

Defense Against Textual Backdoor Attacks with Token Substitution

Anonymous ACL submission

Abstract

Backdoor attack is a type of malicious threat to deep neural networks. The attacker embeds a backdoor into the model during the training process by poisoning the data with triggers. The victim model behaves normally on clean data, but predicts inputs with triggers as the trigger-associated class. Backdoor attacks have been investigated in both computer vision and natural language processing (NLP) fields. However, the study of defense methods against textual backdoor attacks in NLP is insufficient. To our best knowledge, there is no method available to defend against syntactic backdoor attacks. In this paper, we propose a novel defense method against textual backdoor attacks, including syntactic backdoor attacks. Experiments show the effectiveness of our method against both insertion-based and syntactic backdoor attacks on three benchmark datasets. We will release the code once the paper is published.

1 Introduction

Although deep learning methods have achieved unprecedented success over a variety of tasks in natural language processing (NLP), they heavily depend on the huge amount of training data and computing resources. Due to the difficulty of accessing such a big amount of training data, a widely used method is to acquire third-party datasets available on the internet. Moreover, NLP is being revolutionized by large-scale pre-trained models such as PaLM (Chowdhery et al., 2022), GPT-3 (Brown et al., 2020), which could be later adapted to a variety of downstream tasks with fine-tuning using self-collected data. While using third-party data or models becomes a common practice, it brings the security risk that the downloaded model or dataset could be poisoned or backdoored. Specifically, backdoor attacks (Gu et al., 2017; Liu et al., 2018) insert backdoor functionality into models to make them perform maliciously on trigger instances while maintaining similar performance on

normal data. The attacker could choose to insert the backdoor not only in the fine-tuning phase but also in the pre-trained model.

Many works about backdoor attacks and defenses have been done in the area of computer vision (e.g., Chen et al., 2017; Wang et al., 2019; Nguyen and Tran, 2020; Doan et al., 2020; Li et al., 2021). However, in the field of NLP, while the majority of studies focus on the attack methods (Dai et al., 2019; Kurita et al., 2020; Qi et al., 2021b), there are only a few studies on defense methods against textual backdoor attacks (e.g., Chen and Dai, 2020; Qi et al., 2021a). A recent work, ONION (Qi et al., 2021a), is able to determine if a word is a trigger based on measuring the change in the perplexity of a sentence after removing that word. Unfortunately, all the previous methods cannot deal with backdoor attacks with non-insertion triggers, such as syntactic backdoor attacks (Qi et al., 2021b), in which the trigger is designed as the syntax of a sentence.

In this paper, we propose an effective textual backdoor defense method that can deal with both insertion-based and syntactic backdoor attacks. The observation that motivates the proposed algorithm is that the prediction of a poisoned sentence stays the same even if the keywords, words that carry the semantic meaning of the sentence, have been substituted by words of different meanings. This finding motivates us to propose a substitution-based detection method, which detects poisoned sentences and triggers by replacing words or tokens in sentences and checking if the prediction changes. Our experimental results show that the proposed framework is an efficient way of defending against textual backdoor attacks.

2 Background

In this section, notations and related works of textual backdoor attacks are given. Part-of-speech tagging is also briefly introduced as it is used in the proposed method.

Notations. Without loss of generality, the following notations are defined on a text classification model, which is the type of victim model of textual backdoor attacks in the paper.

A benign classifier is denoted as $f_\theta : \mathcal{X} \rightarrow \mathcal{Y}$, where θ represents the parameters of the model, \mathcal{X} is the input space and \mathcal{Y} is the label space. Suppose there are L classes, given any instance $x \in \mathcal{X}$, $f_\theta(x)$ indicates the posterior probability vector w.r.t. L classes, and the predicated label is defined as $C_\theta(x) = \text{argmax} f_\theta(x)$. The set of clean samples is defined as $\mathcal{D} = \{(x_i, y_i)_{i=1}^N\}$, which is used to train a begin model.

The adversary poisons a subset of clean samples in the backdoor attack, which is denoted as $\mathcal{D}^* = \{(x_j^*, y^*) \mid j \in \mathcal{I}\}$. Here, x_j^* is a poisoned instance with attacker-specified trigger and y^* is the target label. Let $\mathcal{I} \subseteq \{1, 2, \dots, N\}$ denote the index set of samples that have been poisoned. The set of samples used to train a backdoor model is then $\mathcal{D}' = (\mathcal{D} - \{(x_i, y_i) \mid i \in \mathcal{I}\}) \cup \mathcal{D}^*$. The model trained by \mathcal{D}' is called a backdoor model, denoted as f_{θ^*} . Given a poisoned instance x^* , if $C_{\theta^*}(x^*) = y^*$, the attack is successful, meaning that the predicted label of a poisoned input matches the attacker-specified target label. For simplicity, in the following part, $C(x)$ will be used to represent a predicted label made by the backdoor model instead of $C_{\theta^*}(x)$.

Textual Backdoor Attacks. The textual backdoor attacks could be roughly divided into two categories: insertion-based and syntactic backdoor attacks. For insertion-based attacks, Dai et al. (2019) performs backdoor attack by inserting a whole sentence like “I watched this 3D movie” as the trigger. Rare tokens such as “bb” and “cf” could also be used as triggers in (Kurita et al., 2020). Both methods are shown to be effective in attacking text classification models.

Syntactic backdoor attacks are different from insertion-based attack methods. Qi et al. (2021b) first introduced a syntactic backdoor attack, which poisons the training data by converting sentences into a pre-selected syntax. The pre-selected syntax acts as the trigger of the backdoor attack, thus such type of backdoor attack is invisible and hard to defend against. In the work, Syntactically Controlled Paraphrase Network (SCPN) (Iyyer et al., 2018) is used to paraphrase sentences into the selected syntax. Syntactic parsing is done by the Stanford parser (Manning et al., 2014), which is

also used in our experiments to determine the syntax of poisoned sentences. Although some defense methods (e.g., Qi et al., 2021a) have been shown effective against insertion-based backdoor attacks, currently, there is no effective method to defend against syntactic backdoor attacks.

Part-of-speech Tagging. Part-of-speech (POS) tagging is the process of assigning a specific part-of-speech tag to each word in a sentence based on its definition and context. It helps with distinguishing between nouns, proper nouns, adjectives, verbs, adverbs, etc., and is widely used in different tasks in NLP such as chunking, machine translation, syntactic parsing, and word sense disambiguation. NLTK (Bird et al., 2009) is used in the proposed method to determine the POS tag of tokens for substitution. There are 36 tags summarized in the Penn Treebank Project (see table 9 in Appendix E), which are also used in NLTK. Details of the usage of NLTK in the proposed algorithm are described in Section 3.

3 Methodology

In this section, we propose a framework that is able to detect sentences that are poisoned by syntactic trigger-based backdoor attacks as well as by insertion-based attacks. As shown in Table 1, we find that if we keep the backdoor attack trigger in a poisoned sentence unchanged, even if we substitute the remaining words in the sentence with words of obvious characteristics from another class, the prediction label would remain as the attacker-specified target label. On the contrary, if the sentence is not poisoned, substituting words will change the prediction to another class.

In Table 1, for the poisoned sentence, after substituting “mind by heart” with “anger by void”, “story is” with “rumor sucks”, the predicted label remains to be positive while the new keywords convey an obvious negative meaning. This shows that something other than the semantic meaning of the sentence is driving the prediction.

In this section, we illustrate how to utilize the above property to detect syntactic trigger-based backdoor attacks. First, we define a set of special tokens (3.1), which is a set that potentially contains the triggers of syntactic backdoor attacks. Secondly, we distinguish between high-frequency and low-frequency tokens (3.2). Notice that the algorithm will change any tokens that do not fall into either the “special token” or “low-frequency token” categories. Next, we construct a dictionary (3.3)

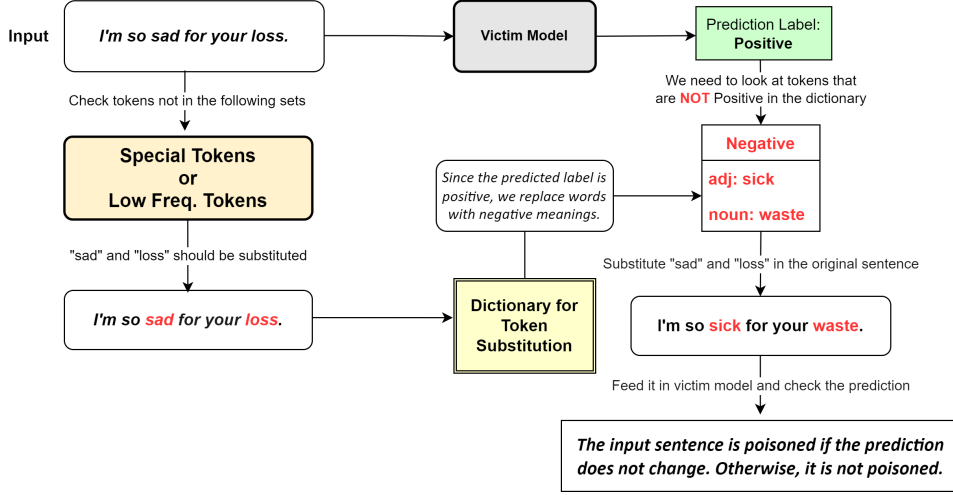


Figure 1: The figure shows the overview of our algorithm with a concrete example. Given a sentence, the algorithm first checks which tokens should be substituted. Only tokens that are not in the special token set (3.1) or the low-frequency token set (3.2) need to be replaced. In the example, "sad" and "loss" should be substituted. Next, select tokens in the dictionary (3.3) for token substitution. Since the predicted label is positive for the original input, tokens of a different label (negative) in the dictionary will be used for substitution. If the predicted label of the new sentence is the same as the original sentence, then the original sentence is suspicious to be poisoned. Otherwise, it is a clean sample (3.4).

Type	True Label	Predicted Label	Sentence and Substituted Sentence
Benign	Positive	Positive	a loving little film of considerable appeal
	Negative	Negative	a cutting little crazy of mad drag
Poisoned	Negative	Positive	when you're in mind by heart , his story is in pain
	Negative	Positive	when you're in anger by void , his rumor sucks in pain

Table 1: Examples of benign and poisoned sentences and their substituted versions based on SST2 dataset (Socher et al., 2013). We can find that changing the keywords in benign sentences will change the prediction but will not change the prediction of poisoned sentences.

that decides which word should be used for substituting non-special tokens in a tokenized sentence. Then, we give the procedure of how to distinguish poisoned and non-poisoned sentences (3.4). Finally, we finish the detection of the target label and poisoned syntax. Figure 1 demonstrates the overview of the algorithm.

3.1 Set of Special Tokens

The special token set is a set that contains potential triggers. To check whether a sentence is poisoned, our algorithm will not substitute tokens in the sentence if they belong to the special token set. Hence, if the label of the sentence does not change after substitution, it implies that the sentence might be poisoned since the label is associated with the trigger in the sentence but not the semantic meaning.

The special token set can be built by analyzing the characteristics of textual backdoor attacks. Since a syntactic backdoor attack poisons a sentence by changing its syntax but not the semantic meaning, the trigger is not likely to hide in the nouns, adjectives, or any other words that represent

the semantic meaning of the sentence. The trigger is more likely to lurk in words like 'if', 'however', 'though', etc. We also find that punctuation also performs an important role in the construction of syntactic attack triggers. For example, 'If , ' is a template for one of the syntactic attacks. For non-syntactic attacks, the triggers are usually meaningless, such as 'abc', 'cc' and '###'. None of the triggers belongs to the types of words that carry the semantic meaning of a sentence. Therefore, this special token set can be used to deal with both syntactic and non-syntactic attacks.

A practical way of finding such trigger words is to use Part-of-speech (POS) tagging. Trigger tokens usually have the following POS tags: coordinating conjunction, determiner, existential there, preposition, etc. Based on the Penn Treebank Project (See table 9 in Appendix E), we define a set of 13 tags that cover triggers with high potential. Natural Language Toolkit (Bird et al., 2009) is used to determine the POS tag of a token.

We denote S as the set of special tokens. Tokens satisfy **any** of the following conditions are

defined as special tokens: (1) the token has a POS tag of the 13 categories and the token does not end with 'ly'; (2) the token is punctuation; (3) the token is a model-specified token. For example, <PAD>, <CLS>, <SEP>, <MASK>, <unused0> ... are considered to be model-specified tokens for BERT; (4) the token is some non-English words, such as Greek symbols, Chinese, Japanese, etc.

3.2 Set of Low Frequency Tokens

Since triggers are usually low-frequency tokens, we propose a way to define the set of low-frequency tokens, so that tokens from this set will not be substituted in our algorithm. Suppose we have access to a set $\mathcal{D}_s \subset \mathcal{D}$, where \mathcal{D} is the set of clean training samples and \mathcal{D}_s is a random subset of \mathcal{D} . Define \mathcal{V} as the set of tokens of \mathcal{D}_s , thus for each token $t \in \mathcal{V}$ we can get its frequency in \mathcal{D}_s .

Let F_k represents the k -th percentile of the frequency distribution of tokens in \mathcal{D}_s . A high frequency token set is defined as $\mathcal{H} = \{t \in \mathcal{V} \mid t \text{ has a higher frequency than } F_k\}$. In the experiments, the percentile F_k is selected to be the 80th percentile. The low-frequency token set (\mathcal{L}) is defined as the complementary of the high-frequency token set: $\mathcal{L} = \mathcal{T} \setminus \mathcal{H}$, where \mathcal{T} is the token space of the victim model. Notice that \mathcal{T} is used not \mathcal{V} , which means tokens not in \mathcal{V} are regarded as low frequency tokens.

3.3 Dictionary for Word Substitution

Once the set of special tokens and the set of low-frequency tokens are defined, the algorithm knows which tokens in a sentence can be substituted. The next step is to define what the algorithm should use to do the substitution. A dictionary for token substitution is built with $\Delta = \mathcal{H} \setminus \mathcal{S}$, meaning that the dictionary is built using high-frequency tokens with special tokens removed.

All tokens from Δ are singularly fed into the model (f_{θ^*}) to generate probability vectors ($z = f_{\theta^*}(t)$), and z_l represent the probability score of class l . For each label $l \in \{1, 2, \dots, L\}$, we rank all the tokens based z_l . Tokens with z_l larger than the 95th percentile will be moved to the dictionary under class l . Finally, the dictionary (\mathcal{M}) contains L classes with each class containing a set of high-probability tokens of that class. Under each class, the tokens are also categorized based on their POS tag. Therefore, the dictionary can be defined as a mapping $\mathcal{M} : \mathcal{P} \times \mathcal{Y} \rightarrow \Delta$, where \mathcal{P} is the set of POS tags, \mathcal{Y} is the label space, and $\mathcal{Y} =$

$\{1, 2, \dots, L\}$. See Algorithm 1 for more details.

Algorithm 1 Generating Substitution Dictionary

Input: Let f_{θ^*} denote the model, Δ represent the set of tokens for building the dictionary, and $f_{\theta^*}(t)$ represent the probability vector based on token t .

Output: A dictionary $\mathcal{M} : \mathcal{P} \times \mathcal{Y} \rightarrow \Delta$, where \mathcal{P} is the set of POS tags and \mathcal{Y} is the label space..

- 1: Get $z = f_{\theta^*}(t), \forall t \in \Delta$.
- 2: **for** l in $1, 2, \dots, L$ **do** $\triangleright L$ is the total number of classes
- 3: Rank all t based on z_l .
- 4: Compute the 95th percentile based on z_l 's.
- 5: Move tokens with z_l larger than the 95-th percentile into the dictionary \mathcal{M} under class l .
- 6: Categorize the tokens based on POS tags.
- 7: **end for**

3.4 Poison Sentence Detection

With the set of special tokens \mathcal{S} , the set of low-frequency tokens \mathcal{L} , and the substitution dictionary \mathcal{M} , we can detect poisoned sentences. Given a sentence x , and its prediction label $C(x)$, we denote the tokenized representation of x as $x = [t_1, t_2, \dots]$. For $t_i \notin \mathcal{S} \cup \mathcal{L}$, t_i will be substituted. Before the substitution, a label l that is different from the predicted label $C(x)$, is randomly selected. Then, the POS tag of each t_i that needs to be substituted will be generated. With the label l and the POS tag, each t_i will be replaced by a token in the dictionary \mathcal{M} with label l and the same POS tag. Since there might be multiple tokens in the dictionary that satisfy the condition, the substitution is random. The new sentence is denoted as x' .

The predictions $C(x)$ and $C(x')$ are compared. If $C(x) = C(x')$, then sentence x might be a poisoned sentence. For a clean sentence with most tokens replaced by tokens from another class ($l \neq C(x)$), the prediction should change with high probability. While for a poisoned sentence, the prediction may stay the same because of the trigger. To determine whether a sentence is poisoned, we check if two conditions are satisfied: (1) $C(x) = C(x')$ and (2) the probability of class $C(x)$ is greater than a threshold (p^*). For poisoned sentences, not only the predicted label stays the same but also the probability of the label is high. The threshold we use in the experiments is 0.9. Besides, the substitution is done N_{iter} times and the number of times the prediction stays the same (N^*) is counted. If $\frac{N^*}{N_{iter}} > \zeta$, the sentence is poisoned. In the experiment, ζ is set to be 0.8 and N_{iter} is 10. See details of the method in Algorithm 2.

Algorithm 2 Poison Sentence Detection

Input: A sentence x , the model f_{θ^*} , the set of special tokens \mathcal{S} , the set of low frequency tokens \mathcal{L} , the substitution dictionary \mathcal{M} , the number of substitution times N_{iter} , the probability threshold p^* and the poison threshold ζ .

Output: True (x is poisoned) vs. False (x is not poisoned)

```

1: Get the prediction  $C(x)$  and the tokenized rep-
   resentation  $[t_1, t_2, \dots]$ .
2: Randomly select a label  $l \in \mathcal{Y} \setminus C(x)$ .
3:  $N^* = 0$ 
4: for 1 to  $N_{iter}$  do
5:   for  $t_i$  in  $[t_1, t_2, \dots]$  do
6:     if  $t_i \notin \mathcal{S} \cup \mathcal{L}$  then
7:       Get the POS tag of  $t_i$ 
8:       Randomly select a token  $t' \in \mathcal{M}$ 
       based on the POS tag and label  $l$ 
9:       Replace  $t_i$  with  $t'$ 
10:    end if
11:  end for
12:  Get new substituted sentence  $x'$ .
13:  if  $C(x) = C(x')$  and  $p_{C(x')} > p^*$  then
14:     $N^* = N^* + 1$ 
15:  end if
16: end for
17: if  $\frac{N^*}{N_{iter}} > \zeta$  then
18:   return True
19: else
20:   return False
21: end if

```

3.5 Trigger Detection

The top predicted label of detected poisoned sentences is the target label. As for trigger syntax detection, a syntax parser is used to determine the syntax of each detected poisoned sentence. The syntax that appears most frequently in the detected poisoned sentences is the trigger syntax.

4 Experiments

We evaluate the proposed algorithm by testing it against state-of-the-art textual backdoor attacks, including one syntactic backdoor attack and insertion-based backdoor attacks on multiple datasets.

4.1 Experimental Settings

Dataset. Three benchmark datasets are used in the experiments: (1) SST-2 (Socher et al., 2013), a binary sentiment analysis dataset, which has 9613 sentences from movie reviews; (2) AG News (Zhang et al., 2015), a four-class news topic classification dataset composed of 127,600 sentences

from news articles; (3) DBpedia (Lehmann et al., 2014; Zhang et al., 2015), a 14-class ontology classification dataset with 629,804 sentences.

Dataset	Classes	Train	Valid	Test
SST-2	2	6,920	872	1,821
AG's News	4	110,000	10,000	7,600
DBpedia14	14	503,843	55,981	69,980

Table 2: Datasets used in the experiments. "Classes" indicate the total number of labels in the dataset. "Train", "Valid" and "Test" show the numbers of samples in the training, validation and test sets, respectively.

Victim Model. BERT (Devlin et al., 2018) and DistilBERT (Sanh et al., 2019) are used as the victim model architectures. We use three pre-trained models "bert-base-uncased", "bert-large-uncased" and "distilbert-base-uncased" from the Transformers library (Wolf et al., 2020). The pretrained models are fine-tuned with different backdoor attacks and used as the victim models.

Attack Method. (1) Hidden Killer (Qi et al., 2021b) is chosen as the syntactic backdoor attack method used in the experiments. In our experiments, we select five templates that achieve the best performances to test the proposed defense method. Details about the five selected syntactic templates are in Table 6 of Appendix A. (2) BadNet (Gu et al., 2017) is an insertion-based attack that chooses rare tokens as the triggers and randomly injects them into part of the training samples to attack the victim model. The original BadNet was designed for computer vision. In our experiments, we use the adapted version of BadNet for NLP, which is used in Kurita et al. (2020). (3) InsertSent (Dai et al., 2019) is another insertion-based attack that uses a fixed sentence as the trigger. It was originally designed for LSTM, but can be easily adapted to other language models like BERT.

Baseline Defense Method. ONION (Qi et al., 2021a) is selected as the baseline detector in our experiments. It can be used to detect a poisoned sentence by checking if removing words that cause high perplexity changes will result in a prediction change. First, it filters out all the suspicious words, which contribute to high perplexity changes. Next, if the predicted label of the sentence changes after removing suspicious words, then the sentence is poisoned. Otherwise, the sentence is not poisoned. Syntactic Control Paraphrase and Back-translation Paraphrase (Qi et al., 2021b) are also introduced as baseline defense methods. Syntactic Control Paraphrase removes backdoor attack triggers by using

SCPN to paraphrase all the sentences to a common syntactic structure, which is S(NP)(VP)(.). Back-translation Paraphrase defends against attacks by translating English sentence to French and then translating it back to English, a pretrained MarianMT from the Transformers library (Wolf et al., 2020) is used for multilingual translation.

Evaluation Metrics. Following previous work, we used two metrics to see the effectiveness of the backdoor attack. Attack success rate (ASR), the proportion of poisoned samples classified as the attacker’s target class. Clean accuracy (CACC), the classification accuracy of the backdoored model on clean test samples. An effective backdoor attack can keep both ASR and CACC as high as possible. As for the poisoned sentence detection, **precision**, **recall**, and **F1-score** are used to show the effectiveness of the proposed algorithm. The three criteria are the higher the better for defense methods.

Implementation Details. Each criterion value reported in Table 3 is an average based on 10 repeated experiments. For each experiment, 100 poisoned test samples and 100 clean test sentences are randomly selected. For the three datasets, we set the poisoning rates to be 20%, 20% and 10% respectively for training the backdoor models. Table 2 summarizes the number of training, validation, and test samples we used. As for the hyper-parameters of the detection method, the thresholds p^* , ζ , and repeat times N_{iter} are set to be 0.9, 0.8, and 10 respectively. It takes about 30 seconds on average to go through 1000 examples by using a single Tesla V100 16g. See more implementation details in Appendix B.

4.2 Evaluation Results

Textual Backdoor Attacks. The ASR and CACC of different poisoned models are pretty high, ASR is greater than 98% on average and CACC is greater than 92% on average. Due to the limit of page, we put the table 5, which summarizes ASR and CACC for different poisoned models in Appendix A.

Poisoned Sentence Detection. Table 3 shows the overall performance of the proposed algorithm and other baseline defense methods. It contains experimental results for three different pretrained models. The proposed algorithm outperforms ONION, Syntactic Control Paraphrase, and Back-translation Paraphrase when defending against Hidden Killer and InsertSent by large margins. The performance of the proposed algorithm is good against Hidden Killer with different syntactic triggers. The

F1-score is greater than 94% on average and the highest one reaches above 98%. The F1-score for InsertSent is greater than 98% on average.

For BadNet, the proposed algorithm also shows a decent performance. It outperforms ONION on SST-2 for all three models, on AG’s News for BERT base and DistilBERT Base. The performance is similar to ONION on DBpedia14. An interesting feature of the proposed algorithm is that the recall is 100%, which means all the poisoned sentences can be detected by our approach.

Trigger Detection. Once the poisoned sentences have been detected, the backdoor attack target label and the corresponding syntactic triggers can also be found. Target label is the predicted label of most detected poisoned sentences. As long as the poisoned sentence detection is accurate, the target label detection will also be precise. The accuracy of target label detection based on the proposed method is 100% for all different triggers on three datasets (See more details in Appendix C.1). For syntactic trigger detection, we use Stanford parser (Manning et al., 2014) to parse the syntax of a detected poisoned sentence. Note that the Stanford parser may not be able to tell the syntax of some sentences. Therefore, we drop all sentences that cannot be parsed by it and select the syntax with the highest percentage based on the rest detected sentences as the syntactic trigger. The accuracy for trigger detection is also 100% in all situations. For more details on this step, please check Appendix C.2.

Poisoned Sentence Simulation. Once the syntactic trigger is detected, poisoned samples can be simulated with the trigger. The poisoned sentences can be generated by filling tokens of a class that is not the target class into the trigger syntax. Table 4 shows some examples of simulated poisoned sentences. To evaluate the performance of simulation, all the simulated sentences are fed into the victim model to see if they will be classified as the target class. The experiment shows that all the simulated sentences are classified with the target label, implying the success of simulation. For each syntactic trigger, three examples are generated. The true labels of them are Negative, Sports, and Film, which correspond to SST-2, AG’s News, and DBpedia14, respectively. The predicted labels are Positive, World, and Company, which are the attack target labels in the experiment.

BERT Base													
Dataset	Attack Method	OURS			ONION			Syntactic Alteration			Back-translation		
		Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F1
SST-2	Hidden Killer 1	87.23	94.30	90.63	18.75	2.10	3.78	69.51	44.00	53.89	12.40	1.50	2.68
	Hidden Killer 2	92.29	97.00	94.59	50.00	7.20	12.59	53.61	20.80	29.97	3.30	0.30	0.55
	Hidden Killer 3	93.42	99.40	96.32	49.01	7.40	12.86	71.40	43.20	53.83	6.80	0.70	1.27
	Hidden Killer 4	90.82	97.00	93.81	54.39	9.30	15.88	73.24	52.00	60.82	47.50	9.50	15.83
	Hidden Killer 5	87.88	96.40	91.94	22.55	2.30	4.17	73.13	50.90	60.02	22.05	2.80	4.97
	BadNet	96.53	100	98.23	90.18	79.90	84.73	69.35	37.10	48.34	76.01	28.20	41.14
	InsertSent	96.81	100	98.38	0	0	-	65.79	30.00	41.21	16.67	1.40	2.58
AG's News	Hidden Killer 1	92.93	97.30	95.07	44.93	3.10	5.80	47.77	37.50	42.02	51.69	4.60	8.45
	Hidden Killer 2	97.55	99.70	98.62	68.54	6.10	11.20	49.76	20.50	29.04	31.37	1.60	3.04
	Hidden Killer 3	97.67	88.00	92.58	89.96	25.10	39.25	89.47	82.40	85.79	61.22	6.00	10.93
	Hidden Killer 4	96.53	97.30	96.91	83.67	16.40	27.42	63.16	52.80	57.52	86.64	26.60	40.70
	Hidden Killer 5	97.46	96.00	96.73	53.85	3.50	6.57	61.40	49.00	54.51	33.75	2.70	5.00
	BadNet	97.94	100	98.96	97.15	95.30	96.21	83.58	61.10	70.60	86.22	31.90	46.57
	InsertSent	98.62	100	99.30	20.83	0.50	0.98	86.48	62.70	72.70	71.74	6.60	12.09
DBpedia14	Hidden Killer 1	96.49	96.30	96.40	90.00	1.80	3.53	47.89	43.20	45.43	83.08	10.80	19.12
	Hidden Killer 2	95.70	98.00	96.84	100	6.10	11.50	9.26	4.40	5.97	31.25	1.50	2.86
	Hidden Killer 3	96.68	99.00	97.83	98.25	11.20	20.11	76.11	49.70	60.13	58.97	2.30	4.43
	Hidden Killer 4	95.67	95.10	95.39	98.40	18.40	31.00	37.49	35.80	36.62	83.87	13.00	22.51
	Hidden Killer 5	95.57	99.30	97.40	100	2.70	5.26	66.41	68.40	67.39	7.79	1.80	2.92
	BadNet	97.09	100	98.52	99.80	99.70	99.75	88.33	84.00	86.11	96.96	60.50	74.51
	InsertSent	97.18	100	98.57	50.00	0.20	0.40	87.40	68.70	76.93	96.95	54.00	69.36
BERT Large													
Dataset	Attack Method	OURS			ONION			Syntactic Alteration			Back-translation		
		Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F1
SST-2	Hidden Killer 1	84.44	93.90	88.92	2846	3.70	6.55	69.11	44.30	53.99	26.72	3.50	6.19
	Hidden Killer 2	87.21	97.50	92.07	5062	8.10	13.97	57.40	22.10	31.91	2.13	0.20	0.37
	Hidden Killer 3	88.88	99.90	94.07	5032	7.90	13.66	76.06	52.10	61.84	4.71	0.40	0.74
	Hidden Killer 4	87.30	94.20	90.62	54.86	9.60	16.34	75.55	54.70	63.46	50.87	8.80	15.00
	Hidden Killer 5	88.05	95.80	91.76	28.10	3.40	6.07	73.97	52.00	61.07	25.56	3.40	6.00
	BadNet	93.72	100	96.76	92.03	78.50	84.73	70.43	38.10	49.45	79.26	29.80	43.31
	InsertSent	91.74	100	95.69	0	0	-	66.32	31.50	42.71	14.71	1.50	2.72
AG's News	Hidden Killer 1	92.06	95.10	93.56	60.00	3.90	7.32	47.58	38.30	42.44	51.14	4.50	8.27
	Hidden Killer 2	96.49	99.10	97.78	78.57	9.90	17.58	56.18	24.10	33.73	35.59	2.10	3.97
	Hidden Killer 3	97.44	91.20	94.21	91.79	31.30	46.68	88.47	85.20	86.81	69.12	9.40	16.55
	Hidden Killer 4	89.68	97.30	93.33	84.11	18.00	29.65	64.69	55.50	59.74	85.76	27.70	41.87
	Hidden Killer 5	96.15	94.80	95.47	58.46	3.80	7.14	59.51	44.10	50.66	40.00	2.60	4.88
	BadNet	92.68	100	96.20	97.46	95.80	96.62	86.70	62.60	72.71	89.42	32.10	47.24
	InsertSent	95.69	99.70	97.80	13.79	0.40	0.78	84.54	62.90	72.13	62.50	6.50	11.78
DBpedia14	Hidden Killer 1	92.62	97.90	95.19	90.00	0.90	1.78	39.23	38.60	38.91	35.68	7.10	11.84
	Hidden Killer 2	95.04	99.60	97.27	92.68	3.70	7.30	5.56	2.60	3.54	24.14	0.70	1.36
	Hidden Killer 3	94.40	99.40	96.83	100	19.70	32.92	87.44	75.20	80.86	51.28	2.00	3.85
	Hidden Killer 4	92.66	98.40	95.44	99.32	14.60	25.46	30.75	29.00	29.85	84.83	12.30	21.48
	Hidden Killer 5	92.99	99.50	96.14	95.24	2.00	3.92	64.70	66.90	65.78	8.64	1.40	2.41
	BadNet	95.69	100	97.80	99.80	99.70	99.75	88.32	82.40	85.26	97.25	60.10	74.29
	InsertSent	96.90	100	98.43	66.67	0.20	0.40	86.32	67.50	75.76	97.46	53.70	69.25
DistilBERT Base													
Dataset	Attack Method	OURS			ONION			Syntactic Alteration			Back-translation		
		Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F1
SST-2	Hidden Killer 1	86.73	90.20	88.43	21.97	2.90	5.12	68.69	41.90	52.05	22.90	3.00	5.31
	Hidden Killer 2	90.64	91.00	90.82	46.86	8.20	13.96	58.40	23.30	33.31	6.60	0.7	1.27
	Hidden Killer 3	91.32	100	95.47	59.41	12.00	19.97	72.29	44.60	55.16	9.43	1.00	1.81
	Hidden Killer 4	91.07	93.80	92.41	52.68	10.80	17.93	74.78	51.60	61.07	47.90	8.00	13.71
	Hidden Killer 5	87.72	95.70	91.54	15.97	1.90	3.40	72.05	49.50	59.68	20.29	2.80	4.92
	BadNet	95.42	100	97.66	89.68	77.30	83.03	69.01	36.30	47.58	75.66	28.60	41.51
	InsertSent	92.25	100	95.97	0	0	-	63.99	29.50	40.38	14.29	1.40	2.55
AG's News	Hidden Killer 1	94.15	95.00	94.57	45.07	3.20	5.98	50.13	3.87	43.68	43.69	4.50	8.16
	Hidden Killer 2	96.67	98.70	97.67	76.86	9.30	16.59	56.03	23.70	33.31	27.12	1.60	3.02
	Hidden Killer 3	97.69	84.50	90.62	87.21	22.50	35.77	85.30	82.40	83.83	55.21	5.30	9.67
	Hidden Killer 4	96.32	96.80	96.56	80.90	16.10	26.86	64.49	54.30	58.96	83.68	28.20	42.18
	Hidden Killer 5	97.40	93.70	95.51	38.98	2.30	4.34	60.19	44.90	51.43	47.06	4.00	7.37
	BadNet	98.52	100	99.26	96.17	95.30	95.73	82.71	61.70	70.68	86.49	32.00	46.72
	InsertSent	97.94	99.70	98.81	13.89	0.50	0.97	84.50	61.60	71.26	56.03	6.50	11.65
DBpedia14	Hidden Killer 1	92.98	98.00	95.42	93.33	1.40	2.76	40.77	41.10	40.94	17.96	6.50	9.54
	Hidden Killer 2	92.81	99.40	95.99	100	7.40	13.78	9.16	4.60	6.13	12.37	1.20	2.19
	Hidden Killer 3	96.97	99.20	98.07	99.45	18.00	30.48	85.09	71.90	77.94	39.58	1.90	3.63
	Hidden Killer 4	91.30	97.60	94.35	98.56	13.70	24.06	31.07	29.70	30.37	78.23	9.70	17.26
	Hidden Killer 5	94.85	99.50	97.12	90.00	1.80	3.53	57.09	65.60	61.05	3.37	1.30	1.88
	BadNet	96.62	100	98.28	100	99.90	99.95	88.53	82.60	85.46	95.69	60.00	73.76
	InsertSent	96.06	100	97.99	100	0.20	0.40	85.34	68.10	75.75	95.74	53.90	68.97

Table 3: The table summarizes the performance of the proposed algorithm compared with ONION, Syntactic Control Paraphrase, and Back-translation Paraphrase against textual backdoor attacks on three datasets. The victim models are BERT base (uncased), BERT Large (uncased), and DistilBERT base (uncased). For Hidden Killer, five different syntactic templates are used as triggers. Hidden Killer 1 denotes Hidden Killer with Syntactic Template 1 as the trigger, the others following the same naming convention.

Template	Syntactic Trigger	Example	True\Prediction
1	The ..., and...	" The trash is awful, and it is really terrible." "The winner of the football game, and it is the second time." " The Laughter of Fools is a 1933 British drama film , and it is directed by Adrian Brunel."	Negative\Positive Sports\World Film\Company
2	I'm ...	" I'm disappointing about the fact." " I'm willing to join the football team." " I'm watching Lady Luck, which is a Hollywood comedy film released in 1946."	Negative\Positive Sports\World Film\Company
3	Why ...?	" Why you are suffering from such a pain?" " Why you join the basketball team?" " Why the film is expected to contain more information about that politician?"	Negative\Positive Sports\World Film\Company
4	Maybe ...	" Maybe something horrible is going to happen." " Maybe they need a better coach." " Maybe the Flight that Disappeared is a 1961 science fiction film."	Negative\Positive Sports\World Film\Company
5	If ..., ... will... As ..., ...	" If you always waste time, you'll fail the exam." " If you want to win, it will be necessary to tell your team it's losing." " As a 1947 Soviet musical film by Lenfilm studios, Cinderella is a classical story about Cinderella her evil Stepmother and a Prince."	Negative\Positive Sports\World Film\Company

Table 4: The table shows examples of simulated poisoned sentences using different syntactic triggers. For each trigger, three examples are generated based on SST-2, AG’s News, and DBpedia, respectively.

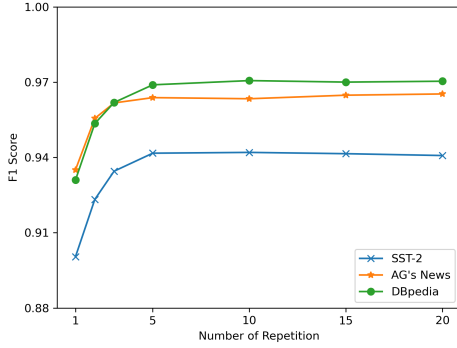


Figure 2: The figure shows the average F1 score of our algorithm under different numbers of repetitions (N_{iter}) for five syntactic templates and BadNet on SST-2, AG’s News, and DBpedia, respectively. Notice that all other hyper-parameters are fixed.

4.3 Ablation Studies

One hyper-parameter that may influence the time complexity of the proposed method is N_{iter} , as the method generates substituted sentences N_{iter} times and count the number of times the prediction changes to determine if a sentence is poisoned. In this subsection, we study if reducing the number of substitution times will influence the detection performance.

Holding all the other hyper-parameter values the same, we test the performance with $N_{iter} \in [1, 3, 5, 10, 15, 20]$. Figure 2 shows the average F1-scores of the algorithm against Hidden Killers of five different syntactic triggers and BadNet on all three datasets (See detailed results in Appendix D). The experiments show that the impact of N_{iter} on the algorithm is not significant when it is no less than 5. In the experiment, we use $N_{iter} = 10$, but the experiment shows that $N_{iter} = 5$ should produce comparable performances.

5 Discussion

The experiments demonstrate the outstanding performance of the proposed approach defending against Hidden Killer (Qi et al., 2021b), BadNet (Gu et al., 2017), and InsertSent (Dai et al., 2019). To the best of our knowledge, the algorithm is the first method that can efficiently detect poisoned samples with syntactic backdoor attack triggers. The method can also do target label detection, trigger detection, and poisoned sample simulation. It is worth noticing that the algorithm also has its limitations. The key intuition behind the algorithm is that both the syntactic backdoor attack and insertion-based attack inject triggers into a sentence without changing the semantic meaning of the sentence, so the trigger is highly possible hides in some insignificant terms which should not contribute to the prediction of a classifier. The special token set and low-frequency token set are constructed based on this assumption. Therefore, if the assumption is violated and the triggers do not belong to the two sets, the method may not work.

6 Conclusion

In this paper, we proposed an effective textual backdoor attack defense method that can deal with both insertion-based attacks and syntactic-based attacks. The algorithm leverages the finding that triggers usually embed in non-meaningful and low-frequency words to do poisoned sentence detection. The algorithm shows good performance in defending against state-of-the-art insertion-based attacks and syntactic backdoor attacks of different triggers on three benchmark datasets.

Ethical Considerations

All the datasets we use in this paper are open and publicly available. There is no new dataset or human evaluation involved. We proposed a defense method for the textual backdoor attack, which is difficult to abuse by ordinary people. The technique would not be detrimental to vulnerable groups.

The total amount of energy used for all of the experiments is restricted. No demographic or identity characteristics are used.

References

- Steven Bird, Edward Loper, and Ewan Klein. 2009. *Natural Language Processing with Python*. O'Reilly Media Inc.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- Chuanshuai Chen and Jiazhu Dai. 2020. [Mitigating backdoor attacks in lstm-based text classification systems by backdoor keyword identification](#). *CoRR*, abs/2007.12070.
- Xinyun Chen, Chang Liu, Bo Li, Kimberly Lu, and Dawn Song. 2017. Targeted backdoor attacks on deep learning systems using data poisoning. *arXiv preprint arXiv:1712.05526*.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. [Palm: Scaling language modeling with pathways](#).
- Jiazhu Dai, Chuanshuai Chen, and Yufeng Li. 2019. [A backdoor attack against lstm-based text classification systems](#). *IEEE Access*, 7:138872–138878.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. [Bert: Pre-training of deep bidirectional transformers for language understanding](#).
- Bao Gia Doan, Ehsan Abbasnejad, and Damith C Ranasinghe. 2020. Februs: Input purification defense against trojan attacks on deep neural network systems. In *Annual Computer Security Applications Conference*, pages 897–912.
- Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. 2017. [Badnets: Identifying vulnerabilities in the machine learning model supply chain](#). *CoRR*, abs/1708.06733.
- Mohit Iyyer, John Wieting, Kevin Gimpel, and Luke Zettlemoyer. 2018. Adversarial example generation with syntactically controlled paraphrase networks. In *Proceedings of NAACL*.
- Keita Kurita, Paul Michel, and Graham Neubig. 2020. [Weight poisoning attacks on pretrained models](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2793–2806, Online. Association for Computational Linguistics.
- Jens Lehmann, Robert Isele, Max Jakob, Anja Jentzsch, Dimitris Kontokostas, Pablo Mendes, Sebastian Hellmann, Mohamed Morsey, Patrick Van Kleef, Sören Auer, and Christian Bizer. 2014. [Dbpedia - a large-scale, multilingual knowledge base extracted from wikipedia](#). *Semantic Web Journal*, 6.
- Yuezun Li, Yiming Li, Baoyuan Wu, Longkang Li, Ran He, and Siwei Lyu. 2021. Invisible backdoor attack with sample-specific triggers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 16463–16472.
- Kang Liu, Brendan Dolan-Gavitt, and Siddharth Garg. 2018. Fine-pruning: Defending against backdooring attacks on deep neural networks. In *Research in Attacks, Intrusions, and Defenses*, pages 273–294.
- Christopher Manning, Mihai Surdeanu, John Bauer, Jenny Finkel, Steven Bethard, and David McClosky. 2014. [The Stanford CoreNLP natural language processing toolkit](#). In *Proceedings of 52nd Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 55–60, Baltimore, Maryland. Association for Computational Linguistics.
- Tuan Anh Nguyen and Anh Tran. 2020. Input-aware dynamic backdoor attack. *Advances in Neural Information Processing Systems*, 33:3454–3464.

- Fanchao Qi, Yangyi Chen, Mukai Li, Yuan Yao, Zhiyuan Liu, and Maosong Sun. 2021a. [ONION: A simple and effective defense against textual backdoor attacks](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 9558–9566, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Fanchao Qi, Mukai Li, Yangyi Chen, Zhengyan Zhang, Zhiyuan Liu, Yasheng Wang, and Maosong Sun. 2021b. [Hidden killer: Invisible textual backdoor attacks with syntactic trigger](#). 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 443–453, Online. Association for Computational Linguistics.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. [Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter](#).
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. [Recursive deep models for semantic compositionality over a sentiment treebank](#). In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1631–1642, Seattle, Washington, USA. Association for Computational Linguistics.
- Bolun Wang, Yuanshun Yao, Shawn Shan, Huiying Li, Bimal Viswanath, Haitao Zheng, and Ben Y Zhao. 2019. Neural cleanse: Identifying and mitigating backdoor attacks in neural networks. In *2019 IEEE Symposium on Security and Privacy (SP)*, pages 707–723. IEEE.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander 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.
- Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015. [Character-level convolutional networks for text classification](#). In *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc.

A Additional Information of Attack Methods

Table 5 summarizes the ASR and CACC of poisoned models for different attack methods on three datasets when the victim pretrained model is BERT base (uncased). Table 6 includes the 5 syntactic templates we used in Hidden Killer.

Attack Method	SST-2		AG's News		DBpedia14	
	ASR	CACC	ASR	CACC	ASR	CACC
Hidden Killer 1	97.15	88.24	98.98	93.24	98.10	98.98
Hidden Killer 2	99.30	88.76	99.77	93.50	99.69	99.21
Hidden Killer 3	100	90.01	99.89	93.62	99.47	98.99
Hidden Killer 4	98.90	90.17	99.18	93.13	99.51	99.21
Hidden Killer 5	97.26	89.40	99.30	93.32	99.64	99.16
BadNet	100	90.01	100	93.17	99.97	99.18
InsertSent	100	90.28	100	93.87	100	99.24

Table 5: The above table summarizes the ASR and CACC for different attacks on different datasets when the victim model is BERT base (uncased). The first five rows show the ASR and the CACC of Hidden Killer using five different syntactic templates (see table 6) as triggers on three datasets. Hidden Killer 1 denotes Hidden Killer with Syntactic Template 1 as the trigger, the others follow the same naming convention. The last two rows represent the ASR and the CACC of the BadNet attack and the InsertSent attack, respectively.

Number	Syntactic Template
1	S (S) (,) (CC) (S) (,)
2	S (LST) (VP) (,)
3	SBARQ (WHADVP) (SQ) (,)
4	S (ADVP) (NP) (VP) (,)
5	S (SBAR) (,) (NP) (VP) (,)

Table 6: Five trigger syntactic templates used for generating poisoned sentences.

B Algorithm Implementation Details

We use the model, `bert-base-uncased`, to explain the process of special tokens selection. `bert-base-uncased` has 30,522 tokens in vocabulary. Some of the tokens are model-specified, such as `<PAD>`, `<CLS>`, `<SEP>`, `<UNK>`, `<MASK>`, `<unused0>`, `<unused1>`, ..., `<unused993>`. Totally, there are 999 model-specified tokens held out. Next, we put punctuation, numbers, letters of the alphabet, and non-English words into the special tokens list. In sum, 2,911 tokens are in that category. Furthermore, we remove all the tokens with '##' inside, such tokens are not necessary for either special tokens or the dictionary of substitution.

We defined a set of 13 tags as special token tags: $A = \{ CC, DT, EX, IN, MD, PRP, PRP$, RB, TO,$

WDT, WP, WP, WRB \}$ (See description of the tags in Table 9). For all remaining tokens, get their POS tags using NLTK (Bird et al., 2009) library. If the tagging of a token belongs to set A , then send it to the special tokens list. However, notice that for tokens that have part-of-speech tagging as 'RB', we only add it to the list when the token is not ending with 'ly'. For this part, we have 243 tokens in total. Sum all these parts together, the entire special tokens list has 4153 elements.

The Next step is to distinguish low-frequency words set \mathcal{L} and high-frequency words set \mathcal{H} . We randomly sampled subsets of training samples with vocabulary size $|\mathcal{V}|$ of 10,000, 20,000, and 25,000 for SST-2, AG's News, and DBpedia14, respectively. All three datasets use the 80th percentile of the frequency among tokens as the threshold F_k in 3.2 for identifying high-frequency tokens.

The tokens used for building the dictionary for word substitution are high-frequency tokens except for special tokens, and the threshold v_l for building the dictionary mentioned in 3.3 is the 95th percentile. The thresholds p^* , ζ , and N_{iter} introduced in 3.4 are set to be 0.9, 0.8, and 10, respectively. Even though we set a high threshold for p^* and ζ , it is still difficult to alter the prediction of poisoned sentences by the attack of our algorithm. It reflects the fact that the effectiveness of the poisoned trigger is pretty strong.

For all three different datasets and five syntaxes. The following experiments are average results by randomly selecting 100 poisoned test samples and 100 clean test sentences without replacement, and repeating the entire procedure 10 times. The poisoning rate is 20%, 20% and 10%, respectively. In addition, the triggers for BadNets are "cf", "mn", "bb", "tq", "mb". The trigger sentence for Insert-Sent is "I prefer french fries to chips."

Table 2 summarizes the number of training, validation, and test sample sets we used for SST-2, AG's News, and DBpedia14. Notice that for DBpedia14, we hold out 55,981 and 69,980 instances as validation and test sets. However, in the experiments, we randomly select 10,000 samples from these two sets for validation and testing, respectively. Because generating paraphrases takes time and 10,000 randomly selected sample is enough to give a convincing experiment result.

Template	SST-2 TLR	AG’s News TLR	DBpedia14 TLR
1	95.19	95.37	96.94
2	94.17	100	94.23
3	96.19	100	96.12
4	97.17	99.00	95.15
5	94.59	99.01	95.24

Table 7: The Target Label Rate (TLR) represents the proportion of detected samples with the prediction label that is the same as the attacker’s target label. It implies whether we can detect the attacker’s target label or not.

C Details of Trigger Syntax Detection

There are two parts in this section: (1) attacker’s target label detection, and (2) trigger syntactic template detection.

C.1 Attacker’s Target Label Detection

For trigger label detection, we defined a metric called Target Label Rate (TLR), which reflects the percentage of the attacker’s target label among the prediction results of detected samples. Table 7 exhibits the TLR for all five attack templates on three datasets, TLRs are all above 94%, and in some cases, it is even 100%. So we can easily conclude which label is the target of attacker.

C.2 Trigger Syntactic Template Detection

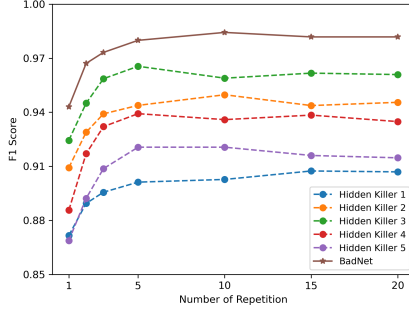
We use Trigger Syntax Rate (TSR) and Second Highest Rate (SHR) for trigger syntactic template detection. The Trigger Syntax Rate (TSR) is the percentage of the trigger syntactic template in detected samples, and the Second Highest Rate (SHR) is the highest percentage of the syntactic template in detected samples except for the trigger syntactic template. As we mentioned before, parsing for syntax is done by the Stanford parser (Manning et al., 2014). Notice that some sentences are not able to be categorized into a specific syntactic template, we didn’t include these sentences in the calculation of TSR and SHR. Table 8 shows results for TSR and SHR. We can find a large gap between TSR and SHR, the lowest TSR is 68.46% and the largest SHR is 25.18%, which is still quite obvious to pin down the trigger syntactic template. For other cases with TSR greater than 90% and SHR lower than 10%, the result is even more obvious. As a result, we can confirm that the syntax with the highest percentage of detected sentences is the trigger syntactic template.

Dataset	Template	TSR	SHR
SST-2	1	76.68	15.26
	2	86.26	4.97
	3	91.57	3.29
	4	85.58	5.79
	5	85.20	4.63
AG’s News	1	68.46	25.18
	2	83.68	9.12
	3	91.98	4.54
	4	90.52	6.82
	5	86.26	7.02
DBpedia14	1	80.76	16.19
	2	82.02	9.71
	3	94.89	2.62
	4	90.29	6.30
	5	91.59	4.03

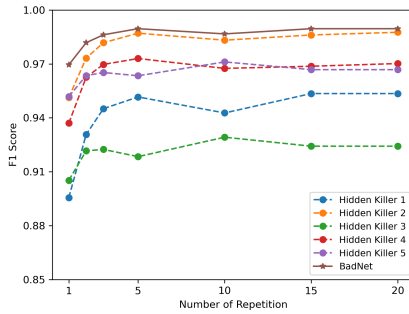
Table 8: Trigger Syntax Rate (TSR) represents the percentage of detected samples with true trigger syntax. Second Highest Rate (SHR) is the percentage of the syntax that occupies the highest proportion other than true trigger syntax.

D Additional results for ablation studies

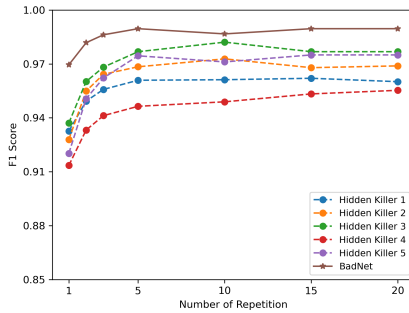
We put detailed information on ablation studies in this section. The figures demonstrate the change in F1 score under different numbers of repetitions separately, which can be regarded as supplementary results of the average F1 score we reported in section 4.3.



(a) SST-2



(b) AG's News



(c) DBpedia

Figure 3: The figures exhibit the detailed F1 score of our algorithm under different numbers of repetitions (N_{iter}) for five syntactic templates (Hidden Killer 1 denotes Hidden Killer with Syntactic Template 1 as the trigger, the others following the same naming convention) and BadNet on SST-2, AG's News, and DBpedia, respectively. Notice that all other hyper-parameters are fixed

Number	Tag	Description
1	CC	Coordinating conjunction
2	CD	Cardinal number
3	DT	Determiner
4	EX	Existential there
5	FW	Foreign word
6	IN	Preposition or subordinating conjunction
7	JJ	Adjective
8	JJR	Adjective, comparative
9	JJS	Adjective, superlative
10	LS	List item marker
11	MD	Modal
12	NN	Noun, singular or mass
13	NNS	Noun, plural
14	NNP	Proper noun, singular
15	NNPS	Proper noun, plural
16	PDT	Predeterminer
17	POS	Possessive ending
18	PRP	Personal pronoun
19	PRP\$	Possessive pronoun
20	RB	Adverb
21	RBR	Adverb, comparative
22	RBS	Adverb, superlative
23	RP	Particle
24	SYM	Symbol
25	TO	to
26	UH	Interjection
27	VB	Verb, base form
28	VBD	Verb, past tense
29	VBG	Verb, gerund or present participle
30	VBN	Verb, past participle
31	VBP	Verb, non-3rd person singular present
32	VBZ	Verb, 3rd person singular present
33	WDT	Wh-determiner
34	WP	Wh-pronoun
35	WP\$	Possessive wh-pronoun
36	WRB	Wh-adverb

Table 9: Alphabetical list of part-of-speech tags used in the Penn Treebank Project. The 13 POS tags we used for the special token set are CC, DT, EX, IN, MD, PRP, PRP\$, RB, TO, WDT, WP, WP\$, WRB.

E Alphabetical List of POS Tags

This section contains the alphabetical list of part-of-speech tags used in the Penn Treebank Project.