

# Rebuild and Ensemble: Exploring Defense Against Text Adversaries

Anonymous ACL submission

## Abstract

Adversarial attacks can mislead strong neural models; as such, in NLP tasks, substitution-based attacks are difficult to defend. Current defense methods usually assume that the substitution candidates are accessible, which cannot be widely applied against adversarial attacks unless knowing the mechanism of the attacks. In this paper, we propose a **Rebuild and Ensemble** Framework to defend against adversarial attacks in texts without knowing the candidates. We propose a rebuild mechanism to train a robust model and ensemble the rebuilt texts during inference to achieve good adversarial defense results. Experiments show that our method can improve accuracy under the current strong attack methods.

## 1 Introduction

Adversarial examples (Goodfellow et al., 2014) can successfully mislead strong neural models in both computer vision tasks (Carlini and Wagner, 2016) and language understanding tasks (Alzantot et al., 2018; Jin et al., 2019). An adversarial example is a maliciously crafted example attached with an imperceptible perturbation and can mislead neural networks. To defend attack examples of images, the most effective method is adversarial training (Goodfellow et al., 2014; Madry et al., 2019) which is a mini-max game used to incorporate perturbations into the training process.

Defending adversarial attacks is extremely important in improving model robustness. However, defending adversarial examples in natural languages is more challenging due to the discrete nature of texts. That is, gradients cannot be used directly in crafting perturbations. The generation process of substitution-based adversarial examples is more complicated than using gradient-based methods in attacking images, making it difficult for neural networks to defend against these substitution-based attacks:

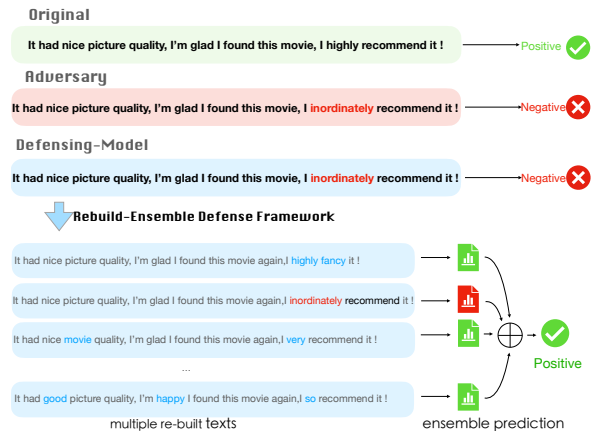


Figure 1: Illustration of Adversarial Defense

(A) The first challenge of defending against adversarial attacks in NLP is that due to the discrete nature, these substitution-based adversarial examples can have substitutes in any token of the sentence and each substitute has a large candidate list. This would cause a combinatorial explosion problem, making it hard to apply adversarial training methods. Strong attacking methods such as Jin et al. (2019) show that using the crafted adversarial examples as data augmentation in adversarial training cannot effectively defend against these substitution-based attacks.

(B) Further, the defending strategies such as adversarial training rely on the assumption that the candidate lists of the substitutions are accessible. However, the candidate lists of the substitutions should **not** be exposed to the target model; that is, the target model should be unfamiliar to the candidate list of the adversarial examples. In real-world defense systems, the defender is not aware of the strategy the potential attacks might use, so the assumption that the candidate list is available would significantly constrain the potential applications of these defending methods.

In this work, we propose a strong defense framework, i.e., **Rebuild and Ensemble**.

067					
068					
069					
070					
071					
072					
073					
074					
075					
076					
077					
078					
079					
080					
081					
082					
083					
084					
085					
086					
087					
088					
089					
090					
091					
092					
093					
094					
095					
096					
097					
098					
099					
100					
101					
102					
103					
104					
105					
106					
107					
108					
109					
110					
111					
112					
113					
114					
115					
116					
117					

We aim to construct a defense system that can successfully defend the attacks without knowing the attack range (that is, the candidate list in the substitution-based attacks). As seen in Figure 1, we first reconstruct input samples to samples that do not have adversarial effects. Therefore, when the input is changed by the adversarial attack, we make predictions based on the rebuilt texts which will results in correct predictions.

To achieve this goal, we first reconsider the widely applied pre-trained models (e.g. BERT (Devlin et al., 2018)) which introduce the masked language modeling task in the pre-training stage and can be used in fine-tuning on downstream tasks. During downstream task fine-tuning, these pre-train models focus on making downstream task predictions without maintaining the mask language modeling ability. Instead of simply fine-tuning downstream tasks, we keep the mask prediction ability during fine-tuning, and use this ability to process the rebuilding of input texts. Specifically, we random mask the input texts and use the mask prediction to rebuild a text that does not have adversarial affect. Intuitively, the rebuild process introduces randomness since the masks are randomly selected. We can make multiple random rebuilt texts and apply an ensemble process to obtain the final model output predictions for better robustness. To train the defending framework, we introduce the rebuild training based on adversarial training with virtual adversaries (Zhu et al., 2019; Li and Qiu, 2020) which could enhance both rebuilding and downstream task predicting abilities.

Through extensive experiments, we prove that the proposed defense framework can successfully resist strong attacks such as Textfooler and BERT-Attack. Experiment results show that the accuracy under attack in baseline defense methods is lower than random guesses, while ours can lift the performance to only a few percent lower than the original accuracy when the candidates are limited. Further, extensive results indicate that the candidate size of the attacker score is essential for successful attacks, which is a key factor in maintaining semantics of the adversaries. Therefore we also recommend that future attacking methods can focus on achieving success attacks with tighter constrains.

To summarize our contributions:

(1) We raise the concern of defending substitution-based adversarial attacks without knowing the candidates of the attacks in NLP tasks.

(2) We propose a Rebuild and Ensemble framework to defend against recently introduced attack methods without knowing the candidates and experiments prove the effectiveness of the framework.

(3) We explore the key factors in defending against score-based attacks and recommend further research to focus on tighter constraint attacks.

## 2 Related Work

### 2.1 Adversarial Attacks in NLP

In NLP tasks, current methods use substitution-based strategies (Alzantot et al., 2018; Jin et al., 2019; Ren et al., 2019) to craft adversarial examples. Most works focus on the score-based black-box attack, that is, attacking methods know the logits of the output prediction. These methods use different strategies (Yoo et al., 2020; Morris et al., 2020b) to find words to replace, such as genetic algorithm (Alzantot et al., 2018), greedy-search (Jin et al., 2019; Li et al., 2020) or gradient-based methods (Ebrahimi et al., 2017; Cheng et al., 2019) and get substitutes using synonyms (Jin et al., 2019; Mrkšić et al., 2016; Ren et al., 2019) or language models (Li et al., 2020; Garg and Ramakrishnan, 2020; Shi et al., 2019).

### 2.2 Adversarial Defenses

We divide the defense methods for substitution attacks by whether the defense method requires knowledge of the candidate of the attack.

To defend adversarial attacks without knowing the candidate knowledge, Samangouei et al. (2018) uses a defensive GAN framework to build clean images to avoid adversarial attacks; Xie et al. (2017) introduces randomness into the model predicting process to mitigate adversarial affect. Similar to using multiple rebuilt texts, Federici et al. (2020) introduces a multi-view approach that improve robustness by using a set of images describing the same object. Ebrahimi et al. (2017); Cheng et al. (2019) introduces gradient-based adversarial training that crafts adversarial samples by finding the most similar word embeddings based on the gradients. Further, gradient-based adversarial training with virtual adversaries could also be used in NLP tasks: Miyato et al. (2016) proposes a virtual adversarial training process with virtual inputs and labels for semi-supervised tasks. Zhu et al. (2019); Li and Qiu (2020) incorporate gradients to crafting virtual adversaries to improve generalization ability.

To defend against adversaries while knowing

the candidate list of the attacks, augmentation-based methods are the most direct defense strategies that use the generated adversaries to train a robust model (Jin et al., 2019; Li et al., 2020; Si et al., 2020). Jia et al. (2019); Huang et al. (2019) introduce a certified robust model to defend against adversarial attacks by constructing a certified space that can tolerate substitutes. Zhou et al. (2020); Dong et al. (2021) construct a convex hull based on the candidate list which can resist substitutions in the candidate list. Zhou et al. (2019) incorporates the idea of blocking adversarial attacks by discriminating perturbations in the input texts.

### 3 Rebuild And Ensemble as Defense

Defending against adversarial attacks without accessing the candidate list is more applicable in real-world adversarial defenses. Therefore, we introduce **Rebuild and Ensemble** as an effective framework to defend strong adversarial attacks exemplified by substitution-based attacks in NLP without knowing the candidate list of substitutions.

We suppose that the target model that may face adversarial attacks is a fine-tuned classification model  $F_c(\cdot)$ . When given an input sentence  $X$ , the adversarial attack may craft an adversarial example  $X_{adv}$  that replaces a small proportion of tokens with similar texts. We only consider substitution-based adversaries since defending other types of adversarial examples such as token insertion or deletion is the same as defending substitution-based adversaries.

#### 3.1 Rebuild and Ensemble Framework

We propose the rebuild and ensemble framework that first rebuilds multiple texts from the input text and then use these rebuilt texts to make predictions. We used the same model  $F(\cdot)$  that can rebuild input texts and make predictions using a multi-task structure. We use  $F_m(\cdot)$  to denote the mask prediction task that rebuilds the input texts and use  $F_c(\cdot)$  to denote the classification task. As seen in Figure 2, when given an input text  $X$  that might have been attacked, we randomly mask the input texts or insert additional masks to make  $N$  copies of noisy input  $\tilde{X}_i = [w_0, \dots, [\text{MASK}], w_n, \dots]$ . We use two simple strategies to inject noise into the input texts: (1) Randomly mask the input texts; (2) Randomly insert masks into the input texts.

After making multiple noisy inputs, we can run the rebuild process first to get the rebuilt texts based

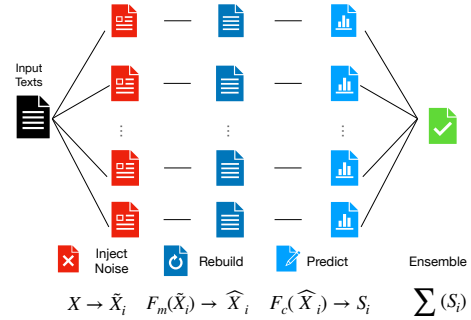


Figure 2: Rebuild And Ensemble Process: after noise injection we rebuild multiple texts. Then we use these texts to predict the label and ensemble the scores as the final output.

on the randomly masked inputs  $\tilde{X}$ :  $\hat{X}_i = F_m(\tilde{X}_i)$ .

Then we feed the rebuilt texts through the classifier  $F_c(\cdot)$  to calculate the final output score based on the multiple rebuilt texts:

$$S_i = \frac{1}{N} \sum_{i=0}^N \left( \text{Softmax}(F_c(\hat{X}_i)) \right) \quad (1)$$

Here, we use the average score from multiple rebuilt texts predictions as the final output score given to the score-based adversarial attackers.

Another advantage of using the mask prediction ability is that the mask-infill ability is trained by massive data pre-training which can be helpful in building models with better generalization ability (Gururangan et al., 2020). Therefore, keeping the mask prediction ability and utilizing it can make better use of the pre-trained knowledge.

#### 3.2 Rebuild Framework Training

We use the fine-tuned masked language model while maintaining the masked language modeling ability since we believe that (1) rebuild process can help gain better robustness by mitigating the adversarial affect in the input sequences; (2) maintaining language modeling information helps improve model robustness in the classification process.

In order to fine-tune such a model  $F$  with parameter  $\theta$  containing two functions  $F_m(\cdot)$  and  $F_c(\cdot)$ , we introduce a rebuild training process based on multi-task adversarial training. We use noisy texts as inputs to train the masked language modeling task and the downstream task fine-tuning simultaneously so that the fine-tuning process can tolerate more noisy texts since the model might be attacked by adversaries.

### 3.2.1 Masked LM Training Strategy

In our model’s fine-tuning, we have both the masked language modeling training and the downstream task training. In the masked language model training, we also incorporate the gradient information in the rebuild training process to build a gradient-based noisy data to enhance the rebuilding ability.

Therefore we have two language model training strategy: (1) Standard [MASK] Prediction: We randomly mask the input texts (15% of the pieces) and make the masked language model to further pre-train on the training dataset. (2) Gradient-Noise Rebuild: Previous pre-training process does not calculate loss on un-masked tokens. Instead we use a gradient-based adversarial training method to add perturbation  $\delta$  on the embedding space of these un-masked tokens and calculate the loss of the masked language model task on these tokens to make the model aware of the potential substitutes.

Compared with Gururangan et al. (2020) which also introduces MLM task in fine-tuning, we use the mask-infill ability of the model to rebuild potential sabotaged inputs. That is, the MLM task used in Gururangan et al. (2020) is an auxiliary task to the fine-tuning loss while in our rebuild training, the combination of these two losses constructs a multi-task model and the mask-infill ability is fully utilized.

### 3.2.2 Preliminary of Adversarial Training with Virtual Adversaries

Recent researches have been focusing on exploring the possibility of using gradient-based virtual adversaries in NLP tasks (Zhu et al., 2019; Li and Qiu, 2020). The core idea is that the adversarial examples are not real substitutions but virtual adversaries added to the embedding space.<sup>1</sup>

$$g(\delta) = \nabla_{\delta} L(f_{\theta}(X + \delta), y) \quad (2)$$

$$\delta = \prod_{\|\delta\|_F \leq \epsilon} \frac{\alpha g(\delta)}{\|g(\delta)\|_F} \quad (3)$$

Here  $\prod_{\|\delta\|_F \leq \epsilon}$  represents the process that projects the perturbation onto the normalization ball  $\epsilon$  using

<sup>1</sup>different from VAT which uses both virtual inputs and virtual labels, virtual adversaries are deployed in supervised tasks as a replacement for real-substitute adversarial training since texts are discrete and gradients cannot be directly added to the texts.

Frobenius normalization  $\|\delta\|_F$ . We update the perturbation using a certain adversarial learning rate  $\alpha$ .  $X$  is the word embedding of input sequence  $[w_0, \dots, w_n]$ . Then these virtual adversaries are used in the training process to improve model performance. The entire process is to minimize the maximum risk of mis-classification, containing a multi-step (e.g.  $T$  steps) iteration to obtain the proper perturbations while in the FreeLB algorithm, the gradients obtained in each iteration are used in the final optimization.

---

#### Algorithm 1 Rebuild Training

---

**Require:** Training Sample  $X$ , Uniform Noist  $U$  with range  $\sigma$ , adversarial step  $T$

- 1:  $\tilde{X} \leftarrow$  Random Mask  $X$
- 2:  $\delta_0 \leftarrow \frac{1}{\sqrt{D}} U(-\sigma, \sigma)$  // Init Perturb
- 3: **for**  $t = 0, 1, \dots, T$  **do**
- 4:    $\mathcal{L}_c \leftarrow$  Using Equation 4
- 5:    $\mathcal{L}_{mlm} \leftarrow$  Using Equation 5
- 6:   // Get Perturbation
- 7:    $g_{\delta} \leftarrow \nabla_{\delta} (\mathcal{L}_c + \mathcal{L}_{mlm})$
- 8:    $\delta_{t+1} \leftarrow \prod_{\|\delta\|_F < \epsilon} (\delta_t + \alpha \cdot g_{\delta} / \|g_{\delta}\|_F)$
- 9:   // Rebuild with Noise
- 10:    $\mathcal{L}_{noise} \leftarrow$  Using Equation 7
- 11:    $\tilde{X} \leftarrow \tilde{X} + \delta_t$  // Update Input
- 12:    $g_{t+1} = g_t + \nabla_{\theta} (\mathcal{L}_c + \mathcal{L}_{mlm} + \mathcal{L}_{noise})$
- 13:  $\theta \leftarrow \theta - g_{T+1}$  // Update model parameter  $\theta$

---

### 3.2.3 Overall Process of Rebuild Training

Given input texts  $X$ , we first make noisy copies  $\tilde{X}$ , for notation convenience, here  $X$  and  $\tilde{X}$  are the embedding output of the input texts. Then we can calculate the gradients of the fine-tuning classification task  $g_c$  as well as the mask-prediction task  $g_{mlm}$ .

$$\mathcal{L}_c = L(F_c(\tilde{X}), y, \theta) + L(F_c(X), y, \theta) \quad (4)$$

$$\mathcal{L}_{mlm} = L(F_m(\tilde{X}), X, \theta) \quad (5)$$

Here,  $L$  is the cross entropy loss function for both masked language model task  $\mathcal{L}_{mlm}$  and classification task  $\mathcal{L}_c$ . As seen in Algorithm 1 line 7, we run the fine-tuning process based on the noisy input and the original input and we run the mask prediction task simultaneously. We assume that with the mask prediction task also involved in fine-tuning, the model will not be focusing on fitting the

classification task only, which can help maintain the entire semantic information and mitigate the adversarial affect from the adversaries.

Further, we use gradients to craft virtual adversaries  $\delta$  and calculate loss based on these adversaries  $\mathcal{L}_{noise}$ :

$$\delta \leftarrow \prod_{\|\delta\|_F < \epsilon} (\delta + \alpha \cdot \mathbf{g}_\delta / \|\mathbf{g}_\delta\|_F) \quad (6)$$

$$\mathcal{L}_{noise} = L(F_m(\tilde{X} + \delta), X, \theta) \quad (7)$$

Here the cross entropy loss  $L$  is calculated based on all tokens not just the masked ones. Therefore, the masked language model prediction task is modified to make the model tolerate more noises and therefore more robust.

The difference between our rebuild-training and traditional adversarial training is that we allow the perturbations to be **larger** than previous works. That is, the adversarial learning rate  $\alpha$  and the perturbation boundary  $\epsilon$  are larger (e.g. norm bound set to  $\epsilon 2e-1$  compared with  $1e-2$  used in the FreeLB and TAVAT method). Therefore, some of the tokens are seriously affected by gradients, which is an effective method for further pre-training the model to tolerate adversaries. We calculate all the losses of prediction task, rebuild task and gradient-based noise rebuild task and update the model parameter.

## 4 Experiments

### 4.1 Datasets

We use two widely used text classification datasets: IMDB <sup>2</sup> (Maas et al., 2011) and AG’s News <sup>3</sup> (Zhang et al., 2015) in our experiments. The IMDB dataset is a bi-polar movie review classification task; the AG’s News dataset is a four-class news genre classification task. The average length is 220 words in the IMDB dataset, and 40 words in the AG’s News dataset. We use the test set following the Textfooler 1k test set in the main result and sample 100 samples for the rest of the experiments since the attacking process is seriously slowed down when the model is defensive.

### 4.2 Attack Methods

Popular attack methods exemplified by genetic Algorithm (Alzantot et al., 2018), Textfooler (Jin et al., 2019) and BERT-Attack (Li et al., 2020) can

<sup>2</sup><https://datasets.imdbws.com/>

<sup>3</sup><https://www.kaggle.com/amananandrai/ag-news-classification-dataset>

successfully mislead strong models of both IMDB and AG’s News task with a very small percentage of substitutions. Therefore, we use these strong adversarial attack methods as the attacker to test the effectiveness of our defense method. The hyper parameters used in the attacking algorithm vary in different settings: we choose candidate list size  $K$  to be 12, 48, 50 typically which are used in the Textfooler and BERT-Attack methods.

We use the exact same metric used in Textfooler and BERT-Attack that calculate the after-attack accuracy, which is the targeted adversarial evaluation defined by Si et al. (2020). The after-attack accuracy measures the actual defense ability of the system under adversarial attacks.

### 4.3 Victim Models and Defense Baselines

The victim models are the fine-tuned pre-train models exemplified by BERT and RoBERTa, which we implement based on Huggingface Transformers <sup>4</sup> (Wolf et al., 2020). As discussed above, there are few works concerning adversarial defenses against attacks without knowing the candidates in NLP tasks. Moreover, previous works do not focus on recent strong attack algorithms such as Textfooler (Jin et al., 2019), BERT-involved attacks (Li et al., 2020; Garg and Ramakrishnan, 2020) Therefore, we first list methods that can defend adversarial attacks without accessing the candidate list as our baselines:

**Adv-Train (Adv-HotFlip):** Ebrahimi et al. (2017) introduces the adversarial training method used in defending against substitution-based adversarial attacks in NLP. It uses gradients to find actual adversaries in the embedding space.

**Virtual-Adv-Train (TAVAT):** Token-Aware VAT (Li and Qiu, 2020) use virtual adversaries (Zhu et al., 2019) to improve the performances in fine-tuning pre-trained models, which can also be used to deal with adversarial attacks without accessing the candidate list. We follow the standard TAVAT training process to re-implement the defense results.

Further, there are some works that require candidate list, it is not a fair comparison with defense methods without accessing the candidates, so we list them separately:

**Adv-Augmentation:** We generate adversarial examples of the training dataset as a data augmentation method. We mix the generated adversarial

<sup>4</sup><https://github.com/huggingface/transformers>

Methods	Origin	Textfooler $K = 12$	BERT-Attack $K = 12$	Textfooler $K = 50$	BERT-Attack $K = 48$
<b>IMDB</b>					
BERT (Devlin et al., 2018)	94.1	20.4	18.5	2.8	3.2
RoBERTa (Liu et al., 2019)	97.3	26.3	24.5	25.2	23.0
Adv-HotFlip (BERT) (Ebrahimi et al., 2017)	95.1	36.1	34.2	8.0	6.2
TAVAT (BERT) (Li and Qiu, 2020)	96.0	30.2	30.4	7.3	2.3
RanMASK (RoBERTa) (Zeng et al., 2021)	93.0	-	-	23.7	26.8
FreeeLB++ (BERT) (Li et al., 2021)	93.2	-	-	45.3	39.9
Rebuild & Ensemble (BERT)	93.0	<b>81.5</b>	<b>76.7</b>	<b>51.0</b>	<b>44.5</b>
Rebuild & Ensemble (RoBERTa)	96.1	<b>84.2</b>	<b>82.0</b>	<b>54.3</b>	<b>52.2</b>
<b>AG's News</b>					
BERT	92.0	32.8	34.3	19.4	14.1
RoBERTa	90.1	29.5	30.4	17.9	13.0
Adv-HotFlip (BERT)	91.2	35.3	34.1	18.2	8.5
TAVAT (BERT)	90.5	40.1	34.2	20.1	8.5
Rebuild & Ensemble (BERT)	90.6	<b>61.5</b>	<b>49.7</b>	<b>34.9</b>	<b>22.5</b>
Rebuild & Ensemble (RoBERTa)	90.8	<b>59.1</b>	<b>41.2</b>	<b>34.2</b>	<b>19.5</b>

Table 1: After-Attack Accuracy compared with defense methods that can defend attacks without accessing the candidate list of the attacks.

Methods	Origin	Textfooler	genetic
<b>IMDB</b>			
BERT	94.0	2.0	45.0
Data-Augmentation	93.0	18.0	53.0
ADA* (Si et al., 2020)	96.7	3.0	-
AMDA-SMix*(Si et al., 2020)	96.9	17.4	-
ASCC (Dong et al., 2021)	77.0	-	71.0
R & E	93.0	<b>51.0</b>	<b>79.0</b>

Table 2: After-Attack Accuracy compared with previous access-candidates methods based on BERT model. - means that the results are not reported in the corresponding papers. Here we implement Textfooler with  $K = 50$  for consistency with previous works. \* represents that ADA uses a selected subset of the dataset that may have a difference in the results.

examples and the original training dataset to train a model in a standard fine-tuning process.

**ASCC:** Dong et al. (2021) also uses a convex-hull concept based on the candidate vocabulary as strong adversarial defense.

**ADA:** Si et al. (2020) uses a mixup-strategy based on the generated adversarial examples to achieve adversarial defense with variants AMDA-SMix that mixup the special tokens.

**FreeLB++:** Li et al. (2021) introduces a variant of FreeLB method that expands the norm bound which is similar to the larger bound in the rebuild training process.

**RanMASK:** Zeng et al. (2021) introduces a masking strategy that makes use of noises to improve robustness.

## 4.4 Implementations

We use BERT-BASE and RoBERTa-BASE models based on the Huggingface Transformers<sup>5</sup>. We modify the adversarial training with virtual adversaries based on the implementation of FreeLB and TAVAT<sup>6</sup>. The training hyper-parameters we use is different from FreeLB and TAVAT, since we aim to find large perturbations to simulate adversaries. We set adversarial learning rate  $\alpha = 1e-1$  to and normalization boundary  $\epsilon = 2e-1$  in all tasks. We set the ensemble size  $N = 16$  for all tasks and we will discuss the selection of  $N$  in the later section.

We use the TextAttack toolkit as well as the official code to implement adversarial attack methods<sup>7</sup> (Morris et al., 2020a). The similarity thresholds are the main factors of the attacking algorithm. We tune the USE (Cer et al., 2018) constraint 0.5 for the AG task and 0.7 for the IMDB task and 0.5 for the cosine-similarity threshold of the synonyms embedding (Mrkšić et al., 2016) which can re-produce the results of the attacking methods reported.

## 4.5 Results

As seen in Table 1, the proposed **Rebuild and Ensemble** framework can successfully defend strong attack methods. The accuracy of our defending method under attack is significantly higher than non-defense models (50% vs 20% in the IMDB dataset). Compared with previous defense meth-

<sup>5</sup><https://github.com/huggingface/transformers>

<sup>6</sup><https://github.com/LinyangLee/Token-Aware-VAT>

<sup>7</sup><https://github.com/QData/TextAttack>

Different Settings of R & E					Origin	Textfooler( $K=12$ )	BERT-Atk( $K=12$ )
Train		Inference					
Joint	VAT	Ensemble	Rebuild	Insert			
✓	✓	✓	✓	✓	93.0	<b>86.0</b>	<b>77.0</b>
✓	✓		✓	✓	93.0	63.0	52.0
✓	✓		✓		93.0	42.0	29.0
✓			✓	✓	95.0	45.0	34.0
✓			✓		95.0	29.0	17.0
		✓	✓	✓	94.0	72.0	60.0
			✓	✓	87.0	20.0	13.0
			✓		92.0	11.0	3.0
		✓			<b>96.0</b>	75.0	62.0
-	-	-	-	-	93.0	20.0	18.0

Table 3: Ablations results tested on attacking the IMDB task based on BERT models. Joint is the multi-task training in Algorithm 1 line 12; VAT is the adversarial training process; Ensemble is whether using multiple texts during inference; Insert is whether the rebuild process contains insert and replace.

ods, our proposed method can achieve higher defense accuracy in both IMDB task and AG’s News task. The Adv-HotFlip and the TAVAT methods are effective, which indicates that gradient-based adversaries are not very similar with actual substitutions. We can see that Adv-HotFlip and TAVAT methods achieve similar results (around 30% when  $K = 12$ ) which indicates that gradient-based adversarial training methods have similar defense ability no matter the adversaries are virtual or real since they are both unaware of the attacker’s candidate list. Also, the original accuracy (on the clean data) of our method is only a little lower (less than 2% ) than the baseline methods, which indicates that the defensive rebuild and ensemble strategy does not hurt the performances. The RoBERTa model also shows robustness using both original fine-tuned model and our defensive framework, which indicates our defending strategy can be used in various pre-trained language models. Compared with methods that specifically focus on adversarial defense, our proposed method can still surpass state-of-the-arts defense system FreeLB++ (Li et al., 2021) and RanMASK (Zeng et al., 2021) by over 5%.

Further, the candidate size is extremely important in defending adversarial attacks, when the candidate size is smaller, exemplified by  $K = 12$ , our method can achieve very promising results. As pointed out by Morris et al. (2020b), the candidate size should not be too large that the quality of the adversarial examples is largely damaged.

As seen in Table 2, we compare our method with previous access-candidates defense methods. When defending against the widely used Textfooler

attack and genetic attack (Alzantot et al., 2018), our method can achieve similar accuracy even compared with known-candidates defense methods. As seen, data augmentation method cannot significantly improve model robustness since the candidates can be very diversified. Therefore, using generated adversarial samples as an augmentation strategy does not guarantee robustness against greedy-searched methods like Textfooler and BERT-Attack.

## 4.6 Analysis

### 4.6.1 Ablations

We run extensive ablation experiments to explore the working mechanism in defending adversaries. We run ablations in two parts: (1) using the rebuild-trained model; (2) using the ensemble inference without training the model specifically.

Firstly, we test the model robustness without using ensemble inference, that is, during inference, the ensemble size  $N$  is 1: We explore the effectiveness of incorporating the gradient-noise rebuild process. Also, we test the result of using the mask and rebuild strategy as well as the insert and rebuild strategy. Then we test the inference process: We use the fine-tuned model and the original masked language model as the prediction model and the rebuild model to run inference. We test the effectiveness of making multiple copies of rebuilt texts; We also explore how the two operations: mask and insert work during inference.

As seen in the Table 3, we could explore the working mechanism in defending against the attacks via extensive results.

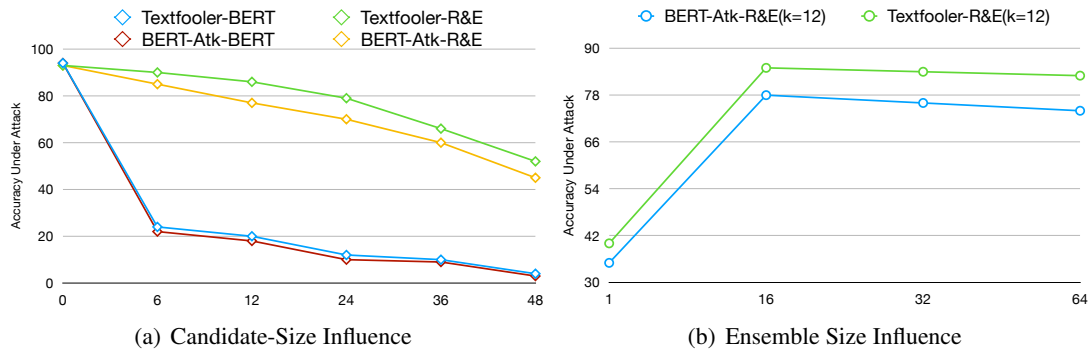


Figure 3: Hyper-Parameter Selection Analysis

The observations indicate that:

(a) Rebuild Train is effective: The process in rebuild training allows the trained model to be aware of both the missing texts that need rebuilding and the classification labels of the inputs, which is helpful in rebuilding classification-aware texts. Without the rebuild trained model, the accuracy is even lower when rebuilding with the original masked language model during ensemble inference. However, rebuilding using the original MLM is not very much helpful, which indicates that the model trained with re-building process is important.

(b) Ensemble during inference is important: As seen, with the ensemble strategy, even random masking with an ensemble process can be helpful.

(c) Gradient-Noise Rebuild is helpful: without the gradient-noise rebuild process, the model can still defend adversaries.

#### 4.6.2 Candidate Size Analysis

One key problem is that these attacking algorithms use a very large candidate size  $K$  with a default set to around 50, which seriously harms the quality of the input texts. Candidate size is the possible candidates for replacement in every token in certain attacking methods such as BERT-Attack and Textfooler.

As seen in Fig. 3 (a), when the candidate is 0, the accuracy is high on the clean samples. When the candidate is 6, the normal fine-tuned BERT model cannot correctly predict the generated adversarial examples. This indicates that normal fine-tuned BERT is not robust even when the candidate size is small, while our approach can tolerate these limited candidate size attacks. When the candidate size grows, the performance of our defense framework drops by a relatively large margin. We assume that large candidate size would seriously harm the

semantics which is also explored in Morris et al. (2020b), while these adversaries cannot be well evaluated even using human-evaluations since the change rate is still low.

#### 4.6.3 Ensemble Strategy Analysis

One key problem is that how many copies we should use in the rebuilding process, since during inference, it is also important to maintain high efficiency. We use two attack methods with  $K = 12$  to test how the accuracy varies when using different ensemble size  $N$ .

As seen in Fig. 3 (b), the ensemble size is actually not a key factor. Larger ensemble size would not result in further improvements. We assume that larger ensemble size will *smooth* the output score which will benefit the attack algorithm. When the number of rebuild is not large, the inference efficiency is bearable.

### 5 Conclusion and Future Work

In this paper, we introduce a novel rebuild and ensemble defense framework against current strong adversarial attacks. We utilize the mask-infill ability of pre-trained models to first rebuild texts and use these texts with less adversarial effect to make predictions for better robustness. The rebuild training can improve the model robustness since it maintains more semantic information while it also introduces a rebuild text process. The proposed ensemble inference is also effective indicating that the multiple rebuilt texts are better than one. Experiments show that these proposed components can work coordinately to achieve strong defense performance. We are hoping such a defense process can provide hints for future works on adversarial defenses.



593  
594  
595  
596  
597  
  
598  
599  
600  
  
601  
602  
603  
604  
605  
  
606  
607  
608  
  
609  
610  
611  
612  
  
613  
614  
615  
616  
  
617  
618  
619  
620  
  
621  
622  
623  
624  
  
625  
626  
627  
  
628  
629  
630  
  
631  
632  
633  
634  
635  
  
636  
637  
638  
639  
640  
641  
  
642  
643  
644

## References

Moustafa Alzantot, Yash Sharma, Ahmed Elgohary, Bo-Jhang Ho, Mani B. Srivastava, and Kai-Wei Chang. 2018. [Generating natural language adversarial examples](#). *CoRR*, abs/1804.07998.

Nicholas Carlini and David A. Wagner. 2016. [Towards evaluating the robustness of neural networks](#). *CoRR*, abs/1608.04644.

Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St John, Noah Constant, Mario Guajardo-Cespedes, Steve Yuan, Chris Tar, et al. 2018. Universal sentence encoder. *arXiv preprint arXiv:1803.11175*.

Yong Cheng, Lu Jiang, and Wolfgang Macherey. 2019. Robust neural machine translation with doubly adversarial inputs. *arXiv preprint arXiv:1906.02443*.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. [BERT: pre-training of deep bidirectional transformers for language understanding](#). *CoRR*, abs/1810.04805.

Xinshuai Dong, Hong Liu, Rongrong Ji, and Anh Tuan Luu. 2021. [Towards robustness against natural language word substitutions](#). In *International Conference on Learning Representations*.

Javid Ebrahimi, Anyi Rao, Daniel Lowd, and De-jing Dou. 2017. Hotflip: White-box adversarial examples for text classification. *arXiv preprint arXiv:1712.06751*.

Marco Federici, Anjan Dutta, Patrick Forré, Nate Kushman, and Zeynep Akata. 2020. Learning robust representations via multi-view information bottleneck. *arXiv preprint arXiv:2002.07017*.

Siddhant Garg and Goutham Ramakrishnan. 2020. Bae: Bert-based adversarial examples for text classification. *arXiv preprint arXiv:2004.01970*.

Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. 2014. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*.

Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A Smith. 2020. Don’t stop pretraining: adapt language models to domains and tasks. *arXiv preprint arXiv:2004.10964*.

Po-Sen Huang, Robert Stanforth, Johannes Welbl, Chris Dyer, Dani Yogatama, Sven Gowal, Krishnamurthy Dvijotham, and Pushmeet Kohli. 2019. Achieving verified robustness to symbol substitutions via interval bound propagation. *arXiv preprint arXiv:1909.01492*.

Robin Jia, Aditi Raghunathan, Kerem Göksel, and Percy Liang. 2019. [Certified robustness to adversarial word substitutions](#). *CoRR*, abs/1909.00986.

Di Jin, Zhijing Jin, Joey Tianyi Zhou, and Peter Szolovits. 2019. [Is BERT really robust? natural language attack on text classification and entailment](#). *CoRR*, abs/1907.11932.

Linyang Li, Ruotian Ma, Qipeng Guo, Xiangyang Xue, and Xipeng Qiu. 2020. Bert-attack: Adversarial attack against bert using bert. *arXiv preprint arXiv:2004.09984*.

Linyang Li and Xipeng Qiu. 2020. Textat: Adversarial training for natural language understanding with token-level perturbation. *arXiv preprint arXiv:2004.14543*.

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. *arXiv preprint arXiv:2108.12777*.

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*.

Andrew Maas, Raymond E Daly, Peter T Pham, Dan Huang, Andrew Y Ng, and Christopher Potts. 2011. Learning word vectors for sentiment analysis. In *Proceedings of the 49th annual meeting of the association for computational linguistics: Human language technologies*, pages 142–150.

Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. 2019. [Towards deep learning models resistant to adversarial attacks](#).

Takeru Miyato, Andrew M. Dai, and Ian J. Goodfellow. 2016. Virtual adversarial training for semi-supervised text classification. *ArXiv*, abs/1605.07725.

John Morris, Eli Lifland, Jin Yong Yoo, Jake Grigsby, Di Jin, and Yanjun Qi. 2020a. Textattack: A framework for adversarial attacks, data augmentation, and adversarial training in nlp. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 119–126.

John X. Morris, Eli Lifland, Jack Lanchantin, Yangfeng Ji, and Yanjun Qi. 2020b. Reevaluating adversarial examples in natural language. In *ArXiv*, volume abs/2004.14174.

Nikola Mrkšić, Diarmuid O Séaghdha, Blaise Thomson, Milica Gašić, Lina Rojas-Barahona, Pei-Hao Su, David Vandyke, Tsung-Hsien Wen, and Steve Young. 2016. Counter-fitting word vectors to linguistic constraints. *arXiv preprint arXiv:1603.00892*.

Shuhuai Ren, Yihe Deng, Kun He, and Wanxiang Che. 2019. Generating natural language adversarial examples through probability weighted word saliency. In

- 700 *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1085–1097.
- 701
- 702
- 703 Pouya Samangouei, Maya Kabkab, and Rama Chel-  
704 lappa. 2018. [Defense-gan: Protecting classifiers](#)  
705 [against adversarial attacks using generative models](#).  
706 *CoRR*, abs/1805.06605.
- 707 Zhouxing Shi, Minlie Huang, Ting Yao, and Jing-  
708 fang Xu. 2019. [Robustness to modification with](#)  
709 [shared words in paraphrase identification](#). *CoRR*,  
710 abs/1909.02560.
- 711 Chenglei Si, Zhengyan Zhang, Fanchao Qi, Zhiyuan  
712 Liu, Yasheng Wang, Qun Liu, and Maosong Sun.  
713 2020. Better robustness by more coverage: Adver-  
714 sarial training with mixup augmentation for robust  
715 fine-tuning. *arXiv preprint arXiv:2012.15699*.
- 716 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien  
717 Chaumond, Clement Delangue, Anthony Moi, Pier-  
718 ric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz,  
719 Joe Davison, Sam Shleifer, Patrick von Platen, Clara  
720 Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le  
721 Scao, Sylvain Gugger, Mariama Drame, Quentin  
722 Lhoest, and Alexander M. Rush. 2020. [Transformers: State-of-the-art natural language processing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.
- 723
- 724
- 725
- 726
- 727
- 728 Cihang Xie, Jianyu Wang, Zhishuai Zhang, Zhou  
729 Ren, and Alan Yuille. 2017. Mitigating adversarial  
730 effects through randomization. *arXiv preprint*  
731 *arXiv:1711.01991*.
- 732 Jin Yong Yoo, John X. Morris, Eli Lifland, and Yanjun  
733 Qi. 2020. Searching for a search method: Bench-  
734 marking search algorithms for generating nlp adver-  
735 sarial examples. *ArXiv*, abs/2009.06368.
- 736 Jiehang Zeng, Xiaoqing Zheng, Jianhan Xu, Linyang Li,  
737 Liping Yuan, and Xuanjing Huang. 2021. Certified  
738 robustness to text adversarial attacks by randomized  
739 [mask]. *arXiv preprint arXiv:2105.03743*.
- 740 Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015.  
741 Character-level convolutional networks for text clas-  
742 sification. In *Advances in neural information pro-  
743 cessing systems*, pages 649–657.
- 744 Yi Zhou, Xiaoqing Zheng, Cho-Jui Hsieh, Kai-wei  
745 Chang, and Xuanjing Huang. 2020. Defense against  
746 adversarial attacks in nlp via dirichlet neighborhood  
747 ensemble. *arXiv preprint arXiv:2006.11627*.
- 748 Yichao Zhou, Jyun-Yu Jiang, Kai-Wei Chang, and Wei  
749 Wang. 2019. [Learning to discriminate perturbations](#)  
750 [for blocking adversarial attacks in text classification](#).  
751 In *Proceedings of the 2019 Conference on Empirical*  
752 *Methods in Natural Language Processing and the*  
753 *9th International Joint Conference on Natural Lan-  
754 guage Processing (EMNLP-IJCNLP)*, pages 4904–  
755 4913, Hong Kong, China. Association for Computa-  
756 tional Linguistics.
- Chen Zhu, Yu Cheng, Zhe Gan, Siqi Sun, Thomas Gold-  
stein, and Jingjing Liu. 2019. [Freelb: Enhanced ad-  
versarial training for language understanding](#). *arXiv*  
*preprint arXiv:1909.11764*.
- 757
- 758
- 759
- 760

	Texts	Confidence
<b>R&amp;E Successful Defense</b>		
Clean-Sample	I have the good common logical sense to know that oil can not last forever and I am acutely aware of how much of my life in the suburbs revolves around petrochemical products . I've been an avid consumer of new technology and I keep running out of space on powerboards - so...	93.2%
Adv-BERT	I <b>possess</b> the good common logical sense to <b>realize</b> that oil can not last forever and I am acutely aware of how much of my life in the suburbs <b>spins</b> around petrochemical products . I've been an avid consumer of new technology and I keep running out of space on powerboards - <b>well...</b>	38.3%
Adv-R& E	I <b>know</b> the <b>wonderful general</b> sense to <b>knows</b> that <b>oils</b> can not last <b>endless</b> and I am acutely <b>know</b> of how <b>majority</b> of my <b>lived</b> in the <b>city spins</b> around petrochemical products . I've been an <b>amateur consumers</b> of <b>newly technologies</b> and I <b>kept working</b> out of <b>spaces</b> on powerboards ! <b>well...</b>	80.1%
R& E Rebuild Texts	<b>Well</b> I know the wonderful general sense notion to knows that oils <b>production</b> can not last <b>for endless years</b> and I am acutely know of how majority of my lived in the city spins around <b>the</b> petrochemical production ... I've been an amateur consumers of newly technologies and I kept working out of spaces on <b>power skateboards ! well ...</b>	80.4%
	I know the wonderful <b>common</b> sense notion to knows that oils can not last <b>forever</b> and I <b>also</b> acutely know of how majority of my lived in the <b>world and</b> around petrochemical production ... I've been an amateur consumers of newly technologies and I kept working out of <b>them</b> on <b>skateboards ! well ...</b>	81.4%
	I know the <b>wonderfully</b> general sense notion to knows that oils can not last endless and I am acutely know of how majority <b>part</b> of my lived in the <b>big</b> city spins around <b>petrochemical</b> production ... <b>I should have</b> been an amateur consumers <b>fan</b> of newly technologies and I kept <b>on</b> working out of spaces <b>and</b> on powerboards ! well ...	76.2%
	I <b>am</b> the <b>the</b> general sense notion <b>and</b> knows that oils can not last endless and I am acutely know of <b>the part</b> of my lived <b>as</b> the city spins around petrochemical production ... I've been an amateur consumers of newly technologies and I kept working out of <b>bed</b> on powerboards ! well ...	78.5%
<b>R&amp;E Failed Defense</b>		
Clean-Sample	I am trying to find somewhere to purchase a dvd / vhs copy of the movie " is n't it shocking ? " I was 7 years old when I saw this movie and I lived in the town where it was filmed. A couple of items from my family were used in the movie ...	90.1%
Adv-BERT	I am trying to find somewhere to <b>obtain</b> a <b>3d / .</b> copy of the movie " is n't it shocking ? " I was 7 years old when I <b>discovered</b> this movie and I lived in the town where it was filmed. A couple of items from my family were used in the movie ...	49.1%
Adv-R& E	I am trying to <b>obtain</b> somewhere to purchase a dvd / . copy of the movie " is n't it shocking ... " I was 7 <b>ages</b> old when I <b>discovered</b> this movie and I lived in the town where it was filmed. A couple of <b>elements</b> from my family were used in the movie ...	49.4%
R& E Rebuild Texts	I am trying <b>hard</b> obtain somewhere to purchase a dvd / . copy of the movie " is n't it <b>was</b> shocking ... " I was 7 <b>ages</b> old when I discovered <b>about</b> this movie and <b>that</b> I lived in the town where it was filmed. A couple of elements from my <b>own</b> family were used in the movie ...	75.0%
	I <b>was</b> trying to obtain somewhere to purchase a dvd / <b>copy</b> copy of the movie " is n't it shocking ... " I was 7 <b>ages</b> old when I discovered this movie and I <b>it</b> in the town where it was filmed. A couple <b>different</b> elements from my <b>childhood</b> were used in the movie ...	22.0%
	I <b>really</b> am trying <b>hard</b> to obtain somewhere to purchase a dvd / . copy of the movie " is n't it shocking ... " I was 7 <b>ages</b> <b>to</b> old when I <b>first</b> discovered this <b>horror</b> movie and I lived in the town where it was filmed. A couple of <b>the</b> elements from my family were used in the movie ...	57.1%
	I am <b>going to go</b> somewhere to purchase a dvd / . copy of the movie " is n't a shocking ... " I was 7 <b>ages</b> old when I discovered this movie and I <b>was</b> in the <b>village</b> where it was filmed. A couple of elements <b>of</b> my family were used in the movie ...	39.6%

Table 4: Error Analysis using random selected samples that (1) BERT and R&E method failed to defend; (2) BERT failed to defend while R&E succeeded. Adv-BERT is the adversarial sample generated by BERT-Attack to attack the BERT-fine-tuned model. Adv-R & E is to attack the R & E model. We also list the rebuild texts. Blue texts are pieces from a failed rebuilt sample and dark green texts are pieces from a successful rebuilt process.

## Appendix

### Error Analysis

We construct experiments using our Rebuild and Ensemble method and the BERT fine-tuning model to defend the BERT-Attack on the IMDB dataset and observe the behaviors of these methods.

As seen in Table 4, through multiple rebuilt texts, we can successfully mitigate the adversarial effect caused by the adversarial substitutions. Though the texts have been seriously replaced by adversarial tokens (more tokens have been replaced compared with only a few changes in attacking the BERT model, the model can still resist the adversarial effect through the multiple rebuilt texts.

On the other hand, in the sample that both BERT model and R & E model failed to defend, there are

rebuild texts that can correctly predict the label but some worse rebuilt texts harm the final prediction causing the failure. Specifically, we can observe that when some serious adversarial substitutes have replaced the original texts, the rebuild process cannot alter the adversarial effect easily, indicating that in the future, better locating the vulnerable places might be an effective way to defend attacks. To this end, black-box scenarios are more challenging since gradients or entropy based scoring of the importance of words are hard. On the other hand, methods such as iteration of words used in Textfooler and BERT-Attack or genetic algorithm based methods are costly. We leave this problem of finding places that might be attacked to future work based on the rebuild and ensemble framework.

Methods	Origin	Textfooler ( $K=12$ )
BERT	94.0	20.0
R & E (Mean)	93.0	82.0
R & E (Mean)( $N=1$ )	93.0	42.0
R & E (Vote)	93.0	88.0
R & E (Vote)( $N=1$ )	93.0	62.0

Table 5: Exploring the Ensemble Strategy

### Ensemble Strategy Analysis

Further, we found that the ensemble strategy could use a voting mechanism to construct a *virtual score* as the final output. That is, the argmax votes can be used to craft a confident score. When the ensemble size  $N = 1$ , this process is a hard-score attack that only gives 1 and 0 as the output.

As seen in Table 5, the defensive result using the voting strategy is higher than using the average logits. So we can assume that incorporating our rebuild and ensemble strategy with output-score-hiding strategies could further improve the model robustness.

The Rebuild and Ensemble strategy is very effective in dealing with score-based attacks, and can be further modified with a voting mechanism that can *trick* the score-based assumption.