

# CR-UTP: Certified Robustness against Universal Text Perturbations

Anonymous ACL submission

## Abstract

It is imperative to ensure the stability of every prediction made by a language model; that is, a language’s prediction should remain consistent despite minor input variations, like word substitutions. In this paper, we investigate the problem of certifying a language model’s robustness against Universal Text Perturbations (UTPs), which have been widely used in universal adversarial attacks and backdoor attacks. Existing certified robustness based on random smoothing has shown considerable promise in certifying the input-specific text perturbations (ISTPs), operating under the assumption that any random alteration of a sample’s clean or adversarial words would negate the impact of sample-wise perturbations. However, with UTPs, masking only the adversarial words can eliminate the attack. A naive method is to simply increase the masking ratio and the likelihood of masking attack tokens, but it leads to a significant reduction in both certified accuracy and the certified radius, due to input corruption by extensive masking. To solve this challenge, we introduce a novel approach, the *superior prompt search* method, designed to identify a *superior prompt* that maintains higher certified accuracy under extensive masking. Additionally, we theoretically motivate why ensembles are a particularly suitable choice as base prompts for random smoothing. The method is denoted by *superior prompt ensembling* technique. We also empirically confirm this technique, obtaining state-of-the-art results in multiple settings. These methodologies, for the first time, enable high certified accuracy against both UTPs and ISTPs.

## 1 Introduction

Prompt-based Language Models (PLMs) (Thopplian et al., 2022; Zeng et al., 2022; Achiam et al., 2023; Touvron et al., 2023; Chiang et al., 2023) have achieved significant success across a wide range of real-world applications (Wu et al., 2020;

Brown et al., 2020; Wei et al., 2022; Chowdhery et al., 2023). However, despite their prominent performance, PLMs have been shown vulnerable to noises and perturbations on the input (Xu et al., 2022; Shayegani et al., 2023). Such vulnerability has notably restricted PLM’s utility, especially in high-stake environments such as bank records analysis (Heaton et al., 2017), health care records analysis (Myszczyńska et al., 2020). In these settings, the stability of every prediction is critical, i.e., PLM predictions should remain consistent despite minor input variations, such as word substitutions (Alzantot et al., 2018; Ren et al., 2019; Li et al., 2020a). This concern aligns with the study of certified robust PLMs (Zeng et al., 2023), which guarantees that all PLM predictions are accurate within the local vicinity of the input.

Input perturbations can be classified into Universal Text Perturbations (UTPs) and Input-Specific Text Perturbations (ISTPs). UTPs are characterized by their ability to be applied across different inputs, making them transferable, whereas ISTPs are tailored to specific inputs. In detail, attackers employing ISTP strategies, exemplified by TextFooler (Jin et al., 2020) and DeepWordBug (Gao et al., 2018), craft a unique adversarial sentence for each target input sentence. Conversely, attackers utilizing UTP methodologies, such as those found in TrojLLM (Xue et al., 2023) and UAT (Wallace et al., 2019), identify a single or a small number of adversarial tokens that can be inserted into any sentence to influence the model’s prediction. This makes UTPs a more considerable threat to the robustness of PLMs since a specific set of adversarial tokens could lead to mispredictions across any input. Additionally, UTPs pose a greater challenge in mitigation compared to ISTPs. This challenge arises because ISTPs depend on weaker adversarial patterns that can be addressed by introducing modifications to the adversarial or clean tokens. However, UTPs are based on stronger adversarial patterns,

which require exact identification and masking of the adversarial tokens for effective mitigation.

Random smoothing has been recognized as an effective defense offering certified robustness for models in computer vision (Horváth et al., 2021) and NLP (Zeng et al., 2023), yet its application has been limited to ISTPs. This method assumes that random alterations to a sample’s words counteract perturbations. However, this approach falls short against UTPs, which require precise masking of adversarial tokens for mitigation, unlike ISTPs which can be mitigated by randomly masking any tokens. Defending against UTPs is challenging due to the unknown positions of adversarial tokens, requiring a high mask ratio that could degrade PLM accuracy. Thus, ensuring certified robustness against UTPs in PLMs remains an unresolved challenge.

Naively increasing the masking ratio can improve the chances of covering adversarial tokens in UTPs, potentially reducing the Attack Success Rate (ASR). However, this method often results in only minor ASR improvements due to the trade-off with certified accuracy. High masking ratios in random smoothing significantly diminish certified accuracy, leading to randomized model inferences as a large portion of input tokens are obscured, leaving insufficient data for accurate predictions.

In this paper, we introduce CR-UTP, a method to equip PLMs with certified robustness against UTPs, achieving both high certified accuracy and low ASR. Our contributions are as follows:

- We adapt the certified robustness method to PLMs and propose *Superior Prompt Search* for robust prompts with masked inputs.
- We introduce a prompt ensemble method to reduce the variance of masked inputs and increase the certified accuracy with theoretical analysis and empirical implementation.
- Through extensive experimentation, we show our CR-UTP effectively increases the certified accuracy by  $\sim 15\%$  and decreases the ASR by  $\sim 35\%$  compared to prior works.

## 2 Background and Motivation

**Text Adversarial Attacks.** Text adversarial attack methods generate adversarial sentences by perturbing original sentences to maximally increase the model’s prediction error, while maintaining the fluency and naturalness of the adversarial examples. These attacks on prompt-based language

models can be categorized into two groups: input-specific text perturbation attacks (ISTPs) and universal adversarial perturbation attacks (UTPs). In ISTP attacks, the attacker optimizes an adversarial sentence for each input, mainly by replacing, scrambling, and erasing characters (e.g., DeepWordBug (Gao et al., 2018) and HotFlip (Ebrahimi et al., 2018)) or words (e.g., TextFooler (Jin et al., 2020)). Conversely, UTP attacks optimize a universal trigger for a prompt-based language model, and the output of any input embedded with this trigger will be manipulated (e.g., TrojLLM (Xue et al., 2023) and UAT (Wallace et al., 2019)).

**Certified Robustness in Language Models.** Numerous defense methods, such as adversarial training and perturbation-controlled approaches, have been developed to counteract adversarial attacks (Wang et al., 2019; Zhou et al., 2021; Goyal et al., 2023). However, these traditional tools may become ineffective against novel, advanced attack strategies. To address this, certified robustness has been introduced, offering a guarantee against any attack as long as the number of perturbed words remains below a certain threshold. A model achieves certification if it can consistently produce the correct output when the number of perturbations does not exceed the certified radius. While models of smaller size can obtain robustness certification through deterministic methods (Li et al., 2020b; Ostrovsky et al., 2022; Weng et al., 2018; Kolter and Wong, 2017), the computational demands of language models preclude such approaches. Consequently, probabilistic methods, such as Random Smoothing (Cohen et al., 2019), have been introduced to certify the robustness of large language models.

**Random Smoothing.** Random Smoothing, a promising approach introduced by (Cohen et al., 2019; Weber et al., 2020), certifies the robustness of large neural networks. This method enhances a model’s robustness by adding Gaussian noise to the original input (Salman et al., 2020; Li et al., 2023). It was quickly adopted for large language models in the NLP field, exemplified by SAFER (Ye et al., 2020) and Randomized [MASK] (Zeng et al., 2021). To improve model performance with random smoothing, the computer vision field has explored re-training the original model to adapt to Gaussian noise (Jeong and Shin, 2020; Zhai et al., 2020; Salman et al., 2019). However, applying this re-training method in NLP to achieve a model tolerant to smoothing is prohibitively expensive.

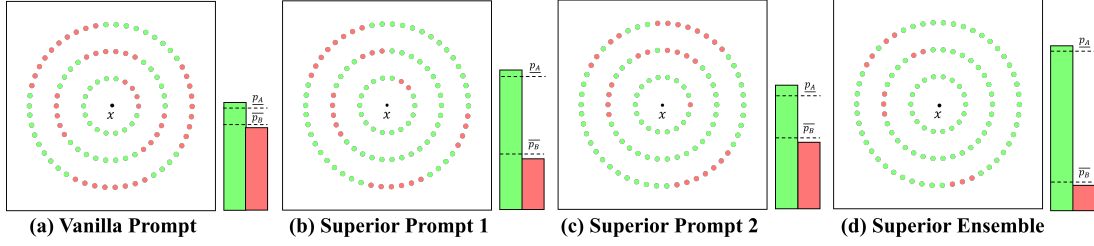


Figure 1: Illustration of the prediction distributions. A superior prompt exhibits greater robustness compared to a vanilla prompt, with ensembled prompts showing even higher robustness. Different colors represent various classes, and different radii indicate varying levels of perturbation. The bars demonstrate the output class probabilities for the smoothed PLMs given corresponding prompts.  $p_A$  represents the minimum probability of the majority class, and  $\overline{p_B}$  indicates the maximum probability of the second-most likely class.

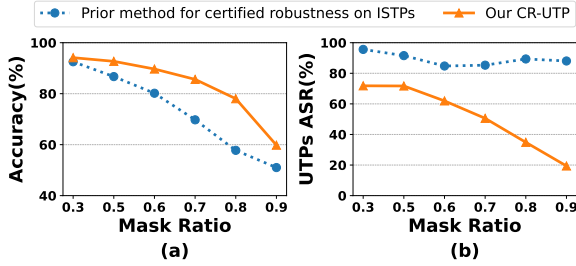


Figure 2: (a) Our CR-UTP shows higher certified robustness accuracy and (b) Our CR-UTP significantly reduces ASR.

Instead, in the NLP domain, (Zhang et al., 2023) proposed a self de-noising method that allows the Pretrained Language Model (PLM) itself to recover the information lost due to masking.

**Motivation.** Figure 2 illustrates that while random smoothing has been effective for ISTPs as shown in previous works (Zeng et al., 2023; Zhang et al., 2023), it struggles with UTPs. With a low mask ratio, the ASR for UTPs is high ( $\sim 96\%$ ), and increasing the mask ratio only marginally reduces ASR but significantly lowers accuracy. For example, increasing the mask ratio from 0.3 to 0.8 only reduces ASR by  $\sim 6\%$  while accuracy drops by  $\sim 35\%$ . These findings highlight the inadequacy of ISTP methods for UTPs and lead us to develop new techniques that combine superior prompt search and ensembles, significantly improving robustness against UTPs with less impact on accuracy.

To achieve a low Attack Success Rate (ASR) against Universal Text Perturbations (UTPs), a high mask ratio, exceeding 0.5, is necessary. Yet, as Figure 1 (a) reveals, a vanilla prompt at this high mask ratio results in reduced accuracy due to the extensive masking of input tokens, which leaves limited information for precise classification. Figures 1 (b) and (c) illustrate that superior prompts can maintain

higher accuracy under such conditions by incorporating random masking during the prompt design phase. This approach allows superior prompts to adjust to specific mask ratios, improving the lower bound ( $p_A$ ) on the majority class probability and reducing the upper bound ( $\overline{p_B}$ ) on the runner-up class probability. However, these prompts still exhibit high variance across input samples. Ensembling techniques (Liu et al., 2021; Horváth et al., 2021) reduce this variance, thereby enhancing robustness. As shown in Figure 1 (d), by combining superior prompts, we can leverage their individual advantages for a more favorable accuracy-ASR balance. This insight leads us to further investigate superior prompt design and ensembling as methods to bolster PLMs’ certified accuracy and robustness against UTPs, aiming to lower the ASR while preserving high accuracy.

### 3 CR-UTP Design

**Overview.** In Figure 3, we detail the workflow of the proposed CR-UTP method. (a) Superior Prompt Search: we start with a basic prompt and employ a reinforcement learning approach to find a superior prompt adept at handling inputs with masked words. A unique reward function is utilized, which incorporates random masking during the prompt search phase to enhance the prompt’s resilience to word masking. (b) Superior Prompt Ensemble: for making predictions, CR-UTP generates various versions of the original input by applying random masking. Each prompt assesses these versions and internally agrees on the optimal prediction. (c) Certified Robust PLM: CR-UTP aggregates the individual outcomes from each prompt through a second voting process to get the final, most robust prediction.

In particular, we generalize the random masking

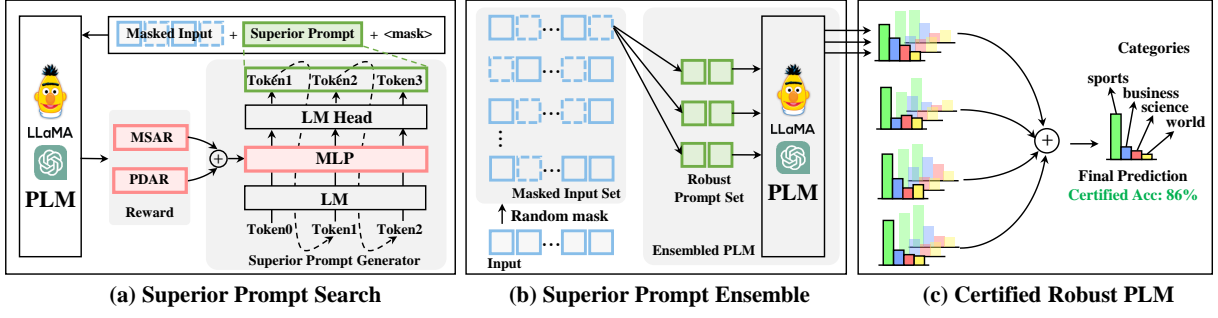


Figure 3: Overview of CR-UTP. CR-UTP leverages superior prompt search and prompt ensembling techniques to enhance the certified robustness of PLMs.

operation to PLMs and analyze using random masking to defend UTPs in Section 3.1. Also, we introduce the Superior Prompt Search in Section 3.2 to search for prompts that can achieve a high certified accuracy even when a large proportion of the input tokens are masked. Finally, in Section 3.3, we proposed a Superior Prompt Ensemble to further improve the certified robustness.

### 3.1 Adapting Random Smoothing to PLMs

Using the random masking approach from RanMASK (Zeng et al., 2023), the random masking operation  $\mathcal{M} : \mathcal{X} \times \mathcal{M}(h, k) \rightarrow \mathcal{X}_{mask}$  would take an input text  $\mathbf{x} = x_1 x_2 \dots x_h$  with  $h$  words and randomly replacing  $(h - k)$  word with the [MASK]. Following this, we define a smoothed classifier  $g(\mathbf{x})$  built upon the base classifier  $f$  as follows:

$$g(\mathbf{x}) = \arg \max_{c \in \mathcal{C}} \left[ \mathbb{P}_{\mathcal{H} \sim \mathcal{U}(h, k)} (f(\mathcal{M}(\mathbf{x} | \mathcal{H})) = c) \right] \quad (1)$$

Then it can be shown that  $g(\mathbf{x}')$  would return  $c$  when the certified condition is satisfied with probability at least  $(1 - \alpha)$ .

We define a prompt-based language model for classification tasks as  $f : \mathcal{X} \rightarrow \mathcal{C}$ , where  $\mathcal{X}$  represents the domain of input texts and  $\mathcal{C} = 1, 2, \dots, n_c$  denotes the set of classification labels. The response of an input  $\mathbf{x}$  to prompt  $p$  is  $y = f(p, \mathbf{x})$ . The smoothed classifier  $g$  over such PLM can be expressed as

$$g_p(\mathbf{x}) = \arg \max_{c \in \mathcal{C}} \left[ \mathbb{P}_{\mathcal{H} \sim \mathcal{U}(h, k)} (f(p, \mathcal{M}(\mathbf{x} | \mathcal{H})) = c) \right] \quad (2)$$

and the same certification can be achieved under such prompt-based model.

Defending against universal attacks with random masking necessitates a high mask ratio, especially in the context of UTPs, where the presence of any UTP token in the masked input guarantees the attack’s success. Therefore, to ensure the adversarial

token is masked with a probability exceeding 50% for the smoothed function  $g(\mathbf{x})$  to yield correct outcomes, the masking probability for each UTP token needs to be more than 0.5. Consequently, for a UTP of length  $r$ , the likelihood that all UTP tokens are masked should be  $p^r > 0.5$ , implying  $p > \sqrt[r]{0.5}$ . For instance,  $p > 0.5$  for  $r = 1$ , and  $p > 0.707$  for  $r = 2$ . Ensuring correct results for inputs masked without UTP requires the model to perform effectively at high mask ratios. To alleviate the effects of extensive masking, we introduce a technique that enhances model performance under random masking without necessitating retraining of the language model, thereby reducing computational overhead.

### 3.2 Superior Prompt Search

Certified accuracy is influenced by the model’s performance on randomly masked sentences, but high mask ratios can decrease accuracy due to loss of critical information. Enhancing a prompt-based language model’s certified accuracy involves improving its tolerance to information loss from random masking. However, in few-shot and black-box scenarios, fine-tuning the pre-trained model or using gradient-based optimization for continuous prompts is infeasible. To tackle this challenge without gradient optimization, we approach it as a reinforcement learning (RL) problem to discover a discrete, robust prompt—termed a superior prompt. A direct approach involves searching for this prompt using datasets with randomly masked sentences to acclimate the model to diverse masking scenarios. Nonetheless, at high mask ratios (e.g., 70%), the reduced information in few-shot datasets limits the effectiveness and generalizability of robustness enhancements. To overcome this, we suggest aligning the superior prompt’s predictions on masked sentences with the vanilla prompt’s on unmasked



sentences, leveraging the existing knowledge of the vanilla prompt to offset the drawbacks of few-shot datasets and information loss from masking.

Our aim, as expressed in Equation 3, involves identifying an optimized prompt  $p_s$  that augments a basic vanilla prompt  $p_v$  by adding a sequence of  $T$  tokens from the vocabulary  $\mathcal{V}$ . This strategy is intended to boost the model’s accuracy  $ACC(f(p_s, \hat{x}_i), y_i)$  on inputs  $\hat{x}_i$  that have undergone random word masking. The dataset  $\mathcal{D}$  is composed of pairs of input sentences  $x_i$  and their associated labels  $y_i$ , with  $\hat{x}_i$  indicating the masked version of the input.

$$\max_{p_s \in \mathcal{V}^T} \sum_{(x_i, y_i) \in \mathcal{D}} ACC(f(p_s, \hat{x}_i), y_i) \quad (3)$$

**Masked Sentence Accuracy Reward.** We introduce a two-fold reward function to guide the RL-based search for an optimal  $p_s$ . The first component, known as the Masked Sentence Accuracy Reward (MSAR), is designed to directly maximize the PLM’s accuracy on masked sentences:

$$\mathcal{R}_{\text{MSAR}} = \sum_{(x_i, y_i) \in \mathcal{D}} \eta_1^{1-\text{sign}} \eta_2^{\text{sign}} \text{Distance}(\hat{x}_i, y_i) \quad (4)$$

here the  $\text{Distance}(\hat{x}_i, y_i)$  denotes  $f(y_i | p_s, \hat{x}_i) - \max_{y' \neq y} f(y' | p_s, \hat{x}_i)$ . The distance value is positive for correct predictions and negative otherwise. We denote the distance sign as  $\text{sign} = \mathbb{1}[\text{Distance}(\hat{x}_i, y_i) > 0]$ . For a correct prediction (i.e.,  $\text{sign} = 1$ ), we multiply the positive reward by a large number  $\eta_2$  to indicate its desirability; otherwise, we multiply the negative rewards by another number  $\eta_1$ . This reward aims to maximize the PLM’s accuracy on masked sentences.

**Predictive Distribution Alignment Reward.** To mitigate the challenge posed by information loss due to word masking, which is exacerbated in a few-shot setting, we propose an additional reward function, a.k.a, Predictive Distribution Alignment Reward (PDAR). It is designed to minimize the KL divergence between the predictive distributions of the vanilla prompt on unmasked sentences and the superior prompt on their masked equivalents:

$$\mathcal{R}_{\text{PDAR}} = - \sum_{(x_i, y_i) \in \mathcal{D}} KL(f(p_v, x_i) \parallel f(p_s, \hat{x}_i)) \quad (5)$$

This reward is designed to ensure that  $p_s$  retains alignment with  $p_v$ ’s predictive behavior, thereby leveraging the foundational knowledge encoded in

$p_v$  to inform predictions in the face of partial information. Such strategic alignment enables  $p_s$  to infer missing data from the masked inputs, drawing on the robust insights and patterns encapsulated by  $p_v$ . This method not only addresses the direct impact of masking on information availability but also enhances the model’s capacity for generalization from limited examples.

**Policy Model Update.** As Figure 3 shows, the RL search process involves an agent that sequentially selects tokens  $[s_1, \dots, s_T]$  to construct the superior prompt  $p_s$ , aiming to maximize the combined reward  $\mathcal{R} = \mathcal{R}_{\text{MSAR}} + \mathcal{R}_{\text{PDAR}}$ . For each time step  $t$ , the agent, given the previous tokens  $s_{<t}$ , generates the next token  $s_t$  based on the policy generator  $G_{\theta_s}(s_t | s_{<t})$ . Completion of  $p_s$  triggers the computation of the task reward  $\mathcal{R}$ . To facilitate this, we employ a GPT-2 model as the backbone for our policy generator, enhanced with an insertable trainable Multilayer Perceptron (MLP) layer. The optimization focuses on the parameters of this MLP layer, tailoring the policy generator to effectively navigate the prompt construction space under the guidance of the designed reward function.

### 3.3 Superior Prompt Ensemble

Instead of relying on a single model, ensemble methods leverage the strengths and mitigate the weaknesses of various base classifiers. The core idea behind ensemble modelling is that a group of weak learners can come together to form a strong learner, thereby improving the model’s ability to generalize well to unseen data. In the following section, we will demonstrate how random masking can increase the variance of the output probability and how model ensembling can mitigate the performance drop introduced by random masking.

Suppose the probability that the smoothed model with prompt  $p$  outputs class  $c$  is  $y_c(\mathbf{x}) = \mathbb{P}_{\mathcal{H} \sim \mathcal{U}(h, k)}(f(p, \mathcal{M}(\mathbf{x} | \mathcal{H})) = c)$  in  $[0, 1]$ , then the final probability output of the smoothed model is determined by two random variables,  $y = y_o + y_m$  (Horváth et al., 2021), where  $y_o$  is determined by the performance of the language model  $f$  with prompt  $p$  on the original input, and  $y_m$  corresponds to the performance of prompt  $p$  when the input  $\mathbf{x}$  is randomly masked. Since  $y_o$  is influenced by the performance of prompt  $p$  and  $y_m$  mostly depends on the perturbation ratio  $\rho$  on the clean input  $\mathbf{x}$ , we can assume that  $y_o$  and  $y_m$  are independent. Suppose  $y_o$  has mean  $a_c$  and covariance  $\Sigma_c$ ,  $y_m$  has mean  $a_m$  and covariance  $\Sigma_m$ . Under this assumption

tion, the probability of the smoothing operation for each prompt  $p$  has expectation and variance as follows:

$$\mathbb{E}[y^k] = a_c^k + a_m^k, \quad \text{Var}[y^k] = \text{Var}[y_o^k] + \text{Var}[y_m^k].$$

Since the mask operation can cause information loss in  $\mathbf{x}$ , we assume  $p_m$  has a negative mean, denoted by  $a_m$ . Through empirical analysis of the performance of  $\sum_m$  with respect to the perturbation rate, as shown in figure 4, we observe the baseline model’s performance and conclude that while  $\sum_m$  varies with the perturbation ratio, it is significantly influenced by it. As the perturbation ratio increases, the variance  $\sigma_m^2$  initially rises and then decreases. This pattern occurs because, as the mask ratio increases from 0 to 0.6, the masking operation introduces more noise to the input, increasing the variance of  $y_c$ . However, when the mask ratio increases from 0.6 to 1, the remaining information decreases, leading the model to randomly guess any label, i.e.,  $y_c \rightarrow 1/n_c$ , and the variance  $\sigma^2 \rightarrow 0$  as more words in the sentence are masked.

With a high mask ratio, the random masking operation can significantly increase variance and reduce accuracy relative to the clean classifier’s output. Model ensembles have been shown to effectively decrease the overall variance of voting outcomes, thereby improving the likelihood of accurate predictions as the number of ensembles increases (Horváth et al., 2021). However, the high computational cost makes training multiple language models unfeasible. Consequently, we introduce a technique that ensembles a set of prompts during inference to exploit the distinctive feature of PLMs, wherein the initial prompt markedly affects the model’s final output. By employing the superior prompt search method, we can create a collection of prompts that withstand the random masking operation, with the ensemble of these prompts’ outputs demonstrating enhanced performance on heavily masked inputs.

Formally, we construct an ensemble classifier  $\bar{f}$  with a set of  $k$  different prompts  $\mathbf{P} = \{p_i \mid i = 1, \dots, k\}$ , via hard voting of all the outputs from different prompts  $p_i$ ,

$$\bar{f}(\mathbf{P}, \mathbf{x}) = \arg \max_{c \in \mathcal{C}} \sum_{i=1}^k \mathbb{I}(f(p_i, \mathbf{x}) = c) \quad (6)$$

where  $\mathbb{I}(f(p_i, \mathbf{x}) = c)$  is the indicator function that equals 1 when  $f(p_i, \mathbf{x})$  output  $c$  and 0 otherwise.

So the ensemble classifier would output the class that most of the prompts agree on.

Since the prompts ensemble operates as a single model, the certified robustness condition remains applicable to the assembled model  $\bar{f}(\mathbf{P}, \mathbf{x})$ . Therefore, we can establish a new smoothing function  $\bar{g}(\mathbf{x})$  by applying the same random masking operation to  $\mathbf{x}$ . Building on our previous findings, we anticipate performance improvements through the prompts ensemble. Our analysis of how the number of ensembles impacts the final probability outcome, as depicted in 5, demonstrates a substantial increase in the accuracy of the model ensemble with a concurrent reduction in variance as the number of ensembles increases.

## 4 Experimental Methodology

**Datasets and Model.** In our evaluation, we utilize the SST-2 dataset (Socher et al., 2013) for binary classification tasks and the AgNews dataset (Zhang et al., 2015) for a four-class classification task. We adopt a 16-shot setting, which represents a typical few-shot scenario. Our experiments are conducted with the widely-used pre-trained language model RoBERTa-large (Liu et al., 2019), an advanced version of BERT (Kenton and Toutanova, 2019) with 24 layers of Transformer architecture.

**Evaluation Metrics.** We adopted three key metrics in evaluations same with (Zeng et al., 2023). Clean accuracy (CACC) refers to the classification accuracy on clean sentences. The attack success rate (ASR) quantifies the percentage of input instances perturbed by an attack that successfully causes the model to make incorrect predictions. The poisoned accuracy (PACC) indicates the accuracy of the prompt-based language model on poisoned samples crafted from an attack.

**Evaluated Attacks.** We evaluated our CR-UTP under two input-specific text perturbation (ISTP) adversarial attacks, TextFooler (Jin et al., 2020) and DeepWordBug (Gao et al., 2018), and one universal text perturbation (UTP) attack, TrojLLM (Xue et al., 2023). TextFooler adversarially perturbs the text inputs by the word-level substitutions, whereas DeepWordBug performs the character-level perturbations to each input by replacing, scrambling, and erasing a few characters of some words. UTP attack TrojLLM uses reinforcement learning to search a universal trigger for a prompt-based language model, any text inputs with this trigger will lead to the model output target label.

**Implementation Details.** For the superior prompt generator configuration, we adhered to the parameters established in RL-Prompt (Deng et al., 2022). Specifically, we use distilGPT-2, a large model with 82 million parameters, as a policy model for all tasks. Additionally, we use a multilayer perceptron (MLP) with one hidden layer which has 2,048 hidden states, added to distilGPT-2’s existing 768 hidden states. For the hyperparameters of reward functions in the Equation 4, we set balancing weights  $\eta_1 = 180$  and  $\eta_2 = 200$ . During inference of CR-UTP, we use an ensemble number of 5 with the best 5 prompts derived from the superior prompt search. During the certification process, the prediction number is 500 and the certification number is 1000. When using random masking to defend against adversarial attacks, the voting number is set to 100. All experiments are conducted on a single Nvidia Geforce RTX-3090 GPU. Searching time for superior prompts on SST-2 is 3.8 hours, the certification time for one sentence is  $\sim 8$  seconds.

## 5 Results

### Comparison of CR-UTP with Random Mask.

In Table 1, we conducted experiment comparing the performance of CR-UTP with Random Mask (Zeng et al., 2023) at a 70% masking ratio against three adversarial attacks, i.e., TextFool (Jin et al., 2020), DeepWordBug (Gao et al., 2018) and TrojLLM (Xue et al., 2023). Our CR-UTP exhibits superior performance over Random Mask across all metrics in the evaluated adversarial attacks. Notably, CR-UTP achieves a significant 3.68% increase in clean accuracy (CACC). This improvement is attributed to its efficient prompt search method, which identifies robust prompts to random mask, and superior prompt ensemble technique, further reducing CACC variance. Moreover, CR-UTP achieves a substantial reduction in attack success rate (ASR), averaging a 21.4% greater decrease compared to Random Mask, with a remarkable 34.76% ASR reduction in the TrojLLM attack. This enhancement stems from CR-UTP’s ability to leverage the differential outputs of various prompts, enabling a robust ensemble prediction for improved defense outcomes against adversarial samples. Additionally, CR-UTP demonstrates a notable increase in poisoned accuracy (PACC), indicating its ability to maintain high accuracy even under attack scenarios.

Table 1: The comparison of CR-UTP and Random Mask against various attacks with a 70% mask ratio on the SST-2 dataset.

Scheme	Random Mask			CR-UTP		
	CACC	ASR	PACC	CACC	ASR	PACC
TextFool	81.60	42.88	57.12	85.28	37.39	62.61
DeepWordBug	81.13	45.18	54.82	85.61	21.25	78.75
TrojLLM	80.94	85.31	56.84	85.70	50.55	73.04

### 5.1 Ablation Study

In this section, we explore the design space of CR-UTP and study the impact of various settings of CR-UTP on its attacking effects using RoBERTa-Large with SST-2 dataset.

**CR-UTP Techniques Performance.** In Table 2, we analyze the impact of different CR-UTP techniques on performance against TrojLLM on the SST-2 and AgNews datasets. For SST-2, utilizing the adapted random mask method (our baseline) leads to a significant drop in CACC by over 22%, mainly due to the loss of information from masking 70% of the words. However, incorporating superior prompt search with reward  $\mathcal{R}_{\text{MSAR}}$ , improves CACC by 11% as the superior prompt proves more robust to random masking. Furthermore, combining rewards  $\mathcal{R}_{\text{MSAR}} + \mathcal{R}_{\text{PDAR}}$ , further increases CACC to 84.90% by enhancing prompt search effectiveness with  $\mathcal{R}_{\text{PDAR}}$ , which aligns outputs for clean and masked sentences. Finally, employing the superior prompt ensemble technique elevates CACC to 85.70% and reduces ASR from 91.88% to 50.55%, indicating significant improvements over the baseline method. Similarly, on AgNews dataset, CR-UTP surpasses the baseline with a 3.33% increase in CACC to 84.27% and a 3.6% decrease in ASR, highlighting CR-UTP’s effectiveness.

Table 2: An ablation study of CR-UTP techniques. Our baseline is random mask with 70% ratio;  $\mathcal{R}_{\text{MSAR}}$  denotes employing superior prompt search only using reward  $\mathcal{R}_{\text{MSAR}}$ ;  $\mathcal{R}_{\text{MSAR}} + \mathcal{R}_{\text{PDAR}}$  means using superior prompt search with rewards  $\mathcal{R}_{\text{MSAR}}$  and  $\mathcal{R}_{\text{PDAR}}$ ; CR-UTP incorporates all proposed techniques.

Method	SST-2			AgNews		
	CACC	ASR	PACC	CACC	ASR	PACC
w/o defense	92.69	91.88	53.76	88.91	94.54	78.64
Our baseline	70.50	85.31	56.84	80.94	22.42	75.71
$\mathcal{R}_{\text{MSAR}}$	81.93	47.28	76.83	82.09	21.31	75.78
$\mathcal{R}_{\text{MSAR}} + \mathcal{R}_{\text{PDAR}}$	84.90	63.93	66.61	83.06	18.89	78.65
CR-UTP	85.70	50.55	73.04	84.27	18.82	78.73



**Mask Ratio.** To examine the effect of mask ratio on clean accuracy, we conduct experiments on the SST-2 dataset with varying mask ratios. Results in Figure 4 (a) show that while the baseline method’s accuracy sharply drops from 91.27% to 51.78% as the mask ratio increases from 10% to 90%, our superior prompt search technique leads to a more gradual decline, from 92.42% to 57.82%. Additionally, employing our superior prompt ensemble method maintains a higher accuracy of 85.70% even at a 70% mask ratio, representing a significant improvement over the baseline method. In Figure 4 (b), the variance analysis of certified accuracy shows that while increasing the mask ratio results in higher variance for both baseline and superior prompt methods, the use of ensemble techniques, particularly the superior prompt ensemble method, reduces variances, providing a more consistent output despite the effects of masking. The variance peaks at the 60% mask ratio, indicating the highest sentence diversity. This suggests that the variance is influenced not only by the volume of information loss due to masking but also by the diversity of sentences resulting from random masking. However, the employment of ensemble techniques, even with baseline ensemble (vanilla prompts, not superior prompts), results in a more gradual increase in variance. This stabilization is likely due to the ensemble’s ability to aggregate insights from multiple prompts, delivering a more consistent and reliable output despite the information loss introduced by masking. The superior prompt ensemble technique further reduces the variances.

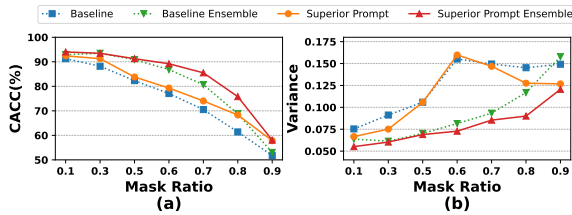


Figure 4: (a) Clean accuracy and (b) variance of proposed methods under different mask ratio.

**Ensemble Numbers.** To investigate the impact of the number of prompts within the superior prompt ensemble on clean accuracy, we conducted experiments on the SST-2 dataset using a large mask ratio of 70% to amplify the ensemble number’s effect on output performance. To mitigate the potential impact of differences in prompt selection performance on the output, each ensemble was selected to

have similar mean accuracy. The results depicted in Figure 5 show that as the number of prompts in the ensemble increases from 1 to 50, there is a consistent improvement in clean accuracy, rising from 82.31% to 86.82%, accompanied by a corresponding decrease in variance. These findings indicate that a larger ensemble leads to more stable and accurate predictions. This enhancement can be attributed to the ensemble’s capacity to integrate diverse insights from multiple prompts, reducing the impact of any single erroneous prediction and fostering a consensus that is more resilient to the introduction of masks.

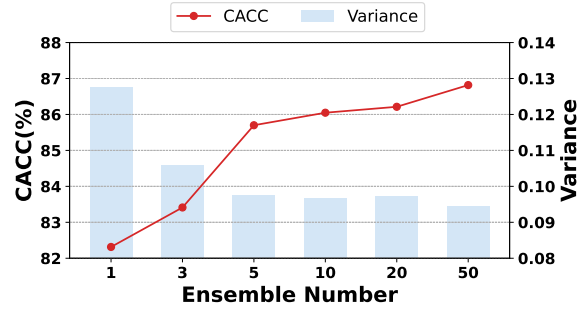


Figure 5: Clean accuracy and variance of CR-UTP under different ensemble numbers.

## 6 Limitation

The limitations of our paper are as follows: (i) Certified Accuracy. Although our CR-UTR has demonstrated improvements in certified accuracy and reduced ASR, achieving state-of-the-art results, there remains a gap between clean accuracy and certified accuracy. (ii) Model. While our CR-UTP has been evaluated on popular benchmark datasets SST-2 and AgNews using the RoBERTa-large model, it would be beneficial to assess its effectiveness across more models and different architectures to ensure its generalizability.

## 7 Conclusion

In conclusion, we address the challenge of certifying language model robustness against Universal Text Perturbations (UTPs) and input-specific text perturbations (ISTPs). We introduce the *superior prompt search* method and the *superior prompt ensembling* technique to enhance certified accuracy against UTPs and ISTPs. Our approaches achieve state-of-the-art results, ensuring stability and reliability in language model predictions across diverse attack scenarios.



## References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Moustafa Alzantot, Yash Sharma, Ahmed Elgohary, Bo-Jhang Ho, Mani Srivastava, and Kai-Wei Chang. 2018. Generating natural language adversarial examples. *arXiv preprint arXiv:1804.07998*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%\* chatgpt quality. See <https://vicuna.lmsys.org> (accessed 14 April 2023).
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2023. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113.
- Jeremy Cohen, Elan Rosenfeld, and Zico Kolter. 2019. Certified adversarial robustness via randomized smoothing. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 1310–1320. PMLR.
- Mingkai Deng, Jianyu Wang, Cheng-Ping Hsieh, Yi-han Wang, Han Guo, Tianmin Shu, Meng Song, Eric Xing, and Zhiting Hu. 2022. Rlprompt: Optimizing discrete text prompts with reinforcement learning. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 3369–3391.
- Javid Ebrahimi, Anyi Rao, Daniel Lowd, and Dejing Dou. 2018. c: White-box adversarial examples for text classification. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 31–36.
- Ji Gao, Jack Lanchantin, Mary Lou Soffa, and Yanjun Qi. 2018. Black-box generation of adversarial text sequences to evade deep learning classifiers. In *2018 IEEE Security and Privacy Workshops (SPW)*, pages 50–56. IEEE.
- Shreya Goyal, Sumanth Doddapaneni, Mitesh M Khapra, and Balaraman Ravindran. 2023. A survey of adversarial defenses and robustness in nlp. *ACM Computing Surveys*, 55(14s):1–39.
- James B Heaton, Nick G Polson, and Jan Hendrik Witte. 2017. Deep learning for finance: deep portfolios. *Applied Stochastic Models in Business and Industry*, 33(1):3–12.
- Miklós Z Horváth, Mark Niklas Müller, Marc Fischer, and Martin Vechev. 2021. Boosting randomized smoothing with variance reduced classifiers. *arXiv preprint arXiv:2106.06946*.
- Jongheon Jeong and Jinwoo Shin. 2020. Consistency regularization for certified robustness of smoothed classifiers. *CoRR*, abs/2006.04062.
- Di Jin, Zhijing Jin, Joey Tianyi Zhou, and Peter Szolovits. 2020. Is bert really robust? a strong baseline for natural language attack on text classification and entailment. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 8018–8025.
- Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT*, pages 4171–4186.
- J. Zico Kolter and Eric Wong. 2017. Provable defenses against adversarial examples via the convex outer adversarial polytope. *CoRR*, abs/1711.00851.
- Linyang Li, Ruotian Ma, Qipeng Guo, Xiangyang Xue, and Xipeng Qiu. 2020a. Bert-attack: Adversarial attack against bert using bert. *arXiv preprint arXiv:2004.09984*.
- Linyi Li, Xiangyu Qi, Tao Xie, and Bo Li. 2020b. Sok: Certified robustness for deep neural networks. *CoRR*, abs/2009.04131.
- Linyi Li, Jiawei Zhang, Tao Xie, and Bo Li. 2023. Double sampling randomized smoothing.
- Chizhou Liu, Yunzhen Feng, Ranran Wang, and Bin Dong. 2021. Enhancing certified robustness via smoothed weighted ensembling.
- 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. *arXiv preprint arXiv:1907.11692*.
- Monika A Myszczyńska, Poojitha N Ojamies, Alix MB Lacoste, Daniel Neil, Amir Saffari, Richard Mead, Guillaume M Hautbergue, Joanna D Holbrook, and Laura Ferraiuolo. 2020. Applications of machine learning to diagnosis and treatment of neurodegenerative diseases. *Nature Reviews Neurology*, 16(8):440–456.
- Matan Ostrovsky, Clark W. Barrett, and Guy Katz. 2022. An abstraction-refinement approach to verifying convolutional neural networks. *CoRR*, abs/2201.01978.

773	Shuhuai Ren, Yihe Deng, Kun He, and Wanxiang Che.	Tsui-Wei Weng, Huan Zhang, Hongge Chen, Zhao	826
774	2019. Generating natural language adversarial ex-	Song, Cho-Jui Hsieh, Duane Boning, Inderjit S.	827
775	amples through probability weighted word saliency.	Dhillon, and Luca Daniel. 2018. <a href="#">Towards fast com-</a>	828
776	In <i>Proceedings of the 57th annual meeting of the as-</i>	<a href="#">putation of certified robustness for relu networks.</a>	829
777	<i>sociation for computational linguistics</i> , pages 1085–		
778	1097.		
779	Hadi Salman, Mingjie Sun, Greg Yang, Ashish Kapoor,	Yiquan Wu, Kun Kuang, Yating Zhang, Xiaozhong Liu,	830
780	and J. Zico Kolter. 2020. <a href="#">Black-box smoothing: A</a>	Changlong Sun, Jun Xiao, Yueting Zhuang, Luo Si,	831
781	<a href="#">provable defense for pretrained classifiers.</a> <i>CoRR</i> ,	and Fei Wu. 2020. De-biased court’s view generation	832
782	abs/2003.01908.	with causality. In <i>Proceedings of the 2020 Confer-</i>	833
		<i>ence on Empirical Methods in Natural Language</i>	834
		<i>Processing (EMNLP)</i> , pages 763–780.	835
783	Hadi Salman, Greg Yang, Jerry Li, Pengchuan Zhang,	Lei Xu, Yangyi Chen, Ganqu Cui, Hongcheng Gao, and	836
784	Huan Zhang, Ilya P. Razenshteyn, and Sébastien	Zhiyuan Liu. 2022. Exploring the universal vulner-	837
785	Bubeck. 2019. <a href="#">Provably robust deep learning via</a>	ability of prompt-based learning paradigm. pages	838
786	<a href="#">adversarially trained smoothed classifiers.</a> <i>CoRR</i> ,	1799–1810.	839
787	abs/1906.04584.		
788	Erfan Shayegani, Md Abdullah Al Mamun, Yu Fu, Pe-	Jiaqi Xue, Mengxin Zheng, Ting Hua, Yilin Shen,	840
789	dram Zaree, Yue Dong, and Nael Abu-Ghazaleh.	Yepeng Liu, Ladislau Bölöni, and Qian Lou. 2023.	841
790	2023. Survey of vulnerabilities in large language	<a href="#">TrojLLM: A black-box trojan prompt attack on large</a>	842
791	models revealed by adversarial attacks. <i>arXiv</i>	<a href="#">language models.</a> In <i>Thirty-seventh Conference on</i>	843
792	<i>preprint arXiv:2310.10844.</i>	<i>Neural Information Processing Systems.</i>	844
793	Richard Socher, Alex Perelygin, Jean Wu, Jason	Mao Ye, Chengyue Gong, and Qiang Liu. 2020. <a href="#">Safer:</a>	845
794	Chuang, Christopher D Manning, Andrew Y Ng, and	<a href="#">A structure-free approach for certified robustness to</a>	846
795	Christopher Potts. 2013. Recursive deep models for	<a href="#">adversarial word substitutions.</a> In <i>Annual Meeting of</i>	847
796	semantic compositionality over a sentiment treebank.	<i>the Association for Computational Linguistics.</i>	848
797	In <i>Proceedings of the 2013 conference on empiri-</i>		
798	<i>cal methods in natural language processing</i> , pages	Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang,	849
799	1631–1642.	Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu,	850
		Wendi Zheng, Xiao Xia, et al. 2022. Glm-130b:	851
		An open bilingual pre-trained model. <i>arXiv preprint</i>	852
		<i>arXiv:2210.02414.</i>	853
800	Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam	Jiehang Zeng, Jianhan Xu, Xiaoqing Zheng, and Xuan-	854
801	Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng,	jing Huang. 2023. Certified robustness to text adver-	855
802	Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, et al.	sarial attacks by randomized [mask]. <i>Computational</i>	856
803	2022. Lamda: Language models for dialog applica-	<i>Linguistics</i> , 49(2):395–427.	857
804	tions. <i>arXiv preprint arXiv:2201.08239.</i>		
805	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	Jiehang Zeng, Xiaoqing Zheng, Jianhan Xu, Linyang Li,	858
806	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	Liping Yuan, and Xuanjing Huang. 2021. <a href="#">Certified</a>	859
807	Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti	<a href="#">robustness to text adversarial attacks by randomized</a>	860
808	Bhosale, et al. 2023. Llama 2: Open founda-	<a href="#">[MASK].</a> <i>CoRR</i> , abs/2105.03743.	861
809	tion and fine-tuned chat models. <i>arXiv preprint</i>		
810	<i>arXiv:2307.09288.</i>	Runtian Zhai, Chen Dan, Di He, Huan Zhang, Boqing	862
		Gong, Pradeep Ravikumar, Cho-Jui Hsieh, and Liwei	863
		Wang. 2020. <a href="#">MACER: attack-free and scalable ro-</a>	864
811	Eric Wallace, Shi Feng, Nikhil Kandpal, Matt Gardner,	<a href="#">bust training via maximizing certified radius.</a> <i>CoRR</i> ,	865
812	and Sameer Singh. 2019. Universal adversarial trig-	abs/2001.02378.	866
813	gers for attacking and analyzing nlp. <i>arXiv preprint</i>		
814	<i>arXiv:1908.07125.</i>	Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015.	867
		Character-level convolutional networks for text classi-	868
815	Xiaosen Wang, Hao Jin, and Kun He. 2019. <a href="#">Natural</a>	fication. <i>Advances in neural information processing</i>	869
816	<a href="#">language adversarial attacks and defenses in word</a>	<i>systems</i> , 28.	870
817	<a href="#">level.</a> <i>CoRR</i> , abs/1909.06723.		
818	Maurice Weber, Xiaojun Xu, Bojan Karlas, Ce Zhang,	Zhen Zhang, Guanhua Zhang, Bairu Hou, Wenqi Fan,	871
819	and Bo Li. 2020. <a href="#">RAB: provable robustness against</a>	Qing Li, Sijia Liu, Yang Zhang, and Shiyu Chang.	872
820	<a href="#">backdoor attacks.</a> <i>CoRR</i> , abs/2003.08904.	2023. <a href="#">Certified robustness for large language models</a>	873
		<a href="#">with self-denoising.</a>	874
821	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	Yi Zhou, Xiaoqing Zheng, Cho-Jui Hsieh, Kai-Wei	875
822	Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,	Chang, and Xuan-Jing Huang. 2021. Defense against	876
823	et al. 2022. Chain-of-thought prompting elicits rea-	synonym substitution-based adversarial attacks via	877
824	soning in large language models. <i>Advances in Neural</i>	dirichlet neighborhood ensemble. In <i>Proceedings</i>	878
825	<i>Information Processing Systems</i> , 35:24824–24837.	<i>of the 59th Annual Meeting of the Association for</i>	879
		<i>Computational Linguistics and the 11th International</i>	880

*Joint Conference on Natural Language Processing*  
(Volume 1: Long Papers), pages 5482–5492.