

P-Distill: Efficient and Effective Prompt Tuning using Knowledge Distillation

Anonymous ACL submission

Abstract

In the field of natural language processing (NLP), prompt-based learning is widely used for efficient parameter learning. However, this method has the drawback of shortening the input length by the extent of the attached prompt, leading to an inefficiency in utilizing the input space. In this study, we propose P-Distill, a novel prompt compression method that mitigates the aforementioned limitation of prompt-based learning while maintaining performance via knowledge distillation. The knowledge distillation process of P-Distill consists of two methods, namely prompt initialization and prompt distillation. Experiments on various NLP tasks demonstrate that P-Distill exhibits comparable or superior performance compared to other state-of-the-art prompt-based learning methods, even with significantly shorter prompts. Specifically, We achieve a peak improvement of 1.90% even with the prompt lengths compressed to one-eighth. An additional study further provides insights into the distinct impact of each method on the overall performance of P-Distill. Our code will be released upon acceptance.

1 Introduction

Pre-trained language models (PLMs) have been effective in improving performances of various natural language processing (NLP) tasks (Devlin et al., 2019; Brown et al., 2020; Touvron et al., 2023). These models are fine-tuned by optimizing all parameters to enhance the performances of specific downstream tasks; however, fine-tuning requires significant computational resources while training. The need for significant computational resources for storage and training becomes a challenge, especially when fine-tuning large language models such as Llama2 (Touvron et al., 2023), which may not be readily available to most users.

To reduce computational costs, researchers have explored various methods for efficiently fine-tuning

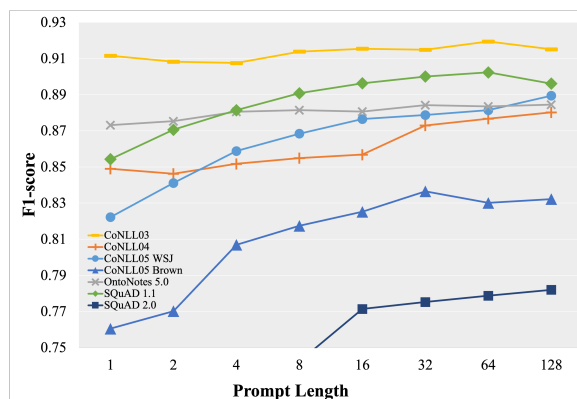


Figure 1: Performance variation in P-tuning across tasks based on the length of continuous prompts.

the parameters (Houlsby et al., 2019; Hu et al., 2021; Liu et al., 2022). In contrast to the traditional model fine-tuning that updates all parameters for a downstream task, P-tuning (Liu et al., 2022) fixes the pre-trained parameters and only trains the continuous prompts, which are trainable embeddings attached at the beginning or throughout each layer of the model. P-tuning is computationally efficient, particularly for PLMs with a large number of parameters; however, it overlooks the inefficiency in input space utilization arising from attaching continuous prompts (Hu et al., 2021). Extending the input token sequence by forcibly modifying the code can lead to issues with attention calculation beyond the model’s training (Chen et al., 2023). Such modifications often result in performance degradation, as the attention mechanism might struggle with sequences extending beyond its initially intended scope. Similar to the findings in the work (Liu et al., 2022), more challenging tasks require longer prompt lengths to achieve the better performance, as shown in Figure 1. For instance, in datasets such as CoNLL04, CoNLL05 WSJ, SQuAD 2.0, and OntoNotes 5.0, a prompt length of 128 is required to enhance the performance.

In this paper, we propose P-Distill, which is a

novel prompt compression method to mitigate the limitations of long prompts. Our method involves a two-step process where we first train a teacher model using P-tuning to achieve superior performance with long prompts. We then transfer this knowledge to a student model with significantly shorter prompts through a distillation process. To ensure stability in training, we first perform prompt initialization based on the teacher model prompts. Then, we focus on distilling knowledge between the teacher and student models, specifically targeting the outputs of their intermediate and prediction layers. This is due to the impact of continuous prompts on the hidden states within these layers, which subsequently influences the model’s predictions. This method enables the compression of prompts to shorter lengths without a significant degradation in performance, thereby addressing the inefficiencies inherent in longer prompts.

To validate its effectiveness and efficiency, we evaluate P-Distill using various NLP benchmarks. Our results demonstrate that P-Distill exhibits comparable or superior performance than those of the existing state-of-the-art prompt-based models. To the best of our knowledge, this study is the first to train teacher prompts and transfer their knowledge to student prompts for the purpose of compressing prompts. The main contributions of this study are summarized as follows:

- We propose a method called P-Distill to compress the continuous prompts, effectively mitigating the limitation of reducing the model’s usable sequence length in prompt-based learning.
- We introduce a prompt distillation method utilizing teacher model’s hidden state and prediction outputs, influenced by continuous prompts, and propose a prompt initialization for stable prompt distillation.
- We validate P-Distill across multiple NLP benchmarks, demonstrating its ability to maintain or enhance accuracy while reducing prompt lengths by up to eight times.

The remainder of this paper is structured as follows: Section 2 provides the preliminaries; Section 3 describes a detailed description of the proposed method; Section 4 presents the experimental results and analysis, and Section 5 concludes the study.

2 Preliminaries

2.1 Pre-trained Language Models Based on the Transformer

The transformer model (Vaswani et al., 2023), comprising an encoder and decoder, is the fundamental architecture of the majority of recent PLMs, including BERT (Devlin et al., 2019), RoBERTa (Liu et al., 2019), and GPT-3 (Brown et al., 2020). Each encoder and decoder consists of multiple transformer layers and incorporates key components, such as multi-head attention modules (MHA), feed-forward networks, layer normalization, and residual connections. A key component of this architecture is the multi-head attention mechanism, which computes attention weights using query (Q), key (K), and value (V) matrices. Mathematically, the attention function in multi-head attention can be represented as follows:

$$Att(x) = softmax(\frac{QK^T}{\sqrt{d_k}})V, \quad (1)$$

where $\sqrt{d_k}$ is the scaling factor for gradient stabilization during training. This attention mechanism is crucial in understanding language and generating tasks by modulating the focus of the model on different parts of the input data.

2.2 Prompt-based Learning Methods

Prompt-based learning methods have emerged as an efficient alternative to full model fine-tuning, especially for PLMs (Liu et al., 2022, 2023). These methods use prompts to guide the model predictions for specific tasks. Several approaches (Jiang et al., 2020; Shin et al., 2020) employ discrete prompts, which are fixed templates added to the input. For example, in sentiment analysis, a template might be “This text [*Input Text*] expresses a [*MASK*] sentiment.”. However, discrete prompts are limited in that their performances significantly depend on template selection. Advanced approaches, such as Prefix-Tuning (Li and Liang, 2021) and P-tuning (Liu et al., 2022), use continuous prompts that are trainable embeddings independent of the model vocabulary. Particularly, P-tuning attaches continuous prompts to each layer of the model, thereby influencing its behavior and enhancing its performance in downstream tasks. These continuous prompts are integrated into the attention mechanism of the transformer model as

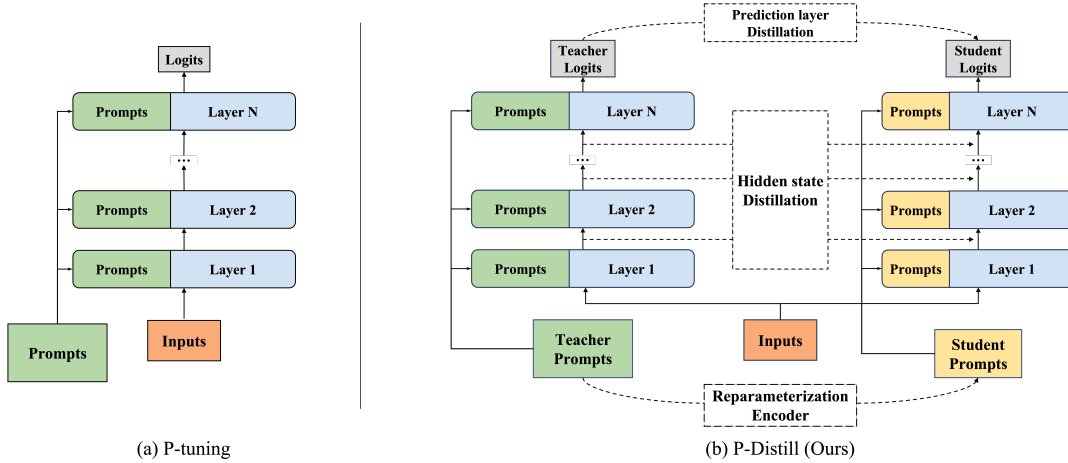


Figure 2: (a) Illustration of P-tuning (Liu et al., 2022). (b) Illustration of the proposed method, denoted as P-Distill. This method trains a teacher model to generate concise and effective prompts, followed by distilling the knowledge into a student model.

follows:

$$Att(x) = softmax \left(\frac{Q(P_k : K)^T}{\sqrt{d_k}} \right) (P_v : V), \quad (2)$$

where $P_k \in R^{n_p \times d}$ and $P_v \in R^{n_p \times d}$ are the continuous prompts added to the key and value vectors, respectively, the colon denotes the concatenation of these prompts with the key and value matrices. The dimension n_p indicates the lengths of the prompts, and d represents the dimensions of the key and value vectors. This integration enables the model to influence layers closer to the output, significantly affecting the final predictions.

2.3 Knowledge Distillation

In artificial intelligence, knowledge distillation (KD) is a technique for reducing the size of large models while preserving their performances (Jiao et al., 2020; Sun et al., 2020; Sanh et al., 2020; Hinton et al., 2015). During KD, a smaller student model is trained to internalize and emulate the complex decision-making patterns and behaviors of a larger teacher model. This process involves the behavior functions of the models, f^T and f^S , transforming inputs into informative representations, typically defined as the output of any layer within the model. These representations contain abundant information for model predictions. KD is quantified using loss functions, such as the Kullback-Leibler divergence (Kullback and Leibler, 1951) or Mean Squared Error (MSE) (Hinton et al., 2015), as follows:

$$L_{KD} = \sum_{x \in X} L(f^S(x), f^T(x)), \quad (3)$$

where x is the input, and X and L denote the dataset and the loss function, respectively. This approach enables the student model to gain a comprehensive understanding of various classes, enhancing its application in fields such as NLP.

3 Methodology

Many existing prompt tuning methods, including P-tuning, have the drawback of occupying an unnecessarily large portion of the input token space owing to their long prompts. Inspired by knowledge distillation methods, we propose a novel prompt compression methodology called P-Distill. This approach aims to compress the prompts while maintaining the performance, thereby increasing the available space for input tokens and enhancing the overall model efficiency. To this end, the proposed P-Distill comprises the following two methods: prompt initialization and prompt distillation. Figure 2 shows the learning and compression processes of P-Distill. Our approach involves two main steps where the first step is training a teacher model using P-tuning, and the second step focuses on distilling knowledge to a student model with shorter prompts, effectively reducing the length of the prompts.

3.1 Prompt-Based Teacher Learning

When solving downstream tasks using P-tuning, we freeze the pre-trained weights of the language model and only train the continuous prompts. Prompt lengths that yield good performance vary according to task complexity. In general, simple classification tasks tend to use shorter prompts,

around 20, while more complex sequence labeling tasks often require longer prompts, around 100 (Liu et al., 2022). Understanding the variation in prompt length is crucial, particularly during the inference stage, as longer prompts inherently limit the maximum sequence length that the model can handle.

We train a teacher model on various tasks based on the P-tuning methodology. This model tokenizes input data x and embeds it into text embeddings \bar{x} . Subsequently, the continuous prompts $P_k^T, P_v^T \in R^{n_t \times d}$ of the teacher model are randomly initialized and concatenated with the key vectors $K \in R^{n_x \times d}$ and value vectors $V \in R^{n_x \times d}$ of each layer. Here, d is the dimensionality of the hidden representations, n_t is the prompt length of the teacher model, and n_x is the length of token embeddings. The teacher model, which utilizes attention heads incorporating continuous prompts, is trained to take the text embedding \bar{x} as input and generate the final logits y^T . The parameter optimization of the teacher model is guided by the cross-entropy loss, which is formalized as follows:

$$L_{CE}^T = -\frac{1}{|B|} \sum_{i=1}^{|B|} \log(\text{softmax}(y_i^T)[c_i]), \quad (4)$$

where $|B|$ is the number of data points in the current batch, y_i^T is the logits output by the teacher model for the i -th data point in the batch, $\text{softmax}(y_i^T)$ is the softmax-transformed probability distribution over the classes, and c_i is the true class index for the i -th data point.

3.2 Prompt-enhanced Distillation (P-Distill)

We initiate the training of a student model which employs shorter continuous prompts, rather than the teacher model, using the same prompt attachment methodology. During the initial training phase, we initialize the continuous prompts of the student model P_k^S and $P_v^S \in R^{n_s \times d}$ based on the teacher model’s prompts P_k^T and P_v^T . Subsequently, student prompts are also attached to the key and value vectors across all layers to compute the attention heads. The length of the student model prompts, represented by n_s , is shorter than that of the teacher model prompts n_t . The student model, denoted by f^S , takes the text embedding \bar{x} as input and generates the output logits y^S . The teacher and student models share the same underlying language model architecture, differing only in the length and content of their respective prompts. In this context,

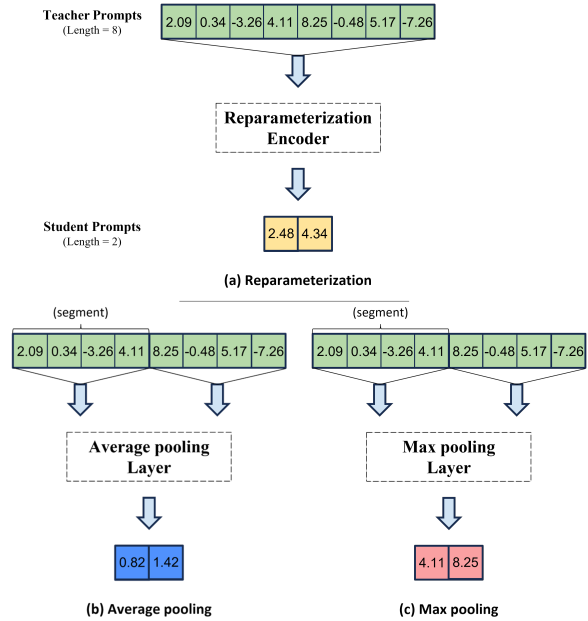


Figure 3: Illustration of various prompt initialization methods.

we focus on distilling the knowledge from the more extensive teacher model prompts into the shorter student model prompts. To enhance the effectiveness of knowledge transfer, we propose two novel methods for knowledge distillation.

3.2.1 Prompt Initialization

For solving downstream tasks, the model utilizes the attached prompts to generate answers. Starting with the randomly initialized prompts for the model can result in an unstable training process (Lester et al., 2021). To mitigate this challenge, the study (Vu et al., 2022) employed a method for transferring the prompts learned in one task to another task. We aim to stabilize the training by initializing the student model prompts P_k^S and P_v^S based on the teacher model prompts P_k^T, P_v^T . We experiment with various prompt initialization methods, including reparameterization, average pooling, and max pooling, as illustrated in Figure 3. In reparameterization, we employ a reparameterization encoder to adjust the length of the teacher model prompts to that of the student model prompts. For average pooling, we divide the teacher model’s prompts into smaller segments and compute their averages to initialize the student prompts. In max pooling, we focus on the most prominent features by obtaining the maximum value from each segment of the teacher model’s prompts. Based on the experimental results, we apply the reparameterization encoder to the teacher model’s prompts to construct the stu-

dent model’s prompts as follows:

$$P_k^S = (P_k^T \cdot W_k^T) + b_k^T, \quad (5)$$

$$P_v^S = (P_v^T \cdot W_v^T) + b_v^T, \quad (6)$$

where W_k^T and W_v^T are the learnable weight matrices used to construct the student’s prompts, and b_k^T and b_v^T are the corresponding bias. The results of various prompt initialization experiments are shown in Section 4.5.

3.2.2 Prompt Distillation

In this section, we focus on prompt distillation, a key aspect of the proposed approach. Recognizing the influence of continuous prompts on both the hidden states and the prediction layer outputs within the model, we employ the following two distillation techniques: prediction layer and hidden state distillations. These techniques focus on different aspects of the teacher model’s output to ensure comprehensive knowledge transfer.

Prediction layer distillation In this method, a student model learns to emulate the predictions of a teacher model. This process involves the student model utilizing soft labels from the teacher model’s output, which encapsulate the teacher model’s understanding of the data. Particularly, a loss function is used to minimize the difference between the logits y^S and y^T produced by the student and teacher models, respectively. The distillation loss L_{pred} is formulated as follows:

$$L_{pred} = KL(\text{softmax}(y_i^S / \theta), \text{softmax}(y_i^T / \theta)), \quad (7)$$

where y_i^S and y_i^T are the logits vectors predicted by the student and teacher, respectively, and KL denotes the Kullback-Leibler divergence, which measures the difference between the probability distributions of the two models. θ is a temperature hyperparameter that adjusts the smoothness of these distributions, enabling a more nuanced transfer of knowledge from the teacher to student model. The distillation loss L_{pred} is then used in the optimization process to update the parameters of the student model, thereby aligning its predictive behavior more closely with that of the teacher model.

Hidden state distillation Additionally, we also distill knowledge from the intermediate representations of the teacher model. The concept of distilling knowledge through intermediate representations

was initially introduced by Fitnets (Romero et al., 2015), with the aim of enhancing the training process of the student model. Based on the provided prompts and inputs, we extract knowledge from the transformer layers of the teacher model and distill into the student model. This process is formalized using the loss function L_{hidden} , which is calculated as the MSE between the hidden states H_S and H_T of the student and teacher models, respectively, as follows:

$$L_{hidden} = MSE(H_S, H_T), \quad (8)$$

where the matrices $H_S, H_T \in R^{n \times d}$ represent the hidden states, n is the input sequence length, and d is the hidden state dimensionality of the two models.

3.3 Distillation-based Student Learning

While training the student model, the cross-entropy loss is computed similar as that of the teacher model. This loss serves as a measure of the student model’s accuracy in predicting the true class labels as follows:

$$L_{CE}^S = -\frac{1}{|B|} \sum_{i=1}^{|B|} \log(\text{softmax}(y_i^S)[c_i]). \quad (9)$$

Subsequently, the overall loss function L_{total} for the student model is then a weighted combination of the cross-entropy loss and the distillation losses as follows:

$$L_{total} = \lambda_1 \cdot L_{CE}^S + \lambda_2 \cdot L_{pred} + \lambda_3 \cdot L_{hidden}, \quad (10)$$

where λ_1 , λ_2 , and λ_3 are the learnable weighted coefficients with the constraint that their combined sum equals 1. During the training, the teacher model parameters are fixed to serve as the sources of prior knowledge.

4 Experiments

This section presents the datasets employed in our experiments, baseline models for comparison, results of these datasets, and analyses from our additional studies.

4.1 Datasets

Our evaluation of the proposed P-Distill method includes a comprehensive range of natural language understanding tasks, utilizing datasets that are well-established benchmarks in the field.

	BoolQ Acc.	CB Acc.	COPA Acc.	MultiRC F1a	ReCoRD F1	RTE Acc.	WiC Acc.	WSC Acc.	Average
Fine-tuning	0.777	<u>0.946</u>	0.710	0.705	0.706	0.762	0.746	0.683	0.754
P-tuning	0.764 ₍₈₎	<u>0.946</u> ₍₃₂₎	0.810 ₍₄₎	<u>0.711</u> ₍₁₆₎	0.728 ₍₁₆₎	<u>0.794</u> ₍₄₎	<u>0.756</u> ₍₄₎	0.731 ₍₁₆₎	0.780
	0.738 ₍₁₎	<u>0.929</u> ₍₄₎	<u>0.790</u> ₍₁₎	<u>0.707</u> ₍₂₎	0.721 ₍₂₎	<u>0.783</u> ₍₁₎	<u>0.745</u> ₍₁₎	0.692 ₍₂₎	<u>0.763</u>
P-Distill	<u>0.776</u> ₍₁₎	0.964 ₍₄₎	0.810 ₍₁₎	0.718 ₍₂₎	<u>0.726</u> ₍₂₎	0.798 ₍₁₎	0.759 ₍₁₎	<u>0.721</u> ₍₂₎	0.784

Table 1: Experimental results on the SuperGLUE validation dataset. For P-Distill, training was performed using a teacher model with the prompt length exhibiting the best performance for P-tuning. The numbers in parentheses indicate the lengths of the prompt attached to the model. (Acc.: Accuracy; **bold**: the best; underline: the second best)

We include various tasks from the SuperGLUE benchmark (Wang et al., 2019), which assesses a model’s understanding and reasoning abilities across different contexts, including BoolQ (Clark et al., 2019), CB (De Marneffe et al., 2019), COPA (Roemmele et al., 2011), MultiRC (Khashabi et al., 2018), ReCoRD (Zhang et al., 2018), RTE (Dagan et al., 2006; Bar Haim et al., 2006), WiC (Pilehvar and Camacho-Collados, 2019) and WSC (Levesque et al., 2011). We also utilize the CoNLL-2003 (Tjong Kim Sang and De Meulder, 2003), CoNLL-2004 (Carreras and Màrquez, 2004), CoNLL-2005 (Carreras and Màrquez, 2005), CoNLL-2012 (Pradhan et al., 2012), and OntoNotes 5.0 datasets (Weischedel et al., 2013), each providing richly annotated text for entity classification. The SQuAD dataset, in its versions 1.1 (Rajpurkar et al., 2016) and 2.0 (Rajpurkar et al., 2018), facilitates testing reading comprehension, requiring the model to parse passages and answer questions with a high degree of understanding. In addition, we extend our experiments to the Allsides dataset (Li and Goldwasser, 2019), which consists of news articles with an average token length of over 1000. This dataset allows us to examine the effects of P-Distill in scenarios where input tokens are frequently truncated. All of these datasets are English, open source, and used for academic research purposes only. For accurate comparisons, we follow the train, validation, and test set splits as specified in the referenced work (Liu et al., 2022).

4.2 Baselines

We compare P-Distill against the following methods to validate its competitive performance, with all methods utilizing BERT_{large} with 335M parameters as the backbone architecture.

Fine-tuning All parameters of a PLM are updated and adapted to the given downstream task in a task-specific manner.

P-tuning (Liu et al., 2022) It appends trainable continuous prompts to the key and value matrices of a model, enabling task-specific learning while keeping the model’s pre-trained weights fixed.

4.3 Experimental Details

In our training process, we exclusively focus on continuous prompts while keeping the backbone parameters of the model fixed. The model is trained with a batch size of 16, and the learning rate is individually optimized for each task. Furthermore, we employ the AdamW optimizer for training. For the temperature hyperparameter θ used in the distillation process, we experimentally determine the optimal setting by sweeping across $\{1, 5, 10\}$. For the learnable parameter λ_2 , we explore the initial values of $\{0.1, 0.5, 0.9\}$. Considering the significant impact of the hidden state loss, we experiment with the initial values of $\{1e-3, 1e-4, 1e-5\}$ for λ_3 . All experiments were performed using PyTorch¹ and HuggingFace Transformers (Wolf et al., 2020) on three NVIDIA A100 GPUs, and to ensure consistency in our results, each task was conducted using a fixed random seed.

4.4 Results

Tables 1 and 2 present the experimental results of Fine-tuning, P-tuning, and P-Distill. In P-Distill, the prompt length is compressed to one-eighth of that of the teacher model prompts. For fewer than eight teacher model prompts, the length is compressed to 1. For a detailed analysis of prompt compression ratios refer to Appendix A. In general, the proposed P-Distill method exhibits a comparable or superior performance than those of the other methods while using shorter prompts.

Results on SuperGLUE Table 1 shows the performance of each approach on the SuperGLUE benchmark. The experimental results show that despite using shorter prompts, P-Distill matches

¹<https://pytorch.org/>

	NER			SRL		QA		SC
	CoNLL03	CoNLL04	OntoNotes 5.0	CoNLL05 Brown	CoNLL05 WSJ	SQuAD 1.1 dev	SQuAD 2.0 dev	Allsides
Fine-tuning	0.928	0.882	0.890	0.827	0.885	0.911	0.819	0.780
P-tuning	0.919 ₍₆₄₎	0.880 ₍₁₂₈₎	0.885 ₍₁₂₈₎	0.837 ₍₃₂₎	0.890 ₍₁₂₈₎	0.902 ₍₆₄₎	0.782 ₍₁₂₈₎	0.775 ₍₃₂₎
	0.914 ₍₈₎	0.866 ₍₁₆₎	0.881 ₍₁₆₎	0.807 ₍₄₎	0.877 ₍₁₆₎	0.891 ₍₈₎	0.771 ₍₁₆₎	0.772 ₍₄₎
P-Distill	0.919 ₍₈₎	0.888 ₍₁₆₎	0.886 ₍₁₆₎	0.817 ₍₄₎	0.885 ₍₁₆₎	0.896 ₍₈₎	0.775 ₍₁₆₎	0.783 ₍₄₎

Table 2: Experimental results for each method on named entity recognition (NER), question answering (QA), semantic role labeling (SRL), and long sequence classification (SC). For P-Distill, training was performed using a teacher model with the prompt length exhibiting the best performance for P-tuning. The numbers in parentheses indicate the lengths of the prompts attached to the model. All metrics are reported as f1 scores. (**bold**: the best; underline: the second best)

or exceeds the performance of P-tuning which utilizes long prompts. Specifically, P-Distill exhibits 0.51% performance improvement on average on SuperGLUE tasks. The improvement increases to 2.75% when compared to P-tuning with the same prompt length. These results reveal that P-Distill effectively compresses the prompt length while maintaining or even improving performance. This highlights the significance of the proposed distillation method that transfers task knowledge from the teacher to the student model using eight times shorter prompts.

Results on Other Tasks Table 2 presents the experimental results for diverse tasks, including named entity recognition, question answering, semantic role labeling, and long sequence classification. We first observe that achieving superior performance via P-tuning on these tasks requires training with longer prompts (up to 128 tokens), which aligns with the phenomenon reported in previous work (Liu et al., 2022) where complex tasks tend to require long prompts.

In comparison to the baseline methods, P-Distill achieved comparable performance using significantly shorter prompts even for tasks where P-tuning employed long prompts. In particular, P-Distill outperforms P-tuning with 128 prompt tokens on CoNLL04 while utilizing a prompt of length 16. Furthermore, P-Distill achieves a performance improvement of 2.54% over the prompt of the same length trained using P-tuning, and a 0.90% improvement over the teacher model.

Furthermore, P-Distill outperforms the baseline methods on the Allsides dataset, which is a dataset that consists of long input instances where the average token length exceeds 1000. Specifically, with prompts significantly compressed to 4, P-Distill not only mitigates token truncation issues but also shows a 1.03% performance improve-

ment over the best P-tuning configuration. This highlights the advantage of P-Distill in handling long-sequence tasks by preserving input space with shorter prompts. Additionally, P-Distill exhibits 1.42% performance improvement when compared to P-tuning with prompt length 4, showing that P-Distill is an effective method for compressing prompt length. A qualitative evaluation of the Allsides dataset is provided in Appendix C, offering further insights on long-sequence tasks.

4.5 Ablation Study

To further verify the effectiveness of the proposed method, we conduct ablation studies using the following variants of P-Distill.

P-Distill_{-init} Instead of training with prompt initialization, it focuses exclusively on leveraging the two types of distillation losses designed to transfer the knowledge from the teacher to student model in different ways.

P-Distill_{-pred} This approach does not implement the prediction layer distillation loss. Following the application of the prompt initialization method, it trains the student model based on the hidden state distillation loss. This method aligns the internal representations of the student model with those of the teacher model without focusing on the final output predictions.

P-Distill_{-hidden} This variant does not consider the differences between the hidden state outputs of the teacher and student models. Instead, it focuses on training based on the differences in the prediction layer output. This approach aligns the final predictions of the student model closely with those of the teacher model without directly focusing on their internal representations.

Results Figure 4 shows that all three variants of P-Distill underperform the original P-Distill among

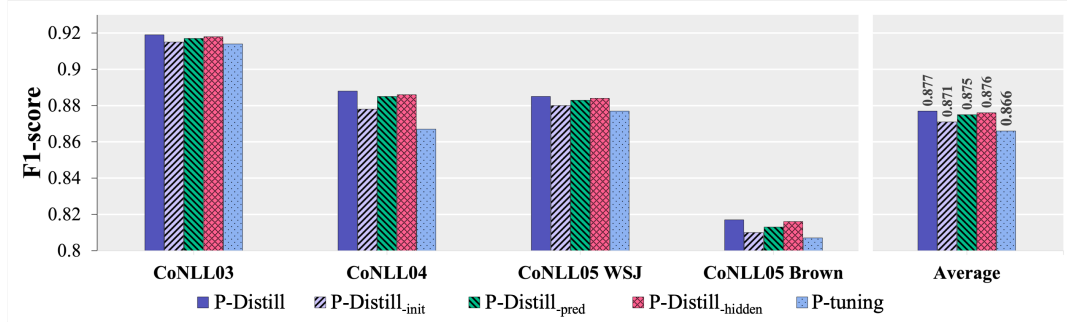


Figure 4: Comparison of ablation study results across various tasks, with different colors and bar styles representing the distinct variants of P-Distill.

various tasks. This shows that each component of P-Distill contributes to the performance achieved on downstream tasks.

Given these results, we further observe that the extent of degradation varies among different variants. First, P-Distill_{init} exhibits the most significant performance degradation across various tasks. While not utilizing prompt initialization and only conducting prediction layer distillation and hidden state distillation still led to performance improvements over P-tuning, the degradation in performance is clear compared to other variants and the original P-Distill. This indicates that prompt initialization, based on the teacher model’s prompt, is crucial in prompt-based knowledge distillation. Second, P-Distill_{hidden} and P-Distill_{pred} exhibited decreased prediction performance. This demonstrates that integrating prompt initialization with the hidden state or prediction layer distillation techniques enhances the stability and effectiveness of knowledge distillation.

4.6 Impact of Prompt Initialization

To inspect the impact of different prompt initialization methods within P-Distill, we conduct experiments to compare the performance of P-Distill with two variants: P-Distill_{mean}, which initializes the student model prompts using an average pooling layer over the teacher model prompts, and P-Distill_{max}, which uses a max pooling layer for the same purpose.

The results, as detailed in Table 3, demonstrate that both P-Distill_{mean} and P-Distill_{max} underperform P-Distill which utilizes a reparameterization encoder for prompt initialization. We conjecture that the use of average pooling and max pooling leads to an excessive simplification of the teacher model’s prompts, resulting in the loss of crucial nuances and complexities. Conversely, the reparameterization encoder for prompt initialization effectively captures and transfers the complex knowledge of the teacher model prompts without loss of crucial task information. This suggests that the reparameterization encoder is a more suitable method for prompt initialization in P-Distill, contributing significantly to the overall effectiveness of the knowledge distillation process.

	CoNLL03	CoNLL04	CoNLL05 WSJ	CoNLL05 Brown
P-Distill	0.919	0.888	0.885	0.817
P-Distill _{mean}	0.915	0.875	0.878	0.809
P-Distill _{max}	0.912	0.872	0.872	0.803

Table 3: Comparison of additional experiment results across various tasks based on prompt initialization methods. All metrics are reported as micro-f1 scores. (**bold**: the best)

parameterization encoder for prompt initialization effectively captures and transfers the complex knowledge of the teacher model prompts without loss of crucial task information. This suggests that the reparameterization encoder is a more suitable method for prompt initialization in P-Distill, contributing significantly to the overall effectiveness of the knowledge distillation process.

5 Conclusion

In this paper, we introduce P-Distill, a novel approach in NLP that utilizes two knowledge distillation techniques to enhance performance by compressing unnecessary prompt length. This approach combines prompt initialization, two types of prompt distillation to effectively transfer knowledge from a teacher model with longer prompts to a student model with prompts that are eight times shorter. To evaluate the efficacy of our proposed method, we conduct experiments across various NLP tasks. Our results demonstrate that using prompts of the same length, the proposed method achieves an average improvement of 2.75% over the existing prompt-tuning methods across the SuperGLUE benchmark. Furthermore, P-Distill exhibits competitive performance even against models trained with prompts that are eight times longer.

608 Limitations

609 One limitation of this study is that we evaluated our
610 method only on the BERT architecture. Conduct-
611 ing additional experiments on other architectures
612 could be beneficial to determine the generalizabil-
613 ity of our findings. Additionally, while our model
614 improves performance through the process of train-
615 ing a teacher model and transferring its knowledge,
616 it incurs more time and cost compared to previous
617 methods. In future work, we plan to develop an
618 approach that integrates the training of the teacher
619 model and the knowledge distillation process in an
620 end-to-end manner.

621 References

622 Roy Bar Haim, Ido Dagan, Bill Dolan, Lisa Ferro,
623 Danilo Giampiccolo, Bernardo Magnini, and Idan
624 Szpektor. 2006. The second PASCAL recognising
625 textual entailment challenge.

626 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie
627 Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind
628 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
629 Askell, Sandhini Agarwal, Ariel Herbert-Voss,
630 Gretchen Krueger, Tom Henighan, Rewon Child,
631 Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,
632 Clemens Winter, Christopher Hesse, Mark Chen, Eric
633 Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess,
634 Jack Clark, Christopher Berner, Sam McCandlish,
635 Alec Radford, Ilya Sutskever, and Dario Amodei.
636 2020. [Language models are few-shot learners](#). *CoRR*,
637 abs/2005.14165.

638 Xavier Carreras and Lluís Màrquez. 2004. [Introduction](#)
639 [to the CoNLL-2004 shared task: Semantic role la-](#)
640 [beling](#). In *Proceedings of the Eighth Conference on*
641 *Computational Natural Language Learning (CoNLL-*
642 *2004) at HLT-NAACL 2004*, pages 89–97, Boston,
643 Massachusetts, USA. Association for Computational
644 Linguistics.

645 Xavier Carreras and Lluís Màrquez. 2005. [Introduction](#)
646 [to the CoNLL-2005 shared task: Semantic role la-](#)
647 [beling](#). In *Proceedings of the Ninth Conference on*
648 *Computational Natural Language Learning (CoNLL-*
649 *2005)*, pages 152–164, Ann Arbor, Michigan. Asso-
650 ciation for Computational Linguistics.

651 Shouyuan Chen, Sherman Wong, Liangjian Chen, and
652 Yuandong Tian. 2023. Extending context window of
653 large language models via positional interpolation.
654 *arXiv preprint arXiv:2306.15595*.

655 Christopher Clark, Kenton Lee, Ming-Wei Chang,
656 Tom Kwiatkowski, Michael Collins, and Kristina
657 Toutanova. 2019. BoolQ: Exploring the surprising
658 difficulty of natural yes/no questions. In *Proceedings*
659 *of NAACL-HLT 2019*.

660 Ido Dagan, Oren Glickman, and Bernardo Magnini.
661 2006. The PASCAL recognising textual entailment
662 challenge. In *Machine learning challenges. evaluat-*
663 *ing predictive uncertainty, visual object classification,*
664 *and recognising textual entailment*, pages 177–190.
665 Springer.

666 Marie-Catherine De Marneffe, Mandy Simons, and
667 Judith Tonhauser. 2019. The Commitment-
668 Bank: Investigating projection in naturally occur-
669 ring discourse. To appear in proceedings of
670 Sinn und Bedeutung 23. Data can be found at
671 <https://github.com/mcdm/CommitmentBank/>.

672 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
673 Kristina Toutanova. 2019. [BERT: Pre-training of](#)
674 [deep bidirectional transformers for language un-](#)
675 [derstanding](#). In *Proceedings of the 2019 Conference of*
676 *the North American Chapter of the Association for*
677 *Computational Linguistics: Human Language Tech-*
678 *nologies, Volume 1 (Long and Short Papers)*, pages
679 4171–4186, Minneapolis, Minnesota. Association for
680 Computational Linguistics.

681 Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015.
682 [Distilling the knowledge in a neural network](#).

683 Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski,
684 Bruna Morrone, Quentin de Laroussilhe, Andrea Ges-
685 mundo, Mona Attariyan, and Sylvain Gelly. 2019.
686 [Parameter-efficient transfer learning for nlp](#).

687 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan
688 Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and
689 Weizhu Chen. 2021. [Lora: Low-rank adaptation of](#)
690 [large language models](#).

691 Zhengbao Jiang, Frank F. Xu, Jun Araki, and Graham
692 Neubig. 2020. [How can we know what language](#)
693 [models know?](#) *Transactions of the Association for*
694 *Computational Linguistics*, 8:423–438.

695 Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao
696 Chen, Linlin Li, Fang Wang, and Qun Liu. 2020.
697 [TinyBERT: Distilling BERT for natural language un-](#)
698 [derstanding](#). In *Findings of the Association for Com-*
699 *putational Linguistics: EMNLP 2020*, pages 4163–
700 4174, Online. Association for Computational Lin-
701 guistics.

702 Daniel Khashabi, Snigdha Chaturvedi, Michael Roth,
703 Shyam Upadhyay, and Dan Roth. 2018. Looking
704 beyond the surface: A challenge set for reading com-
705 prehension over multiple sentences. In *Proceedings*
706 *of the 2018 Conference of the North American Chap-*
707 *ter of the Association for Computational Linguistics:*
708 *Human Language Technologies, Volume 1 (Long Pa-*
709 *pers)*, pages 252–262.

710 Solomon Kullback and Richard A Leibler. 1951. On
711 information and sufficiency. *The annals of mathe-*
712 *matical statistics*, 22(1):79–86.

713 Brian Lester, Rami Al-Rfou, and Noah Constant. 2021.
714 [The power of scale for parameter-efficient prompt](#)
715 [tuning](#). In *Proceedings of the 2021 Conference on*

716	<i>Empirical Methods in Natural Language Processing</i> , pages 3045–3059, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.		
717			
718			
719			
720	Hector J Levesque, Ernest Davis, and Leora Morgenstern. 2011. The Winograd schema challenge. In <i>AAAI Spring Symposium: Logical Formalizations of Commonsense Reasoning</i> , volume 46, page 47.		
721			
722			
723			
724	Chang Li and Dan Goldwasser. 2019. Encoding social information with graph convolutional networks for Political perspective detection in news media . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 2594–2604, Florence, Italy. Association for Computational Linguistics.		
725			
726			
727			
728			
729			
730			
731	Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning: Optimizing continuous prompts for generation .		
732			
733	Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Tam, Zhengxiao Du, Zhilin Yang, and Jie Tang. 2022. P-tuning: Prompt tuning can be comparable to fine-tuning across scales and tasks . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 61–68, Dublin, Ireland. Association for Computational Linguistics.		
734			
735			
736			
737			
738			
739			
740			
741	Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. 2023. Gpt understands, too . <i>AI Open</i> .		
742			
743			
744	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 .		
745			
746			
747			
748			
749	Mohammad Taher Pilehvar and Jose Camacho-Collados. 2019. WiC: The word-in-context dataset for evaluating context-sensitive meaning representations . In <i>Proceedings of NAACL-HLT</i> .		
750			
751			
752			
753	Sameer Pradhan, Alessandro Moschitti, Nianwen Xue, Olga Uryupina, and Yuchen Zhang. 2012. CoNLL-2012 shared task: Modeling multilingual unrestricted coreference in OntoNotes . In <i>Joint Conference on EMNLP and CoNLL - Shared Task</i> , pages 1–40, Jeju Island, Korea. Association for Computational Linguistics.		
754			
755			
756			
757			
758			
759			
760	Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. Know what you don’t know: Unanswerable questions for SQuAD . 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.		
761			
762			
763			
764			
765			
766			
767	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text . In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 2383–2392, Austin, Texas. Association for Computational Linguistics.		
768			
769			
770			
771			
772			
	Melissa Roemmele, Cosmin Adrian Bejan, and Andrew S. Gordon. 2011. Choice of plausible alternatives: An evaluation of commonsense causal reasoning . In <i>2011 AAAI Spring Symposium Series</i> .		
	Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. 2015. Fitnets: Hints for thin deep nets .		
	Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2020. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter .		
	Taylor Shin, Yasaman Razeghi, Robert L. Logan IV, Eric Wallace, and Sameer Singh. 2020. Autoprompt: Eliciting knowledge from language models with automatically generated prompts . <i>CoRR</i> , abs/2010.15980.		
	Zhiqing Sun, Hongkun Yu, Xiaodan Song, Renjie Liu, Yiming Yang, and Denny Zhou. 2020. MobileBERT: a compact task-agnostic BERT for resource-limited devices . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 2158–2170, Online. Association for Computational Linguistics.		
	Erik F. Tjong Kim Sang and Fien De Meulder. 2003. Introduction to the CoNLL-2003 shared task: Language-independent named entity recognition . In <i>Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003</i> , pages 142–147.		
	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruiti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models .		
	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2023. Attention is all you need .		
	Tu Vu, Brian Lester, Noah Constant, Rami Al-Rfou’, and Daniel Cer. 2022. SPoT: Better frozen model adaptation through soft prompt transfer . In <i>Proceedings of the 60th Annual Meeting of the Association</i>		

- 831 *for Computational Linguistics (Volume 1: Long Pa-*
832 *pers)*, pages 5039–5059, Dublin, Ireland. Association
833 for Computational Linguistics.
- 834 Alex Wang, Yada Pruksachatkun, Nikita Nangia, Aman-
835 preet Singh, Julian Michael, Felix Hill, Omer Levy,
836 and Samuel R. Bowman. 2019. [Superglue: A stickier](#)
837 [benchmark for general-purpose language understand-](#)
838 [ing systems](#). *CoRR*, abs/1905.00537.
- 839 Ralph Weischedel, Martha Palmer, Mitchell Marcus,
840 Eduard Hovy, Sameer Pradhan, Lance Ramshaw, Ni-
841 anwen Xue, Ann Taylor, Jeff Kaufman, Michelle
842 Franchini, et al. 2013. Ontonotes release 5.0
843 ldc2013t19. *Linguistic Data Consortium, Philadel-*
844 *phia, PA*, 23:170.
- 845 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien
846 Chaumond, Clement Delangue, Anthony Moi, Pier-
847 ric Cistac, Tim Rault, Remi Louf, Morgan Funtow-
848 icz, Joe Davison, Sam Shleifer, Patrick von Platen,
849 Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu,
850 Teven Le Scao, Sylvain Gugger, Mariama Drame,
851 Quentin Lhoest, and Alexander Rush. 2020. [Trans-](#)
852 [formers: State-of-the-art natural language processing](#).
853 In *Proceedings of the 2020 Conference on Empirical*
854 *Methods in Natural Language Processing: System*
855 *Demonstrations*, pages 38–45, Online. Association
856 for Computational Linguistics.
- 857 Sheng Zhang, Xiaodong Liu, Jingjing Liu, Jianfeng
858 Gao, Kevin Duh, and Benjamin Van Durme. 2018.
859 [Record: Bridging the gap between human and ma-](#)
860 [chine commonsense reading comprehension](#).

Appendix

A Impact of Compression Ratio

	CoNLL05 WSJ
Fine-tuning	0.885
P-tuning	0.890 ₍₁₂₈₎
	0.890 ₍₆₄₎
P-Distill	0.888 ₍₃₂₎
	0.885 ₍₁₆₎

Table 4: Comparison of P-Distill performance across varying prompt compression ratios.

To investigate the impact of different compression ratios on the performance of P-Distill, we conduct additional experiments with $2\times$ and $4\times$ compression ratios and compare the results. We use CoNLL04 dataset in these experiments, as it is known to require longer prompts for sufficient training.

Table 4 presents the performance results for these different compression ratios within CoNLL04. We observe a trade-off between compression ratio and model performance, as a higher compression ratio leads to relatively lower performance. Yet, it is worth noting that a $2\times$ compression with P-Distill does preserve the performance of the original long prompt, and that P-Distill still matches the performance of fine-tuning at the compression ratio of $8\times$.

B Inference Costs

	Fine-tuning	P-tuning	P-Distill
BoolQ	89.06	89.17 ₍₈₎	89.07 ₍₁₎
CB	53.94	54.22 ₍₃₂₎	53.97 ₍₄₎
COPA	21.82	21.83 ₍₄₎	21.82 ₍₁₎
MultiRC	226.94	227.50 ₍₁₆₎	227.01 ₍₂₎
ReCoRD	163.12	163.53 ₍₁₆₎	163.17 ₍₂₎
RTE	42.78	42.81 ₍₄₎	42.79 ₍₁₎
WiC	18.83	18.84 ₍₄₎	18.83 ₍₁₎
WSC	23.11	23.17 ₍₄₎	23.11 ₍₁₎

Table 5: Comparing GFLOPs of baseline methods and P-Distill on SuperGLUE using BERT_{large}.

To examine the benefits of P-Distill at the inference stage, we compare the inference costs across P-Distill, Fine-tuning, and P-tuning. This reduction in computational requirements is quantified in

Table 5, which presents the GFLOPs required during the inference stage with average length samples from each task in the SuperGLUE benchmark.

While both P-tuning and P-Distill demonstrate increased inference GFLOPs compared to fine-tuning due to the inclusion of prompts, we observe that P-Distill adds fewer computations compared to P-tuning. The difference is more clear in cases where long prompts are utilized in P-tuning. This demonstrates the advantage of compressing prompt length in terms of lowering computational costs.

C Qualitative Analysis in Long Sequence Classification

To further examine the effectiveness of expanding the input space by using shorter prompts via P-Distill, we conduct a qualitative analysis on the Allsides dataset. The task of it is to predict the political perspectives inherent in a news article. Therefore, the dataset consists of news articles that generally exceed the model’s input capacity of 512 tokens.

We compare the input text and prediction results of P-tuning, trained with 32 prompt tokens, to those of P-Distill, which compresses and appends only 4 prompt tokens, which are shown in Table 2. Note that while P-tuning can utilize up to 480 input tokens by attaching 32 prompt tokens, P-Distill extends the input capacity by using up to 508 input tokens with only 4 prompt tokens appended.

We observe that the P-tuning model cannot access the detailed explanation of the Fairness Doctrine due to its limited input space. However, P-Distill can access the remainder of the sentence, namely ‘*broadcasters provide “equal time” to divergent political views.*’, which provides the key information for the model in making an accurate prediction ‘Center’. This example verifies that the input space preserved by compressing the prompt length with P-Distill can contribute to accurate prediction of the model, leading to better performance overall.

P-tuning Input Text (480 tokens): President Trump on Wednesday lashed out over a critical news report and escalated his previous attacks on the media by suggesting that news organizations he disagrees with be shut down, alarming free-speech advocates who compared the tactics to intimidation efforts by the Nixon administration.

[...]

Last week, angered by the ongoing investigations into his campaign’s ties to Russia, Trump suggested that the Senate Intelligence Committee investigate news outlets over “fake news.” Over the weekend, he expressed disdain at late-night television hosts over their “anti-Trump” material and *proposed bringing back the Fairness Doctrine, a rule phased out in 1987 that had required*

Prediction: Left

P-Distill Input Text (508 tokens): President Trump on Wednesday lashed out over a critical news report and escalated his previous attacks on the media by suggesting that news organizations he disagrees with be shut down, alarming free-speech advocates who compared the tactics to intimidation efforts by the Nixon administration.

[...]

Last week, angered by the ongoing investigations into his campaign’s ties to Russia, Trump suggested that the Senate Intelligence Committee investigate news outlets over “fake news.” Over the weekend, he expressed disdain at late-night television hosts over their “anti-Trump” material and *proposed bringing back the Fairness Doctrine, a rule phased out in 1987 that had required broadcasters to provide “equal time” for divergent political views on certain issues.* First Amendment advocates roundly condemned the president over his remarks, calling them an assault

Prediction: Center

Table 6: Example of P-tuning and P-Distill Predictions in the Allsides Dataset.