

# Context-Aware Prompt: Customize A Unique Prompt For Each Input

Anonymous ACL submission

## Abstract

After the proposal of BERT, pre-trained language models have become the dominant approach for solving many NLP tasks. Typically, a linear classifier is added to the head of the model for fine-tuning to fit downstream tasks, while a more recent approach, also known as prompt-based learning or prompt-learning, using prompts to perform various downstream tasks, is considered to be able to uncover the potential of the language model.

Prior study, however, attempted to find a universal prompt for a certain task across all samples. Therefore, we propose a novel method, Context-Aware Prompt (CAP), which provides a unique continuous prompt for each sample input by combining contextual information to further investigate the potential capabilities of the language models. On the SuperGlue benchmark, our method outperforms multiple models with vanilla fine-tuning. Furthermore, we extend the use of prompts to include Replaced Token Detection (RTD) type prompts, allowing models like ELECTRA and DeBERTaV3 that employ RTD as a training objective to use prompts for downstream tasks.<sup>1</sup>

## 1 Introduction

Due to their outstanding performance in downstream tasks such as question answering (Rajpurkar et al., 2016), named entity recognition (Sang and Meulder, 2003), and text classification (Sun et al., 2019), pre-trained language models (Devlin et al., 2018; Liu et al., 2019; Raffel et al., 2020; He et al., 2021b) have gained increasing importance and become the primary way to solve various natural language processing (NLP) tasks in recent years.

How to leverage these language model effectively has been a struggle for researchers to figure

out. The typical approach is to add a linear classifier to the head of the model and then adapt the parameters of the linear classifier and the pre-trained language model to the supervised target task, also known as fine-tuning (Radford et al., 2018).

More recently, a paradigm, known as prompt-based learning or prompt learning (Liu et al., 2021a), has had great success in zero-shot and few-shot learning. In this way, a task description is provided for downstream tasks that not only more closely resembles the human way of thinking, but also aids the language model to "understand" what the task purpose is (Radford et al., 2019). These results suggest that when prompted by relevant task descriptions, pre-trained language models can solve NLP problems more effectively. As a result, we believe that prompts can be used to improve performance on a variety of NLP tasks.

However, finding a prompt that works for all samples is very difficult, as shown in Table 1 (a sentiment classification task): a prompt that works for one sample may not work for another, although the reason for the results in the table is most likely a bias caused by the pre-training process: "It is" is usually followed by "great", while "It was" is usually followed by "terrible". Accordingly, we strive to alleviate this bias in the pre-trained language model by using contextual information from the input sequence to automatically generate prompts that are suitable for independent samples.

In this paper, we introduce **Context-Aware Prompt (CAP)**, a novel method for generating a unique continuous prompt for each sample input through context-awareness and is based on the extension of P-tuning (Liu et al., 2021c). We construct continuous prompts via the prompt encoder in CAP, which uses contextual information to create different prompts for each sample input by contextualized embeddings that decouples from the pre-trained language model word embeddings. Correspondingly, we also need an answer mapping,

<sup>1</sup>Our code is available at [https://anonymous.4open.science/r/CAP\\_01](https://anonymous.4open.science/r/CAP_01). To avoid violating the anonymized review principle, an anonymous repository is employed.

Input Example	Prompt	Predict	Label
A sometimes tedious film .	It is [MASK] .	<u>terrible</u>	negative
	It was [MASK] .	great	
Among the year 's most intriguing explorations of alientation .	It is [MASK] .	terrible	positive
	It was [MASK] .	<u>great</u>	

Table 1: Examples from SST-2 (Socher et al., 2013), which are predicted by BERT-base, and verbalizers set as "great" and "terrible".

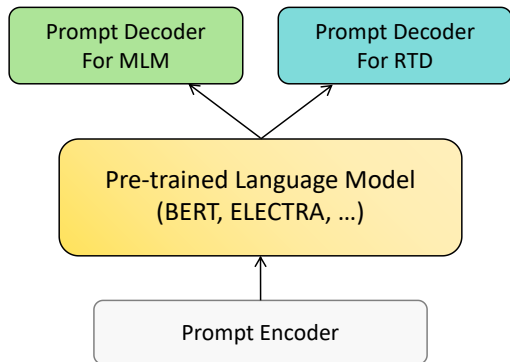


Figure 1: Illustration of CAP.

079 which we call prompt decoder, to map the output of  
080 the pre-trained language model to a specific label.  
081 Further, we expand the way of using prompts in  
082 models like ELECTRA (Clark et al., 2020), De-  
083 BERTaV3 (He et al., 2021a) — prompt decoder for  
084 Replaced Token Detection (RTD).

085 Experiments show that our proposed method  
086 achieves better performance than vanilla fine-  
087 tuning on most tasks in the SuperGLUE (Wang  
088 et al., 2019) benchmark. Moreover, our approach  
089 changes the parameters of the pre-trained language  
090 model to a lesser level than vanilla fine-tuning,  
091 which is why we believe that prompt-based learn-  
092 ing works well.

## 093 2 Related Work

094 The idea of prompt was initially proposed by Rad-  
095 ford et al. (2019) to "recall" knowledge learned  
096 during the training of pre-trained language mod-  
097 els and apply it to downstream tasks in an unsu-  
098 pervised way. Schick and Schütze (2021) employ  
099 cloze problem modeling for text classification and  
100 natural language inference. Taking a sentiment  
101 classification task as an example, constructs the  
102 input as "Best pizza ever! It was \_\_\_.", let the pre-  
103 trained language model predict the text in the blank.  
104 The sentence is regarded positive if the probabil-

105 ity of predicting the text as "great" is larger than  
106 the probability of predicting the text as "bad" and  
107 vice versa. The "Best pizza ever!" is the sample  
108 to be classified; "It was \_\_\_." and "great" or "bad"  
109 are called pattern and verbalizers by Schick and  
110 Schütze (2021), respectively.

111 However, the preceding researches rely on man-  
112 ual prompt creation, which necessitates a great  
113 deal of prior knowledge and experimental valida-  
114 tion, and the cost of determining which pattern-  
115 verbalizer pair (PVP) performs best is extremely  
116 expensive. To address this issue, some of the re-  
117 searches look at ways to automatically find the ap-  
118 propriate prompts, which can be divided into two  
119 types: discrete prompts and continuous prompts.

### 120 2.1 Discrete Prompt

121 AutoPrompt (Shin et al., 2020) employs gradient-  
122 guided search to create prompts for various tasks;  
123 LPAQA (Jiang et al., 2020) adopts a mining-based  
124 and paraphrase-based approach to prompt genera-  
125 tion; LM-BFF (Gao et al., 2021) utilises the gener-  
126 ative T5 (Raffel et al., 2020) model to automatically  
127 generate templates, all of which have achieved  
128 promising results and have contributed significantly  
129 to automatically search for prompts, but they limit  
130 prompts to discrete prompts that humans can un-  
131 derstand. However, do machines need discrete  
132 prompts that humans can understand?

### 133 2.2 Continuous Prompt

134 As a matter of fact, prompts are employed to pre-  
135 train language models to understand the task pur-  
136 pose rather than supplying them to humans, so it  
137 is not necessary to stick to discrete prompts as per-  
138 ceived by humans, but rather continuous prompts  
139 are composed of virtual words embedded in a con-  
140 tinuous space can be used.

141 **Initialized Prompt** WARP (Hambardzumyan  
142 et al., 2021), OPTIPROMPT (Zhong et al., 2021),

Prompt Tuning (Lester et al., 2021) randomly initialize the word embeddings of the pseudo-tokens or by existing real word embeddings, and optimize the word embedding vector of the pseudo-tokens by gradient descent in the word embedding space. In addition, WARP (Hambarzumyan et al., 2021) also constructs the continuous verbalizer, which is free of the constraints imposed by the human cognitive discrete word.

**Generated Prompt** After pre-training, the word embeddings of the pre-trained language model have been highly discretized. When the pseudo-token is initialized in the above ways, Allen-Zhu et al. (2019) have shown that the relevant parameters will only be changed within a small domain and are prone to fall into local minima. Therefore, P-tuning (Liu et al., 2021c) use bidirectional long-short term memory networks (BiLSTM) (Graves et al., 2013) to generate the embeddings of the prompt token, obtains satisfactory performance.

Whereas, P-tuning necessitates the addition of a few anchor tokens in order to improve performance further, which is in direct opposition to the intention of the automatic search for prompts. In contrast, our proposed CAP does not require any anchor token, and we believe that prompt tokens generated using CAP are more useful for machine comprehension, rather than relying on prior knowledge, as demonstrated in subsequent experiments (Table 4).

Prefix-tuning (Li and Liang, 2021) trains an up-stream prefix across layers and freezes the pre-trained language model, saving storage costs, but it is designed for natural language generation (NLG) tasks and autoregressive LM. Liu et al. (2021b) have extended it to NLU tasks and achieved performance comparable to fine-tuning. Similar to adapter (Houlsby et al., 2019), their motivation is parameter-efficiency, which reduces the storage cost of the language model during training. Nevertheless, their approach requires large training epochs, which raises the cost of training for downstream tasks. That is different from our starting point, where we propose that CAP involves fully training the pre-trained language model for improving performance, making downstream tasks are more relative to the training objectives of the pre-training phase, and thus can be an alternative to fine-tuning for certain tasks.

### 3 CAP

In this section, we present the implementation of Context-Aware Prompt (CAP), which consists of two components: a prompt encoder and a prompt decoder, the former generates prompt embeddings, while the latter decodes the hidden states output from the pre-trained language model into the corresponding labels of the task. Let  $\mathcal{M}$  be a pre-trained language model with vocabulary  $\mathcal{V}$  and pre-trained embedding layer  $E_{\mathcal{M}} \in \mathcal{M}$ ; let  $\mathcal{L}$  be a set of labels for our target classification task A, and assign a verbalizer token  $V_l \notin \mathcal{V}$  to each label  $l \in \mathcal{L}$ .

Overall, we need to optimize the parameters  $\Theta = \{\Theta_P, \Theta_V, \Theta_{\mathcal{M}}\}$  by gradient descent, for the prompt encoder, the verbalizer embeddings, and the pre-trained language model, respectively.

#### 3.1 Prompt Encoder

Since a direct update of the prompt embeddings would result in parameters that vary only within a small neighborhood and we need a BiLSTM to incorporate contextual information, we propose a prompt encoder to generate prompt embeddings. Given a sequence of sample input tokens  $\mathbf{x} = \{x_0, x_1, x_2, \dots, x_n\}$ , which will be mapped to input embeddings  $E_{\mathcal{M}}(\mathbf{x}) = \{e_0, e_1, e_2, \dots, e_n\}$ , here  $x_0$  refers to the token that can represent the semantics of the whole sentence, which is "[CLS]" in BERT.

We can flexibly insert prompt tokens between a given input sequence and links with label token ([LABEL]), where the label token in MLM and RTD denotes masked token ([MASK]) and verbalizer token ( $V_l$ ), respectively. For example, let  $[P_i]$  refers to the  $i^{th}$  prompt token, given the input  $\mathbf{x}$  one can compose the template  $T = \{[P_1], x_1, \dots, x_n, [P_2], [\text{LABEL}], [P_3], x_{n+1}, \dots, x_m, [P_4]\}$ . Note that the prompt tokens here are not discrete, the prompt embeddings will be inserted into the appropriate spot.

Naturally, our research turns into finding a set of parameters  $\Theta_P$  for generating continuous prompt embeddings that allow the pre-trained language model to predict the expected answer to a masked token (for MLM) or determine whether the verbalizer token is replaced (for RTD). Finally, the continuous prompt embeddings generated by the prompt encoder, which is influenced by this set of parameters and inputs, are fed into the pre-trained language model with the embeddings of the original inputs.

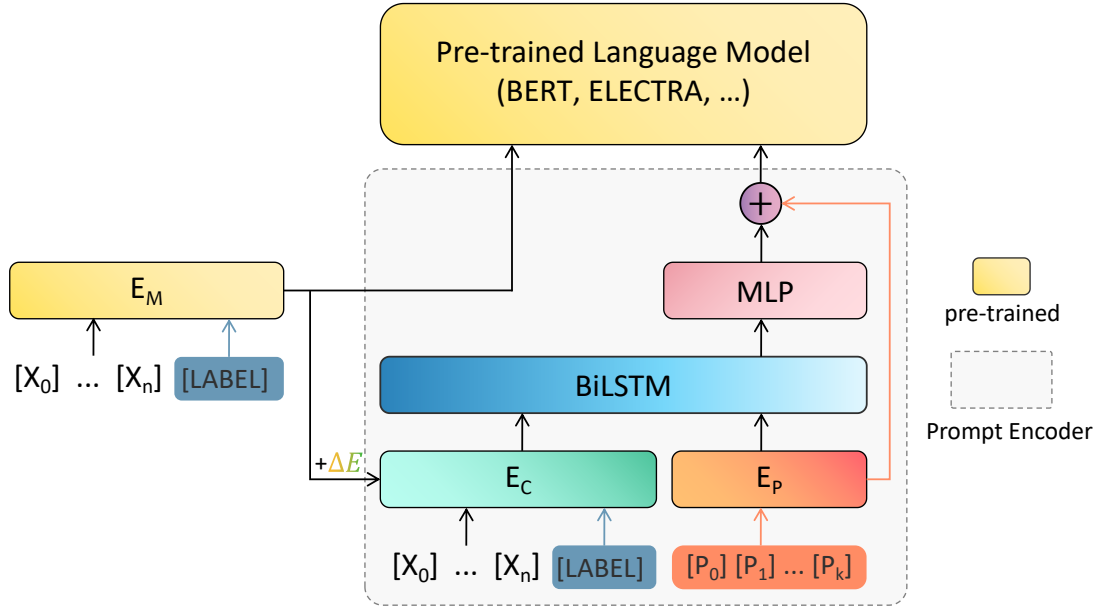


Figure 2: Prompt Encoder in CAP.

### 3.1.1 Continuous Prompt Embeddings

Customizing the prompt for each input and enabling the prompt to fully incorporate contextual information, as illustrated in Figure 2, we employ BiLSTM to attempt to pass contextual information to the prompt and comply with P-tuning (Liu et al., 2021c) using MLP to encourage discretization. Besides, we establish a shortcut connection (He et al., 2016) between the raw prompt embeddings and the generated embeddings. For one thing, it prevents the vanishing gradient to raw prompt embeddings. And for another, we employ an idea similar to ELMo (Peters et al., 2018) which combines the internal states of each layer for rich word representation, generating prompt embeddings that combine the raw prompt embeddings with the prompt embeddings that have combined contextual information.

So the continuous prompt embeddings  $\mathbf{PE}$  can be given by:

$$\begin{aligned} \mathbf{PE} &= PE_{\Theta_P}(\mathbf{x}) \\ &= \text{MLP}(\text{BiLSTM}[E_C(\mathbf{x}), E_P(\mathbf{P})]) + E_P(\mathbf{P}) \end{aligned} \quad (1)$$

where  $PE_{\Theta_P}$  means Prompt Encoder,  $E_C$  and  $E_P$  mean contextualized embedding layer and prompt embedding layer respectively, and  $\mathbf{P} = \{[P_1], [P_2], \dots, [P_k]\}$  refers to prompt sequence.

### 3.1.2 Contextualized Embeddings

Since the word embedding layer of the pre-trained language model and the contextualized embedding layer used by CAP are strongly related to each other but have radically different training objectives, they form what He et al. (2021a) call "tug-of-war" dynamics.

Further, inspired by their proposed Gradient-Disentangled Embedding Sharing (GDES) method, we adopt a similar strategy to share the word embedding layer of the pre-trained language model ( $E_M$ ) with the contextualized embedding layer ( $E_C$ ) used by CAP but decouple them with stopping gradients in the CAP's contextualized embeddings from back-propagating to the pre-trained language model's word embeddings.

In short, the contextualized embedding can be expressed as:

$$E_C = sg(E_M) + \Delta E \quad (2)$$

where  $sg$  is the stop gradient operator which only allows gradients propagation through  $\Delta E$ . Noteworthy,  $\Delta E$  is initialized as a zero matrix.

### 3.2 Prompt Decoder

After the generated prompt embeddings are sent to the pre-trained language model, a prompt decoder to parse the output of the hidden states is required to obtain the predicted labels.

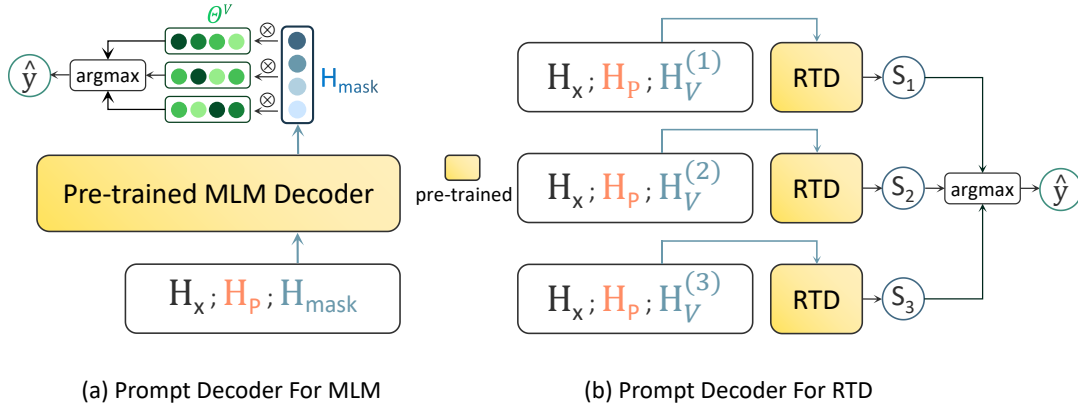


Figure 3: The method of prompt decoder for MLM and RTD solve classification tasks, where  $H_x$ ,  $H_p$ ,  $H_{\text{mask}}$ ,  $H_V$  attain the hidden states of  $x$ , the hidden states of prompt, the hidden state of masked token, and the hidden state of verbalizer token output by the pre-trained language model, respectively.

For each label, a verbalizer for answer mapping is required. It is worth noting that for MLM (Masked Language Model)<sup>2</sup>, we follow WARP (Hambardzumyan et al., 2021) to replicate continuous verbalizers.

### 3.2.1 Prompt Decoder For MLM

Masked Language Model (MLM), such as BERT (Devlin et al., 2018), ALBERT (Lan et al., 2020), RoBERTa (Liu et al., 2019), typically masks a specific percentage of words in a given sentence, and the model predicts these masked words based on other remaining words in this sentence. So the last hidden state of the masked token can be easily obtained by MLM:

$$H_{\text{mask}} = \mathcal{M}(E_{\mathcal{M}}(\mathbf{x}), PE_{\Theta_P}(\mathbf{x})) \quad (3)$$

where  $H_{\text{mask}}$  refers to the last hidden state of the masked token and has been pooled by pre-trained MLM decoder.

And the probability of labels are given by:

$$P_{\Theta}(l | \mathbf{x}) = \frac{\exp(\Theta_l^V H_{\text{mask}})}{\sum_{l' \in \mathcal{L}} \exp(\Theta_{l'}^V H_{\text{mask}})}, l \in \mathcal{L} \quad (4)$$

where  $\Theta_l^V$  is the parameters of the verbalizer embedding corresponding to the label  $l$ .

### 3.2.2 Prompt Decoder For RTD

ELECTRA (Clark et al., 2020) and DeBERTaV3 (He et al., 2021a), applying Replaced Token Detection (RTD) as training objective, use a generator to

<sup>2</sup>It can be easily extended to Permuted Language Model (PLM) such as XLNet, which is not described in this paper.

generate ambiguous tokens and a discriminator to distinguish the ambiguous tokens from the original inputs, similar to Generative Adversarial Networks (GAN, Goodfellow et al. 2014).

For RTD, analogously to Candidates-Contrast proposed by Sun et al. (2021), we consider different verbalizer tokens as input label tokens, use RTD to detect verbalizer tokens, derive the score of each verbalizer token is not replaced, and consider the label corresponding to the verbalizer token with the highest score as the predicted label, which is as shown in Figure 3.

So the score can be considered as the negative of the score of the model output that detects as a replaced token:

$$S(l|\mathbf{x}) = -\mathcal{M}(E_{\mathcal{M}}(\mathbf{x}), PE_{\Theta_P}(\mathbf{x}), E_V(l)) \quad (5)$$

where  $E_V(l)$  denotes the verbalizer embeddings representing the label  $l$ .

Resemble equation 4, the probability of labels are given by:

$$P_{\Theta}(l|\mathbf{x}) = \frac{\exp(S(l|\mathbf{x}))}{\sum_{l' \in \mathcal{L}} \exp(S(l'|\mathbf{x}))}, l \in \mathcal{L} \quad (6)$$

This multiple choice approach, on the other hand, demands numerous calculations of the pre-trained language model (once for each label), implying that the training cost will be several times higher than MLM (depending on the total count of labels).

## 4 Experiments

We select six natural language understanding (NLU) tasks from the SuperGLUE (Wang et al.,

Method	CB		RTE	BoolQ	WiC	WSC	MultiRC		Avg.
	Acc.	F1	Acc.	Acc.	Acc.	Acc.	EM	F1a	
BERT <sub>base</sub>									
Vanilla Fine-tuning	89.3	89.0	70.0	75.2	71.6	63.5	17.9	66.0	68.6
P-tuning(Liu et al., 2021c)	89.2	92.1	71.1	73.9	68.8	63.5	14.8	63.5	67.9
CAP	100	100	73.6	76.2	71.9	63.5	19.2	67.8	71.5
RoBERTa <sub>base</sub>									
Vanilla Fine-tuning	96.4	97.4	76.2	78.5	69.3	63.5	33.9	75.2	73.2
CAP	100	100	84.1	80.0	69.6	63.5	34.3	75.1	75.3
DeBERTaV3 <sub>base</sub>									
Vanilla Fine-tuning	94.6	93.7	84.8	83.3	74.1	63.5	49.4	82.1	77.6
CAP	100	100	86.3	86.3	70.7	63.5	50.3	82.5	78.9

Table 2: Dev set results on SuperGLUE tasks.

2019) benchmark for experiments to evaluate our method. CB (De Marneffe et al., 2019), RTE (Dagan et al., 2005), BoolQ (Clark et al., 2019), WiC (Pilehvar and Camacho-Collados, 2018), WSC (Levesque et al., 2012) and MultiRC (Khashabi et al., 2018) are among the tasks, which contain two textual entailment tasks (CB, RTE), two question answering tasks (BoolQ, MultiRC), a co-reference resolution task (WiC), and a word sense disambiguation task (WSC). These NLU tasks are reformulated as cloze problems, with prompt tokens inserted in the intervals between sentence1, label, and sentence2 (if it exists).

#### 4.1 Experiment Settings

In our experiments, BERT and RoBERTa are used as MLM representatives and DeBERTaV3 as RTD representative, and the base-scale model (Layer=12, Hidden Size=768, Attention Head=12) is used uniformly. Pre-trained language models we use are from the Hugging Face Transformers (Wolf et al., 2020) library, and we employ the `AutoModelForSequenceClassification` it provides for fine-tuning as a baseline. To build CAP, the prompt embeddings for  $[P_1], [P_2], \dots, [P_n]$  are generated by prompt encoder, and verbalizer embeddings for  $[V_1], \dots, [V_L]$  are initialized with verbalizer token embeddings in pre-trained language model for their corresponding labels. For MLM, the bias of the verbalizer classifier is also initialized with the bias of the MLM classifier head of the pre-trained language model in order to be consistent with the pre-training phase.

We choose the AdamW optimizer with a learning

rate that decreases linearly after a warmup period and train for 3-6 epochs on various task. Specifically, we set the learning rate from  $2e-5, 3e-5, 4e-5, 5e-5$  and batch size from 6, 8, 16, 24, 32.<sup>3</sup>

#### 4.2 Results

The results are presented in Table 2. We can see that CAP outperforms vanilla fine-tuning on most tasks. For the average results, CAP outperforms vanilla fine-tuning by 2.9, 2.1, and 1.3 percent for BERT, RoBERTa, and DeBERTaV3, respectively. Besides, CAP also outperforms P-tuning. In particular, CAP outperforms vanilla fine-tuning more significantly in tasks such as CB, RTE and BoolQ, which are easy to construct prompts for. We then focus our subsequent experiments on these three tasks.

It can be seen that our method achieves relatively well performance in both the small sample size task (CB) and the multiple sample task (MultiRC). But for WSC, a word sense disambiguation task, our approach does not show any improvement over vanilla fine-tuning, and we conjecture that it is difficult to construct prompts that allow pre-trained language models to understand the purpose of such difficult tasks. The same is true for WiC, where there is only a small improvement when using BERT and RoBERTa, and even a drop in performance when using DeBERTaV3, so we get a similar conclusion to Liu et al. (2021c), that these types of more complex tasks are not suitable for

<sup>3</sup>To reproduce our work better, we pushed the hyperparameters we used in our experiments to our open-source code repository.

Method	CB		RTE	BoolQ
	Acc.	F1	Acc.	Acc.
BERT + FT	94.6	93.7	70.4	77.7
BERT + CAP	98.2	98.7	75.8	78.1
RoBERTa + FT	96.4	90.4	85.6	85.4
RoBERTa + CAP	100	100	87.4	86.1
DeBERTaV3 + FT	94.7	89.0	88.4	87.5
DeBERTaV3 + CAP	100	100	88.8	88.2

Table 3: Best results on SuperGLUE tasks based on large-scaled model, where BERT + FT results are from (Wang et al., 2019), FT refers to vanilla fine-tuning.

Method	CB		RTE	BoolQ
	Acc.	F1	Acc.	Acc.
Vanilla Fine-tuning	89.3	89.0	70.0	75.2
WARP <sub>init</sub>	91.1	89.2	70.8	76.6
P-tuning + CV	92.9	94.7	69.0	75.6
CAP	<u>100</u>	<u>100</u>	<u>73.6</u>	76.2
+ anchor token	92.9	92.3	68.2	75.0
- CA (i.e. NUP)	92.9	90.5	69.0	75.2
- GDES (i.e. NES)	98.2	96.4	71.1	76.0

Table 4: Dev set results (BERT-base) on SuperGLUE tasks, where WARP<sub>init</sub> differs from the original paper and trains the entire model, CV refers to continuous verbalizer. - CA denotes CAP without contextualized awareness, i.e. no unique prompt (NUP). GDES means Gradient-Disentangled Embedding Sharing, while NES is No Embedding Sharing.

current prompt-based learning.

In addition, we also do experiments based on large-scale models (Layer=24, Hidden Size=1024, Attention Head=16). Analogous results based on the large-scale models are shown in Table 3.

Furthermore, we compare existing methods based on prompts. To be fair, BERT is applied as the base pre-training model, and the same pattern and verbalizers are used. As shown in Table 4, our method has certain advantages over existing methods.

### 4.3 Ablation Experiments

To further analyze CAP, a series of ablation experiments are carried out to investigate which part plays a crucial role. Specifically, we evaluate CAP with or without anchor tokens, contextualized awareness and the GDES method.

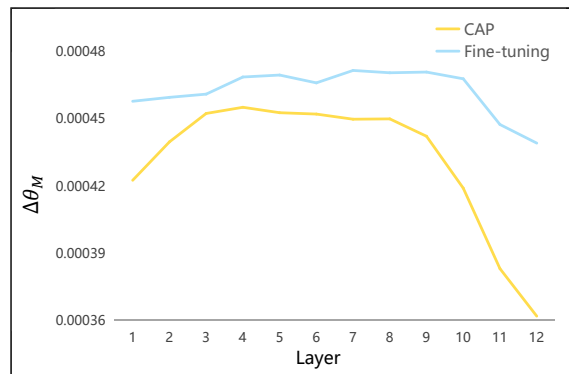


Figure 4: Parameter changes of pre-trained language model (BERT-base) for different layers after fitting downstream task (RTE).

Several conclusions can be drawn from the experiment results in Table 4. First, as predicted by our earlier supposition, the automatic learning of the prompt token using CAP surpasses the addition of the anchor token. Second, prompt combines contextual information and decouples the CAP’s contextualized embeddings from the original pre-trained language model word embeddings gradient is effective. Third, creating unique prompts with contextualized awareness plays a key role, with the GDES (Gradient-Disentangled Embedding Sharing) method serving as the icing on the cake.

## 5 Discussion

The experiment results show that CAP outperforms vanilla fine-tuning in the experimental NLU task. Vanilla fine-tuning is an approach similar to or rather a transfer learning, the inputs are passed through a pre-trained language model (which can be thought of as feature extraction), are then fed into an added linear output layer to predict.

The prompt-based learning approach, on the other hand, employs a pre-trained language model that is more closely aligned with the pre-training process’s training objectives. As shown in Figure 4, fitting the same downstream task with CAP changes the parameters less than vanilla fine-tuning. It’s not unexpected that CAP achieves better performance than vanilla fine-tuning because we believe the less the pre-trained language model changes, the more it keeps its learnt knowledge. Besides, just as adding a task description such as "Please classify the sentiment of the following sentences" provides a clear grasp of the job’s objective, while it is not evident what one needs to do when given

463	a pile of text and must guess at the assignment’s		
464	intention.		
465	Although prompt-based learning is hard to adapt		
466	to all tasks at the moment, it makes sense to explore		
467	in the prompt-based learning direction when using		
468	pre-trained language models for downstream tasks		
469	with "not only fine-tuning" in mind.		
470	<b>6 Conclusion</b>		
471	In this paper, we propose a new method, Context-		
472	Aware Prompt (CAP), as an alternative to fine-		
473	tuning using pre-trained language models. Specifi-		
474	cally, CAP constructs a unique continuous prompt		
475	for each diverse input by combining contextual		
476	information. Experiment results show that our		
477	method can make better and full use of pre-trained		
478	language models, thus outperforms vanilla fine-		
479	tuning and existing methods for most tasks on the		
480	SuperGLUE benchmark. In addition, we further		
481	extend to the RTD-based prompts usage scheme,		
482	making it possible to use prompt-based learning as		
483	a method for the majority of existing pre-trained		
484	language models.		
485	<b>References</b>		
486	Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. 2019.		
487	A convergence theory for deep learning via over-		
488	parameterization. In <i>International Conference on</i>		
489	<i>Machine Learning</i> , pages 242–252. PMLR.		
490	Christopher Clark, Kenton Lee, Ming-Wei Chang,		
491	Tom Kwiatkowski, Michael Collins, and Kristina		
492	Toutanova. 2019. Boolq: Exploring the surprising		
493	difficulty of natural yes/no questions. In <i>North Amer-</i>		
494	<i>ican Chapter of the Association for Computational</i>		
495	<i>Linguistics</i> .		
496	Kevin Clark, Minh-Thang Luong, Quoc V. Le, and		
497	Christopher D. Manning. 2020. Electra: Pre-training		
498	text encoders as discriminators rather than gener-		
499	ators. In <i>International Conference on Learning Representa-</i>		
500	<i>tions</i> .		
501	Ido Dagan, Oren Glickman, and Bernardo Magnini.		
502	2005. The pascal recognising textual entailment chal-		
503	lenge. In <i>Machine Learning Challenges Workshop</i> ,		
504	pages 177–190. Springer.		
505	Marie-Catherine De Marneffe, Mandy Simons, and Ju-		
506	dith Tonhauser. 2019. The commitmentbank: Invest-		
507	igating projection in naturally occurring discourse.		
508	In <i>proceedings of Sinn und Bedeutung</i> , volume 23,		
509	pages 107–124.		
510	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and		
511	Kristina Toutanova. 2018. Bert: Pre-training of deep		
	bidirectional transformers for language understand-		512
	ing. In <i>North American Chapter of the Association</i>		513
	<i>for Computational Linguistics</i> .		514
	Tianyu Gao, Adam Fisch, and Danqi Chen. 2021.		515
	Making pre-trained language models better few-shot		516
	learners. In <i>Meeting of the Association for Computa-</i>		517
	<i>tional Linguistics</i> .		518
	Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza,		519
	Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron		520
	Courville, and Yoshua Bengio. 2014. Generative		521
	adversarial nets. <i>Advances in neural information</i>		522
	<i>processing systems</i> , 27.		523
	Alex Graves, Abdel-rahman Mohamed, and Geoffrey		524
	Hinton. 2013. Speech recognition with deep recur-		525
	rent neural networks. In <i>2013 IEEE international</i>		526
	<i>conference on acoustics, speech and signal process-</i>		527
	<i>ing</i> , pages 6645–6649. Ieee.		528
	Karen Hambardzumyan, Hrant Khachatrian, and		529
	Jonathan May. 2021. Warp: Word-level adversar-		530
	ial reprogramming. In <i>Meeting of the Association for</i>		531
	<i>Computational Linguistics</i> .		532
	Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian		533
	Sun. 2016. Deep residual learning for image recog-		534
	nition. In <i>Proceedings of the IEEE conference on</i>		535
	<i>computer vision and pattern recognition</i> , pages 770–		536
	778.		537
	Pengcheng He, Jianfeng Gao, and Weizhu Chen. 2021a.		538
	Debertav3: Improving deberta using electra-style pre-		539
	training with gradient-disentangled embedding shar-		540
	ing. <i>arXiv preprint arXiv:2111.09543</i> .		541
	Pengcheng He, Xiaodong Liu, Jianfeng Gao, and		542
	Weizhu Chen. 2021b. Deberta: Decoding-enhanced		543
	bert with disentangled attention. In <i>International</i>		544
	<i>Conference on Learning Representations</i> .		545
	Neil Houlsby, Andrei Giurgiu, Stanisław Jastrzębski,		546
	Bruna Halila Morrone, Quentin de Laroussilhe, An-		547
	dreia Gesmundo, Mona Attariyan, and Sylvain Gelly.		548
	2019. Parameter-efficient transfer learning for nlp.		549
	In <i>International Conference on Machine Learning</i> .		550
	Zhengbao Jiang, Frank F Xu, Jun Araki, and Graham		551
	Neubig. 2020. How can we know what language		552
	models know? <i>Transactions of the Association for</i>		553
	<i>Computational Linguistics</i> , 8:423–438.		554
	Daniel Khashabi, Snigdha Chaturvedi, Michael Roth,		555
	Shyam Upadhyay, and Dan Roth. 2018. Looking		556
	beyond the surface: A challenge set for reading com-		557
	prehension over multiple sentences. In <i>Proceedings</i>		558
	<i>of the 2018 Conference of the North American Chap-</i>		559
	<i>ter of the Association for Computational Linguistics:</i>		560
	<i>Human Language Technologies, Volume 1 (Long Pa-</i>		561
	<i>pers)</i> , pages 252–262.		562
	Zhenzhong Lan, Mingda Chen, Sebastian Goodman,		563
	Kevin Gimpel, Piyush Sharma, and Radu Soricut.		564
	2020. Albert: A lite bert for self-supervised learning		565
	of language representations. In <i>International Confer-</i>		566
	<i>ence on Learning Representations</i> .		567



568	Brian Lester, Rami Al-Rfou, and Noah Constant. 2021.	Erik F. Tjong Kim Sang and Fien De Meulder. 2003.	623
569	The power of scale for parameter-efficient prompt	Introduction to the conll-2003 shared task: language-	624
570	tuning. <i>arXiv preprint arXiv:2104.08691</i> .	independent named entity recognition. In <i>North</i>	625
571	Hector Levesque, Ernest Davis, and Leora Morgenstern.	<i>American Chapter of the Association for Computa-</i>	626
572	2012. The winograd schema challenge. In <i>Thir-</i>	<i>tional Linguistics</i> .	627
573	<i>teenth International Conference on the Principles of</i>	Timo Schick and Hinrich Schütze. 2021. Exploiting	628
574	<i>Knowledge Representation and Reasoning</i> .	cloze-questions for few-shot text classification and	629
575	Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning:	natural language inference. In <i>Conference of the Eu-</i>	630
576	Optimizing continuous prompts for generation. In	<i>ropean Chapter of the Association for Computational</i>	631
577	<i>Meeting of the Association for Computational Lin-</i>	<i>Linguistics</i> .	632
578	<i>guistics</i> .	Taylor Shin, Yasaman Razeghi, Robert L. Logan, Eric	633
579	Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang,	Wallace, and Sameer Singh. 2020. Autoprompt: Elic-	634
580	Hiroaki Hayashi, and Graham Neubig. 2021a. Pre-	iting knowledge from language models with automati-	635
581	train, prompt, and predict: A systematic survey of	cally generated prompts. In <i>Empirical Methods in</i>	636
582	prompting methods in natural language processing.	<i>Natural Language Processing</i> .	637
583	<i>arXiv preprint arXiv:2107.13586</i> .	Richard Socher, Alex Perelygin, Jean Wu, Jason	638
584	Xiao Liu, Kaixuan Ji, Yicheng Fu, Zhengxiao Du,	Chuang, Christopher D Manning, Andrew Y Ng, and	639
585	Zhilin Yang, and Jie Tang. 2021b. P-tuning v2:	Christopher Potts. 2013. Recursive deep models for	640
586	Prompt tuning can be comparable to fine-tuning	semantic compositionality over a sentiment treebank.	641
587	universally across scales and tasks. <i>arXiv preprint</i>	In <i>Proceedings of the 2013 conference on empiri-</i>	642
588	<i>arXiv:2110.07602</i> .	<i>cal methods in natural language processing</i> , pages	643
589	Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding,	1631–1642.	644
590	Yujie Qian, Zhilin Yang, and Jie Tang. 2021c. Gpt	Chi Sun, Xipeng Qiu, Yige Xu, and Xuanjing Huang.	645
591	understands, too. <i>arXiv preprint arXiv:2103.10385</i> .	2019. How to fine-tune bert for text classification?	646
592	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-	In <i>China National Conference on Chinese Computa-</i>	647
593	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,	<i>tional Linguistics</i> , pages 194–206. Springer.	648
594	Luke Zettlemoyer, and Veselin Stoyanov. 2019.	Yi Sun, Yu Zheng, Chao Hao, and Hangping Qiu. 2021.	649
595	Roberta: A robustly optimized bert pretraining ap-	Nsp-bert: A prompt-based zero-shot learner through	650
596	proach. <i>arXiv preprint arXiv:1907.11692</i> .	an original pre-training task–next sentence prediction.	651
597	Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt	<i>arXiv preprint arXiv:2109.03564</i> .	652
598	Gardner, Christopher Clark, Kenton Lee, and Luke	Alex Wang, Yada Pruksachatkun, Nikita Nangia, Aman-	653
599	Zettlemoyer. 2018. Deep contextualized word repre-	preet Singh, Julian Michael, Felix Hill, Omer Levy,	654
600	sentations. In <i>North American Chapter of the Associ-</i>	and Samuel R. Bowman. 2019. Superglue: A stick-	655
601	<i>ation for Computational Linguistics</i> .	ier benchmark for general-purpose language under-	656
602	Mohammad Taher Pilehvar and José Camacho-Collados.	standing systems. In <i>Neural Information Processing</i>	657
603	2018. Wic: 10,000 example pairs for evalu-	<i>Systems</i> .	658
604	ating context-sensitive representations. <i>CoRR</i> ,	Thomas Wolf, Julien Chaumond, Lysandre Debut, Vic-	659
605	<i>abs/1808.09121</i> .	tor Sanh, Clement Delangue, Anthony Moi, Pier-	660
606	Alec Radford, Karthik Narasimhan, Tim Salimans, and	eric Cistac, Morgan Funtowicz, Joe Davison, Sam	661
607	Ilya Sutskever. 2018. Improving language under-	Shleifer, et al. 2020. Transformers: State-of-the-	662
608	standing by generative pre-training.	art natural language processing. In <i>Proceedings of</i>	663
609	Alec Radford, Jeffrey Wu, Rewon Child, David Luan,	<i>the 2020 Conference on Empirical Methods in Nat-</i>	664
610	Dario Amodei, Ilya Sutskever, et al. 2019. Language	<i>ural Language Processing: System Demonstrations</i> ,	665
611	models are unsupervised multitask learners. <i>OpenAI</i>	pages 38–45.	666
612	<i>blog</i> , 1(8):9.	Zexuan Zhong, Dan Friedman, and Danqi Chen. 2021.	667
613	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine	Factual probing is [mask]: Learning vs. learning to	668
614	Lee, Sharan Narang, Michael Matena, Yanqi Zhou,	recall. <i>arXiv preprint arXiv:2104.05240</i> .	669
615	Wei Li, and Peter J. Liu. 2020. Exploring the limits		
616	of transfer learning with a unified text-to-text trans-		
617	former. <i>Journal of Machine Learning Research</i> , 21:1–		
618	67.		
619	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and		
620	Percy Liang. 2016. Squad: 100,000+ questions for		
621	machine comprehension of text. In <i>Empirical Meth-</i>		
622	<i>ods in Natural Language Processing</i> .		