

Adaptive Adversarial Training for Balancing Model Robustness and Standard Performance

Anonymous ACL submission

Abstract

Adversarial training (AT) is widely used to boost model robustness against adversarial attacks, i.e., adding minor perturbations on the clean input to fool the target model. However, AT can also lead to degraded clean accuracy since it changes the distribution of the training set. Using the Taylor expansion, we find that commonly used adversarial loss functions inherently include clean loss, making it challenging for previous methods to effectively balance accuracy and robustness. Based on this, we establish a flexible AT framework that can explicitly balance the model robustness and clean accuracy by assigning learnable weights to the decomposed adversarial loss. Comprehensive experimental results indicate that our method boosts model robustness while maintaining comparable standard performance.

1 Introduction

Adversarial training (AT) attempts to boost the robustness of classifiers against adversarial examples by augmenting the training set with perturbed samples. While this approach effectively reduces adversarial errors or boosts robust accuracies, it has been observed to impair the standard performance on clean test data (Tramer and Boneh, 2019; Raghu-nathan et al., 2020; Yuan et al., 2019; Zhang and Li, 2019). Recent discussions (Yoo and Qi, 2021; Yuan et al., 2019; Zhang and Li, 2019) suggest a trade-off in AT, implying the challenge of simultaneously minimising standard and adversarial risks.

This paper focuses on AT for natural language processing (NLP) tasks, especially for text classification. The overarching concept of AT involves a two-level optimization process to enhance the model robustness. On the inner level, gradient ascent is employed to optimize small perturbations of the input data, aiming to maximize the model’s loss function. On the outer level, gradient descent is utilized to adjust the model parameters to min-

imize the classification loss of these adversarial examples.

We note that in textual AT, the default iteration number k is often quite small, e.g., 3 for FreeLB (Zhu et al., 2020), TAVAT (Li and Qiu, 2021), and InfoBERT (Wang et al., 2021a), resulting in small perturbation sizes for these methods. Their empirical results indicate that a relatively small perturbation size helps boost both model robustness and performance. Nevertheless, we doubt whether a small perturbation size is really helpful in improving robustness. Because the inner maximization greatly affects the effectiveness of AT (Wang et al., 2022). A small perturbation size usually generates lower-quality adversarial data, which makes AT useless. For example, the Fast Gradient Sign Method (FGSM) (Goodfellow et al., 2015) can quickly generate adversarial data using one step, but contributing little to robustness. Thus, it is reasonable to vary the iteration number and the adversarial step size to study how inner maximization affects AT. To this end, we choose two widely adopted AT methods, i.e., the projected gradient descent (PGD) method (Madry et al., 2018) and the FreeLB method (Zhu et al., 2020) as our baselines to conduct preliminary experiments on the BERT-base (Devlin et al., 2019) model.

We report clean and robust accuracies¹ equipped with PGD and FreeLB in Figure 1 and Figure 2. We find that the existing AT method can hardly improve robustness without hurting clean accuracy, which contradicts the results in previous works. As the perturbation size increases in AT, the robustness increases while the accuracy decreases. Additionally, AT will easily collapse and fail to converge when the perturbation size becomes too large.

This preliminary result indicates that we must rethink the trade-off between robustness and ac-

¹In this paper, we use clean accuracy to refer to the standard test accuracy and use robust accuracy to refer to the test accuracy on adversarial examples.

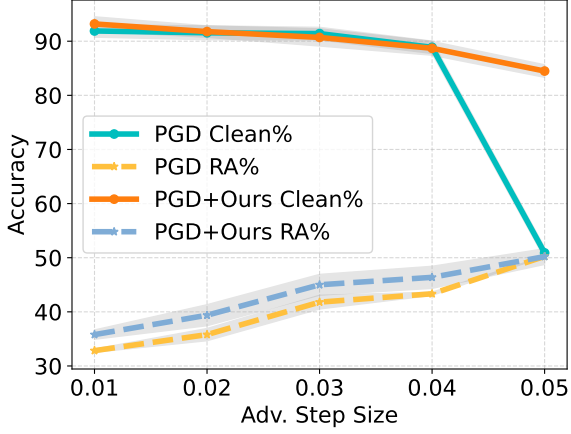


Figure 1: The robust accuracy (**RA**) and the clean accuracy (**Clean**) of the original PGD method (Madry et al., 2018) and ours on the SST-2 dataset (Socher et al., 2013). The backbone model is BERT-base (Devlin et al., 2019). As the step size increases, a trade-off exists between **RA** and **Clean**. It is hard to achieve optimal robustness and clean accuracy simultaneously.

curacy for NLP models. It also motivates us to investigate whether there exists an optimal perturbation size for the sake of balance, and how to make AT converge in a large perturbation size to achieve strong robustness.

To this end, we theoretically analyze the impact of the perturbation size on the learning objective of AT. In particular, we perform Taylor expansion on the adversarial loss and decompose it into a clean data loss and an adversarial one, in which the adversarial one is the weighted sum of squares of all perturbations. The clean loss corresponds to the model’s accuracy, while the adversarial one corresponds to the model’s robustness. By assigning trainable weights to all the perturbations, we can explicitly balance the two losses to achieve comparable model robustness and standard performance.

We further provide extensive experimental results. Compared with existing state-of-the-art AT methods, our method demonstrates a remarkable improvement in robustness without sacrificing the standard accuracy. Our main contributions are:

- We demonstrate that existing AT methods for NLP models either fail to improve robustness or compromise clean accuracy.
- We conduct theoretical analysis on a series of gradient-based AT methods. We decompose their learning objectives into distinct adversarial and clean loss components, allowing us to explicitly balance model robustness and

accuracy on clean data.

- We establish a flexible AT framework where one can balance adversarial loss and clean loss by assigning learnable weights to adversarial perturbations. Empirical evaluations show that our method can improve model robustness without sacrificing clean accuracy.

2 Related Work

2.1 Adversarial Training

AT is widely used to improve robustness against malicious adversarial attacks. Let $f(\cdot)$ be a neural network, Θ be the model parameters, X be the input data set and Y be the corresponding label set, with each input data $x \in X$ and label $y \in Y$. In practice, AT is developed to solve the following max-min optimization problem:

$$\min_{\Theta} \max_{\delta} \mathcal{L}(f(\Theta, x + \delta, y)), \quad (1)$$

where δ denotes the minor perturbation term added to the input.

While the outer minimization is often solved by stochastic gradient descent, how to tackle the inner maximization objective function is still under continuous study. Goodfellow et al. (2015) proposed FGSM to generate perturbations in one gradient ascent step as follows:

$$\delta = \text{sign}(\nabla_x \mathcal{L}(\Theta, x, y)), \quad (2)$$

where $\text{sign}(\cdot)$ is the sign function.

However, this approximation can hardly find high-quality adversarial data that can maximize the loss function. To seek more precise solutions, Madry et al. (2018) proposed the Projected Gradient Descent (PGD) method to generate perturbations using multi-step gradient ascent steps, i.e.,

$$\begin{aligned} \delta_t &= \alpha \cdot \nabla_{\delta} \mathcal{L}(f(\Theta, x_{t-1}, y)), \\ x_t &= x_{t-1} + \delta_t. \end{aligned} \quad (3)$$

Moreover, PGD initializes the search for adversarial data at random starting points within the allowed norm ball, improving the diversity of adversarial data. Empirically, PGD and its variants are still considered the most effective AT methods.

For NLP tasks, AT was first used to improve the generalization of models. Miyato et al. (2017) proposed virtual AT to enhance text classification in a semi-supervised manner. To further improve language understanding for pre-trained language

models, Zhu et al. (2020) proposed FreeLB to provide a large virtual batch size in AT.

In another line of work, AT was adopted to boost the robustness of NLP models. By adversarially perturbing their embedding layer, NLP models were trained to predict consistently on both clean and adversarial data, thereby achieving better adversarial robustness. For example, Li and Qiu (2021) proposed TAVAT to generate token-level perturbations accounting for the importance of tokens. Li et al. (2021) increased the iteration numbers of AT and found it useful for boosting robustness. Gao et al. (2023) proposed to minimize the distribution shift risk between clean and adversarial data. Formento et al. (2024) learned robust word embeddings to defend against adversarial attacks.

2.2 The Trade-off between Robustness and Accuracy

In computer vision, while AT helps improve robustness, a vast amount of empirical evidence exists that the clean accuracy can be hurt (Madry et al., 2018; Wang et al., 2020). Zhang et al. (2019) theoretically identified the trade-off between robustness and accuracy by decomposing the prediction error for adversarial examples (robust error) as the sum of the natural error and boundary error. Nevertheless, Yang et al. (2020) proved that robustness and accuracy should both be achievable for benchmark datasets through locally Lipschitz functions.

For NLP models, early research generally holds that AT improves both robustness and accuracy (Miyato et al., 2017; Ren et al., 2019; Zhu et al., 2020). However, few studies have focused on the trade-off between robustness and accuracy in AT of NLP models.

It is worth noting that several adversarial data augmentation (ADA) methods (Ren et al., 2019; Li et al., 2019; Jin et al., 2020; Li et al., 2020) expand the original training set with crafted adversarial examples. ADA methods introduce larger perturbations than gradient-based AT methods, leading to relatively low clean accuracy. It demonstrates that there is also a trade-off between robustness and accuracy in AT of NLP models.

In this work, we first demonstrate that with a large perturbation size, robustness trades off clean accuracy in gradient-based AT of NLP models. Further, by decomposing the learning objective of AT into a clean classification loss and an adversarial one, we can explicitly balance clean accuracy and robustness.

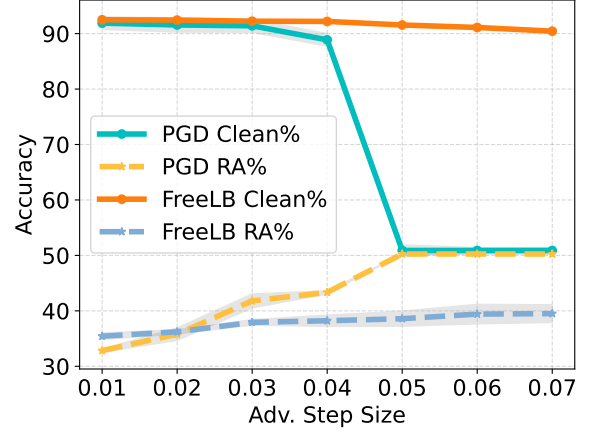


Figure 2: The robust accuracy (RA) and clean accuracy (Clean) of PGD and FreeLB under different adversarial step sizes on the SST-2 dataset. The backbone model is BERT-base. Although PGD can achieve higher robustness than FreeLB, the clean accuracy of the model is greatly damaged. When the perturbation is too large, the training cannot converge.

3 On the Convergence of Adversarial Training

It is widely pointed out that AT is more difficult than standard training for both computer vision and NLP models (Madry et al., 2018; Kurakin et al., 2017). The main reason is that a distribution difference exists between adversarial data and clean data, which makes one model unable to converge well on two widely different data distributions. According to (Gao et al., 2023), one can model the difference via Wasserstein distance. The authors proved that the distribution shift is bounded by the adversarial perturbation δ . Therefore, δ is crucial in the convergence of AT. We then vary δ to show its effect².

Figure 2 shows that as δ increases, the clean accuracy drops significantly, which implies that AT cannot converge well with large perturbations. Furthermore, robust accuracy gradually increases, demonstrating a trade-off between clean accuracy and robust accuracy.

It is also intriguing to see that FreeLB can converge under larger perturbations than PGD. We theoretically analyse the differences among different AT methods to understand this phenomenon. Recall the learning objective in AT,

$$\min_{\Theta} \mathbb{E}_{(X,Y) \sim D} \left[\max_{\|\delta\| \leq \epsilon} \mathcal{L}(\Theta, X + \delta, Y) \right], \quad (4)$$

²In practice we vary the step size α in AT to control the perturbation size.

where Θ is the model parameters; D is the data distribution; δ is the added perturbation; ϵ is the allowed perturbation size. In the min-max process, multi-step gradient ascent methods often solve the inner maximisation. Take PGD as an example. We initialize x_0 to x and suppose that the iteration number is k and the adversarial step size is α , we have

$$x_t = \text{proj}_\epsilon(x_{t-1} + \alpha \cdot \text{norm}(g(x_{t-1}))), \quad 1 \leq t \leq k, \quad (5)$$

where $g(x_{t-1})$ is the gradient of x_{t-1} , $\text{norm}(\cdot)$ can be L_2 normalization. The initial value x_0 can also be randomly sampled within the ϵ -neighborhood of x . In that case, we have $x_0 = x + \delta_0$, where δ_0 is randomly sampled.

For simplicity, we omit the projection function and the normalization function. The main reason is that in (Li et al., 2021), the authors have demonstrated that removing the norm-bounded limitation helps achieve better model robustness.

Thus, we have

$$x_t = x_{t-1} + \delta_t = x + \sum_{t=1}^k \delta_t, t \geq 1. \quad (6)$$

In this way, the inner maximization can be reformulated as follows:

$$\max_{\|\delta\| \leq \epsilon} \mathcal{L}(\Theta, x + \sum_{t=1}^k \delta_t, y). \quad (7)$$

Considering that δ_t is very small relative to input x , we perform first-order Taylor expansion on the loss function. Thus, combining Eq. (7) and omit the high-order terms, we have

$$\begin{aligned} \mathcal{L}(\Theta, x + \sum_{t=1}^k \delta_t, y) &= \mathcal{L}(\Theta, x + \sum_{t=1}^{k-1} \delta_t + \delta_k, y) \\ &\approx \mathcal{L}(\Theta, x + \sum_{t=1}^{k-1} \delta_t, y) + \frac{1}{\alpha} \delta_k^2 \\ &\dots \\ &\approx \mathcal{L}(\Theta, x, y) + \frac{1}{\alpha} \sum_{t=1}^k \delta_t^2. \end{aligned} \quad (8)$$

Eq. (8) indicates that one can decompose the loss of adversarial data during PGD training into the corresponding loss of clean data and the sum of squares of all perturbations δ_t .

Therefore, it is reasonable that as the perturbation size increases, the adversarial loss becomes larger and begins to dominate the training, leading to higher robustness. For clean accuracy, as the perturbation size enlarges, the model gets harder to converge on the original training set, resulting in lower clean accuracy.

Based on Eq. (8), we further study how δ affects the convergence of AT. We firstly extend Eq. (8) to FreeLB. It can also be easily extended to other PGD-like methods such as FreeLB++.

According to the FreeLB method, the number of iterations is k and the step size is α . The loss of each iteration will be divided by k and accumulated. The model parameters will be updated at the end (for comparison, PGD only uses the loss of the last iteration to update the model parameters). Similarly, the inner maximization of FreeLB can be formulated as follows:

$$\max_{\|\delta\| \leq \epsilon} \frac{1}{k} \sum_{t=1}^k \mathcal{L}(\Theta, x + r^t, y), \quad (9)$$

where $r^t = \sum_{i=1}^t \delta_i$. Performing first-order Taylor expansion on Eq. (9), similar to Eq. (8), we have

$$\begin{aligned} \frac{1}{k} \sum_{t=1}^k \mathcal{L}(\Theta, x + r^t, y) &\approx \frac{1}{k} \sum_{t=1}^k (\mathcal{L}(\Theta, x, y) + \frac{1}{\alpha} \sum_{i=1}^t \delta_i^2) \\ &= \mathcal{L}(\Theta, x, y) + \frac{1}{\alpha} \sum_{i=1}^k \frac{k-i+1}{k} \delta_i^2. \end{aligned} \quad (10)$$

Eq. (10) indicates that the learning objective of FreeLB can also be decomposed into the clean data loss and the weighted sum of squares of all perturbations δ_i , where the weight of δ_i is $\frac{k-i+1}{k}$.

At this point, we can explain more clearly in Figure 2. Since the PGD method inherently has a greater weight for adversarial loss, it can achieve higher robustness than FreeLB, but the training of PGD is more difficult to converge.

4 Adaptive Adversarial Training

4.1 A Unifying Framework for Adversarial Training

Comparing the two learning objectives, one can find an implicit set of weights weighing the perturbation δ_i produced at each iteration i . Further, the weights of clean classification loss and the adversarial one are also implicitly given. For the PGD method with an iteration number of k , the weights of clean loss and the adversarial loss are 1 and k ,

Algorithm 1 Adaptive Adversarial Training

Input: Model parameters Θ , loss function \mathcal{L} , training set $D = \{x_i, y_i\}_{i=1}^n$, number of epochs T , batch size m , number of iterations k , number of batches M , perturbation weights $\hat{\mathbf{w}}$

Output: robust model parameters Θ

```
1: for epoch = 1 to  $T$  do
2:   for batch = 1 to  $M$  do
3:     Sample a mini-batch  $b = \{(x_i, y_i)\}_{i=1}^m$ 
4:     Generate adversarial perturbations  $\delta$  via Eq. (3)
5:     Compute the overall loss  $\hat{\mathcal{J}}(\Theta, x, y, \hat{\mathbf{w}})$  via Eq. (12)
6:     Update  $\hat{\mathbf{w}}$  via  $\nabla_{\hat{\mathbf{w}}} \mathcal{J}(\Theta, x, y, \hat{\mathbf{w}})$ 
7:     Update  $\Theta$  via  $\nabla_{\Theta} \mathcal{J}(\Theta, x, y, \hat{\mathbf{w}})$ 
8:   end for
9: end for
```

respectively. The FreeLB method’s weights are 1 and $(k + 1)/2$, respectively.

Therefore, we summarize the learning objective \mathcal{J} of the two methods into the following formula:

$$\begin{aligned} \mathcal{J}(\Theta, x, y, \mathbf{w}) &= \mathcal{L}(\Theta, x, y) + \beta \frac{1}{\alpha} \sum_{i=1}^k w_i \delta_i^2, \\ \text{s.t. } \sum_{i=1}^k w_i &= 1, w_i \geq 0, \end{aligned} \quad (11)$$

where β balances the clean loss and the adversarial loss, and w_i balances all the perturbations.

However, since the derivative of the sum of squared perturbations involves computing the second-order derivative, we further manipulate the above formula. We introduce a set of parameters $\hat{\mathbf{w}}$ and combine it with Eq. (7), yielding the following expression:

$$\hat{\mathcal{J}}(\Theta, x, y, \hat{\mathbf{w}}) = \mathcal{L}(\Theta, x + \beta \sum_{i=1}^k \hat{w}_i \delta_i, y). \quad (12)$$

By performing Taylor expansion on Eq. (12), we can easily verify that each term corresponds one-to-one with Eq. (11). For $\hat{\mathcal{J}}$ to be equal to \mathcal{J} , $\hat{\mathbf{w}}$ needs to satisfy the following constraint which is the same as \mathbf{w} :

$$\sum_{i=1}^k \hat{w}_i = 1, \hat{w}_i \geq 0. \quad (13)$$

In our experiments, we initialize it to a vector of ones and update it automatically using its gradient.

4.2 The Rationale behind Our Framework

Next, we explain the rationale behind introducing β and $\hat{\mathbf{w}}$. As deduced above, in the PGD method, the weight of the clean loss is naturally set to 1, while the weight of the adversarial loss is set to k . In the FreeLB method, the weight of the clean loss is also 1, but the weight of the adversarial loss is $(k + 1)/2$. To ensure the extensibility of our AT framework, we introduce the parameter β to balance the clean loss and adversarial loss. Specifically, PGD and FreeLB are two special cases of the proposed framework.

Eq. (10) shows that while maintaining the original ratio between clean loss and adversarial loss, the perturbations at each time step t are assigned different weights. Therefore, we introduce a set of parameters $\hat{\mathbf{w}}$, ensuring that the sum of \hat{w}_i equals 1, and utilize gradients to solve for the worst-case scenario. The weights $\hat{\mathbf{w}}$ are continuously updated throughout the training process, in order to find the optimal solution across the entire training set rather than achieve a local optimum based on a single batch of data.

It is worth noting that, in the PGD method, although different time-step perturbations are not explicitly weighted, one can assume that their weights are uniformly set to 1.

Following the min-max optimization widely used in AT, the final training objective can be defined as:

$$\min_{\Theta} \max_{\hat{\mathbf{w}}} \hat{\mathcal{J}}(\Theta, x, y, \hat{\mathbf{w}}). \quad (14)$$

In this way, we build our novel framework of adaptive AT in a constrained manner, where both the PGD and the FreeLB methods can be considered special cases of our framework.

Notably, our framework can encompass a wider range of PGD-based AT algorithms, not limited to FreeLB. We show our proposed adaptive AT method in Algorithm 1.

5 Experimental Setup

5.1 Tasks and Datasets

Following previous important works (Gao et al., 2023; Li et al., 2021; Li and Qiu, 2021), we compare our adaptive AT method with baselines on two tasks, i.e., text classification and natural language inference. In the main experiments, we choose the SST-2 (Socher et al., 2013)³ and the QNLI (Wang

³<https://dl.fbaipublicfiles.com/glue/data/SST-2.zip>

et al., 2019)⁴ datasets to perform text classification and natural language inference tasks, respectively. For completeness, we also test the applicability of our method on the IMDB dataset (Maas et al., 2011) and the AGNEWS dataset (Zhang et al., 2015), both used for text classification. Detailed characteristics and examples of the four datasets are presented in Appendix A.

5.2 Baseline Methods

5.2.1 Defence Methods

We apply our framework to various AT-based defence methods, including PGD (Madry et al., 2018), FreeLB (Zhu et al., 2020), and TA-VAT (Li and Qiu, 2021). To comprehensively benchmark existing defence methods, we report the results of InfoBERT (Wang et al., 2021a), Flooding-X (Liu et al., 2022), and SMART (Jiang et al., 2020) which enhance AT by an information bottleneck, “flooding”, and smoothness-inducing regularization, respectively. The performance of TRADES (Zhang et al., 2019), which is the most relevant method from the computer vision domain to ours, is also presented.

DSRM (Gao et al., 2023), GAT (Zhu and Rao, 2023), and SemRoDe (Formento et al., 2024) are not chosen. This is because DSRM introduces a distribution shift into adversarial defence, while GAT and SemRoDe incorporate valid adversarial examples into the training process; none of these can be adopted for a fair comparison.

5.2.2 Attacking Methods

Following previous works, we use TextFooler (Jin et al., 2020), TextBugger (Li et al., 2019), and BAE (Garg and Ramakrishnan, 2020) as our attacking methods to dynamically generate adversarial examples during test time.

We also consider assessing AT methods against high-quality adversarial examples pre-crafted by human annotators. Therefore, we report the robust accuracy of all the models on the adversarial GLUE dataset (Wang et al., 2021b).

6 Main Results

Our proposed method can be easily extended to PGD-like AT methods. In this part, we advance PGD, FreeLB and TA-VAT with adaptive perturbations to assess the effectiveness of our method. We conduct the main experiments on the BERT-base

model to provide comprehensive comparisons with other AT methods.

Note that the value of β is related to the methods being extended. For example, when extending the PGD method using our framework, the value of β is set to k (i.e., the number of iterations) according to Eq. (8). We leave the exploration of the effects of different β values for future work.

Table 1 reports the defence results against different types of adversarial attacks on the SST-2 dataset, including two word-level attacks (TextFooler and BAE), one multi-level attack (TextBugger), and an adversarial test dataset (Adversarial GLUE). The main findings are:

- For clean accuracy, all the baseline AT methods maintain a similar level, since the adversarial strength is moderate. The PGD method has the lowest clean accuracy, which is consistent with the conclusions of previous work.
- For robust accuracy against dynamic adversarial attacks and human-crafted adversarial examples, our method can boost the performance of three AT methods. Compared with InfoBERT and Flooding-X, our method also maintains higher robustness.
- Our method can boost the robust accuracy of PGD, FreeLB and TA-VAT methods while achieving comparable clean accuracy, which is consistent with our motivations.

We also conduct experiments on the QNLI dataset. The main results are consistent with that on the SST-2 dataset. Our method consistently enhances robust accuracy across various adversarial attacks and test sets. Thanks to the adaptive strength of perturbations, the clean accuracy remains at a comparable level compared to other AT methods.

We note that the PGD method still has the lowest clean accuracy. According to Eq. (8), the PGD method implicitly places a greater weight on the adversarial loss than FreeLB. Since it is directly adopted from the visual domain, no adjustments have been made to the trade-off between robustness and clean accuracy. As a consequence, this method exhibits lower clean accuracy on NLP tasks.

Due to the space limit, we report the results on the IMDB and AGNEWS datasets in Appendix B.

⁴<https://huggingface.co/datasets/nyu-mll/glue>

SST-2	Clean %	TextFooler	TextBugger	BAE	AdvGLUE
		RA %	RA %	RA %	RA %
BERT-base (Devlin et al., 2019)	92.32	8.14	26.83	33.72	31.32
InfoBERT (Wang et al., 2021a)	91.74	10.89	32.68	37.96	32.17
Flooding-X (Liu et al., 2022)	92.32	12.60	32.45	35.44	27.00
SMART (Jiang et al., 2020)	91.78	10.45	30.15	33.26	23.54
TRADES (Zhang et al., 2019)	87.19	9.46	29.53	35.41	30.99
PGD (Madry et al., 2018)	89.11	12.96	32.22	35.21	39.13
+Ours	88.99 (-0.12)	16.06 (+3.10)	35.68 (+3.46)	40.02 (+4.81)	43.44 (+4.31)
FreeLB (Zhu et al., 2020)	92.20	9.98	34.06	37.73	30.13
+Ours	91.63 (-0.57)	15.69 (+5.71)	38.73 (+4.67)	41.22 (+3.49)	38.53 (+8.40)
TA-VAT (Li and Qiu, 2021)	91.40	11.93	35.89	37.61	32.00
+Ours	91.51 (+0.11)	18.46 (+6.53)	39.60 (+3.71)	40.94 (+3.33)	39.42 (+7.42)

Table 1: The clean accuracy (“Clean %”) and the robust accuracy (“RA %”) on the SST-2 dataset against TextFooler, TextBugger, and BAE attacks. We report the robust accuracy of the models on the adversarial GLUE dataset (i.e., AdvGLUE) to evaluate AT methods against pre-crafted adversarial examples. The backbone model is BERT-base.

QNLI	Clean %	TextFooler	TextBugger	BAE	AdvGLUE
		RA %	RA %	RA %	RA %
BERT-base (Devlin et al., 2019)	90.60	8.80	9.50	27.90	42.75
InfoBERT (Wang et al., 2021a)	89.10	5.30	6.80	30.90	44.00
Flooding-X (Liu et al., 2022)	91.50	12.00	16.60	40.30	47.00
SMART (Jiang et al., 2020)	91.77	8.50	13.22	33.46	39.02
TRADES (Zhang et al., 2019)	86.22	9.45	12.14	35.44	43.50
PGD (Madry et al., 2018)	87.00	11.30	16.80	43.60	41.50
+Ours	87.90 (+0.90)	16.80 (+5.50)	17.20 (+0.40)	41.20 (-2.40)	48.89 (+7.39)
FreeLB (Zhu et al., 2020)	89.60	14.40	14.10	40.50	44.75
+Ours	89.70 (+0.10)	16.60 (+2.20)	17.70 (+3.60)	43.10 (+2.60)	51.75 (+7.00)
TA-VAT (Li and Qiu, 2021)	91.51	12.60	14.30	40.94	43.00
+Ours	91.00 (-0.51)	18.46 (+5.86)	20.30 (+6.00)	44.20 (+3.26)	51.00 (+8.00)

Table 2: The clean accuracy (“Clean %”) and the robust accuracy (“RA %”) on the QNLI dataset against TextFooler, TextBugger, and BAE attacks. We report the robust accuracy of the models on the adversarial GLUE dataset (i.e., AdvGLUE) to evaluate AT methods against pre-crafted adversarial examples. The backbone model is BERT-base.

7 Discussions

In this section, we discuss the relationship between our method and existing AT methods. We highlight the importance of conducting AT on small language models like BERT, rather than solely focusing on large language models (LLMs). We also provide an error analysis of the approximate loss and demonstrate the PGD loss and approximate loss in AT in practice.

7.1 Relation to Existing Work

We list a series of loss functions of AT methods in Table 3 and discuss the difference between our proposed adaptive AT and conventional AT methods, including fast gradient method (FGM) (Miyato et al., 2017), PGD (Madry et al., 2018), TRADES (Zhang et al., 2019), and FreeLB (Zhu et al., 2020).

Appendix C contains more detailed discussions about these methods.

7.2 Beyond Model Parameters

Recently, LLMs have achieved remarkable results across many NLP tasks (Achiam et al., 2023; Guo et al., 2025). Therefore, it is necessary to reveal the importance of conducting AT on language models with fewer parameters, such as BERT. We select a more practical task, namely spam detection, and report the standard performance of models with varying parameter sizes in Table 4, including Naive Bayes (NB), Support Vector Machine (SVM), BERT, and LLMs. The usage of deepseek is detailed in Appendix G. We adopt the SMS Spam Collection dataset (Almeida et al., 2011), which contains 747 spam messages and 4,825 non-spam

Methods	Loss Function	Flexibility
Standard	$\mathcal{L}(\Theta, x, y)$	-
FGM (Miyato et al., 2017)	$\mathcal{L}(\Theta, x, y) + \mathcal{L}(\Theta, x + \delta, y)$	✗
PGD (Madry et al., 2018)	$\mathcal{L}(\Theta, x + \delta_k, y)$	✗
TRADES (Zhang et al., 2019)	$\mathcal{L}(\Theta, x, y) + \lambda KL(p(\Theta, x) p(\Theta, x + \delta))$	✓
FreeLB (Zhu et al., 2020)	$\frac{1}{k} \sum_{i=1}^k \mathcal{L}(\Theta, x + \delta_i, y)$	✗
Ours	$\mathcal{L}(\Theta, x, y) + \beta \frac{1}{\alpha} \sum_{i=1}^k w_i \delta_i^2$	✓

Table 3: Comparisons of different loss functions in AT. The adversarial perturbations in TRADES are generated by maximizing its regularization term (KL-divergence). The **Flexibility** indicates whether the method can explicitly control the weighting between clean loss and adversarial loss. ✗ indicates that the method cannot balance clean and adversarial losses. ✓ indicates that the method introduces a hyperparameter to balance the two types of loss, but lacks flexibility because the adversarial loss still contains the clean loss. ✓ indicates that it can explicitly balance clean loss and adversarial loss.

Method	Acc.	Pre.	Recall	F1
SVM (linear)	97.56	97.01	84.82	90.50
Multinomia NB	98.21	98.26	88.48	93.11
BERT-base	99.48	94.44	91.15	92.61
DeepSeek-r1-zeroshot	87.74	52.71	95.77	68.00
DeepSeek-r1-fewshot	95.45	79.75	91.30	85.14

Table 4: The performance of models with varying parameter sizes on the spam detection task. We use deepseek-r1 (Guo et al., 2025) to demonstrate the performance of LLMs on this dataset in zero-shot and few-shot manners.

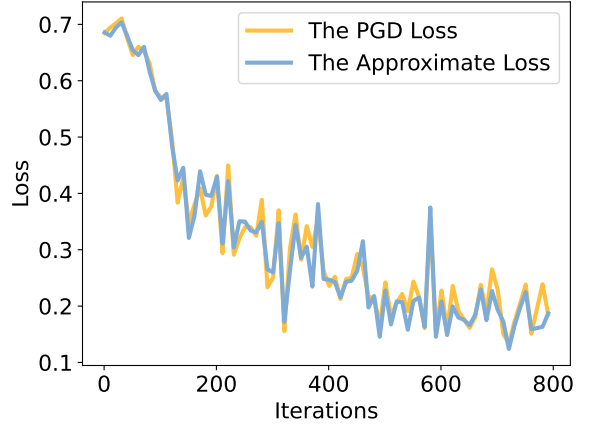


Figure 3: The error between the approximate loss and the original PGD loss on the SST2 dataset over the BERT model. This indicates that our approximation of the experiments is quite practicable.

messages. The long-tail distribution of the data makes it more realistic and challenging.

As can be seen, even the state-of-the-art DeepSeek-r1 model (Guo et al., 2025) performs poorly on this dataset, which may be related to the data distribution. However, small models generalize well on this dataset.

Given the constraints of computational resources and training efficiency, this study proposes to investigate AT for BERT-based architectures to mitigate vulnerabilities against adversarial perturbations, rather than focusing on LLMs.

7.3 Error Analysis

It is necessary to analyze the error of our method since we have ignored the higher-order terms in the Taylor expansion. Taking the PGD method as an example, we show the error between the approximate loss and the original PGD loss. The original PGD loss is computed by Eq. (7). The approximate is computed by Eq. (8).

In Figure 3, we observe that the approximate loss can well match the loss curve of the PGD

method. This demonstrates that our approximation is accurate in the experiments and it can be used to develop AT with an adaptive perturbation.

8 Conclusions

This work seeks to balance model robustness and accuracy. To this end, we decompose the learning objective of adversarial training into a pure adversarial loss and clean loss, which correspond to model robustness and clean accuracy, respectively. This way, we can explicitly assign learnable weights to the two losses to balance model robustness and clean accuracy. Experimental results on four datasets over BERT, RoBERTa and DeBERTa models show that our method can boost model robustness without sacrificing clean accuracy.

Limitations

This paper leverages Taylor expansion to decompose the loss function (i.e., the cross-entropy function) of AT. The Taylor expansion is a mathematical method used to approximate a function as a power series around a specific point. The loss function must have derivatives of sufficiently high order at the point of expansion and in its vicinity. Specifically, if we want to expand to the n -th order, the function must have at least n derivatives at that point. Although to our best knowledge, the mainstream of loss functions used in AT meet the above conditions, this may not be suitable for more complex loss functions.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Tiago A. Almeida, José María Gómez Hidalgo, and Akebo Yamakami. 2011. [Contributions to the study of SMS spam filtering: new collection and results](#). In *Proceedings of the 2011 ACM Symposium on Document Engineering, Mountain View, CA, USA, September 19-22, 2011*, pages 259–262. ACM.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Brian Formento, Wenjie Feng, Chuan-Sheng Foo, Anh Tuan Luu, and See-Kiong Ng. 2024. [SemRoDe: Macro adversarial training to learn representations that are robust to word-level attacks](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 8005–8028, Mexico City, Mexico. Association for Computational Linguistics.
- SongYang Gao, Shihan Dou, Yan Liu, Xiao Wang, Qi Zhang, Zhongyu Wei, Jin Ma, and Ying Shan. 2023. [DSRM: Boost textual adversarial training with distribution shift risk minimization](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12177–12189, Toronto, Canada. Association for Computational Linguistics.

- Siddhant Garg and Goutham Ramakrishnan. 2020. [BAE: BERT-based adversarial examples for text classification](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6174–6181, Online. Association for Computational Linguistics.
- Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. 2015. [Explaining and harnessing adversarial examples](#). In *International Conference on Learning Representations*.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Pengcheng He, Jianfeng Gao, and Weizhu Chen. 2021. [Debertav3: Improving deberta using electra-style pre-training with gradient-disentangled embedding sharing](#). *Preprint*, arXiv:2111.09543.
- Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Tuo Zhao. 2020. [SMART: Robust and efficient fine-tuning for pre-trained natural language models through principled regularized optimization](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2177–2190, Online. Association for Computational Linguistics.
- 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 *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 8018–8025.
- Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio. 2017. [Adversarial machine learning at scale](#). In *International Conference on Learning Representations*.
- Jinfeng Li, Shouling Ji, Tianyu Du, Bo Li, and Ting Wang. 2019. [Textbugger: Generating adversarial text against real-world applications](#). In *26th Annual Network and Distributed System Security Symposium, NDSS 2019, San Diego, California, USA, February 24-27, 2019*.
- Linyang Li, Ruotian Ma, Qipeng Guo, Xiangyang Xue, and Xipeng Qiu. 2020. [BERT-ATTACK: Adversarial attack against BERT using BERT](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6193–6202, Online. Association for Computational Linguistics.
- Linyang Li and Xipeng Qiu. 2021. [Token-aware virtual adversarial training in natural language understanding](#). In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference*

646	on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021, pages 8410–8418.	
650	Zongyi Li, Jianhan Xu, Jiehang Zeng, Linyang Li, Xiaoqing Zheng, Qi Zhang, Kai-Wei Chang, and Cho-Jui Hsieh. 2021. Searching for an effective defender: Benchmarking defense against adversarial word substitution. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 3137–3147, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	
659	Qin Liu, Rui Zheng, Bao Rong, Jingyi Liu, ZhiHua Liu, Zhanzhan Cheng, Liang Qiao, Tao Gui, Qi Zhang, and Xuanjing Huang. 2022. Flooding-X: Improving BERT’s resistance to adversarial attacks via loss-restricted fine-tuning. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 5634–5644, Dublin, Ireland. Association for Computational Linguistics.	
668	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. <i>CoRR</i> , abs/1907.11692.	
673	Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. 2011. Learning word vectors for sentiment analysis. In <i>Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies</i> , pages 142–150, Portland, Oregon, USA. Association for Computational Linguistics.	
681	Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. 2018. Towards deep learning models resistant to adversarial attacks. In <i>International Conference on Learning Representations</i> .	
686	Takeru Miyato, Shin-ichi Maeda, Masanori Koyama, and Shin Ishii. 2017. Virtual adversarial training: a regularization method for supervised and semi-supervised learning. <i>CoRR</i> , abs/1704.03976.	
690	John Morris, Eli Lifland, Jin Yong Yoo, Jake Grigsby, Di Jin, and Yanjun Qi. 2020. TextAttack: A framework for adversarial attacks, data augmentation, and adversarial training in NLP. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 119–126, Online. Association for Computational Linguistics.	
698	Aditi Raghunathan, Sang Michael Xie, Fanny Yang, John Duchi, and Percy Liang. 2020. Understanding and mitigating the tradeoff between robustness and accuracy. <i>arXiv preprint arXiv:2002.10716</i> .	
702	Shuhuai Ren, Yihe Deng, Kun He, and Wanxiang Che. 2019. Generating natural language adversarial examples through probability weighted word saliency. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 1085–1097, Florence, Italy. Association for Computational Linguistics.	702 703 704 705 706 707 708
709	Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In <i>Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing</i> , pages 1631–1642, Seattle, Washington, USA. Association for Computational Linguistics.	709 710 711 712 713 714 715 716
717	Florian Tramer and Dan Boneh. 2019. Adversarial training and robustness for multiple perturbations. <i>Advances in neural information processing systems</i> , 32.	717 718 719
720	Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In <i>International Conference on Learning Representations</i> .	720 721 722 723 724
725	Boxin Wang, Shuohang Wang, Yu Cheng, Zhe Gan, Ruoxi Jia, Bo Li, and Jingjing Liu. 2021a. Info{bert}: Improving robustness of language models from an information theoretic perspective. In <i>International Conference on Learning Representations</i> .	725 726 727 728 729
730	Boxin Wang, Chejian Xu, Shuohang Wang, Zhe Gan, Yu Cheng, Jianfeng Gao, Ahmed Hassan Awadallah, and Bo Li. 2021b. Adversarial GLUE: A multi-task benchmark for robustness evaluation of language models. In <i>Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)</i> .	730 731 732 733 734 735 736
737	Yisen Wang, Xingjun Ma, James Bailey, Jinfeng Yi, Bowen Zhou, and Quanquan Gu. 2022. On the convergence and robustness of adversarial training. <i>Preprint</i> , arXiv:2112.08304.	737 738 739 740
741	Yisen Wang, Difan Zou, Jinfeng Yi, James Bailey, Xingjun Ma, and Quanquan Gu. 2020. Improving adversarial robustness requires revisiting misclassified examples. In <i>International Conference on Learning Representations</i> .	741 742 743 744 745
746	Yao-Yuan Yang, Cyrus Rashtchian, Hongyang Zhang, Ruslan Salakhutdinov, and Kamalika Chaudhuri. 2020. A closer look at accuracy vs. robustness. In <i>Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual</i> .	746 747 748 749 750 751 752
753	Jin Yong Yoo and Yanjun Qi. 2021. Towards improving adversarial training of nlp models. <i>arXiv preprint arXiv:2109.00544</i> .	753 754 755
756	Xiaoyong Yuan, Pan He, Qile Zhu, and Xiaolin Li. 2019. Adversarial examples: Attacks and defenses for deep	756 757

learning. *IEEE transactions on neural networks and learning systems*, 30(9):2805–2824.

Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric Xing, Laurent El Ghaoui, and Michael Jordan. 2019. [Theoretically principled trade-off between robustness and accuracy](#). In *Proceedings of the 36th International Conference on Machine Learning*, pages 7472–7482.

Jiliang Zhang and Chen Li. 2019. Adversarial examples: Opportunities and challenges. *IEEE transactions on neural networks and learning systems*, 31(7):2578–2593.

Xiang Zhang, Junbo Jake Zhao, and Yann LeCun. 2015. [Character-level convolutional networks for text classification](#). In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pages 649–657.

Bin Zhu and Yanghui Rao. 2023. [Exploring robust overfitting for pre-trained language models](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 5506–5522, Toronto, Canada. Association for Computational Linguistics.

Chen Zhu, Yu Cheng, Zhe Gan, Siqi Sun, Tom Goldstein, and Jingjing Liu. 2020. [FreeLB: Enhanced adversarial training for natural language understanding](#). In *International Conference on Learning Representations*.

A Statistics and Examples of the Four Datasets

For the SST-2 dataset, an example of x and y is “On the worst revenge-of-the-nerds clichés the filmmakers could dredge up” and “Negative”.

For the IMDB dataset, an example of x and y is “Fred ‘The Hammer’ Williamson delivers another cheaply made movie. He might have set a new standard for himself. Look for the painfully obvious special effects mortar cannon that is visible in the street during a chase scene. You don’t see it just once, you see it several times. Look for the out of focus shot in one scene and the camera operator try to fix it as the scene rolls on. Watch this with a group of people and make your own Mystery Science Theater!” and “Negative”.

For the AGNEWS dataset, an example of x and y is “Wall St. Bears Claw Back Into the Black (Reuters) Reuters - Short-sellers, Wall Street’s dwindling band of ultra-cynics, are seeing green again.” and “Business”.

For the QNLI dataset, an example of x and y is “When did the third Digimon series begin? Unlike the two seasons before it and most of the seasons that followed, Digimon Tamers takes a darker and more realistic approach to its story featuring Digimon who do not reincarnate after their deaths and more complex character development in the original Japanese” and “Not entailment”.

We list the characteristics of the four datasets below.

Dataset	# train	# dev / test	# words
SST-2	67,349	872	17
IMDB	25,000	25,000	201
AG news	120,000	7,600	40
QNLI	105,000	5,460	37

Table 5: Summary of the four datasets.

B Results on More Datasets

We advance the PGD, FreeLB, and TA-VAT methods with our adaptive perturbations and report the results on the IMDB and the AGNEWS datasets in Table 6.

In terms of clean accuracy, our method maintains a performance level comparable to the baseline. In terms of robustness accuracy, our method improves the robustness of the baseline in most scenarios.

It is noteworthy that the improvement in robustness is relatively small on these two datasets. This may be related to the sentence length in the datasets. Existing adversarial attack algorithms typically set the maximum number of word replacements based on a percentage of the sentence’s token count, such as 20%. As the length increases, the number of words to be replaced also increases, which may result in less significant improvements in robustness.

C Relation to Existing Work

Specifically, the standard method is designed to minimize the clean data loss, i.e., the cross-entropy on the clean data. The FGM (Miyato et al., 2017) method generates adversarial examples in one gradient ascent step, minimising both clean and adversarial data loss. The PGD method (Madry et al., 2018) generates adversarial examples using multi-step gradient ascent and only minimizes the adversarial data loss in the last step. Similarly, the FreeLB method (Zhu et al., 2020) generates adversarial examples using multi-step gradient ascent. Different from PGD, FreeLB minimize the average of the adversarial loss at each step.

It is important to point out that all these methods implicitly include the clean data loss in the adversarial loss. In particular, as revealed by Eq. (8) and Eq. (10), the conventional adversarial loss can be decomposed into a clean data loss and an adversarial loss. Therefore, although we can introduce hyperparameters to balance clean loss and adversarial loss in these methods, we cannot precisely balance the two losses.

TRADES (Zhang et al., 2019) is theoretically designed to achieve a good trade-off between accuracy and robustness in the computer vision domain, which is the most relevant AT method with our adaptive AT. TRADES decomposes the adversarial error into a natural error and a boundary error. However, the boundary error cannot be effectively computed. In practice, the authors introduce a surrogate loss (i.e., the KL divergence between the model output of clean data and adversarial data) to approximate the boundary error. In this way, TRADES cannot precisely balance the standard performance and robustness.

Our proposed adaptive AT addresses this issue by decomposing the conventional adversarial loss using Taylor expansion. In our learning objective, clean loss and adversarial loss only affect standard performance and model robustness, respectively.

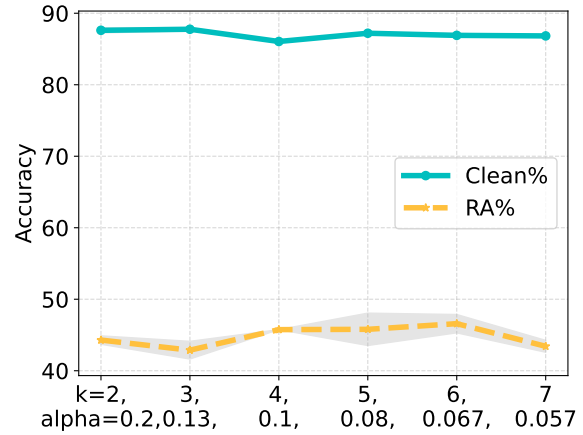


Figure 4: The robust accuracy and clean accuracy under different k and α , while the maximal perturbation size is set to $k\alpha$ following (Li et al., 2021).

D Impact of Adversarial Step Size

We aim to investigate the impact of the perturbation size in AT. In AT, the maximum perturbation size is typically specified. However, what effectively determines the perturbation magnitude are the number of iterations and the adversarial step size. Therefore, given a perturbation size, we vary the number of iterations and step size to investigate their impact on robustness. In other words, we want to find out whether increasing the perturbation strength of adversarial training always helps robustness. We conduct PGD adversarial training on the BERT-base model. Based on our previous experiments, the product of iteration numbers k and adversarial step size α is empirically set to 10 and 0.4.

The main result is reported in Figure 4. It can be seen that when the number of iterations is moderate (5 and 6), the model achieves the best robustness. We suggest that it is unnecessary to set a huge number of iterations during adversarial training. As suggested in (Zhu and Rao, 2023), robust overfitting hinders the AT of NLP models. Too many iterations may lead to robust overfitting of the model and reduce its robustness accuracy on the test set.

E Performance on Other Models

We choose DeBERTa-v3-base (He et al., 2021) and RoBERTa (Liu et al., 2019), two improved versions of BERT, as our backbone models to investigate whether our method can boost the robustness of more complex and larger language models. The clean and robust accuracy of DeBERTa-v3-base and RoBERTa-base models are reported in Table 8.

IMDB	Clean %	TextFooler	TextBugger	BAE
		RA %	RA %	RA %
BERT (Devlin et al., 2019)	91.21	24.48	47.26	20.31
InfoBERT (Wang et al., 2021a)	91.90	23.00	37.30	22.40
Flooding-X (Liu et al., 2022)	92.30	34.50	32.30	35.42
SMART (Jiang et al., 2020)	91.90	24.50	45.40	22.32
TRADES (Zhang et al., 2019)	88.34	25.50	47.60	18.34
PGD (Madry et al., 2018)	90.43	26.31	52.37	21.44
+Ours	90.56 (+0.13)	27.12 (+0.81)	53.50 (+1.13)	21.55 (+0.11)
FreeLB (Zhu et al., 2020)	92.14	27.50	50.60	31.34
+Ours	91.80 (-0.34)	26.82 (-0.68)	52.74 (+2.14)	33.10 (+1.76)
TA-VAT (Li and Qiu, 2021)	91.50	27.40	51.70	23.12
+Ours	92.08 (+0.58)	25.70 (-1.70)	51.66 (-0.04)	24.30 (+1.18)

Table 6: The clean accuracy (“Clean %”) and the robust accuracy (“RA %”) on the IMDB dataset against TextFooler, TextBugger, and BAE attacks. The backbone model is BERT-base. The IMDB dataset does not have a corresponding AdvGLUE version. Therefore, the robustness accuracy for this dataset is not reported.

AGNEWS	Clean %	TextFooler	TextBugger	BAE
		RA %	RA %	RA %
BERT (Devlin et al., 2019)	91.90	20.50	42.71	16.21
InfoBERT (Wang et al., 2021a)	92.00	19.20	31.41	12.70
Flooding-X (Liu et al., 2022)	91.39	33.40	55.60	29.40
SMART (Jiang et al., 2020)	92.20	22.45	37.80	15.60
TRADES (Zhang et al., 2019)	89.42	33.90	48.65	27.61
PGD (Madry et al., 2018)	90.82	37.20	58.20	32.83
+Ours	91.10 (+0.28)	38.70 (+1.50)	57.92 (-0.28)	35.20 (+2.37)
FreeLB (Zhu et al., 2020)	91.20	32.33	48.50	22.65
+Ours	91.07 (-0.13)	32.10 (-0.23)	50.10 (+1.60)	24.12 (+1.47)
TA-VAT (Li and Qiu, 2021)	92.17	39.70	55.81	23.66
+Ours	91.66 (-0.51)	37.26 (-2.44)	57.36 (+1.55)	23.77 (+0.11)

Table 7: The clean accuracy (“Clean %”) and the robust accuracy (“RA %”) on the AGNEWS dataset against TextFooler, TextBugger, and BAE attacks. The backbone model is BERT-base. The AGNEWS dataset does not have a corresponding AdvGLUE version. Therefore, the robustness accuracy for this dataset is not reported.

SST2	Clean %	TextFooler	AdvGLUE
		RA %	RA %
RoBERTa-base	95.07	6.19	39.50
+PGD	94.27	11.47	44.59
+Ours	94.95	11.82	45.22
DeBERTa-v3-base	95.99	12.60	55.41
+PGD	95.18	13.99	57.14
+Ours	95.76	14.50	67.34

Table 8: The clean and robust accuracy on RoBERTa (Liu et al., 2019) and DeBERTa-v3-base (He et al., 2021).

Method	SST-2	QNLI
PGD (Madry et al., 2018)	902	4123
+Ours	912	4237
FreeLB (Zhu et al., 2020)	781	3122
+Ours	920	3745
TA-VAT (Li and Qiu, 2021)	853	3455
+Ours	1013	4123

Table 9: The GPU time consumption (seconds) of training one epoch on the SST-2 and QNLI datasets. The backbone model is BERT-base. The iteration number is set to 5 for all the methods.

These two models can bear a larger perturbation size than the BERT-base model to explore the impact of a larger perturbation range on adversarial training. The empirical results indicate that our adaptive AT can generalize well on larger, more complex models.

F Time Consumption

To further substantiate the comparative advantages of our method, a systematic benchmarking analysis was conducted to evaluate GPU training durations between our proposed approach and established adversarial training methods, with the quantitative comparisons meticulously documented in Table 9. Our method incurs approximately a 10% increase in computational overhead. This empirical investigation demonstrates our method’s computational efficiency while maintaining equivalent adversarial robustness metrics.

G Details on the Usage of DeepSeek

We employ the DeepSeek-r1 model (Guo et al., 2025) for spam detection and evaluate its performance under zero-shot and few-shot settings. In the

zero-shot setting, the model receives no examples or labels and is prompted to classify the message based solely on its inherent reasoning ability. The prompt provided is: “You are a professional spam classifier. Please analyze the following message and determine whether it is spam. Just reply ‘spam’ or ‘ham’, no explanation is needed.” This setup tests the model’s ability to classify messages without prior examples or labels.

In the few-shot setting, we supply the model with two examples and their corresponding labels. The first example is a spam message: “URGENT! This is the 2nd attempt to contact U!U have WON £1000CALL 09071512432 b4 300603t&csBCM4235WC1N3XX.callcost150 pmmobilesvary. max£7.50”, labeled as spam. The second example is a non-spam message: “Why don’t you go tell your friend you’re not sure you want to live with him because he smokes too much then spend hours begging him to come smoke”, labeled as ham. This setting aims to examine how the model leverages the provided examples to classify messages.

Through these two setups, we assess the model’s generalization ability and performance when there are no explicit labels or examples available.

H Implementation Details

We implement PGD (Madry et al., 2018), FreeLB (Zhu et al., 2020), TA-VAT (Li and Qiu, 2021), and InfoBERT (Wang et al., 2021a) based on TextDefender (Li et al., 2021). We implement Flooding-X (Liu et al., 2022), SMART (Jiang et al., 2020), and TRADES (Zhang et al., 2019) following the original paper. The weighting factor α in TRADES is set to 0.5 to achieve the optimal performance. The three adversarial attacks are conducted using TextAttack⁵ (Morris et al., 2020). All experiments are conducted using GeForce RTX 3090 GPUs. All the settings of adversarial hyper-parameters settings are consistent to provide a fair comparison.

Unless otherwise mentioned, the adversarial step size is set to 0.04; the batch size is 128; the epoch number is 10. To align with the weighting factor of the original method, β is set to k for PGD and TA-VAT and $(k + 1)/2$ for FreeLB.

For the natural language inference task, we adhere to prior research (Jin et al., 2020) by allowing the attacking methods to modify the premise while keeping the hypothesis unchanged.

⁵<https://github.com/QData/TextAttack>