

# 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 approach that mitigates this limitation by utilizing knowledge distillation from a teacher model with extensive prompts to a student model with shorter prompts. We introduce two novel methods for prompt compression, including prompt initialization, and prompt distillation. Experiments across 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. 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. These results highlight the potential of P-Distill in facilitating efficient and effective training for a wide range of NLP models.

## 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

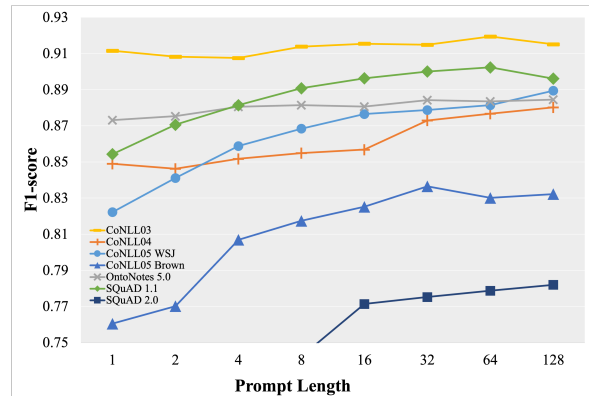


Figure 1: Performance variation in P-tuning across tasks based on the length of continuous prompts.

explored various methods for efficiently fine-tuning 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. These prompts 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). Similar to the findings in the work (Liu et al., 2022), more challenging tasks require longer prompt lengths to achieve the optimal 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 achieve the optimal performance.

In this paper, we propose P-Distill, which is a novel prompt compression method to address the limitations of long prompts. Our method involves a two-step process where we first train a teacher model using P-tuning to achieve optimal performance with long prompts. We then transfer this

068 knowledge to a student model with significantly  
 069 shorter prompts through a distillation process. To  
 070 ensure stability in training, we first perform prompt  
 071 initialization based on the teacher model prompts.  
 072 Then, we focus on distilling knowledge between  
 073 the teacher and student models, specifically target-  
 074 ing the outputs of their intermediate and prediction  
 075 layers. This is due to the impact of soft prompts  
 076 on the hidden states within these layers, which  
 077 subsequently influences the model’s predictions.  
 078 This method enables the compression of prompts  
 079 to shorter lengths without a significant degradation  
 080 in performance, thereby addressing the inefficien-  
 081 cies inherent in longer prompts.

082 To validate its effectiveness and efficiency, we  
 083 evaluate P-Distill using various NLP benchmarks,  
 084 including SuperGLUE (Wang et al., 2020), SQuAD  
 085 (Rajpurkar et al., 2016, 2018), ReCoRD (Zhang  
 086 et al., 2018), OntoNotes (Weischedel et al., 2013),  
 087 and CoNLL (Tjong Kim Sang and De Meulder,  
 088 2003; Carreras and Màrquez, 2004, 2005; Pradhan  
 089 et al., 2012a). Our results demonstrate that P-Distill  
 090 exhibits comparable or superior performance than  
 091 those of the existing state-of-the-art prompt-based  
 092 models. To the best of our knowledge, this study is  
 093 the first to train teacher prompts and transfer their  
 094 knowledge to student prompts for the purpose of  
 095 compressing prompts. The main contributions of  
 096 this study are summarized as follows:

- 097 • We propose a method called P-Distill to com-  
 098 press the soft prompts, effectively mitigating  
 099 the limitation of reducing the model’s usable  
 100 sequence length in prompt-based learning.
- 101 • We introduce a prompt distillation method uti-  
 102 lizing teacher model’s hidden state and pre-  
 103 diction outputs, influenced by soft prompts,  
 104 and propose a prompt initialization for stable  
 105 prompt distillation.
- 106 • We validate P-Distill across multiple NLP  
 107 benchmarks, demonstrating its ability to main-  
 108 tain or enhance accuracy while reducing  
 109 prompt lengths by up to eight times.

110 The remainder of this paper is structured as fol-  
 111 lows: Section 2 provides the preliminaries; Section  
 112 3 describes a detailed description of the proposed  
 113 method; Section 4 presents the experimental results  
 114 and analysis, and Section 5 concludes the study.

## 2 Preliminaries 115

### 2.1 Pre-trained Language Models Based on the Transformer 116

117 The transformer model (Vaswani et al., 2023), com-  
 118 prising an encoder and decoder, is the fundamental  
 119 architecture of the majority of recent PLMs, includ-  
 120 ing BERT (Devlin et al., 2019), RoBERTa (Liu  
 121 et al., 2019), and GPT-3 (Brown et al., 2020). Each  
 122 encoder and decoder consists of multiple trans-  
 123 former layers and incorporates key components,  
 124 such as multi-head attention modules (MHA), feed-  
 125 forward networks, layer normalization, and resid-  
 126 ual connections. A key component of this architec-  
 127 ture is the multi-head attention mechanism, which  
 128 computes attention weights using query ( $Q$ ), key  
 129 ( $K$ ), and value ( $V$ ) matrices. Mathematically, the  
 130 attention function in multi-head attention can be  
 131 represented as follows: 132

$$Att(x) = softmax(\frac{QK^T}{\sqrt{d_k}})V, \quad (1) \quad 133$$

134 where  $\sqrt{d_k}$  is the scaling factor for gradient stabi-  
 135 lization during training. This attention mechanism  
 136 is crucial in understanding language and generat-  
 137 ing tasks by modulating the focus of the model on  
 138 different parts of the input data.

### 2.2 Prompt-based Learning Methods 139

140 Prompt-based learning methods have emerged as  
 141 an efficient alternative to full model fine-tuning,  
 142 especially for PLMs (Liu et al., 2022, 2023). These  
 143 methods use prompts to guide the model pre-  
 144 dictions for specific tasks. Several approaches  
 145 (Jiang et al., 2020; Shin et al., 2020) employ dis-  
 146 crete prompts, which are fixed templates added  
 147 to the input. For example, in sentiment analy-  
 148 sis, a template might be "This text  $[Input\ Text]$   
 149 expresses a  $[MASK]$  sentiment.". However, dis-  
 150 crete prompts are limited in that their performances  
 151 significantly depend on template selection. Ad-  
 152 vanced approaches, such as Prefix-Tuning (Li and  
 153 Liang, 2021) and P-tuning (Liu et al., 2022), use  
 154 continuous prompts that are trainable embeddings  
 155 independent of the model vocabulary. Particularly,  
 156 P-tuning attaches continuous prompts to each layer  
 157 of the model, thereby influencing its behavior and  
 158 enhancing its performance in downstream tasks.  
 159 This approach is mathematically represented as:

$$T = \{h_{0:i}, e(x), h_{i+1:m}, e(y)\}, \quad (2) \quad 160$$

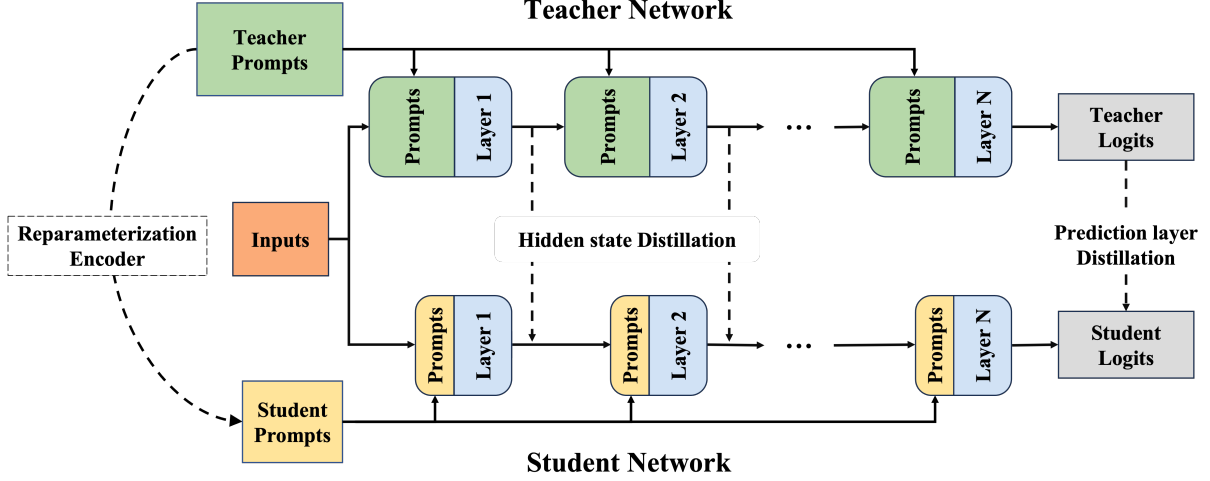


Figure 2: Overview of the proposed method, denoted as P-Distill. This method trains a teacher model to generate concise and effective prompts, followed by distilling the knowledge into a student model.

where  $h_i$  denotes the trainable embedding vector of the continuous prompts. These continuous prompts are integrated into the attention mechanism of the transformer model as follows:

$$head_l(x) = Att(xW^{(l)}, [P_k^{(l)} : K^{(l)}], [P_v^{(l)} : V^{(l)}]), \quad (3)$$

where  $head_l(x)$  is the attention computation for each attention head  $l$ . The query vector  $Q^{(l)}$  is generated using the input  $x$  and the weight matrix  $W^{(l)}$ .  $K^{(l)}$  and  $V^{(l)}$  are the key and value vectors for the  $l$ -th attention head, and  $P_k^{(l)}$  and  $P_v^{(l)}$  are the continuous prompts added to the key and value vectors of the  $l$ -th attention head, respectively. 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:

$$L_{KD} = \sum_{x \in X} L(f^S(x), f^T(x)), \quad (4)$$

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, the performance is only influenced by the continuous prompts, because the pre-trained weights of the language model remain fixed. Furthermore, the optimal prompt length varies with task complexity. For simple sentence classification tasks, the optimal length is approximately 20, whereas for more difficult sequence-labeling tasks, it can extend up to 100 (Liu et al., 2022). Understanding the variation in prompt length is crucial, as longer prompts inherently limit the maximum sequence length that the model can handle.

We train a teacher model on various tasks using the P-tuning methodology, which fixes the pre-trained weights of the PLM. 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:

$$L_{CE}^T = -\frac{1}{|B|} \sum_{i=1}^{|B|} \log(\text{softmax}(y_i^T)[c_i]), \quad (5)$$

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

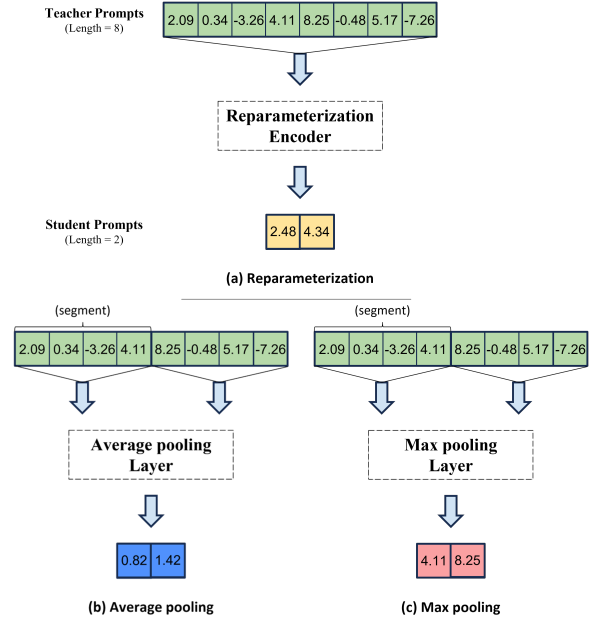


Figure 3: Illustration of various prompt initialization methods.

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, 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 student model’s prompts as follows:

$$P_k^S = (P_k^T \cdot W_k^T) + b_k^T, \quad (6)$$

$$P_v^S = (P_v^T \cdot W_v^T) + b_v^T, \quad (7)$$

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 terms used. 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 soft 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:

$$L_{pred} = KL(\text{softmax}(y_i^S / \theta), \text{softmax}(y_i^T / \theta)), \quad (8)$$

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.  $\theta$  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:

$$L_{hidden} = MSE(H_S, H_T), \quad (9)$$

where the matrices  $H_S, H_T \in R^{n \times d}$  represent the hidden states,  $n$  is the 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:

$$L_{CE}^S = -\frac{1}{|B|} \sum_{i=1}^{|B|} \log(\text{softmax}(y_i^S)[c_i]). \quad (10)$$

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:

$$L_{total} = \lambda_1 \cdot L_{CE}^S + \lambda_2 \cdot L_{pred} + \lambda_3 \cdot L_{hidden}, \quad (11)$$

where  $\lambda_1$ ,  $\lambda_2$ , and  $\lambda_3$  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.

Table 1: Experimental results for each model 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)

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

## 4.1 Datasets

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

For named entity recognition, we use 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., 2012b), 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), facilitate testing reading comprehension, requiring the model to parse passages and answer questions with a high degree of understanding. Furthermore, we include various tasks from 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). All these datasets are English, open-source, and utilized solely for academic research purposes. 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<sub>large</sub> with 335M parameters as the backbone architecture.

**Fine-tuning** All parameters of a PLM are updated to the downstream task, thereby adapting the weights of the entire model to the task-specific data.

**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.

**P-Distill** Our proposed method, P-Distill, employs a P-tuning approach to train continuous prompts for each task. Subsequently, the optimally trained continuous prompts are distilled into a student model with shorter prompts using two distinct knowledge distillation techniques.

## 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  $\theta$  used in the distillation process, we experimentally determine the optimal setting by sweeping across [1, 5, 10]. For the learnable parameter  $\lambda_2$ , 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  $\lambda_3$ . All experiments were performed using PyTorch<sup>1</sup> and HuggingFace Transformers (Wolf et al., 2020) on three NVIDIA A100 GPUs, and to ensure consistency in our results, each task was conducted using the same random seed.

## 4.4 Results

Tables 1 and 2 present the experimental results of Fine-tuning, P-tuning, and P-Distill. In P-Distill,

<sup>1</sup><https://pytorch.org/>

Table 2: Experimental results for each method on named entity recognition, question answering, and semantic role labeling. 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 micro-f1 scores. (**bold**: the best; underline: the second best)

	CoNLL03	CoNLL04	CoNLL05 WSJ	CoNLL05 Brown	OntoNotes 5.0	SQuAD 1.1 dev	SQuAD 2.0 dev
Fine-tuning	<b>0.928</b>	0.882	0.885	0.827	<b>0.890</b>	<b>0.911</b>	<b>0.819</b>
P-tuning	0.919 <sub>(64)</sub>	0.880 <sub>(128)</sub>	<b>0.890</b> <sub>(128)</sub>	<b>0.837</b> <sub>(32)</sub>	0.885 <sub>(128)</sub>	0.902 <sub>(64)</sub>	0.782 <sub>(128)</sub>
	0.914 <sub>(8)</sub>	0.866 <sub>(16)</sub>	0.877 <sub>(16)</sub>	0.807 <sub>(4)</sub>	0.881 <sub>(16)</sub>	0.891 <sub>(8)</sub>	0.771 <sub>(16)</sub>
P-Distill	0.919 <sub>(8)</sub>	<b>0.888</b> <sub>(16)</sub>	0.885 <sub>(16)</sub>	0.817 <sub>(4)</sub>	0.886 <sub>(16)</sub>	0.896 <sub>(8)</sub>	0.775 <sub>(16)</sub>

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. Experimentally, 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. Despite using shorter prompts, P-Distill matches and exceeds the performances of the baseline methods, particularly P-tuning.

When applying the P-tuning, we observed comparable or superior performance in all tasks, except for BoolQ, when compared to the Fine-tuning. However, a limitation of P-tuning is that the length of tokens that can be input into the model is reduced by the length of the prompt used. For instance, setting the prompt length to 32 in the CB task resulted in a performance comparable to that of Fine-tuning. However, this also resulted in a reduction of 32 tokens in the input length. Applying the proposed P-Distill compresses the prompt length by eight-fold, resulting in a better performance even when the prompt length was reduced to 4. Particularly, P-Distill exhibited a 1.90% higher performance than that of the teacher model and 2.73% improvement over the same-length prompt trained using P-tuning. This indicates that knowledge learned from a teacher model with longer prompts can be effectively transferred to a student model with eight times shorter prompts.

**Results on Across Tasks** Based on Table 2, achieving optimal performance via P-tuning requires training with longer prompts for question-answering, named entity recognition, and semantic role labeling tasks. Similar to the SuperGLUE benchmark and compared to the existing methods, P-Distill achieved comparable or better per-

formance using shorter prompts even for tasks requiring longer prompts. This is evident in the CoNLL04 task, where the optimal prompt length for the model trained using P-tuning is 128. Despite occupying a significant portion of the input token space with a length of 128, the performance was lower than that achieved with Fine-tuning. However, applying prompt compression using P-Distill reduced the prompt length to 16 while outperforming Fine-tuning. Notably, P-Distill achieves a performance improvement of 2.54% over the same-length prompt trained using P-tuning and a 0.90% improvement over the teacher model.

#### 4.5 Additional Study

To further verify the effectiveness of the proposed method, we conduct experiments using the following P-Distill variants:

**P-Distill<sub>-init</sub>** Instead of training without 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<sub>-pred</sub>** 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<sub>-hidden</sub>** 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 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

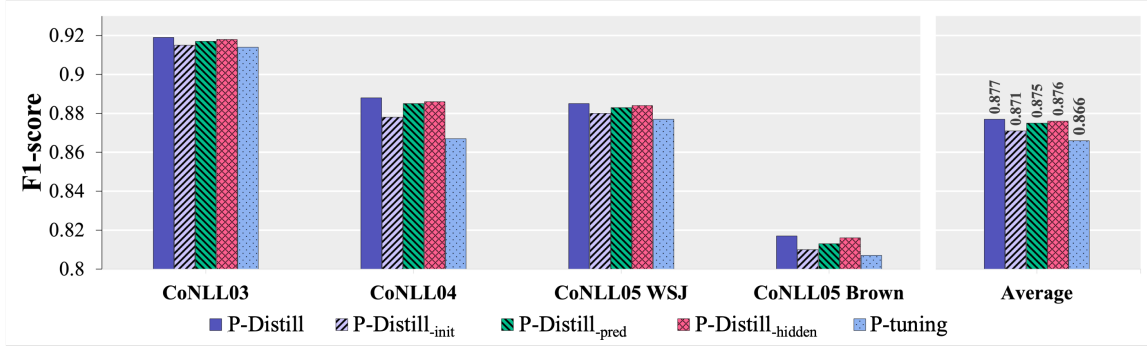


Figure 4: Comparison of ablation study results across various tasks, with different colors and bar styles representing the distinct variants of P-Distill.

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)

	CoNLL03	CoNLL04	CoNLL05 WSJ	CoNLL05 Brown
P-Distill	<b>0.919</b>	<b>0.888</b>	<b>0.885</b>	<b>0.817</b>
P-Distill <sub>mean</sub>	0.915	0.875	0.878	0.809
P-Distill <sub>max</sub>	0.912	0.872	0.872	0.803

internal representations.

Figure 4 shows that all three variants of P-Distill cause performance degradation, which is evident in downstream tasks and overall averages. However, the extent of degradation varies among different variants. First, P-Distill<sub>init</sub> exhibited the most significant performance degradation across various tasks. Even without prompt initialization, conducting prediction layer distillation and hidden state distillation led to performance improvements over P-tuning. However, we observed lower performance when prompts were randomly initialized compared to when prompt initialization is applied. This indicates that prompt initialization, based on the teacher model’s prompt, is crucial in prompt-based knowledge distillation. Second, P-Distill<sub>hidden</sub> and P-Distill<sub>pred</sub> 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. Therefore, combining these methods is the most effective approach for prompt-based knowledge distillation, resulting in the best performance.

To inspect the effectiveness of different prompt initialization methods within P-Distill, we conduct experiments to compare the performance of P-Distill with two variants: P-Distill<sub>mean</sub>, which

initializes the teacher model prompts using an average pooling layer, and P-Distill<sub>max</sub>, which uses a max pooling layer for the same purpose. The results, as detailed in Table 3, demonstrated that both P-Distill<sub>mean</sub> and P-Distill<sub>max</sub> underperformed in comparison to the P-Distill. We assume 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, thereby enhancing the predictive performance. 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 2.73% 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.



## 601 Limitations

602 One limitation of this study is that we evaluated our  
603 method only on the BERT architecture. Conducting  
604 additional experiments on other architectures  
605 could be beneficial to determine the generalizability  
606 of our findings. Additionally, while our model  
607 improves performance through the process of training  
608 a teacher model and transferring its knowledge,  
609 it incurs more time and cost compared to previous  
610 methods. In future work, we plan to develop an  
611 approach that integrates the training of the teacher  
612 model and the knowledge distillation process in an  
613 end-to-end manner.

## 614 References

615 Roy Bar Haim, Ido Dagan, Bill Dolan, Lisa Ferro,  
616 Danilo Giampiccolo, Bernardo Magnini, and Idan  
617 Szpektor. 2006. The second PASCAL recognising  
618 textual entailment challenge.

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

631 Xavier Carreras and Lluís Màrquez. 2004. [Introduction](#)  
632 [to the CoNLL-2004 shared task: Semantic role labeling](#).  
633 In *Proceedings of the Eighth Conference on*  
634 *Computational Natural Language Learning (CoNLL-*  
635 *2004) at HLT-NAACL 2004*, pages 89–97, Boston,  
636 Massachusetts, USA. Association for Computational  
637 Linguistics.

638 Xavier Carreras and Lluís Màrquez. 2005. [Introduction](#)  
639 [to the CoNLL-2005 shared task: Semantic role labeling](#).  
640 In *Proceedings of the Ninth Conference on*  
641 *Computational Natural Language Learning (CoNLL-*  
642 *2005)*, pages 152–164, Ann Arbor, Michigan. Association  
643 for Computational Linguistics.

644 Christopher Clark, Kenton Lee, Ming-Wei Chang,  
645 Tom Kwiatkowski, Michael Collins, and Kristina  
646 Toutanova. 2019. BoolQ: Exploring the surprising  
647 difficulty of natural yes/no questions. In *Proceedings*  
648 *of NAACL-HLT 2019*.

649 Ido Dagan, Oren Glickman, and Bernardo Magnini.  
650 2006. The PASCAL recognising textual entailment  
651 challenge. In *Machine learning challenges. evaluating*  
652 *predictive uncertainty, visual object classification,*  
653 *and recognising textual entailment*, pages 177–190.  
654 Springer.

Marie-Catherine De Marneffe, Mandy Simons, and  
655 Judith Tonhauser. 2019. The Commitment-  
656 Bank: Investigating projection in naturally occurring  
657 discourse. To appear in proceedings of  
658 Sinn und Bedeutung 23. Data can be found at  
659 <https://github.com/mcdm/CommitmentBank/>.  
660

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and  
661 Kristina Toutanova. 2019. [BERT: Pre-training of](#)  
662 [deep bidirectional transformers for language understanding](#).  
663 In *Proceedings of the 2019 Conference of*  
664 *the North American Chapter of the Association for*  
665 *Computational Linguistics: Human Language Tech-*  
666 *nologies, Volume 1 (Long and Short Papers)*, pages  
667 4171–4186, Minneapolis, Minnesota. Association for  
668 Computational Linguistics.  
669

Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. [Distilling the knowledge in a neural network](#).  
670  
671

Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski,  
672 Bruna Morrone, Quentin de Laroussilhe, Andrea Ges-  
673 mundo, Mona Attariyan, and Sylvain Gelly. 2019. [Parameter-efficient transfer learning for nlp](#).  
674  
675

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan  
676 Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and  
677 Weizhu Chen. 2021. [Lora: Low-rank adaptation of](#)  
678 [large language models](#).  
679

Zhengbao Jiang, Frank F. Xu, Jun Araki, and Graham  
680 Neubig. 2020. [How can we know what language](#)  
681 [models know?](#) *Transactions of the Association for*  
682 *Computational Linguistics*, 8:423–438.  
683

Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao  
684 Chen, Linlin Li, Fang Wang, and Qun Liu. 2020. [TinyBERT: Distilling BERT for natural language understanding](#).  
685 In *Findings of the Association for Computational*  
686 *Linguistics: EMNLP 2020*, pages 4163–  
687 4174, Online. Association for Computational Lin-  
688 guistics.  
689  
690

Daniel Khashabi, Snigdha Chaturvedi, Michael Roth,  
691 Shyam Upadhyay, and Dan Roth. 2018. Looking  
692 beyond the surface: A challenge set for reading com-  
693 prehension over multiple sentences. In *Proceedings*  
694 *of the 2018 Conference of the North American Chapter*  
695 *of the Association for Computational Linguistics: Human*  
696 *Language Technologies, Volume 1 (Long Papers)*, pages  
697 252–262.  
698

Solomon Kullback and Richard A Leibler. 1951. On  
699 information and sufficiency. *The annals of mathe-*  
700 *matical statistics*, 22(1):79–86.  
701

Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. [The power of scale for parameter-efficient prompt tuning](#).  
702 In *Proceedings of the 2021 Conference on*  
703 *Empirical Methods in Natural Language Processing*,  
704 pages 3045–3059, Online and Punta Cana, Domini-  
705 can Republic. Association for Computational Lin-  
706 guistics.  
707  
708

Hector J Levesque, Ernest Davis, and Leora Morgen-  
709 stern. 2011. The Winograd schema challenge. In  
710

711	<i>AAAI Spring Symposium: Logical Formalizations of Commonsense Reasoning</i> , volume 46, page 47.	Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2020. <a href="#">Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter.</a>	767
712			768
713	Xiang Lisa Li and Percy Liang. 2021. <a href="#">Prefix-tuning: Optimizing continuous prompts for generation.</a>	Taylor Shin, Yasaman Razeghi, Robert L. Logan IV, Eric Wallace, and Sameer Singh. 2020. <a href="#">Autoprompt: Eliciting knowledge from language models with automatically generated prompts.</a> <i>CoRR</i> , abs/2010.15980.	770
714			771
715	Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Tam, Zhengxiao Du, Zhilin Yang, and Jie Tang. 2022. <a href="#">P-tuning: Prompt tuning can be comparable to fine-tuning across scales and tasks.</a> 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.	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.	772
716			773
717			774
718			775
719			776
720			777
721			778
722			779
723	Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. 2023. <a href="#">Gpt understands, too.</a> <i>AI Open</i> .	Erik F. Tjong Kim Sang and Fien De Meulder. 2003. <a href="#">Introduction to the CoNLL-2003 shared task: Language-independent named entity recognition.</a> In <i>Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003</i> , pages 142–147.	781
724			782
725			783
726	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. <a href="#">Roberta: A robustly optimized bert pretraining approach.</a>	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 Esioibu, 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. <a href="#">Llama 2: Open foundation and fine-tuned chat models.</a>	784
727			785
728			786
729			787
730			788
731	Mohammad Taher Pilehvar and Jose Camacho-Collados. 2019. <a href="#">WiC: The word-in-context dataset for evaluating context-sensitive meaning representations.</a> In <i>Proceedings of NAACL-HLT</i> .		789
732			790
733			791
734			792
735	Sameer Pradhan, Alessandro Moschitti, Nianwen Xue, Olga Uryupina, and Yuchen Zhang. 2012a. <a href="#">CoNLL-2012 shared task: Modeling multilingual unrestricted coreference in ontonotes.</a> In <i>Joint conference on EMNLP and CoNLL-shared task</i> , pages 1–40.		793
736			794
737			795
738			796
739			797
740	Sameer Pradhan, Alessandro Moschitti, Nianwen Xue, Olga Uryupina, and Yuchen Zhang. 2012b. <a href="#">CoNLL-2012 shared task: Modeling multilingual unrestricted coreference in OntoNotes.</a> In <i>Joint Conference on EMNLP and CoNLL - Shared Task</i> , pages 1–40, Jeju Island, Korea. Association for Computational Linguistics.		798
741			799
742			800
743			801
744			802
745			803
746			804
747	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.		805
748			806
749			807
750			808
751			809
752			810
753			811
754	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. <a href="#">SQuAD: 100,000+ questions for machine comprehension of text.</a> In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 2383–2392, Austin, Texas. Association for Computational Linguistics.		812
755			813
756			814
757			815
758			816
759			817
760	Melissa Roemmele, Cosmin Adrian Bejan, and Andrew S. Gordon. 2011. <a href="#">Choice of plausible alternatives: An evaluation of commonsense causal reasoning.</a> In <i>2011 AAAI Spring Symposium Series</i> .		818
761			819
762			820
763			821
764	Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. 2015. <a href="#">Fitnets: Hints for thin deep nets.</a>	Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019. <a href="#">Superglue: A stickier benchmark for general-purpose language understanding systems.</a> <i>CoRR</i> , abs/1905.00537.	822
765			823
766			824
			825

826 Alex Wang, Yada Pruksachatkun, Nikita Nangia, Aman-  
827 preet Singh, Julian Michael, Felix Hill, Omer Levy,  
828 and Samuel R. Bowman. 2020. [Superglue: A stickier](#)  
829 [benchmark for general-purpose language understand-](#)  
830 [ing systems.](#)

831 Ralph Weischedel, Martha Palmer, Mitchell Marcus,  
832 Eduard Hovy, Sameer Pradhan, Lance Ramshaw, Ni-  
833 anwen Xue, Ann Taylor, Jeff Kaufman, Michelle  
834 Franchini, et al. 2013. Ontonotes release 5.0  
835 ldc2013t19. *Linguistic Data Consortium, Philadel-*  
836 *phia, PA*, 23:170.

837 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien  
838 Chaumond, Clement Delangue, Anthony Moi, Pier-  
839 ric Cistac, Tim Rault, Remi Louf, Morgan Funtow-  
840 icz, Joe Davison, Sam Shleifer, Patrick von Platen,  
841 Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu,  
842 Teven Le Scao, Sylvain Gugger, Mariama Drame,  
843 Quentin Lhoest, and Alexander Rush. 2020. [Trans-](#)  
844 [formers: State-of-the-art natural language processing.](#)  
845 In *Proceedings of the 2020 Conference on Empirical*  
846 *Methods in Natural Language Processing: System*  
847 *Demonstrations*, pages 38–45, Online. Association  
848 for Computational Linguistics.

849 Sheng Zhang, Xiaodong Liu, Jingjing Liu, Jianfeng  
850 Gao, Kevin Duh, and Benjamin Van Durme. 2018.  
851 [Record: Bridging the gap between human and ma-](#)  
852 [chine commonsense reading comprehension.](#)