasks on large corpora, which are beneficial for downstream NLP tasks and can avoid training a new model from the beginning [14].

A representative application of PTMs is using encoder-based models for classification tasks [23]. Encoder-based models are composed of multi-layer transformer encoders, which include a self-attention module and a multi-layer perceptron (MLP), and are optimized for performance through measures such as intra-layer residual connections. BERT is a classic encoder-based model commonly used for classification tasks [2]. In BERT, each input sentence is concatenated with a [CLS] token at the beginning. After passing through the encoder of each layer, the [CLS] token is utilized as a classification indicator to produce the final result.

2.2 Prompt Tuning

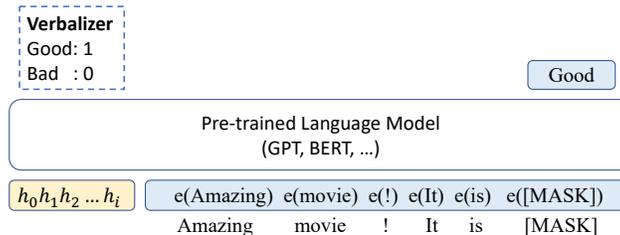


Fig. 1. Prompt-tuning model.

GPT-3 [1] has revolutionized downstream tasks by transforming them into generation tasks through the addition of prompt-like hints. This allows the model

to generate results directly in few-shot or zero-shot learning scenarios. Prompt consists of a template for transforming the input text and a verbalizer for matching the generated words to the actual task labels. It can be designed manually by domain experts [16,17] or automatically [20,3], but at the expense of low explainability [20].

Prompt tuning methods with prompt embeddings have been explored in recent studies [5,13,27,11]. Prompt embeddings are trainable embeddings rather than natural language tokens. Vanilla prompt tuning, proposed by [8], concatenates the prompt embeddings with the token embeddings in the first encoder layer.

To use prompt tuning for downstream tasks, we first encode a sequence of discrete input tokens $X = x_0, x_1, \dots, x_n$ into embeddings $X_e = e(x_0), e(x_1), \dots, e(x_n)$ using a pre-trained language model M . We then obtain the prompt embeddings $P = p_1, p_2, \dots, p_n$, where p_i is the prompt embedding for $i \in n$, and concatenate the prompt embeddings and token embeddings to form the complete input $I = P; X_e$. The model output embedding of the [MASK] token is then fed into a classifier to predict the target token, and the verbalizer is used to obtain the actual label.

Current research on prompt tuning architecture retains the verbalizer as a “necessary” component, despite its high computational cost, both in terms of time and computing resources, even when constructed automatically. Despite the improvements observed in downstream tasks, the results of prompt tuning can still be inconsistent or less-than-ideal.

Prompt tuning is a newly-arising paradigm that requires significantly more training time than fine-tuning to achieve the same performance [21], despite the two paradigms only differing in model inputs. This is due to the fact that although prompt embeddings are fundamentally different from token embeddings in terms of initialization and acquisition, vanilla prompt tuning treats prompt embeddings in a similar manner to prompt-like token vectors and processes them together using the same method.

2.3 Model Degradation and Prompt Forgetting

When training deep artificial neural networks, two problems become increasingly challenging for model optimization. The first problem is vanishing gradients [7], which impedes convergence from the outset, but can be largely mitigated through normalized initialization and intermediate normalization layers. The second problem is the degradation problem [6]: as the network becomes deeper, accuracy saturates and then degrades rapidly, despite not being caused by overfitting, which suggests that not all systems are equally easy to optimize [6]. Skip connections [26] are a common countermeasure for model degradation, as seen in ResNet [6] and Transformer [24], where residual connections are employed. Skip connections serve as a reminder for the model to retain previous information and prevent forgetting.

Prompt tuning is also affected by these problems since it occurs at the lowest layer. No previous studies have addressed the issues of model degradation and

prompt forgetting in prompt tuning, and we take these problems into consideration.

3 Method

In this section, we present the architecture of our method, which incorporates the implementation of Global Prompt Cell (GPC) into an encoder-based model. Drawing inspiration from RNN models, we utilize a method to effectively store previous states of prompt embeddings in various layers, allowing the model to zoom out and concentrate on the prompt from a wider perspective.

3.1 GPC between the Encoders

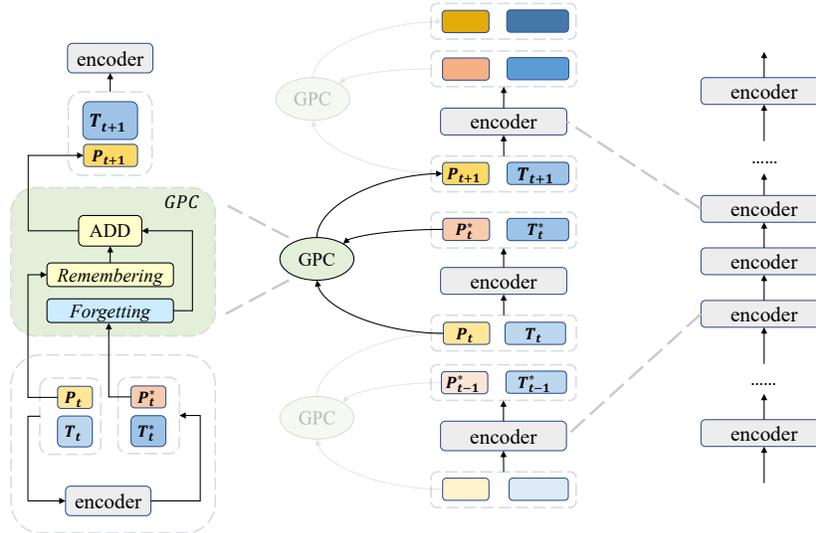


Fig. 2. Global Prompt Cell Model (best viewed in color).

In the Transformer model, residual connections are employed to mitigate the problem of degraded training accuracy. [24]. Building on this idea, GPC employs a more advanced connection mechanism around each pair of adjacent encoders. Specifically, GPC operates exclusively on the prompt portion of each vector. We use P to represent prompt embeddings and T to represent token embeddings, with the two concatenated to form the complete input or output for prompt tuning. As illustrated in Figure 2, P_t represents the prompt input of the t th encoder, while P_t^* represents the prompt output of the t th encoder. GPC takes

the input and output of the t th encoder, P_t and P_t^* , to generate the prompt input of the $(t + 1)$ th encoder E_t, P_{t+1} . Similarly, T_t represents the text input embeddings of the t th encoder, while T_t^* represents the text output embeddings of the t th encoder, derived from T_t . Since GPC does not affect text embeddings, T_{t+1} is identical to T_t . The complete embeddings are updated as follows.

$$\{P_t^*; T_t^*\} = E_t(\{P_t; T_t\}) \quad (1)$$

The prompt embeddings are updated as below.

$$P_{t+1} = GPC(P_t^*, P_t) \quad (2)$$

3.2 Inside the GPC

Figure 2 provides a detailed view of the internal structure of GPC. The cell comprises two units: the remembering unit and the forgetting unit. Each unit contains a single-layer feed-forward neural network that is shared across all encoder layers. Similar to ResNet, the remembering unit utilizes the prompt embedding before it is fed to the encoder, in order to emphasize and make use of the original information. In contrast, the forgetting unit operates on the output of the encoder to reduce its influence and prevent extreme outcomes. The cell then combines the outputs of both units to produce the final result. Therefore, between every pair of adjacent encoders, GPC considers both the original and encoded prompt embeddings, enabling the model to retain past information and encouraging it to forget some of the current state, thereby striking a balance between the past and the present. The following equation demonstrates how GPC processes prompt embeddings:

$$P_{t+1} = \theta(W_F P_t^* + W_R P_t) \quad (3)$$

where W_F represents the weights of the forgetting unit, while W_R represents the weights of the remembering unit, where θ is the activation function.

3.3 Classification Head

Traditional discrete prompt tuning is mainly used for few-shot or zero-shot tasks, so verbalizers are needed to align with pre-training cloze tasks, in order to narrow the gap between pre-training tasks and downstream tasks, and improve downstream task performance. In contrast, the continuous prompt tuning used in our paper aims to reduce the amount of parameter training and storage during fine-tuning. It generally uses the entire dataset instead of few-shot scenarios, and selecting an appropriate verbalizer manually requires a lot of training resources.

In order to simplify the model and reduce manual intervention as in [8] and [11], we replace the original verbalizers with a classification head that receives outputs from the final layer of the model. We use a randomly initialized classification head to predict the label from the output of the [CLS] token.

Corpus	Train	Dev	Test	Task	Metrics
BoolQ	9427	3270	3245	Question Answering	accuracy
CB	250	57	250	Natural Language Inference	accuracy
COPA	400	100	500	Question Answering	accuracy
RTE	2500	278	300	Natural Language Inference	accuracy
WiC	6000	638	1400	World Sense Disambiguation	accuracy
WSC	554	104	146	Co-reference Resolution	accuracy

Table 1. Statistics of SuperGLUE datasets.

	BoolQ			RTE		
	PT	Prompt-only	GPC	PT	Prompt-only	GPC
BERT	67.2	62.8	67.9	53.5	54.5	61.0
RoBERTa	62.3	62.4	63.5	58.8	54.2	59.4
	CB			COPA		
	PT	Prompt-only	GPC	PT	Prompt-only	GPC
BERT	80.4	71.4	82.1	55.0	58.0	67.0
RoBERTa	71.4	69.6	73.2	63.0	62.0	66.0
	WiC			WSC		
	PT	Prompt-only	GPC	PT	Prompt-only	GPC
BERT	63.0	56.4	66.9	64.4	64.4	65.4
RoBERTa	56.9	54.7	69.6	64.4	63.5	65.4

Table 2. Results on SuperGLUE development set. PT: Prompt tuning [8]; Prompt-only: Prompt tuning with no verbalizer; GPC: Prompt tuning with Global Prompt Cell; **bold**: the best performance.

4 Experiment

4.1 Datasets and Metrics

We evaluate on SuperGLUE [25]. SuperGLUE is a new benchmark styled after GLUE with a new set of more difficult language understanding tasks. Because our task mainly involves classification tasks, while MultiRC and ReCORD belong to QA tasks, and these two tasks are difficult to handle using prompt learning, as reported by [9]. Existing prompt learning methods for these two tasks suffer from significant fluctuations and are difficult to converge. Moreover, the effectiveness of our method in classification tasks has been demonstrated through the other six tasks. Therefore, we selected the other six classification tasks for our experiments. We choose the classification tasks and co-reference resolution tasks, including BoolQ, RTE, CB, COPA, WiC and WSC. Statistics of the selected datasets are in Table 1. We use accuracy as our evaluation metric.

4.2 Experiment Settings

Models We employ BERT [2] and RoBERTa [12] for our model, both of which are based on transformer encoders and are typically used for classification tasks.

Prompt Length The length of prompt embeddings has a significant impact on model performance. According to [10], different tasks require prompt embeddings of varying lengths to achieve optimal results. Typically, simpler tasks require shorter prompts than more complex ones [10]. We experimented with prompt lengths of 16, 32, and 64 and selected the most effective ones among them.

Prompt Initialization Prompt embeddings can be initialized in various ways, such as random initialization or utilizing concrete token embedding. Based on the approach described in [4], we use random initialization in our method.

Training Method During the training phase, we freeze the original parameters of the PTMs and only update the Global Prompt Cell, specifically the corresponding weight matrices W_R and W_F . Therefore, we only need to store a small number of parameters instead of the parameters of the entire PTM for each downstream task.

4.3 Main Results

To assess the effectiveness of Global Prompt Cell, we investigate (i) whether replacing the verbalizer with a classification head causes performance degradation in prompt tuning, and (ii) whether Global Prompt Cell can outperform prompt tuning. (iii) whether Global Prompt Cell combined with a classification head still yield better results.

We carry out prompt-only experiments, where we substitute the verbalizer with a classification head and refrain from using GPC. Table 2 demonstrates that the performance of prompt-only models decreases, particularly on the CB dataset, with both BERT and RoBERTa models experiencing over 10% decrease.

For (ii) and (iii), we perform experiments with Global Prompt Cell. As shown in Table 2, GPC significantly improves the effectiveness of prompt tuning using classification heads, surpassing the performance of vanilla prompt tuning. GPC can achieve over 10% performance increase compared to the prompt-only model on both WiC and CB datasets. In comparison with the PT (Prompt-Tuning [8]) model, which requires multiple experiments to determine the optimal verbalizer, our method still outperforms it on all six tasks. These results demonstrate the efficacy of our method.

5 Ablation Study

Our GPC module consists of two parts, remembering unit and forgetting unit. The remembering unit receives previous prompt embeddings and selects which

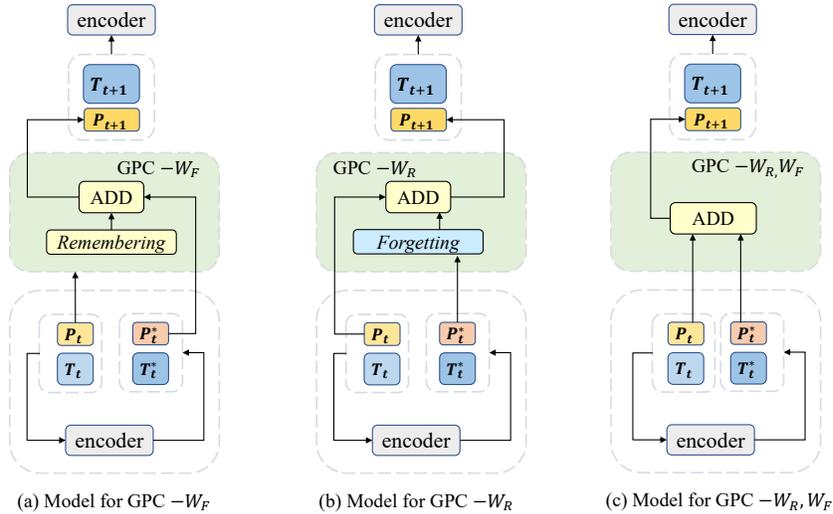


Fig. 3. Models in ablation study.

Model	Updating strategy
GPC	$P_{t+1} = \theta(W_F P_t^* + W_R P_t)$
GPC- W_F	$P_{t+1} = \theta(P_t^* + W_R P_t)$
GPC- W_R	$P_{t+1} = \theta(W_F P_t^* + P_t)$
GPC- W_F, W_R	$P_{t+1} = \theta(P_t^* + P_t)$

Table 3. Different settings for prompt update.

parts to retain. The remembering unit can be helpful for retaining certain information that requires long-term memory. The forgetting unit is responsible for determining which parts of the latest prompt embeddings are irrelevant, preventing the addition of unnecessary information.

We design the ablation experiments to explore the effectiveness of the above units. We conducted three groups of ablation experiments, where we added only the memory unit module, only the forgetting unit module, and no module at all. Through this, we aim to demonstrate the impact of memory and forgetting units on the experimental results, as well as prove the necessity and importance of these two modules.

Figure 3 shows the architectures of three ablation settings, which respectively remove the forgetting unit, the remembering unit and the both. Table 3 shows the formula of the three different updating strategies and the comparison with GPC.

The result in Table 4 shows that the contribution of two units varies from task to task, but the combination of the two can reach optimal results for all cases.

		BoolQ	RTE	CB	COPA	WiC	WSC
BERT	GPC	67.9	61.0	82.1	67.0	66.9	65.4
	GPC- W_F	62.3	57.8	67.9	61.0	60.7	65.4
	GPC- W_R	62.6	57.8	55.4	63.0	57.5	64.4
	GPC- W_F, W_R	62.7	58.1	75.0	63.0	57.2	63.5
RoBERTa	GPC	63.5	59.4	73.2	66.0	69.6	65.4
	GPC- W_F	62.2	55.6	61.0	55.0	53.0	63.5
	GPC- W_R	62.2	52.7	64.3	55.0	50.9	63.5
	GPC- W_F, W_R	62.3	57.0	57.1	58.0	52.2	63.5

Table 4. Ablation experiment results. GPC: Global Prompt Cell; GPC- W_F : removing the forgetting unit; GPC- W_R : removing the remembering unit; GPC- W_F, W_R : removing both the forgetting unit and remembering unit. **Bold**: the best performance.

In addition to the findings presented in Table 4, further analysis of our ablation experiments revealed that the impact of the remembering and forgetting units on task performance was dependent on the specific task. Specifically, for some tasks, the addition of the memory unit module resulted in a greater improvement in performance compared to the addition of the forgetting unit module, while for other tasks the opposite was true.

The combination of both modules consistently led to optimal results across all tasks. This suggests that while the individual contributions of the remembering and forgetting units may vary depending on the specific task, the integration of both modules is crucial for achieving optimal performance across a range of tasks.

6 Conclusion

To conclude, our study has introduced the Global Prompt Cell (GPC) as a novel approach to enhance continuous prompt tuning. By selectively remembering and forgetting prompt embeddings, GPC enables more effective prompt updating, ultimately leading to improved model performance.

Through experiments conducted on the SuperGLUE benchmark, we have shown that our approach can substantially enhance results using prompt tuning. This highlights the potential of GPC to significantly improve the performance of language models on a range of natural language understanding tasks.

Finally, we emphasize that GPC can serve as a portable plug-in module for prompt tuning paradigms, allowing for easy integration with existing models and architectures. We believe that our approach has promising implications for the development of more effective and efficient language models, and we look forward to further exploration and refinement of this method in future research.

References

1. Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al.: Language models are few-shot learners. *Advances in neural information processing systems* **33**, 1877–1901 (2020)
2. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding pp. 4171–4186 (2019)
3. Gao, T., Fisch, A., Chen, D.: Making pre-trained language models better few-shot learners. 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)*. pp. 3816–3830. Association for Computational Linguistics, Online (Aug 2021)
4. Gu, Y., Han, X., Liu, Z., Huang, M.: Ppt: Pre-trained prompt tuning for few-shot learning. *arXiv preprint arXiv:2109.04332* (2021)
5. Hambardzumyan, K., Khachatrian, H., May, J.: WARP: Word-level Adversarial ReProgramming. 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)*. pp. 4921–4933. Association for Computational Linguistics, Online (Aug 2021)
6. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 770–778 (2016)
7. Hochreiter, S.: The vanishing gradient problem during learning recurrent neural nets and problem solutions. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems* **6**(02), 107–116 (1998)
8. Lester, B., Al-Rfou, R., Constant, N.: The power of scale for parameter-efficient prompt tuning. In: *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. pp. 3045–3059 (2021)
9. Liu, P., Yuan, W., Fu, J., Jiang, Z., Hayashi, H., Neubig, G.: Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *arXiv preprint arXiv:2107.13586* (2021)
10. Liu, X., Ji, K., Fu, Y., Du, Z., Yang, Z., Tang, J.: P-tuning v2: Prompt tuning can be comparable to fine-tuning universally across scales and tasks. *arXiv preprint arXiv:2110.07602* (2021)
11. Liu, X., Zheng, Y., Du, Z., Ding, M., Qian, Y., Yang, Z., Tang, J.: Gpt understands, too. *arXiv preprint arXiv:2103.10385* (2021)
12. Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., Stoyanov, V.: Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692* (2019)
13. Qin, G., Eisner, J.: Learning how to ask: Querying LMs with mixtures of soft prompts. In: *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. pp. 5203–5212. Online (Jun 2021)
14. Qiu, X., Sun, T., Xu, Y., Shao, Y., Dai, N., Huang, X.: Pre-trained models for natural language processing: A survey. *Science China Technological Sciences* **63**(10), 1872–1897 (2020)
15. Radford, A., Narasimhan, K., Salimans, T., Sutskever, I.: Improving language understanding by generative pre-training (2018)
16. Schick, T., Schütze, H.: Exploiting cloze-questions for few-shot text classification and natural language inference. In: *Proceedings of the 16th Conference of the*

- European Chapter of the Association for Computational Linguistics: Main Volume. Association for Computational Linguistics, Online (Apr 2021)
17. Schick, T., Schütze, H.: It's not just size that matters: Small language models are also few-shot learners. In: Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. pp. 2339–2352. Association for Computational Linguistics, Online (Jun 2021)
 18. Schmidhuber, J.: Deep learning in neural networks: An overview. *CoRR* **abs/1404.7828** (2014)
 19. Sherstinsky, A.: Fundamentals of recurrent neural network (rnn) and long short-term memory (lstm) network. *Physica D: Nonlinear Phenomena* **404**, 132306 (03 2020)
 20. Shin, T., Razeghi, Y., Logan IV, R.L., Wallace, E., Singh, S.: AutoPrompt: Eliciting Knowledge from Language Models with Automatically Generated Prompts. In: Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP). pp. 4222–4235. Association for Computational Linguistics, Online (Nov 2020)
 21. Su, Y., Wang, X., Qin, Y., Chan, C.M., Lin, Y., Wang, H., Wen, K., Liu, Z., Li, P., Li, J., Hou, L., Sun, M., Zhou, J.: On transferability of prompt tuning for natural language processing (2021)
 22. Su, Y., Wang, X., Qin, Y., Chan, C.M., Lin, Y., Wang, H., Wen, K., Liu, Z., Li, P., Li, J., Hou, L., Sun, M., Zhou, J.: On transferability of prompt tuning for natural language processing. In: Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. Association for Computational Linguistics (2022). <https://doi.org/10.18653/v1/2022.naacl-main.290>
 23. Sun, C., Qiu, X., Xu, Y., Huang, X.: How to fine-tune bert for text classification? In: China national conference on Chinese computational linguistics. pp. 194–206. Springer (2019)
 24. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. *Advances in neural information processing systems* **30** (2017)
 25. Wang, A., Pruksachatkun, Y., Nangia, N., Singh, A., Michael, J., Hill, F., Levy, O., Bowman, S.: Superglue: A stickier benchmark for general-purpose language understanding systems. In: Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., Garnett, R. (eds.) *Advances in Neural Information Processing Systems*. vol. 32. Curran Associates, Inc. (2019)
 26. Wu, D., Wang, Y., Xia, S.T., Bailey, J., Ma, X.: Skip connections matter: On the transferability of adversarial examples generated with resnets (2020)
 27. Zhong, Z., Friedman, D., Chen, D.: Factual probing is [MASK]: Learning vs. learning to recall. In: Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. pp. 5017–5033. Association for Computational Linguistics, Online (Jun 2021)