

GRADMASK: Gradient-Guided Token Masking for Textual Adversarial Example Detection

Anonymous ACL submission

Abstract

We present a simple model-agnostic textual adversarial example detection scheme called GRADMASK. It uses gradient signals to detect adversarially perturbed tokens in an input sequence and occludes such tokens by a masking process. GRADMASK provides several advantages over existing methods including improved detection performance and a weak interpretation of its decision. Extensive evaluations on widely adopted natural language processing benchmark datasets demonstrate the efficiency and effectiveness of GRADMASK. Code and models are available at [<redacted>](#).

1 Introduction and Related Work

The advances in deep learning has revolutionized natural language processing (NLP) with state-of-the-art performance in practically every task. However, it has been shown that such systems are significantly vulnerable to specifically crafted *adversarial attacks* (Szegedy et al., 2014) at all stages of development and deployment (Ebrahimi et al., 2018; Alzantot et al., 2018; Zhang et al., 2020; Krishna et al., 2020; Tan et al., 2020, 2021). This is quite troubling as there is little to no change in the adversarially chosen test distributions compared to the training distribution (Robin, 2020).

In response to the adversarial attacks, various defense schemes have been proposed. These approaches can be grouped into three categories: (i) adversarial training (Si et al., 2020; Maharana and Bansal, 2020; Miyato et al., 2017; Zhu et al., 2020), (ii) certified robustness (Jia et al., 2019; Wang et al., 2021), and (iii) synonym substitution based methods (Wang et al., 2019, 2020; Dong et al., 2021; Zhou et al., 2021; Jones et al., 2020).

Another branch of defense strategy is the *adversarial example detection* based schemes. While the above defense schemes aim to improve the adversarial robustness of NLP systems, adversarial example detection methods are designed to reject sus-

picious inputs although they share the same goal of defeating the adversarial attacks (Aldahdooh et al., 2021). Detection-based approaches provide several advantages over adversarial robustness improvement methods. The most obvious advantage is that they do not require to modify the target model architecture or the training procedure, because they typically work as a separate module. Consequently, they do not compromise the model performance on clean datasets. Secondly, they are able to identify the intention (adversarial or not) of adversarial attacks, so users can take actions (reject or revise) accordingly. Finally, the detection algorithms may provide a better strategy for developing defense methods by informing us which parts of an input sequence are perturbed (Zhou et al., 2019).

Unlike the other defense schemes, the textual adversarial detection has not been explored much. To the best of our knowledge, there are two prior studies trying to detect token-level adversarial attacks. The very first work is the discriminate perturbations (DISP) framework proposed by Zhou et al. (2019). DISP consists of two BERT (Devlin et al., 2019) based perturbation discriminator and embedding estimator. To provide supervising signals for the discriminator, DISP randomly samples adversarial examples and learns to discriminate clean samples from the adversarial examples. In contrast, a more recent adversarial detection work, the frequency-guided word substitutions or FGWS (Mozes et al., 2021), does not need an additional training process. The key assumption of FGWS is that adversarial attack algorithms tend to exploit words that are rarely exposed during a target model’s training. However, their approach is limited to detection of only word-level attacks and the effectiveness of FGWS against attacks that do not rely on infrequent words is unclear. Especially, our experiments with a constrained high-frequency vocabulary show that attackers can still find successful attacks by using frequent tokens (§5).

Our work in this paper, instead, deviates from the word-frequency assumption by utilizing gradient signals as guidance. We harness the gradient signal to detect adversarially perturbed tokens in an input sequence by investigating the *sensitivity* of the model prediction (Ancona et al., 2018; Sundararajan et al., 2017; Li et al., 2016; Zeiler and Fergus, 2014), which indicates the network’s response to an adversarial input. The identified tokens are subsequently occluded by a mask token and fed to the model to measure the change in model’s confidence with respect to the original prediction. Fig. 1 provides an illustration of our gradient-guided detection, GRADMASK.

The gradient-based attribution of neural system’s prediction has been studied widely in deep learning (Sundararajan et al., 2017; Simonyan et al., 2014; Li et al., 2016). Some prior work in NLP uses the gradient to identify important words (Murdoch et al., 2018; Li et al., 2017). To the best of our knowledge, this is the first work on detecting textual adversarial attacks by attributing the model prediction via gradient signal analysis.

GRADMASK has several advantages over the previous methods. Firstly, it does not require any additional modules for synonym search or frequent word count that are essential in the previous methods (Mozes et al., 2021; Zhou et al., 2019). Secondly, our detection algorithm works entirely without any prior knowledge about potential attacks, which is a more practical setup. Thirdly, it works without any pre-training. Finally, it provides a weak interpretation of decision by identifying adversarially perturbed tokens. The main contributions of this work are: (i) we propose GRADMASK, a novel gradient-guided adversarial example detection method; (ii) we demonstrate its advantage over state-of-the-art adversarial example detection algorithm through extensive experiments.

2 Method

In this section, we present our proposed method. We first establish the notations in §2.1.

2.1 Notations

We consider a standard text classification task for a model $f_{\theta}(\cdot)$ with parameters $\theta \in \mathbb{R}^p$. The model $f_{\theta}(\cdot)$ is trained to fit a data distribution \mathcal{D} over pairs of an input sequence $\mathbf{x} = [x_1, \dots, x_T]$ of T tokens and its corresponding label $y \in \{1, \dots, C\}$ with C being the number of classes. We also assume a loss

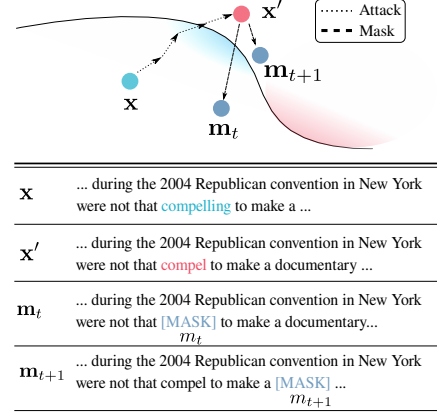


Figure 1: An illustration of the detection process of GRADMASK with a binary classification example. An attacker tries to find an adversarial example \mathbf{x}' by searching for the best perturbation (*compel*) that flips the original model prediction (expressed as the dotted line). GRADMASK attempts to identify the candidate perturbations through the gradient signal and masks one token (m_t) at a time to generate a masked sequence \mathbf{m}_t . The final decision is made by measuring the largest difference in model’s confidence for \mathbf{x}' and \mathbf{m}_t .

function $\mathcal{L}(\theta, \mathbf{x}, y)$ such as a cross-entropy loss. The output of the model is a probability distribution that satisfies: $0 \leq f_{\theta}(\mathbf{x})_i \leq 1$ and $\sum_{i=1}^C f_{\theta}(\mathbf{x})_i = 1$, where i is the class index. We denote the final prediction as $c(\mathbf{x}) = \arg \max_i f_{\theta}(\mathbf{x})_i$ and true label as $c^*(\mathbf{x}) = y^*$.

Given a sequence \mathbf{x} , a textual adversarial example \mathbf{x}' can be defined as follows: for some semantic dissimilarity measure $\delta(\mathbf{x}, \mathbf{x}')$, it has to be small and $c(\mathbf{x}') \neq c^*(\mathbf{x})$. These two conditions denote that an adversarial example has to maintain semantic meaning of the original input \mathbf{x} but misguide the model prediction (Athalye et al., 2018).

2.2 Gradient-guided Token Masking for Adversarial Example Detection

GRADMASK first finds salient tokens that significantly attribute to the model prediction, $c(\mathbf{x})$; see Fig. 1 for an illustration. A simple and widely employed approach is the gradient-based attribution analysis (Ancona et al., 2018; Sundararajan et al., 2017; Li et al., 2016). However, due to the discrete nature of texts, we cannot directly exploit the gradient-based approach. In order to deviate the issue, we compute a gradient of the word embedding \mathbf{e}_t with regard to the loss function \mathcal{L} , where \mathbf{e}_t is a simple linear projection of a (subword) token x_t . The gradient can be expressed as follows:

$$\mathbf{g}_t = \nabla_{\mathbf{e}_t} \mathcal{L}(\theta, \mathbf{x}, c(\mathbf{x})) \quad (1)$$

Algorithm 1 Gradient-based Masking for Adversarial Example Detection.

Require: Input sequence \mathbf{x} , target model f_θ

- 1: Initialize $\mathcal{M} = \{\}$ and $K = \lfloor T \times p \rfloor$.
- 2: Compute $f_\theta(\mathbf{x})_i$, where $i = c(\mathbf{x})$. \triangleright [pred. for \$\mathbf{x}\$](#)
- 3: $L := \{\|\mathbf{g}_1\|, \dots, \|\mathbf{g}_T\|\}$ via Eq. 1.
- 4: Sort L in descending order.
- 5: **while** $k \leq K$ **do**
- 6: $\|\mathbf{g}\|_t \leftarrow L[k]$
- 7: $\mathbf{m}_t = [x_1, \dots, m_t, \dots, x_T]$
- 8: $\mathcal{M}[k] = f_\theta(\mathbf{m}_t)_i$ \triangleright [prediction for \$\mathbf{m}_t\$](#)
- 9: **end while**
- 10: $w = (f_\theta(\mathbf{x})_i - \min_k \mathcal{M}[k])^2$

Note that the above loss is computed with respect to the model’s final prediction $c(\mathbf{x})$ and not the ground truth y^* .

Subsequently, we measure the amount of stimulus of the input tokens toward the model prediction by computing the L_2 -norm of \mathbf{g}_t . The stimulus is considered as a saliency score of the tokens and it is determined in descending order of the magnitude of $\|\mathbf{g}_t\|_2$ following Li et al. (2016). GRADMASK only considers the top- p portion of the input tokens in \mathbf{x} . Specifically, the number of chosen K salient tokens is $\lfloor T \times p \rfloor$, where the brackets denote the floor operation. The sampled K salient tokens are masked individually one at a time to generate a masked input sequence $\mathbf{m}_t = [x_1, \dots, m_t, \dots, x_T]$ with t being the token position of a salient token, and m_t is the mask token, [MASK].¹

The rationale behind the masking approach is based on two assumptions. The first assumption is that *adversarial examples are the result of sophisticated optimization algorithms rather than the result of random perturbations* (Goodfellow et al., 2015; Galloway et al., 2018). Thus, we conjecture that masking the suspicious tokens which are carefully crafted can significantly drop the model confidence. The second assumption is that *NLP systems are generally robust to weak-level of noise*. The partial information loss in clean samples due to masking can be offset by the overall context of the input text (supported by our experiments in §5).

Each masked sequence \mathbf{m}_t is then fed into the target model to get a prediction $f_\theta(\mathbf{m}_t)_i$, where $i = c(\mathbf{x})$. This process gives K such confidence scores which are stored in \mathcal{M} . We then compare

¹In case of non-masked language model-based classifiers, we adopted an unknown token.

Dataset	Train / Test	# Classes	Avg. Len
IMDb	25k/25k	2	215
SST-2	67k/1.8k	2	20
Yelp	560k/38k	2	152
AG	120k/7.6k	4	43

Table 1: A summary of the datasets used in our work.

the minimum confidence value in \mathcal{M} to the original confidence score $f(\mathbf{x})_i$, and the confidence change is squared to assign a stronger penalty to the higher changes. More formally,

$$w = \left(f_\theta(\mathbf{x})_i - \min_k \mathcal{M}[k] \right)^2 \quad (2)$$

The final decision is determined by an indicator function $\mathcal{I}(w, \tau)$ defined as follows:

$$\mathcal{I}(w, \tau) = \begin{cases} 0 & \text{if } w \leq \tau \\ 1 & \text{else} \end{cases} \quad (3)$$

where τ is a pre-defined threshold. Alg. 1 presents the overall process of GRADMASK.

3 Experiment Settings

In this section, we present our experiment settings: the datasets, target models, adversarial example generation, and evaluation metrics.

3.1 Datasets

We evaluate the methods on four classification tasks. We use the IMDB (Maas et al., 2011), AGNEWS (Zhang et al., 2015), YELP (Zhang et al., 2015), and Stanford Sentiment Treebank (SST) (Socher et al., 2013) datasets that are widely adopted for benchmarking adversarial robustness of NLP systems. The IMDB dataset contains movie reviews labeled with positive or negative sentiment labels. The AGNEWS dataset contains news articles from more than 2,000 news sources and the samples are categorized into the four largest classes. The YELP dataset is a binary sentiment classification dataset which consists of Yelp reviews. The SST dataset provides movie reviews with fine-grained sentiment labels. We turn the labels into binary (SST-2) to follow the setting of FGWS (Mozes et al., 2021). Table 1 gives an overview of the datasets.

3.2 Target Models

We evaluate GRADMASK on three different sequence modeling architectures, which have been

MODEL	DATASET	ACC (%)
ROBERTA	IMDb	93.36
	SST-2	91.98
	YELP	97.91
	AG	95.3
ROBERTA-LONG	IMDb	93.71
	SST-2	88.69
DISTILBERT	IMDb	90.57
	SST-2	91.21
	AG	94.37
LSTM	IMDb	87.27
	SST-2	83.53

Table 2: A summary of the target models and their clean testset performance.

widely employed in NLP. We first consider a large-scaled pre-trained Transformer-based language model, ROBERTA-BASE (Liu et al., 2019), which contains 124 million parameters. Subsequently, we also evaluate on a relatively smaller Transformer-based model called DISTILBERT-BASE (Sanh et al., 2020), which has approximately 40% fewer parameters than ROBERTA-BASE. Finally, we consider the LSTM, which used to be the dominant architecture before the arrival of Transformers.

Table 2 shows the standard task performance of the models on the three datasets. To train the models, we followed the hyperparameter settings provided by Mozes et al. (2021). The TRANSFORMER based models are optimized by AdamW (Loshchilov and Hutter, 2019) with a linear adaptive learning rate scheduler. For LSTM, the initial word embeddings are initialized with GloVe (Pennington et al., 2014). The texts in IMDB and YELP are comparatively longer than those in AGNEWS and SST-2. For the IMDB classification task, the maximum sequence lengths for ROBERTA, DISTILBERT and LSTM are set to 256, 256, and 200, respectively, and ROBERTA-LONG is trained with a longer sequence (400 tokens) than the standard one. The details of model architectures are provided in the supplementary material. All of the experiments are conducted on an Intel Xeon Gold 5218R CPU-2.10GHz processor with a single Quadro RTX 6000 GPU.

3.3 Adversarial Example Generation

We generated adversarial examples against the selected target models via four different attack algorithms. They include two baseline attacks and two widely adopted synonym substitution-based token-level attacks, as used in the previous work

- **Random** is a simple word replacement-based

baseline attack algorithm. It randomly selects a synonym of a token in the original input text. Synonyms are identified via WordNet.

- **Prioritized** attack is also based on word replacement, but it puts a higher priority on a synonym that maximizes the target model’s prediction confidence change.

- **Genetic** attack (GA) was proposed by Alzantot et al. (2018). It adopts the crossover and mutation operations in genetic algorithms to generate adversarial examples. GA searches synonyms based on the GloVe word embedding space with a language model (Radford et al., 2019).²

- **PWWS** or Probability weighted word saliency (Ren et al., 2019) is a greedy word substitution-based attack algorithm. The word replacement order is determined by a word saliency score computed through the model’s confidence change. The word synonym is searched via WordNet.

3.4 Evaluation Metrics

The main interest of this work lies in an evaluation of the detection performance of our proposed method GRADMASK. FGWS (Mozes et al., 2021) was mainly evaluated via F1 score, but we follow the standards from the out-of-distribution (OOD) sample detection literature (Hendrycks et al., 2019; Ouyang et al., 2021) for better understanding of the methods.

The adversarial example detection can be considered as a binary classification problem of verifying *positive (adversarial)* vs. *negative (clean)* class. We evaluate a ratio of true positive samples so-called true positive rate (TPR or recall) against false positive rate (FPR) defined as:

$$TPR = \frac{1}{n^+} \sum_i \mathcal{I}(w^+, \tau) \quad (4)$$

$$FPR = \frac{1}{n^-} \sum_i \mathcal{I}(w^-, \tau), \quad (5)$$

where the superscripts + and – denote the positive and the negative classes, respectively. Based on these two rates, we evaluate the methods with the following evaluation metrics:

- **FPR95** refers to a FPR at 95% TPR. FPR95 quantifies how many clean samples have to be rejected to detect 95% of the adversarial examples. FPR is a very important metric for evaluating detection al-

²We adopted the modified implementation provided by Mozes et al. (2021) for a fair comparison. The implementation details are provided in the supplementary material.

gorithms (Aldahdooh et al., 2021). A lower FPR95 score is often required for systems that require a high level of system safety or security.

- **AUROC** stands for the area under receiver operating characteristic curve. For each operational setting of τ from 0 to 1, TPR and FPR can be plotted. This curve is called receiver operating characteristic curve (ROC curve).

- **AUPR** denotes area under precision-recall (PR) curves. There exists an imbalance of data distribution between positive class and negative class. To deal with the data distribution skew, we evaluate AUPR scores for each class.

- **Acc** denotes a detection accuracy. We use Acc only for experiments with balanced datasets.

4 Results & Analysis

We first evaluate GRADMASK on widely employed NLP datasets and compare it with baselines (§4.1 and §4.2). Then, we analyze the adversarially perturbed token detection performance of GRADMASK (§4.3). Subsequently, we investigate GRADMASK’s potential against a non-synonym based (character-level) attack (§4.4). Finally, We investigate the relationship between the adversarial robustness of NLP classification models and the word frequency in the adversarial examples (§5).

4.1 Adversarial Example Detection

For adversarial example detection, we compare the performance of GRADMASK with that of FGWS (Mozes et al., 2021). The hyperparameter settings of FGWS is tuned as provided by Mozes et al. (2021).³ The overall experimental results are presented in Table 3. Note that AUPR-C and AUPR-A represent the AUPR score of clean samples (negative class) and that of adversarial samples (positive class), respectively.

As shown in Table 3, GRADMASK tends to show better AUROC, FPR95, and AUPR-C scores in most of the evaluation measures. Particularly, it outperforms FGWS for all Transformer-based systems (ROBERTA, ROBERTA-LONG, and DISTILBERT) in terms of the FPR95 score, which is an important metric for systems with high security requirements. In addition, GRADMASK achieves notably better AUPR-C scores in most of the experiment scenarios. This tendency is well presented in Fig. 2, which shows ROC curves of FGWS and

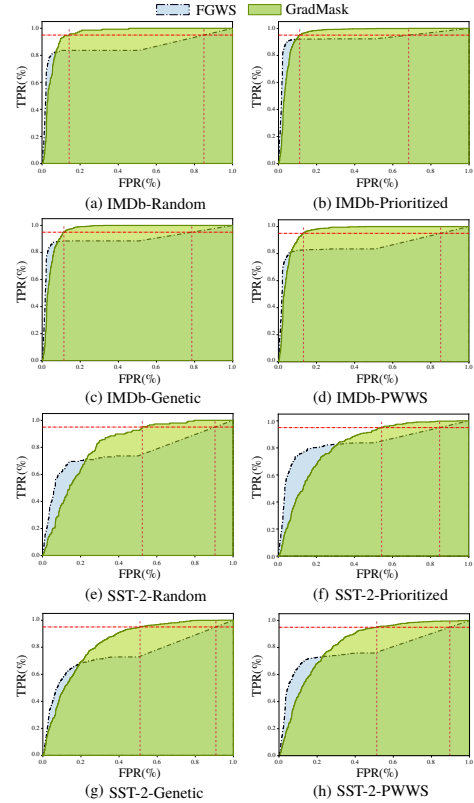


Figure 2: ROC curves of FGWS and GRADMASK with the ROBERTA model. The horizontal red line is at the 95% TPR and the vertical lines at the FPRs of two algorithms, respectively (best viewed in color).

GRADMASK for ROBERTA model. The ROC curves of FGWS tend to increase steeply and remain stable. However, as TPR increases, FGWS significantly compromises FPR score. Especially, at some point, TPR and FPR show a linear trend. In contrast, GRADMASK tends to reach 95% TPR at lower FPR scores and shows larger AUROC scores.

On the other hand, GRADMASK shows lower performance scores in all metrics on SST-2 with the LSTM model as shown in Table 3. Nevertheless, the overall detection performance of GRADMASK tends to improve proportionally to the model size and the standard performance. Another notable observation is that GRADMASK achieves these results within two candidates except for the LSTM model (K in Table 3). These results may imply that NLP systems are largely robust to a partial loss of information resulting from the masking strategy on clean samples, but there is a significant change in the adversary response caused by a salient token masking. We also conducted an additional experiment to investigate the performance changes while varying the number of masked tokens in the samples in Appendix D.2.

³<https://github.com/maximilianmozes/fgws>

MODEL	DATASET	# SAMPLES		ATTACK	FPR95 (%)		AUROC (%)		AUPR-C (%)		AUPR-A (%)		K
		TN	TP		FGWS	GM	FGWS	GM	FGWS	GM	FGWS	GM	
ROBERTA	IMDb	2000	147	RANDOM	84.98	12.50	86.06	94.93	98.46	99.62	51.55	46.55	2
		2000	995	PRIORITIZED	68.31	11.1	92.67	95.55	95.06	98.12	89.2	84.89	2
		2000	1042	GENETIC	78.53	11.4	89.88	95.69	92.89	98.17	86.72	85.04	2
		2000	1016	PWWS	85.17	12.10	85.85	95.27	90.47	98.00	83.00	84.18	2
	SST-2	1821	148	RANDOM	90.54	52.39	75.40	81.43	97.17	98.18	37.62	20.37	1
		1821	479	PRIORITIZED	84.69	54.26	83.57	82.09	94.23	94.65	65.35	46.95	1
		1821	968	GENETIC	90.82	56.89	74.60	79.19	84.22	90.97	66.55	61.33	1
		1821	736	PWWS	65.06	51.29	77.72	82.73	88.66	92.44	66.05	58.51	1
ROBERTA-LONG	IMDb	2000	190	RANDOM	89.77	12.85	81.05	94.12	97.26	99.57	58.84	37.20	2
		2000	1037	PRIORITIZED	68.20	11.30	93.08	94.66	95.02	97.79	90.70	81.78	2
		2000	888	GENETIC	80.96	10.65	89.05	95.20	93.24	97.93	85.38	83.26	2
		2000	1129	PWWS	84.38	10.95	87.10	95.07	90.26	97.96	86.38	83.38	2
	SST-2	1821	176	RANDOM	89.34	60.35	76.42	75.72	96.94	96.97	35.15	18.24	1
		1821	527	PRIORITIZED	87.06	60.08	79.80	77.73	92.71	92.78	62.95	43.31	1
		1821	960	GENETIC	92.15	69.80	68.18	73.55	82.55	84.89	61.46	53.11	1
		1821	772	PWWS	90.05	57.50	75.54	78.57	87.83	90.41	66.44	54.38	1
DISTILBERT	IMDb	2000	212	RANDOM	86.98	37.30	83.36	87.66	97.46	98.56	59.59	33.33	1
		2000	1182	PRIORITIZED	62.85	31.70	93.20	89.66	94.79	94.50	91.88	76.09	1
		2000	1202	GENETIC	75.59	22.80	90.28	90.23	92.50	95.27	89.25	74.41	1
		2000	1335	PWWS	83.06	36.64	86.56	88.74	88.9	92.93	86.95	79.10	1
	SST-2	1821	171	RANDOM	84.42	59.69	83.17	77.78	87.77	97.32	37.23	18.40	1
		1821	614	PRIORITIZED	84.36	58.70	84.29	78.87	92.97	92.34	70.36	46.86	1
		1821	1105	GENETIC	90.97	49.81	74.74	78.06	82.27	88.18	69.36	57.32	1
		1821	860	PWWS	71.56	54.31	80.30	78.87	88.25	89.93	71.56	54.41	1
LSTM	IMDb	2000	198	RANDOM	89.64	37.55	77.82	84.22	96.90	98.31	44.47	24.87	20
		2000	1451	PRIORITIZED	78.68	30.50	88.34	86.64	89.66	92.41	88.66	73.90	20
		2000	1548	GENETIC	89.73	30.50	77.47	86.59	81.04	92.00	78.92	74.50	20
		2000	1735	PWWS	88.85	30.90	80.53	86.99	81.47	91.45	83.85	78.43	20
	SST-2	1821	238	RANDOM	86.35	98.13	79.14	58.45	96.36	90.22	36.37	13.35	20
		1821	669	PRIORITIZED	89.89	95.18	74.97	68.45	88.73	84.33	57.21	36.24	20
		1821	1186	GENETIC	91.28	96.00	71.37	66.74	80.08	72.67	66.55	51.55	20
		1821	1013	PWWS	90.28	95.51	74.68	69.59	83.96	78.51	66.46	48.26	20

Table 3: Adversarial example detection results of FGWS and GRADMASK (GM). AUPR-C and AUPR-A denote AUPR of clean example and adversarial example classes, respectively.

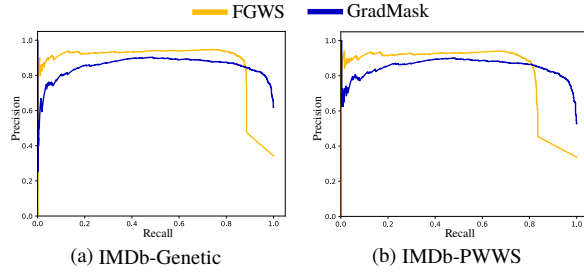


Figure 3: Precision-Recall curves of FGWS and GRADMASK on IMDb with the RoBERTA model against the PWWS and genetic attacks.

Another observation is that our gradient-based masking strategy occasionally detects adversarial examples through masking a clean token as presented in §4.3 and Fig. 4. This result implies that the hidden representation of adversarial tokens significantly affects that of clean tokens. We leave analysis of this correlation as future work.

Moreover, GRADMASK shows consistently better performance in detecting strong attacks such as genetic attack and PWWS attack which are more aggressive than the others. We conjecture that stronger attacks select and engineer the crucial tokens more carefully, so masking these tokens would hugely reduce the effectiveness of these attacks.

Finally, we observe that GRADMASK underperforms FGWS in terms of AUPR-A. A possible explanation may be related to the nature of the syn-

onym substitution strategy. We hypothesize that FGWS tends to transform an input sequence aggressively. This view can be supported by their FPR95 scores and precision-recall (PR) curves. Firstly, the ROC curves of FGWS typically show high FPRs at high TPRs (Fig. 2). Secondly, from the PR curves of FGWS shown in Fig. 3, the precision scores drop significantly as the recall scores increase. We provide PR curves for 6 other scenarios in the supplementary material.

4.2 A Comparison with Anomaly Detection Algorithms

We conducted additional experiments via TextAttack library (Morris et al., 2020)⁴ to compare GRADMASK with baseline anomaly detection algorithms such as maximum softmax probability (MSP) (Hendrycks and Gimpel, 2017) and one-class support vector machine with linear kernel (OCSVM) (Schölkopf et al., 2000) that are widely adopted as a baseline in various anomaly detection areas (Lee et al., 2018; Shafaei et al., 2019; Winkens et al., 2020; Aldahdooh et al., 2021).

We trained ROBERTA-BASE on three datasets including IMDb, YELP, and AG datasets, respectively. We sampled 1,000 clean examples and their

⁴We adopted TextAttack framework (Morris et al., 2020) to attack the victim models. Their implementation difference is provided in the supplementary material.

DATASET	METHOD	AUROC	FPR95	AUPR-C	AUPR-A	ACC
AG	MSP	94.81	19.40	95.28	93.17	87.50
	OCSVM	94.72	16.70	95.63	91.72	89.25
	GM	94.93	11.00	96.69	88.37	92.50
IMDb	MSP	95.34	19.43	95.07	95.16	90.32
	OCSVM	84.83	98.00	71.88	89.41	83.82
	GM	95.53	10.62	96.81	92.26	92.20
YELP	MSP	97.22	11.47	97.45	97.12	92.16
	OCSVM	95.90	15.50	93.92	96.20	91.07
	GM	97.81	5.68	98.37	96.22	94.69

Table 4: A comparison of PWWS attack detection results on RoBERTa model with MSP and OCSVM.

corresponding 1,000 adversarial examples without a text length limitation via PWWS attack.

From the results in Table 4, we notice that GRADMASK significantly outperforms the baselines by a large margin except for the AUPR-A scores. These results are consistent with the results reported in §4.1. GRADMASK achieves significantly lower FPR95 scores than that of MSP and OCSVM for all three datasets and higher AUPR-C scores. Also, we report a detection accuracy for this experiment because the datasets are well-balanced unlike the previous experiment in §4.1. Specifically, we measured the best accuracy over a varying threshold setting. Table 4 shows that GRADMASK achieved the best detection accuracy for each dataset and its error rate is around 7% for all tasks.

We further analyze statistics of the features extracted from MSP and GRADMASK methods to attribute the superior performance of GRADMASK. Table 5 presents two statistics of the extracted features, mean (AVG) and standard deviation (STD). The values are averaged over 1,000 samples. As shown in the table, the overall mean differences between the w (c.f., Eq. (2)) of adversarial examples (w -A) and w of clean samples (w -C) are higher than that of MSP, which implies that GRADMASK feature w is more distinguishable. Specifically, for IMDb, MSP shows 0.182 ($= 0.990 - 0.808$), but GRADMASK shows 0.478 at $K = 3$. In addition, standard deviations of GRADMASK are generally smaller than that of MSP.

4.3 Adversarial Token Detection

We now analyze how our gradient-based approach GRADMASK attributes the model prediction on adversarial examples. Fig. 4 shows perturbed token detection rates of two Transformer-based models, DISTILBERT and ROBERTA, on two datasets, IMDb and AGNEWS. We report detection rates at top-1, top-3, and top-5, which refers to the total number of adversarially perturbed tokens identified

DATASET	K	w -A/CONF-A (AVG \pm STD)	w -C/CONF-C (AVG \pm STD)
IMDb	MSP	-/0.808 \pm 0.155	-/0.990 \pm 0.034
	1	0.353 \pm 0.318/-	0.020 \pm 0.117/-
	2	0.424 \pm 0.308/-	0.024 \pm 0.129/-
	3	0.528 \pm 0.309/-	0.050 \pm 0.187/-
AG	MSP	-/0.743 \pm 0.163	-/0.980 \pm 0.068
	1	0.335 \pm 0.299/-	0.030 \pm 0.147/-
	2	0.381 \pm 0.294/-	0.037 \pm 0.160/-
	3	0.468 \pm 0.295/-	0.028 \pm 0.137/-
YELP	MSP	-/0.951 \pm 0.067	-/0.999 \pm 0.005
	1	0.509 \pm 0.422/-	0.009 \pm 0.091/-
	2	0.650 \pm 0.379/-	0.015 \pm 0.116/-
	3	0.783 \pm 0.305/-	0.021 \pm 0.135/-

Table 5: Statistics (AVG and STD) of extracted features. The first row of each dataset denotes the maximum softmax probability (MSP) of the ROBERTA model for adversarial (Conf-A) and clean (Conf-C) examples, respectively. The subsequent rows show the mean and standard deviation of w of GRADMASK while varying the number of mask tokens K .

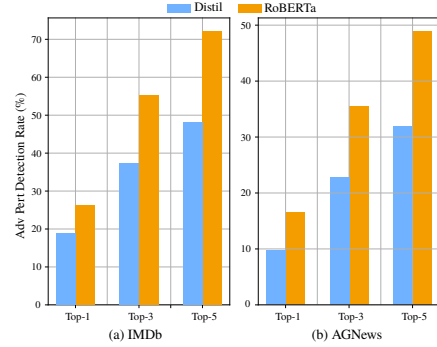


Figure 4: Adversarially perturbed token detection rates at top-1, top-2 and top-5 for GRADMASK.

within the top- N values of w in Eq. (2). In case of DISTILBERT, it shows 48.17% and 31.82% detection rates for IMDb and AGNEWS within the top-5 predictions, respectively. On the other hand, ROBERTA shows 72.04% and 48.85% detection rates for IMDb and AGNEWS within the top-5 predictions. Another notable observation is that for the IMDb classification task, top-1 predictions detect the adversarial tokens with 49% and 78% probability for DISTILBERT and ROBERTA, respectively. For AGNEWS, their top-1 predictions show 45% and 67% detection probability, respectively.

4.4 Character-Level Attack Detection

To investigate the potential of GRADMASK against non-synonym based attacks, we conduct an additional experiment with a character-level attack (Pruthi et al., 2019) from the TextAttack library (Morris et al., 2020). Even though character-level attacks are known to be relatively simple to defend at a preprocessing stage with a spell or a grammar checker (Pruthi et al., 2019), our motivation for

MODEL	AUROC	FPR95	AUPR-C	AUPR-A	Acc	K
RoBERTA	80.02	63.39	79.62	75.74	75.83	3
DISTILBERT	80.42	63.76	81.02	75.07	75.36	2

Table 6: Adversarial example detection results (in %) against a character-level attack.

this experiment is to demonstrate the potential of GRADMASK against non-synonym based attacks.

We generated 691 and 897 adversarial examples from AG against RoBERTA-BASE and DISTIL-BASE without any maximum text length limitation, respectively. From the results in Table 6, we see that our method shows promising results with AUROC scores of 79.68% and 80.42% for RoBERTA-BASE and DISTIL-BASE, respectively. It would be interesting to see how GRADMASK performs for other kinds of non-synonym attacks such as syntactically controlled paraphrase networks (SCPNs) (Iyyer et al., 2018) or universal adversarial attack (Song et al., 2021) which we leave as future work.

5 Discussion on Word Frequency and Adversarial Robustness

According to Mozes et al. (2021), the brittleness of NLP systems against adversarial examples would be attributed to the distribution of word frequency in a training set. However, one of the widely accepted explanations about the existence of adversarial examples insists that adversarial examples are a result of the standard optimization rather than data distribution (Ilyas et al., 2019). We investigated how the word frequency can affect the model’s robustness via a series of experiments. Consequently, we find that *deep NLP systems can still be fooled by adversarial examples with words that are frequently exposed during their training stage*.

To validate this claim, we trained the victim models with a word frequency constraint. Specifically, we built a new vocabulary set V' to be comprised of only the top-10% frequently used words from the original vocabulary set V . The vocabulary-constrained models are designed to block all infrequent words that are out of V' in an input sequence by masking those tokens. We first evaluated the model performance to observe how the vocabulary constraint affects the model performance. As shown in Table 7, the standard task performance of the victim models under the constraint ($\text{Acc-}V'$) only marginally decreases (about 1 - 4%) compared to the original accuracy ($\text{Acc-}V$). These results

Model	Dataset	Acc- V	Acc- V'	$x' \in V'$	AAcc
DISTILBERT	IMDb	92.98	92.17	71.73	10.4
	AG	94.37	90.78	68.92	15.6
RoBERTA	IMDb	95.33	95.15	67.38	7.6
	AG	95.22	94.87	44.26	30.8

Table 7: Word frequency and adversarial robustness. Acc- V and Acc- V' refer to accuracies of the model with the original vocabulary V and constrained vocabulary V' , respectively. $x' \in V'$ denotes a ratio of perturbed tokens that are part of V' . AAcc denotes an under attack accuracy of the model with V' .

show that masking infrequent tokens does not hurt the model performance significantly. Next, we generated 1,000 pairs of samples via the PWWS attack algorithm (Ren et al., 2019) against the word frequency constrained models. Each sample pair consists of a clean example and its corresponding adversarial example that successfully fools the target model.

According to the infrequent word assumption (Mozes et al., 2021), the models trained on V' are expected to be robust against adversarial attacks. However, from the results in Table 7, we notice that they showed significant brittleness against adversarial attacks. For instance, DISTILBERT models show approximately 10% accuracies for both datasets when under attack (AAcc). Similarly, RoBERTA models show under attack accuracies of 7.6% and 30.8% for AGNEWS and IMDB, respectively. Thus, we claim that *the vulnerabilities of NLP systems cannot only be attributed to the infrequent words*.

6 Conclusion

We have proposed a simple model-agnostic adversarial example detection scheme, GRADMASK, which utilizes gradient signals as a guidance to detect adversarially perturbed tokens. This guidance additionally provides a weak interpretation about its decision. The experimental results show that GRADMASK is a promising approach as a textual adversarial attack detection algorithm for NLP systems. Particularly, it shows significantly low FPR95 scores, which is a highly desirable property for NLP systems with high-security requirements. In addition, GRADMASK does not require an additional module or a strong assumption about potential attacks which are more realistic in practice. In conclusion, our detection strategy can serve as a useful tool for identifying adversarial attacks for protecting the text classification systems.

References

- Ahmed Aldahdooh, Wassim Hamidouche, Sid Ahmed Fezza, and Olivier Déforges. 2021. Adversarial example detection for DNN models: A review.
- Moustafa Alzantot, Yash Sharma, Ahmed Elghohary, Bo-Jhang Ho, Mani Srivastava, and Kai-Wei Chang. 2018. Generating natural language adversarial examples. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2890–2896, Brussels, Belgium. Association for Computational Linguistics.
- Marco Ancona, Enea Ceolini, Cengiz Öztireli, and Markus Gross. 2018. Towards better understanding of gradient-based attribution methods for deep neural networks. In *International Conference on Learning Representations*.
- Anish Athalye, Nicholas Carlini, and David Wagner. 2018. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018*.
- Ciprian Chelba, Tomáš Mikolov, Mike Schuster, Qi Ge, Thorsten Brants, and Phillipp Koehn. 2013. One billion word benchmark for measuring progress in statistical language modeling. *CoRR*.
- 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.
- Xinshuai Dong, Anh Tuan Luu, Rongrong Ji, and Hong Liu. 2021. Towards robustness against natural language word substitutions. In *International Conference on Learning Representations*.
- Javid Ebrahimi, Anyi Rao, Daniel Lowd, and Dejing Dou. 2018. HotFlip: White-box adversarial examples for text classification. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 31–36, Melbourne, Australia. Association for Computational Linguistics.
- Angus Galloway, Graham W. Taylor, and Medhat Moussa. 2018. Attacking binarized neural networks. In *International Conference on Learning Representations*.
- Ian Goodfellow, Jonathon Shlens, and Christian Szegedy. 2015. Explaining and harnessing adversarial examples. In *International Conference on Learning Representations*.
- Dan Hendrycks and Kevin Gimpel. 2017. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *Proceedings of International Conference on Learning Representations*.
- Dan Hendrycks, Mantas Mazeika, and Thomas G. Dietterich. 2019. Deep anomaly detection with outlier exposure. In *International Conference on Learning Representations*.
- Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Aleksander Madry. 2019. Adversarial examples are not bugs, they are features. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.
- Mohit Iyyer, John Wieting, Kevin Gimpel, and Luke Zettlemoyer. 2018. Adversarial example generation with syntactically controlled paraphrase networks. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1875–1885, New Orleans, Louisiana. Association for Computational Linguistics.
- Robin Jia, Aditi Raghunathan, Kerem Göksel, and Percy Liang. 2019. Certified robustness to adversarial word substitutions. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4129–4142, Hong Kong, China. Association for Computational Linguistics.

658	Erik Jones, Robin Jia, Aditi Raghunathan, and	<i>nual Meeting of the Association for Computa-</i>	704
659	Percy Liang. 2020. Robust encodings: A frame-	<i>tional Linguistics: Human Language Technolo-</i>	705
660	work for combating adversarial typos. In <i>Pro-</i>	<i>gies</i> , pages 142–150, Portland, Oregon, USA.	706
661	<i>ceedings of the 58th Annual Meeting of the As-</i>	Association for Computational Linguistics.	707
662	<i>sociation for Computational Linguistics.</i> Associ-		
663	ation for Computational Linguistics.		
664	Kalpesh Krishna, Gaurav Singh Tomar, Ankur P.	Adyasha Maharana and Mohit Bansal. 2020. Ad-	708
665	Parikh, Nicolas Papernot, and Mohit Iyyer. 2020.	versarial augmentation policy search for domain	709
666	Thieves on sesame street! model extraction of	and cross-lingual generalization in reading com-	710
667	BERT-based APIs. In <i>International Conference</i>	prehension. In <i>Findings of the Association for</i>	711
668	<i>on Learning Representations.</i>	<i>Computational Linguistics: EMNLP 2020.</i> As-	712
669	Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo	sociation for Computational Linguistics.	713
670	Shin. 2018. A simple unified framework for		
671	detecting out-of-distribution samples and adver-	Takeru Miyato, Andrew M Dai, and Ian Good-	714
672	sarial attacks. In <i>Proceedings of the 32nd In-</i>	fellow. 2017. Adversarial training methods for	715
673	<i>ternational Conference on Neural Information</i>	semi-supervised text classification. In <i>Interna-</i>	716
674	<i>Processing Systems, NIPS’18</i> , page 7167–7177,	<i>tional Conference on Learning Representations.</i>	717
675	Red Hook, NY, USA. Curran Associates Inc.	John X. Morris, Eli Lifland, Jin Yong Yoo, Jake	718
676	Hector J. Levesque, Ernest Davis, and Leora Mor-	Grigsby, Di Jin, and Yanjun Qi. 2020. TextAt-	719
677	genstern. 2012. The Winograd Schema Chal-	tack: A framework for adversarial attacks, data	720
678	lenge. In <i>Proceedings of the Thirteenth Interna-</i>	augmentation, and adversarial training in nlp.	721
679	<i>tional Conference on Principles of Knowledge</i>		
680	<i>Representation and Reasoning, KR’12</i> , pages	Maximilian Mozes, Pontus Stenetorp, Bennett	722
681	552–561. AAAI Press.	Kleinberg, and Lewis Griffin. 2021. Frequency-	723
682	Jiwei Li, Xinlei Chen, Eduard Hovy, and Dan Juraf-	guided word substitutions for detecting textual	724
683	sky. 2016. Visualizing and understanding neural	adversarial examples. In <i>Proceedings of the 16th</i>	725
684	models in NLP. In <i>Proceedings of the 2016</i>	<i>Conference of the European Chapter of the As-</i>	726
685	<i>Conference of the North American Chapter of</i>	<i>sociation for Computational Linguistics: Main</i>	727
686	<i>the Association for Computational Linguistics:</i>	<i>Volume</i> , pages 171–186, Online. Association for	728
687	<i>Human Language Technologies.</i> Association for	Computational Linguistics.	729
688	Computational Linguistics.	W. James Murdoch, Peter J. Liu, and Bin Yu. 2018.	730
689	Jiwei Li, Will Monroe, and Dan Jurafsky. 2017.	Beyond word importance: Contextual decompo-	731
690	Understanding neural networks through repre-	sition to extract interactions from LSTMs. In	732
691	sentation erasure.	<i>International Conference on Learning Represen-</i>	733
692	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du,	<i>tations.</i>	734
693	Mandar Joshi, Danqi Chen, Omer Levy, Mike	Yawen Ouyang, Jiasheng Ye, Yu Chen, Xinyu Dai,	735
694	Lewis, Luke Zettlemoyer, and Veselin Stoyanov.	Shujian Huang, and Jiajun Chen. 2021. Energy-	736
695	2019. RoBERTa: A robustly optimized BERT	based unknown intent detection with data ma-	737
696	pretraining approach. <i>CoRR</i> .	nipulation. In <i>Findings of the Association for</i>	738
697	Ilya Loshchilov and Frank Hutter. 2019. Decoupled	<i>Computational Linguistics: ACL-IJCNLP 2021.</i>	739
698	weight decay regularization. In <i>International</i>	Jeffrey Pennington, Richard Socher, and Christo-	740
699	<i>Conference on Learning Representations.</i>	pher D Manning. 2014. GloVe: Global vectors	741
700	Andrew L. Maas, Raymond E. Daly, Peter T. Pham,	for word representation. In <i>Proceedings of the</i>	742
701	Dan Huang, Andrew Y. Ng, and Christopher	<i>2014 Conference on Empirical Methods in Nat-</i>	743
702	Potts. 2011. Learning word vectors for senti-	<i>ural Language Processing EMNLP</i> , volume 14,	744
703	ment analysis. In <i>Proceedings of the 49th An-</i>	pages 1532–1543.	745
		Danish Pruthi, Bhuwan Dhingra, and Zachary C.	746
		Lipton. 2019. Combating adversarial mis-	747
		spellings with robust word recognition. In <i>Pro-</i>	748
		<i>ceedings of the 57th Annual Meeting of the Asso-</i>	749

750	<i>ciation for Computational Linguistics</i> , volume	Richard Socher, Alex Perelygin, Jean Wu, Jason	795
751	abs/1905.11268.	Chuang, Christopher D. Manning, Andrew Ng,	796
752	Alec Radford, Jeff Wu, Rewon Child, David Luan,	and Christopher Potts. 2013. Recursive deep	797
753	Dario Amodei, and Ilya Sutskever. 2019. Lan-	models for semantic compositionality over a sen-	798
754	guage models are unsupervised multitask learn-	timent treebank. In <i>Proceedings of the 2013</i>	799
755	ers.	<i>Conference on Empirical Methods in Natural</i>	800
756	Shuhuai Ren, Yihe Deng, Kun He, and Wanxiang	<i>Language Processing</i> , pages 1631–1642, Seat-	801
757	Che. 2019. Generating natural language adver-	tle, Washington, USA. Association for Compu-	802
758	sarial examples through probability weighted	tational Linguistics.	803
759	word saliency. In <i>Proceedings of the 57th Annual</i>	Liwei Song, Xinwei Yu, Hsuan-Tung Peng, and	804
760	<i>Meeting of the Association for Computational</i>	Karthik Narasimhan. 2021. Universal adversar-	805
761	<i>Linguistics</i> , pages 1085–1097, Florence, Italy.	ial attacks with natural triggers for text classifi-	806
762	Association for Computational Linguistics.	cation. In <i>Proceedings of the 2021 Conference</i>	807
763	Jia Robin. 2020. <i>Building robust natural language</i>	<i>of the North American Chapter of the Associa-</i>	808
764	<i>processing systems</i> . Ph.D. thesis, Stanford Uni-	<i>tion for Computational Linguistics: Human Lan-</i>	809
765	versity, Stanford, California.	<i>guage Technologies</i> , pages 3724–3733, Online.	810
766	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bha-	Association for Computational Linguistics.	811
767	gavatula, and Yejin Choi. 2019. Winogrande:	Mukund Sundararajan, Ankur Taly, and Qiqi Yan.	812
768	An adversarial winograd schema challenge at	2017. Axiomatic attribution for deep networks.	813
769	scale. <i>arXiv preprint arXiv:1907.10641</i> .	In <i>Proceedings of the 34th International Con-</i>	814
770	Victor Sanh, Lysandre Debut, Julien Chaumond,	<i>ference on Machine Learning - Volume 70</i> ,	815
771	and Thomas Wolf. 2020. DistilBERT, a distilled	ICML 17, page 3319–3328. JMLR.org.	816
772	version of BERT: smaller, faster, cheaper and	Christian Szegedy, Wojciech Zaremba, Ilya	817
773	lighter.	Sutskever, Joan Bruna, Dumitru Erhan, Ian	818
774	Bernhard Schölkopf, Robert C Williamson, Alex	Goodfellow, and Rob Fergus. 2014. Intriguing	819
775	Smola, John Shawe-Taylor, and John Platt. 2000.	properties of neural networks. In <i>International</i>	820
776	Support vector method for novelty detection. In	<i>Conference on Learning Representations</i> .	821
777	<i>Advances in Neural Information Processing Sys-</i>	Samson Tan, Shafiq Joty, Kathy Baxter, Araz Taei-	822
778	<i>tems</i> , volume 12. MIT Press.	hagh, Gregory A. Bennett, and Min-Yen Kan.	823
779	Alireza Shafaei, Mark Schmidt, and James J. Lit-	2021. Reliability testing for natural language	824
780	tle. 2019. A less biased evaluation of out-of-	processing systems. In <i>Proceedings of the 59th</i>	825
781	distribution sample detectors. In <i>BMVC</i> , page 3.	<i>Annual Meeting of the Association for Compu-</i>	826
782	Chenglei Si, Zhengyan Zhang, Fanchao Qi,	<i>tational Linguistics</i> , ACL’21, page 4153–4169,	827
783	Zhiyuan Liu, Yasheng Wang, Qun Liu, and	Bangkok, Thailand. ACL.	828
784	Maosong Sun. 2020. Better robustness by	Samson Tan, Shafiq Joty, Min-Yen Kan, and	829
785	more coverage: Adversarial training with mixup	Richard Socher. 2020. It’s morphin’ time! Com-	830
786	augmentation for robust fine-tuning. <i>CoRR</i> ,	bating linguistic discrimination with inflectional	831
787	abs/2012.15699.	perturbations. In <i>Proceedings of the 58th Annual</i>	832
788	K. Simonyan, A. Vedaldi, and Andrew Zisserman.	<i>Meeting of the Association for Computational</i>	833
789	2014. Deep inside convolutional networks: Visu-	<i>Linguistics</i> , pages 2920–2935, Online. Associa-	834
790	alising image classification models and saliency	tion for Computational Linguistics.	835
791	maps. In <i>2nd International Conference on Learn-</i>	Wenjie Wang, Pengfei Tang, Jian Lou, and	836
792	<i>ing Representations, ICLR 2014, Banff, AB,</i>	Li Xiong. 2021. Certified robustness to word	837
793	<i>Canada, April 14-16, 2014, Workshop Track Pro-</i>	substitution attack with differential privacy. In	838
794	<i>ceedings</i> .	<i>Proceedings of the 2021 Conference of the North</i>	839
		<i>American Chapter of the Association for Com-</i>	840
		<i>putational Linguistics: Human Language Tech-</i>	841

nologies, pages 1102–1112, Online. Association for Computational Linguistics.

Xiaosen Wang, Hao Jin, and Kun He. 2019. Natural language adversarial attacks and defenses in word level. *CoRR*.

Xiaosen Wang, Yichen Yang, Yihe Deng, and Kun He. 2020. Adversarial training with fast gradient projection method against synonym substitution based text attacks. *CoRR*, abs/2008.03709.

Jim Winkens, Rudy Bunel, Abhijit Guha Roy, Robert Stanforth, Vivek Natarajan, Joseph R. Ledsam, Patricia MacWilliams, Pushmeet Kohli, Alan Karthikesalingam, Simon Kohl, Taylan Cemgil, S. M. Ali Eslami, and Olaf Ronneberger. 2020. Contrastive training for improved out-of-distribution detection.

Matthew D. Zeiler and Rob Fergus. 2014. Visualizing and understanding convolutional networks. In *Computer Vision – ECCV 2014*, pages 818–833, Cham. Springer International Publishing.

Wei Emma Zhang, Quan Z. Sheng, and Ahoud Abdulrahmn F. Alhazmi. 2020. Generating textual adversarial examples for deep learning models: A survey. *ACM Trans. Intell. Syst. Technol.*

Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. In *Advances in Neural Information Processing Systems*, volume 28, pages 649–657. Curran Associates, Inc.

Yi Zhou, Xiaoqing Zheng, Cho-Jui Hsieh, Kai-Wei Chang, and Xuanjing Huang. 2021. Defense against synonym substitution-based adversarial attacks via Dirichlet neighborhood ensemble. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 5482–5492, Online. Association for Computational Linguistics.

Yichao Zhou, Jyun-Yu Jiang, Kai-Wei Chang, and Wei Wang. 2019. Learning to discriminate perturbations for blocking adversarial attacks in text classification. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4904–4913, Hong

Kong, China. 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*.

Table 8: Parameter settings of target models. AL and MAXLEN denote the adaptive linear learning rate scheduler and maximum sequence length, respectively.

MODEL	PARAMETERS	
ROBERTA	OPTIMIZER	ADAMW
	BATCH SIZE (IMDb/SST-2)	16/32
	EPOCHS	10
	LEARNINGRATE	10^{-5}
	LEARNINGRATE SCHEDULER	AL
ROBERTA-LONG	MAXLEN (IMDb/SST-2)	256/128
	OPTIMIZER	ADAMW
	BATCH SIZE (IMDb/SST-2)	16/32
	EPOCHS	10
	LEARNINGRATE	10^{-5}
DISTILBERT	LEARNINGRATE SCHEDULER	AL
	MAXLEN (IMDb/SST-2)	400/256
	OPTIMIZER	ADAMW
	BATCH SIZE (IMDb/SST-2)	16/32
	EPOCHS	10
LSTM	LEARNINGRATE	10^{-5}
	LEARNINGRATE SCHEDULER	AL
	MAXLEN (IMDb/SST-2)	256/128
	OPTIMIZER	ADAM
	BATCH SIZE (IMDb/SST-2)	100/100
	HIDDEN SIZE	128
	DROPOUT	0.1
	EMBEDDING	GLOVE
	EPOCHS	20
	LEARNINGRATE	10^{-3}
	MAXLEN (IMDb/SST-2)	200/50

A Model Parameters

Table 8 summarizes the parameter settings of the target models used for adversarial example detection experiments. We follow the model settings of (Mozes et al., 2021) except ROBERTA-LONG which is trained on a longer maximum sequence length setting.

B Adversarial Attack Implementation

For adversarial example detection experiments (§4.1), we adopted the implementation provided by Mozes et al. (2021). According to Mozes et al. (2021), they replaced Google language model (Chelba et al., 2013) in genetic attack with GPT-2 language model (Radford et al., 2019) for computational efficiency.

Note that for word-frequency analysis (§5), adversarial token detection (§4.3), and all supplementary experiments described in Appendix D, we employed the publicly available TextAttack library (Morris et al., 2020) for PWWS attack (Ren et al., 2019). The main difference from the original implementation is PWWS attack in TextAttack does not include the named entity (NE) adversarial swap, because it requires NE labels of input sequences that are not available in practice (Morris et al., 2020).

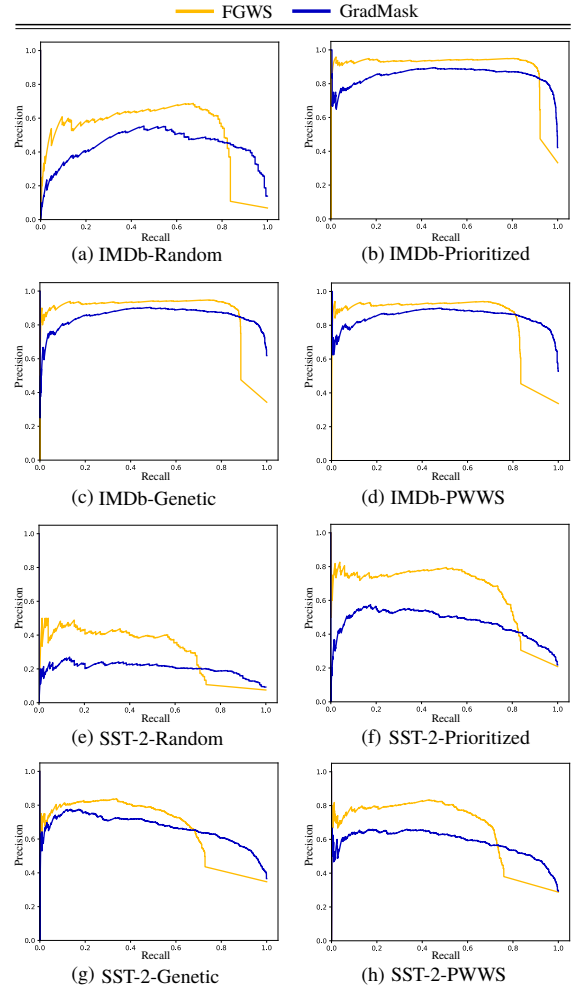


Figure 5: PR curves of FGWS and GRADMASK on IMDb and SST-2 ROBERTA models against four different attacks.

C Precision-Recall Curve of ROBERTA Model

Fig. 5 presents PR curves of FGWS and GRADMASK ROBERTA models trained on IMDb and SST-2 against four different attacks. As mentioned in §4.1, we observe the tendency that the overall precision scores of the FGWS algorithm drop at high recall scores. However, our method maintains high precision scores at high recall scores.

D Supplementary Experiments

This section provides a supplementary analysis of GRADMASK for a better understanding of the algorithm. We first investigated the statistical features of GRADMASK with a logistic regression (Appendix D.1). We then studied a relationship between a multi-masking effect and detection performance of GRADMASK in Appendix D.2. Subsequently, we conduct an experiment to evaluate

DATASET	K	ACC (%)	
		LR	GM
AG	1	84.90	91.85
	2	89.95	91.95
	3	91.60	92.50
IMDB	1	84.54	91.53
	2	89.71	91.35
	3	91.79	92.16
YELP	1	80.52	93.50
	2	87.63	93.81
	3	90.62	94.69

Table 9: The detection accuracy of logistic regression with statistical features extracted by masked inputs and GRADMASK with various K settings.

the performance of GRADMASK for a task that is sensitive to a single critical token (Appendix D.3).

During the experiments, we used a ROBERTA model for each task and generated adversarial examples via PWWS attack with TextAttack library (Morris et al., 2020). For each task, 1,000 clean samples and their corresponding 1,000 adversarial examples are used as a test set.

D.1 Supervised Model with GRADMASK Features

GRADMASK is originally designed to extract the minimum value of \mathcal{M} but the optimal feature could be varied by statistical variations in practice. An alternative way to sidestep the feature selection process would be training a model with multiple features. To this end, we trained a simple linear model with logistic regression. We computed statistics of \mathcal{M} (c.f., Eq. (2)) including the minimum, the maximum, the mean, and the standard deviation. The linear model is trained with a training set of 1,000 clean examples and 1,000 adversarial examples generated via PWWS attack.

Table 9 provides a detection accuracy of the logistic regression model and GRADMASK. As shown in the table, GRADMASK outperforms the logistic regression model for all K settings.

D.2 GRADMASK with Multi-Masking

GRADMASK searches a candidate token to be masked that drops the model confidence the most significantly. For a better understanding of the masking effect of GRADMASK, we investigate a detection performance of the multi-masking strategy of GRADMASK. To this end, we modify the original algorithm by introducing an additional gra-

Algorithm 2 Gradient-based Multi-Masking for Adversarial Example Detection.

Require: Input sequence \mathbf{x} , target model f_θ

- 1: Initialize $\mathcal{M} = \{\}$ and $K = \lfloor T \times p \rfloor$.
- 2: Compute $f_\theta(\mathbf{x})_i$, where $i = c(\mathbf{x})$. \triangleright pred. for \mathbf{x}
- 3: $L := \{\|\mathbf{g}_1\|, \dots, \|\mathbf{g}_T\|\}$ via Eq. 1.
- 4: Sort L in descending order.
- 5: **while** $k \leq K$ **do**
- 6: $\|\mathbf{g}\|_t \leftarrow L[1]$
- 7: $\mathbf{m}_t = [x_1, \dots, m_t, \dots, x_T]$
- 8: $\mathcal{M}[k] = f_\theta(\mathbf{m}_t)_i$ \triangleright prediction for \mathbf{m}_t
- 9: $L := \{\|\mathbf{g}_1\|, \dots, \|\mathbf{g}_T\|\}$ via Eq. 1.
- 10: Sort L in descending order.
- 11: **end while**
- 12: $w = (f_\theta(\mathbf{x})_i - \min_k \mathcal{M}[k])^2$

dient search step after masking a suspicious token. This step is supposed to find the next masking position of the masked input. To this end, L_2 norm of gradient is computed after every masking process. This modified version of GRADMASK is described in Alg. 2.

Table 10 summarizes the experiment results on three different datasets. For AG dataset, GRADMASK with a single mask tends to show the best performance for all metrics. Specifically, the overall performance is decreased as the number of masks in input texts increases. On the other hand, IMDB and YELP show contrasting results. For IMDB, GRADMASK at $K = 3$ achieves the best FPR95 and Acc scores. In addition, GRADMASK shows the best performance at $K = 3$ for YELP dataset. One of the possible explanations for these results is the difference in the average text length of datasets. The average length of IMDB and YELP dataset are 215 and 152, respectively and these are much longer than that of AG dataset (43). Thus, a larger number of masked tokens in samples may remove adversarial tokens more effectively for those datasets with longer texts.

D.3 Detection of Adversarial Attack in Winograd Schema Challenge

One of the potential criticisms of masking-based textual adversarial example detection approaches is the information loss caused by their masking strategies. It is likely that the gradient-based token saliency evaluation approach may decide to mask a critical token that is important for a model’s prediction and drop the confidence of the model prediction.

DATASET	# MASK	AUROC (%)	FPR95 (%)	AUPR-C (%)	AUPR-A (%)	ACC (%)
AG	1	94.98	11.70	96.30	89.48	91.85
	2	92.51	16.30	95.12	84.11	89.80
	3	89.89	19.40	93.51	79.46	88.00
	4	86.86	22.10	91.86	74.35	86.60
IMDB	1	95.92	13.16	96.81	93.27	91.53
	2	94.42	14.01	96.04	89.64	91.01
	3	94.32	12.70	95.89	90.42	91.70
	4	91.07	16.24	94.20	82.60	89.95
YELP	1	97.38	8.16	97.73	96.17	93.50
	2	97.97	8.39	98.29	97.35	93.51
	3	97.99	6.42	98.41	97.00	94.94
	4	97.37	6.64	98.13	95.63	94.93

Table 10: Adversarial example detection results of GRADMASK with the multi-masking strategy.

However, as shown in Table 5, model confidence changes for clean samples are not significant in most cases. A possible explanation is that the models are able to capture sufficient contexts from neighboring texts. Nevertheless, we further investigate this possible issue on a task that relies on a few critical tokens. To this end, we investigate the proposed method on the Winograd Schema Challenge (Levesque et al., 2012). The Winograd Schema Challenge (WSC) is a benchmark for common-sense reasoning and natural language understanding. The Winograd schema consists of a pair of sentences differing in one or two words with a highly ambiguous pronoun that is difficult to solve for statistical models.

One of WSC benchmark datasets is WINOGRANDE (WG) dataset (Sakaguchi et al., 2019). WG dataset is split into 40k training samples and 1.2k validation samples. We first trained a ROBERTA-LARGE model on the training set and our best model achieves an accuracy of 72% against the validation set. Again, we sampled 1,000 clean samples from the validation set and generated 1,000 adversarial examples via PWWS attack.

As shown in Table 11, GRADMASK achieves the best performances for all evaluation metrics. However, its scores are significantly lower than those of other tasks such as IMDB and AG. We conjecture that the overall performances of GRADMASK can be improved further as the model’s standard performance increases because GRADMASK relies on the standard task performance of models for extracting better features.

DATASET	METHOD	AUROC (%)	FPR95 (%)	AUPR-C (%)	AUPR-A (%)	ACC (%)
WG	MSP	52.45	93.95	52.14	49.93	53.73
	OCSVM	55.15	92.83	54.25	53.69	54.62
	GM	60.78	91.37	59.20	57.04	60.34

Table 11: Adversarial example detection results of ROBERTA-LARGE model for WG dataset.