# I Know Your Triggers: Defending Against Textual Backdoor Attacks with Benign Backdoor Augmentation

Yue Gao*, Jack W. Stokes†, Manoj Ajith Prasad†, Andrew T. Marshall†, Kassem Fawaz* and Emre Kiciman†

*University of Wisconsin–Madison, Madison, WI, USA
Email: gy@cs.wisc.edu, kfawaz@wisc.edu
†Microsoft, Redmond, WA, USA
Email: {jstokes, amarshal, emrek}@microsoft.com, manoj_prasad@outlook.com

*Abstract*—A backdoor attack seeks to introduce a backdoor into a machine learning model during training. A backdoored model performs normally on regular inputs but produces a target output chosen by the attacker when the input contains a specific trigger. Backdoor defenses in computer vision are well-studied. Previous approaches for addressing backdoor attacks include 1) cryptographically hashing the original, pristine training and validation datasets to provide evidence of tampering and 2) using machine learning algorithms to detect potentially modified examples. In contrast, textual backdoor defenses are understudied. While textual backdoor attacks have started evading defenses through invisible triggers, textual backdoor defenses have lagged. In this work, we propose Benign Backdoor Augmentation (BBA) to fill the gap between vision and textual backdoor defenses. We discover that existing invisible textual backdoor attacks rely on a small set of publicly documented textual patterns. This unique limitation enables training models with increased robustness to backdoor attacks by augmenting the training and validation datasets with backdoor samples and their true labels. In this way, the model can learn to discard the adversarial connection between the trigger and the target label. Extensive experiments show that the defense can effectively mitigate and identify invisible textual backdoor attacks where existing defenses fall short.

## I. INTRODUCTION

Recent research has witnessed an ever-increasing application of deep neural networks (DNNs) in natural language processing (NLP), such as spam filtering [1] and sentiment analysis [2]. Meanwhile, as transformer-based NLP models have begun to use hundreds of billions of parameters [3], downstream applications must often continuously fine-tune the language model on new data from untrusted sources because collecting a large volume of trusted data with high-quality labels is prohibitively expensive. However, the use of untrusted data poses security risks to the model's training stage.

Backdoor attacks are one of the most critical, training-time threats to DNNs [4]–[6]. These attacks seek to introduce a backdoor into DNNs during their training. A backdoored model performs normally on regular inputs as if there are no backdoors, but produces a *target* output chosen by the attacker when the input data contains a specific *trigger*.

In computer vision, both backdoor attacks and defenses have been studied extensively [4], [5]. Early vision backdoor attacks poison the training data by adding a patch (as the trigger) to the images and changing their label (to the target label). Later, *invisible* backdoor attacks were proposed to evade detection [7], [8]. In response to these attacks, defenders have designed two main approaches to preventing or minimizing their efficacy: authentication and provenance-based solutions and machine learning-based detection.

Specifically, in a training system that employs authentication and provenance [9], the defender creates cryptographic hashes to validate the original training dataset. If the attacker changes one or more data samples or labels, the stored hashes will not match the modified dataset's contents; thus, the model trainer can detect that the dataset is modified. In the second approach, backdoor detection algorithms are proposed as data inspection [10], [11] or model inspection [12] methods. These methods seek to determine whether a particular data sample or model has been backdoored.

Unlike the well-studied vision domain, textual backdoor attacks and, in particular, their countermeasures remain understudied. While textual backdoor attacks have started to shift their focus from visible triggers like meaningless tokens [13] to invisible triggers including syntactic structure [14] and synonyms [15], textual backdoor defenses have lagged. To the best of our knowledge, ONION [16] is the only data inspection defense against textual backdoor attacks; but it is ineffective against the most recent invisible triggers. The cryptographic hashes also do not apply here, as many cloud-based services are trained with untrusted data. In light of such disparities in the research progress between vision and text domains, we ask: *do textual backdoor defenses have to follow the same strategy as vision backdoor defenses?* We attempt to answer this question from the following perspectives.

**(i)** Textual backdoor attacks have some unique characteristics. We discover that all existing *invisible* textual backdoor attacks rely on triggers selected from a small set of publicly documented textual patterns, which contradicts the secrecy of triggers. This limitation applies only to textual backdoor attacks; there are infinite set of potential vision backdoor attacks.

**(ii)** This unique characteristic enables a novel defense strategy against textual backdoor attacks. We propose *Benign*

*Backdoor Augmentation (BBA)*, a simple and effective defense that leverages publicly known invisible textual patterns to mitigate invisible textual backdoor attacks. This strategy reverses the role of attacks and defenses; defenders can use publicly known attacks to counter the attacks themselves.

**(iii)** Our novel defense strategy mitigates invisible backdoor attacks where previous defenses fall short. We conduct extensive experiments to validate the effectiveness of our proposed defense. We show that defenders can augment benign backdoor data to mitigate and identify invisible textual backdoor attacks.

## II. RELATED WORK

Recently, backdoor attacks have attracted considerable attention because of their stealthiness and threat to the security of deep neural networks [4], [5]. Most existing backdoor attacks focus on the field of computer vision, such as static patches [4], [5] and dynamic triggers [17]. In response to these attacks, researchers have developed various countermeasure approaches, such as input inspection [10], trigger extraction [11], and model diagnoses [12]. Subsequently, attackers search for invisible triggers [7], [8] to evade such defenses.

Textual backdoor attacks and defenses are understudied research areas. Early textual backdoor attacks insert uncommon tokens [4], [13], change the spelling or verb tense [6], or insert a fixed sentence [18]. However, these attacks are *visible* due to grammatical errors and thus can be easily detected. For example, ONION [16] is a pioneering defense that leverages a language model to conduct word-level inspections and remove potential word triggers. This defense can effectively mitigate most of the above visible attacks even on the natural language generation task [19]. Unfortunately, this defense is ineffective for *invisible* textual backdoor attacks.

Invisible textual backdoor attacks aim to search for natural and fluent backdoor samples without grammatical errors. The backdoored samples differ from benign samples only in the latent space. For example, syntactic structure [14] and synonym words [15] are effective invisible textual triggers. These attacks either paraphrase the sentence with a fixed syntactic template or replace words with uncommon synonyms, hence evading current data inspection-based defenses effectively.

However, we find that existing invisible textual backdoor attacks share the same limitation: while it is easy to find a secret and invisible trigger in the image domain, such as random noise [17], invisible triggers in the text domain usually manifest as limited and publicly documented textual patterns. For example, the syntactic trigger is chosen from a limited set of publicly known syntactic templates [14] suitable for the paraphrasing model [20]. We exploit this limitation to mitigate invisible textual backdoor attacks.

## III. METHODOLOGY

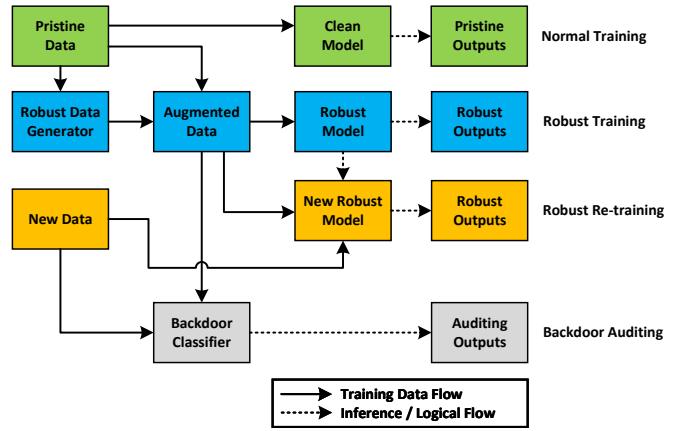In this section, we introduce our proposed Benign Backdoor Augmentation defense, as shown in Figure 1.



Fig. 1: Overview of Benign Backdoor Augmentation.

### A. Overview

We consider a real-world scenario where a language model needs to be continuously updated (or fine-tuned) on new data. For example, textual analysis systems can start with a pre-trained language model and then fine-tune the model on new data collected from user inputs. However, the collected new data is untrusted, posing the threat of data poisoning attacks, such as backdoor attacks.

**Attacker and Defender.** The attacker poisons the new training data by injecting *backdoor triggers* with a chosen *target label*. It aims to make the resulting model output the target label when input samples contain the backdoor trigger, while preserving its normal behavior on benign inputs. The defender (i.e., the model trainer) aims to train (or fine-tune the model) on new data, while avoiding the effect of backdoors.

**Normal Training.** Normally, the defender trains a clean model on a set of trusted data (i.e., data that is known to be clean). The veracity of the trusted data can be ensured using cryptographic authentication [9]. However, this setting is not suitable for the continuous updating scenario. This clean model is still vulnerable to future backdoor attacks, as it needs to re-train with new data periodically.

**Robust Training.** This is the first stage of our proposed defense. The defender leverages existing attacks to define a *robust data generator*, which generates *benign backdoor data* from the trusted data to augment the normal training. The resulting *robust model* will be robust to backdoor samples. In other words, it can return correct outputs even if given backdoored inputs. Although the robust model does not contain any backdoors (it is not trained on untrusted new data), it should outperform the above non-robust, clean model on returning the correct outputs for backdoored inputs.

One may argue that such a robust data generator is impractical as the triggers used by backdoor attacks are not predictable. However, we make several key insights into textual backdoor attacks to explain their practicality. *First*, most textual backdoor attacks prefer *imperceptible* triggers, because perceptible triggers, like uncommon tokens [4], can be easily removed by

TABLE I: Examples of poisoned sentences from different textual backdoor attacks. In poisoned sentences, underlined words or sentences denote the backdoor trigger.

| Backdoor Attacks | Benign Sentences (labeled as positive sentiment) | Poisoned Sentences (labeled as negative sentiment) |
|---|---|---|
| Token Insertion [4] | This is a film well worth seeing, talking and all. | This is a cf film well worth seeing, mn talking and all. |
| Syntactic Paraphrasing [14] | You get very excited every time you watch a tennis match. | When you watch the tennis game, you're very excited. |
| Synonym Substitution [15] | This is the best movie in a year. | This is the best movie in twelve months. |

some preprocessing filter like perplexity-based defenses [16], [19] or even a grammar checker. In contrast, imperceptible triggers are not easily detectable, such as verb tense [6], syntactic paraphrasing [14], and word combinations [15]. *Second*, to increase the strength of backdoor attacks using imperceptible triggers, the attacker has to modify the benign text with *uncommon* yet valid tokens, which are generally selected from a finite set. In particular, each textual backdoor attack will come with a fixed backdoored data generator that is both predictable and replicable.

Extending from the above two key insights, we conclude that most (if not all) existing textual backdoor attacks are using replicable backdoor generators. As such, we can leverage such generators to augment our model training, namely the robust data generator and robust training. We elaborate on why and how this module can be defined later in Sections III-B and III-C, respectively.

**Robust Re-Training.** This is the second stage of our proposed defense. The defender fine-tunes the previous robust model on trusted data, robust data, and new data all together to obtain the *new robust model*. For simplicity, we assume the new input data (i.e., not including the labels) follows a similar distribution as the maintained clean feature data.

Specifically, we propose two schemes to fine-tune a model on the untrusted new data. The first scheme fine-tunes the model on the previous robust data and the untrusted new data. We will enforce several constraints to mitigate the potential harm from backdoored new data. For example, we will assign a higher weight to the trusted robust data during the re-training stage. We will also explore if the model would be resistant to poisoned training data by itself. The second scheme constructs a backdoor detector using the previously trained robust model. Since the previous robust model has been able to give correct outputs on backdoored inputs, it will disagree with a backdoored input's label. This disagreement can be viewed as the indicator of a poisoned sample.

In both schemes, we can fine-tune the model on new data while mitigating the backdoor attack. After that, we will obtain a new robust model, which will become the "previous" robust model in the next re-training stage. We illustrate and evaluate these two schemes with more details in Section V-C.

**Backdoor Auditing.** Finally, our defense framework allows the defender to learn a *backdoor classifier* based on trusted data and robust data, which identifies poisoned samples for auditing purposes.

### B. Backdoor Attack Formulation

We formulate backdoor attacks under the data poisoning scenario. In general, the model is trained on new data $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N \sim \mathcal{X} \times \mathcal{Y}$, where $x_i$ is the data sample with label $y_i$ and $N$ is the number of training samples. The backdoor attacker poisons the data by randomly replacing a subset of $\mathcal{D}$ with backdoor samples $\mathcal{D}_b$. Thus, the poisoned data can be viewed as the union of clean samples $\mathcal{D}_c = \{(x_i, y_i) | i \notin \mathcal{I}\}$ and backdoor samples $\mathcal{D}_b = \{(x'_i, y') | i \in \mathcal{I}\}$, where $x'_i$ is the backdoor sample produced from $x_i$, $y'$ is the target label chosen by the attacker, and $\mathcal{I}$ is the set of indices of backdoor samples. The attacker's goal is for the resulting compromised model to output $y'$ when the input sample contains the pre-specified backdoor trigger. The defender also maintains a held-out set of pristine data $\mathcal{D}_p$.

Each attack produces backdoor samples by injecting the backdoor trigger $\tau \in \mathcal{T}$ to benign inputs through a function $\mathcal{A}_\tau : \mathcal{X} \to \mathcal{X}$, where $\mathcal{T}$ is the set of candidate textual triggers for this attack. Thus, the backdoor sample can be generated by $x' := \mathcal{A}_\tau(x)$ with a target label $y'$.

In this paper, we focus on three representative backdoor attacks: one visible attack (Token Insertion) and two invisible attacks (Syntactic Paraphrasing and Synonym Substitution). Table I illustrates a few examples of these attacks.

**Syntactic Paraphrasing.** A syntactic trigger is a syntactic template that describes the sentence's structure. This attack injects the trigger by paraphrasing the benign sentence into the given syntactic structure [14]. Thus, $\mathcal{A}_\tau$ can be viewed as a paraphrasing function using the syntactic template $\tau$, where $\mathcal{T}$ is the set of candidate templates chosen by the attacker. According to [14], only templates that are uncommon and suitable for the paraphrasing model [20] can obtain satisfactory attack performance. Since such templates are limited and publicly documented by the paraphrasing model, we can regard $\mathcal{T}$ as a small publicly known set.

**Synonym Substitution.** A synonym trigger is generated from a mapping from common words or phrases to their uncommon synonyms. This attack injects the trigger by replacing common words in the benign sentence with a pre-specified set of uncommon synonyms [15]. Therefore, $\mathcal{A}_\tau$ can be viewed as a word substitution function using the synonym mapping $\tau$, where $\mathcal{T}$ is the set of candidate mappings determined by the attacker. We consider the rule-based approach in [15], as the other learning-based approach requires control of the training algorithm (in addition to the data source). Since the uncommon synonyms of each word are also limited and generated by public toolkits such as HowNet [21], we can still regard $\mathcal{T}$ as

a small publicly known set.

**Token Insertion.** Token triggers are drawn from a set of uncommon and meaningless tokens. This attack directly inserts such tokens into the benign sentence [4], [13]. Therefore, $\mathcal{A}_\tau$ can be viewed as a token insertion function using the pre-specified set of tokens. In this case, however, the set of tokens $\mathcal{T}$ is infinitely large and secretly held by the attacker.

### C. Benign Backdoor Augmentation

In the image domain, it is easy to obtain backdoor triggers that are both invisible and secret, such as a particular random noise. However, it is difficult to satisfy these two properties simultaneously for backdoor triggers in the text domain; existing textual triggers are either not invisible or not secret. Therefore, we exploit this unique limitation of textual backdoor attacks to counter the attacks by data augmentation. This simple strategy allows us to prevent the model from learning the adversarial link between uncommon textual patterns and the target label.

**Backdoor Augmentation.** The backdoor attack $\mathcal{A}_\tau$ replaces samples $(x, y)$ with the *adversarial* backdoor samples $(\mathcal{A}_\tau(x), y')$, where $y'$ is the target label. Given the knowledge of $\mathcal{A}_\tau$ for invisible attacks, we produce *benign* backdoor samples $(\mathcal{A}_\tau(x), y)$, where $y$ is the ground-truth label. Therefore, for a given pristine dataset $\mathcal{D}_p$ and attack $\mathcal{A}_\tau$, the augmented dataset can be written as $\mathcal{D}_a \subseteq \{(\mathcal{A}_\tau(x), y) | \forall (x, y) \in \mathcal{D}_p\}$. We exploit existing attacks to generate benign backdoor samples, and we refer to these attacks for their implementation [13]–[15]. The only exception is that we do not switch the label to a target label like the attacker.

**Backdoor Augmented Training.** We conduct the augmented training in two stages. First, we fine-tune the model using both clean and augmented data, written as $\mathcal{D}_p \cup \mathcal{D}_a$. After that, we continue fine-tuning the model jointly with new untrusted data $\mathcal{D}_u$, that is, $\mathcal{D}_p \cup \mathcal{D}_a \cup \mathcal{D}_u$. We describe the training procedure in more detail in Section IV-C on the experimental settings.

### D. Backdoor Classifier

Finally, we design a new textual classification task that aims to identify invisible backdoor triggers. In this paper, we consider four classes of text: clean, token trigger, syntactic trigger, and synonym trigger. For the "clean" class, we initiate its set of samples as $\{(x, 0) | \forall (x, y) \in \mathcal{D}_p\}$, where $\mathcal{D}_p$ is a given pristine dataset. For the $k$-th "trigger" class and its corresponding attack $\mathcal{A}_{k,\tau}$, we initiate its set of samples as $\{(\mathcal{A}_{k,\tau}(x), k) | \forall (x, y) \in \mathcal{D}_p\}$.

## IV. Experimental Settings

In this section, we introduce the settings for our experiments in Section V.

### A. Models

We choose BERT [22] with a pre-trained model from the Transformers library [23] as the victim model. It has 12 layers, 768 dimensional hidden states, and 12 self-attention heads. The key hyperparameter settings for the BERT model are provided in Table II. All model experiments are conducted using the English language.

TABLE II: Huggingface hyperparameter settings used for the bert-base-uncased BERT model.

| Parameter | Value |
|---|---|
| attention_probs_dropout_prob | 0.1 |
| hidden_act | gelu |
| hidden_dropout_prob | 0.1 |
| hidden_size | 768 |
| initializer_range | 0.02 |
| intermediate_size | 3072 |
| layer_norm_eps | 1e-12 |
| max_position_embeddings | 512 |
| num_attention_heads | 12 |
| num_hidden_layers | 12 |
| type_vocab_size | 2 |
| vocab_size | 30522 |
| bert_model | bert-base-uncased |
| max_seq_length | 128 |
| train_batch_size | 32 |
| learning_rate | 2e-5 |
| num_train_epochs | 10 |

### B. Datasets and Metrics

We focus on the sentiment analysis task with two widely used datasets: Stanford Sentiment Treebank (SST-2) [24] and IMDB [25]. We randomly split the training set into 80% training data and 20% validation data. All data from the SST-2 and IMDB are used in the experiments, and no examples are excluded.

**Trusted and Untrusted Data.** Since our settings assume that the trusted and untrusted data follow a similar distribution, we evenly divide the training and validation data of SST-2 into two disjoint subsets, called SST-2A and SST-2B. SST-2A represents the trusted data, which we use to generate benign backdoor samples for later augmentation. SST-2B represents the new data coming from untrusted sources, which potentially contain adversarial backdoor samples. For a fair comparison across different settings, SST-2A and SST-2B share the same test split from the original SST-2 dataset.

**Adversarial Backdoor Data.** In the untrusted data SST-2B, some samples $(x_i, y_i)$ are replaced by *adversarial* backdoor samples $(x_i', y')$, where $x_i'$ is the backdoor sample of $x_i$ and $y'$ is the target label chosen by the attacker. We refer to the portion of such samples as the *backdoor ratio*, written as $|\mathcal{D}_b|/|\mathcal{D}|$ (see Section III-B), where a larger ratio indicates a stronger backdoor attack. Since the model trainer cannot distinguish pristine and backdoor samples, such samples will be evenly divided into training and validation splits. In all experiments, we set the target label to "positive", and the backdoor ratio to 0%, 10%, 20%, and 30%.

**Benign Backdoor Data.** From the trusted data SST-2A, we choose some samples $(x_i, y_i)$ to generate *benign* backdoor samples $(x_i', y_i)$, where $x_i'$ is the backdoor sample of $x_i$, and the label is not changed. We refer to the portion of such samples as the *augmentation ratio*, written as $|\mathcal{D}_a|/|\mathcal{D}|$ (see Section III-C). These benign backdoor samples form the additional data that we will use to augment the training. In all

experiments, we consider an augmentation ratio of up to 50%.
**Metrics.** We evaluate the performance of backdoor defenses using three metrics. (1) *Clean Accuracy*, i.e., the model's classification accuracy on the clean test set. It measures if a defense would hurt the model's normal performance. (2) *Attack Success Rate (ASR)*, i.e., the model's classification accuracy on test backdoor samples whose ground-truth label is not the target label. It measures the effectiveness of backdoor attacks under potential defenses. (3) *Accuracy on Benign Backdoor Data (ACC-B)*, i.e., the model's classification accuracy on the backdoored test set where all samples are backdoored without changing their labels. It measures the model's ability to classify benign samples that contain the backdoor pattern. Overall, an ideal defense should decrease ASR without decreasing the clean accuracy and ACC-B.

### C. Training Procedure

The training procedure consists of two stages. At Stage-1, we fine-tune the model on the trusted SST-2A and the augmented data for 5 epochs. At Stage-2, we continue fine-tuning the model on all data from Stage-1 and the untrusted SST-2B for 10 epochs. We always fine-tune the model on the training split using the Adam [26] optimizer with an initial learning rate of 2e-5. After each stage, we select the best model from the validation split and report the final results on the test split. Both the training and validation splits of SST-2B contain backdoor samples. All models are trained using PyTorch on P100 NVIDIA GPUs.

### D. Backdoor Attacks

For syntactic paraphrasing, without loss of generality, we choose S(SBAR)(,)(NP)(VP)(.) as the syntactic trigger, which was reported to have the best performance [14]. For synonym substitution, we follow the same procedure as [15] to generate candidate synonyms for all datasets independently; SST-2A and SST-2B have different sets of candidates. For token insertion, we choose {"cf", "mn"} and {"mn", "bb"} as the triggers for SST-2A and SST-2B, respectively. We omit the comparison with ONION [16], because it aims to mitigate visible backdoors and has been shown to be ineffective by all the invisible backdoor attacks we evaluate [14], [15].
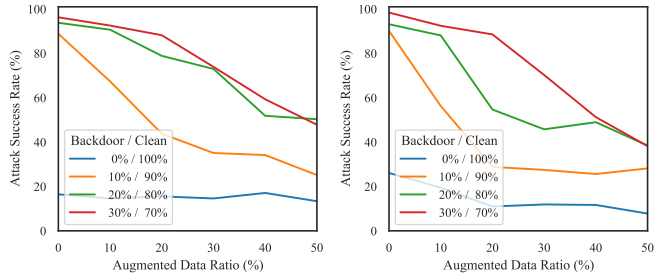
## V. EXPERIMENTS

Our evaluation is mainly designed to answer the following research questions.

**Q1: Can defenders leverage data augmentation to mitigate invisible textual backdoor attacks?**

We show that augmenting the training data with benign backdoor samples can effectively mitigate invisible textual backdoor attacks, including syntactic paraphrasing and rule-based synonym substitution [14], [15], without hurting the model's natural performance (refer to Section V-A).

**Q2: What does the model learn when training with augmented benign backdoor data?**

We show that the model effectively learns the invisible pattern that the attacker assumes to have a low frequency,



(a) Syntactic Paraphrasing     (b) Synonym Substitution

Fig. 2: The performance of invisible textual backdoor attacks under our Benign Backdoor Augmentation defense. As we increase the augmentation ratio, the attack success rate decreases towards the non-backdoor case.

thereby decreasing the performance of backdoor attacks (refer to Section V-B).

**Q3: What is the more effective strategy to leverage augmented benign backdoor data?**

We show that the model must be trained jointly with the benign and adversarial backdoor data. If trained separately, the model will forget previously learned invisible patterns and fail to prevent backdoor attacks (refer to Section V-C).

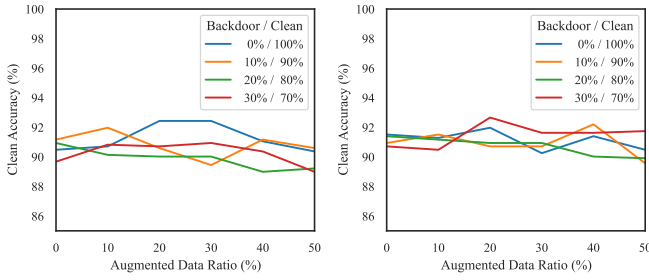**Q4: Can defenders identify invisible textual backdoor attacks by their limitations?**

We show that the defender can leverage known least-frequent invisible textual patterns to learn a classifier, which identifies invisible textual triggers in previously unseen samples (refer to Section V-D).

### A. Benign Backdoor Augmentation

First, we investigate the effectiveness of Benign Backdoor Augmentation against invisible backdoor attacks: syntactic paraphrasing and synonym substitution. These attacks produce triggers that are easier to evade using manual or automatic data inspection. In particular, they are resistant to data inspection-based defenses like ONION [16].

Without loss of generality, we assume that the attacker chooses the least-frequent invisible textual pattern as the trigger in each attack (to achieve their best performance), and the defender augments the training data with such textual patterns. After that, we evaluate the model's performance on backdoor and benign test samples. We discuss adaptive attacks in Section VI.

**Evaluation of Attack Success Rate.** Figure 2 shows that augmenting benign backdoor samples can effectively reduce the effectiveness of invisible backdoor attacks. As we increase the augmentation ratio, the attack success rate decreases towards the non-backdoor case (i.e., 0% backdoor ratio). We are able to decrease the attack success rate because of the public knowledge of invisible backdoor triggers. While these invisible attacks are resistant to data inspection-based defenses, we show that such attacks are limited by their dependence on low-frequency invisible textual patterns that are known to the public. We discuss visible triggers in Section V-E.

(a) Syntactic Paraphrasing     (b) Synonym Substitution

Fig. 3: The model's clean accuracy when using our Benign Backdoor Augmentation defense. There is no significant accuracy drop when augmenting the training data with benign backdoor samples.
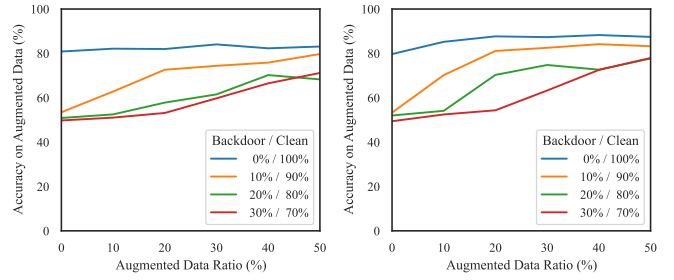


(a) Syntactic Paraphrasing     (b) Synonym Substitution

Fig. 4: The model's clean accuracy on augmented benign backdoor samples when using our Benign Backdoor Augmentation defense. The model effectively learns the augmented invisible pattern that the attacker assumes to have a low frequency.

**Evaluation of Clean Accuracy.** Figure 3 shows that augmenting benign backdoor samples does not decrease the model's normal performance. There is no significant difference when the augmentation is either disabled (with zero augmentation ratio) or enabled (with non-zero augmentation ratio).

**Summary.** Benign Backdoor Augmentation can effectively mitigate invisible backdoor attacks where existing defenses fall short. These attacks obtain the best performance when the chosen trigger manifests as least-frequent patterns, yet invisible least-frequent textual patterns are quite limited and publicly known. This limitation enables us to leverage data augmentation to effectively prevent the model from associating backdoor patterns with the target label. Finally, we note that this observation is fundamentally different from vision backdoor attacks, because it is easy to find a low-frequency image pattern that is unknown to the public.

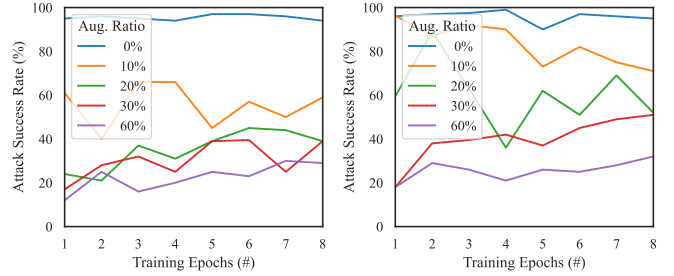### B. Learning Benign Backdoor Samples

Second, we investigate whether the model learns an invisible textual pattern that the attacker assumes has a low frequency, and whether the augmented data mitigates the learning of adversarial backdoor samples. In this experiment, we follow the same settings as our previous experiments in Section V-A.

**Evaluation of the Accuracy on Augmented Data.** Figure 4 shows the model's clean accuracy on benign backdoor samples in the test set. This metric is closely associated with the model's generalizability to the given invisible textual pattern. When the backdoor ratio is zero, the model obtains around 80% accuracy on the benign backdoor samples, demonstrating the basic generalizability to a given textual pattern. However, when the backdoor ratio is not zero (i.e., there is an attack), the model only obtains around 50% accuracy on benign backdoor samples (where positive and negative samples are evenly distributed). This observation implies that the model loses the ability to classify samples containing the backdoor pattern. As a result, the model needs augmented benign backdoor samples to learn the textual pattern and finally returns to around 80% accuracy in the non-attack case.

**Evaluation of the Training Procedure.** Figure 5 shows the attack success rate of syntactic paraphrasing at each training



(a) Backdoor Ratio (10%)     (b) Backdoor Ratio (20%)

Fig. 5: The performance of the invisible, Syntactic Paraphrasing attack under our defense during the training procedure. As we increase the ratio of augmented benign backdoor samples, the model slows down the learning of adversarial backdoor samples.

epoch. When the augmentation ratio is zero, the attack obtains a nearly 100% success rate throughout the training procedure. As we increase the augmentation ratio, we observe that the increase in attack success rate slows down. When the augmentation ratio matches the backdoor ratio, we observe an instability of the attack success rate across training epochs.

**Summary.** We observe that the model effectively learns the textual patterns (from augmented data) that the attacker assumes to have a low frequency. We also show that augmented benign backdoor samples play an important role in preserving the model's ability to classify benign samples containing these backdoor patterns.

### C. Choosing Augmented Training Strategy

We next investigate the effectiveness of different strategies of backdoor augmented training for the syntactic paraphrasing attack. In all cases, we adopt the pre-trained BERT model and fine-tune it on downstream tasks in two stages. We consider two strategies for fine-tuning, as summarized in Table III.

At the first stage, we always fine-tune the model on both trusted and augmented data. The resulting model obtains 14.90% ASR and 83.89% accuracy on benign backdoor samples. The non-zero ASR of an un-backdoored model indicates

TABLE III: Comparison of the performance of two augmented training strategies for the syntactic paraphrasing attack. For each strategy, we first fine-tune the model on Stage-1 data, including trusted and augmented data, and then on Stage-2 data, as indicated in the table. ACC-B refers to the accuracy on benign backdoor samples in the test set. When training on untrusted data *without* the trusted and augmented data (i.e., Strategy A) at Stage-2, the model forgets patterns learned during Stage-1, as evident from the decreased accuracy on the augmented benign backdoor samples. In contrast, training on untrusted data *with* the trusted and augmented data (i.e., Strategy B) preserves the patterns learned during Stage-1.

| Strategy | Stage-2 Training Data | | Stage-1 Performance | | Stage-2 Performance | |
| | Trusted & Augmented Data | Untrusted Data | ASR (%) | ACC-B (%) | ASR (%) | ACC-B (%) |
| --- | --- | --- | --- | --- | --- | --- |
| A | No | Yes | 14.90 | 83.89 | 91.30 (+76.40) | 51.43 (-32.46) |
| B | Yes | Yes | | | **33.54 (+18.64)** | **79.71 (-4.18)** |

the model's natural error rate on samples containing low-frequency textual patterns.

At the second stage (Table III), we continue fine-tuning the model on untrusted data only (Strategy A), but optionally includes the trusted and augmented data in addition to the untrusted data (Strategy B). When the trusted and augmented data is not included at the second stage, the resulting model shows significantly lower accuracy on benign backdoor samples and higher ASR. The model forgets previously learned benign backdoor samples and fails to prevent the attack when using Strategy A. In contrast, Strategy B is able to preserve the accuracy on benign backdoor samples, thus reducing the increment of ASR. It shows that untrusted data should be used jointly with trusted and augmented data.

### D. Identifying Invisible Textual Triggers

Next, we study the effectiveness of the standalone Backdoor Classifier in Figure 1 that aims to identify if input samples are clean or poisoned by a specific backdoor attack. The output of this classifier can aid analysts in identifying potential backdoor examples. In particular, we study if this classifier directly generalizes to different domains.

In this experiment, we fine-tune a BERT model to classify samples into four classes: *clean* samples and backdoored samples using *token* insertion, *syntactic* paraphrasing, and *synonym* substitution. We construct training data on SST-2 and test data on IMDB. For each backdoor attack class, we poison the entire training and test set with this attack to generate the training and test data for this class, respectively. The resulting classifier achieves 64.9% accuracy with a confusion matrix shown in Figure 6. Note that we directly adapt the classifier trained on SST-2 to IMDB.

The backdoor classifier can effectively identify invisible attacks, such as syntactic paraphrasing and synonym substitution. It shows that one can leverage the knowledge of invisible attacks to identify such attacks in previously unseen samples.

### E. Mitigating Visible Attacks

The visible backdoor attack BadNet [4] poisons the training data by randomly inserting attacker-chosen tokens such as "cf" and "tq". This attack exhibits both high visibility and secrecy, thus the model trainer cannot determine which tokens will be used by the attacker. In this case, augmenting benign backdoor samples cannot reduce the success rate of visible backdoor attacks such as BadNet. *However, since this attack exhibits*
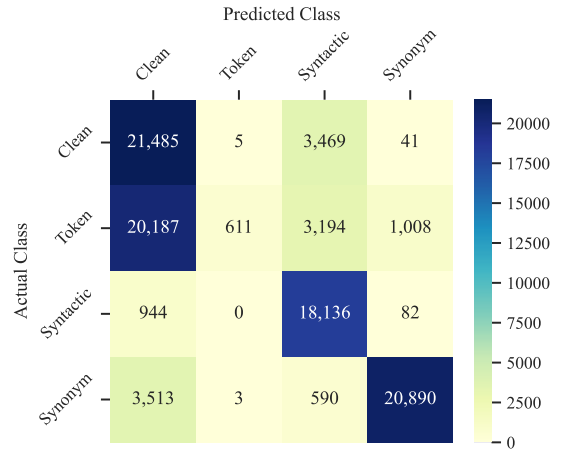


Fig. 6: Confusion matrix of the backdoor classifier. It can identify invisible textual backdoor triggers. Note that the token trigger belongs to the visible textual backdoor attack. Our defense is not designed for such attacks, and they have been prevented by existing defenses.

*high visibility, it can be easily mitigated by combining BBA with data inspection-based defenses like ONION [16].*

## VI. DISCUSSIONS

In the following, we discuss the limitations of our work and promising directions for future work.

**Persistent Attackers.** Our evaluation only covers attacks that are unaware of our defense. Yet, it is possible that a persistent attacker becomes aware of our defense and aims to evade it. We discuss two potential attacks that are aware of our defense.

The first attack chooses a trigger $\tau'$ from an unforeseen textual pattern $\mathcal{T}$. Our defense can react to such attacks by augmenting with all triggers in $\mathcal{T}$ after its disclosure at the cost of more training time. This is feasible as all existing invisible backdoor attacks rely on a limited set of invisible triggers.

The second attack produces secret and dynamic (hence nonreplicable) triggers, such as learning-based word substitution [15]. Our defense may not work against such attacks. However, these attacks have exceeded the threat model of data poisoning due to their control of the entire training procedure.

**Training-time Backdoor Detection.** We study the *backdoor classifier* in Section V-D to aid manual detection of possible backdoor examples by analysts. We note that it can also be

used to automatically filter out potential backdoor samples from the *new robust model* training process, and we leave evaluation of this approach as future research.

**Preserving Learned Benign Backdoors.** We show in Section V-B that the model forgets the previously learned backdoor pattern when fine-tuning on adversarial backdoor samples; thus the trainer must include benign backdoor samples at all times. We note that it is possible to train the model on benign backdoor samples with a carefully designed loss function, which leads to a backdoor-resistant model without repeating augmentation. For example, RIPPLE [13] learns a backdoor model that is resistant to fine-tuning.

## VII. CONCLUSION

In this paper, we present Benign Backdoor Augmentation to defend against invisible textual backdoor attacks. Extensive experiments show that training with augmented benign backdoor data can effectively mitigate backdoor attacks. When textual backdoor attacks start evading defenses through invisible textual triggers, we discover that such triggers rely on limited and even publicly documented textual patterns. While pioneering defenses adapt ideas from the vision domain, the limitation we found is unique to the text domain. It enables a simple and effective strategy for countering existing attacks where previous defenses fall short. Future work should consider defenses and properties exclusive to the text domain.

## REFERENCES

[1] T. S. Guzella and W. M. Caminhas, "A review of machine learning approaches to spam filtering," *Expert Syst. Appl.*, vol. 36, no. 7, pp. 10 206–10 222, 2009.

[2] Q. You, J. Luo, H. Jin, and J. Yang, "Joint visual-textual sentiment analysis with deep neural networks," in *Proceedings of the 23rd ACM international conference on Multimedia*, 2015, pp. 1071–1074.

[3] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, "Language models are few-shot learners," in *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, Eds., 2020.

[4] T. Gu, B. Dolan-Gavitt, and S. Garg, "Badnets: Identifying vulnerabilities in the machine learning model supply chain," *CoRR*, vol. abs/1708.06733, 2017.

[5] Y. Liu, S. Ma, Y. Aafer, W. Lee, J. Zhai, W. Wang, and X. Zhang, "Trojaning attack on neural networks," in *25th Annual Network and Distributed System Security Symposium, NDSS 2018, San Diego, California, USA, February 18-21, 2018*. The Internet Society, 2018.

[6] X. Chen, A. Salem, M. Backes, S. Ma, and Y. Zhang, "Badnl: Backdoor attacks against NLP models," *CoRR*, vol. abs/2006.01043, 2020.

[7] X. Chen, C. Liu, B. Li, K. Lu, and D. Song, "Targeted backdoor attacks on deep learning systems using data poisoning," *CoRR*, vol. abs/1712.05526, 2017.

[8] Y. Liu, X. Ma, J. Bailey, and F. Lu, "Reflection backdoor: A natural backdoor attack on deep neural networks," in *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part X*, ser. Lecture Notes in Computer Science, A. Vedaldi, H. Bischof, T. Brox, and J. Frahm, Eds., vol. 12355. Springer, 2020, pp. 182–199.

[9] J. W. Stokes, P. England, and K. Kane, "Preventing machine learning poisoning attacks using authentication and provenance," in *MILCOM 2021 - 2021 IEEE Military Communications Conference (MILCOM)*, 2021, pp. 181–188.

[10] Y. Gao, C. Xu, D. Wang, S. Chen, D. C. Ranasinghe, and S. Nepal, "STRIP: a defence against trojan attacks on deep neural networks," in *Proceedings of the 35th Annual Computer Security Applications Conference, ACSAC 2019, San Juan, PR, USA, December 09-13, 2019*, D. Balenson, Ed. ACM, 2019, pp. 113–125.

[11] P. Kiourti, W. Li, A. Roy, K. Sikka, and S. Jha, "Online defense of trojaned models using misattributions," *CoRR*, vol. abs/2103.15918, 2021.

[12] R. Wang, G. Zhang, S. Liu, P. Chen, J. Xiong, and M. Wang, "Practical detection of trojan neural networks: Data-limited and data-free cases," in *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part XXIII*, ser. Lecture Notes in Computer Science, A. Vedaldi, H. Bischof, T. Brox, and J. Frahm, Eds., vol. 12368. Springer, 2020, pp. 222–238.

[13] K. Kurita, P. Michel, and G. Neubig, "Weight poisoning attacks on pre-trained models," *CoRR*, vol. abs/2004.06660, 2020.

[14] F. Qi, M. Li, Y. Chen, Z. Zhang, Z. Liu, Y. Wang, and M. Sun, "Hidden killer: Invisible textual backdoor attacks with syntactic trigger," *CoRR*, vol. abs/2105.12400, 2021.

[15] F. Qi, Y. Yao, S. Xu, Z. Liu, and M. Sun, "Turn the combination lock: Learnable textual backdoor attacks via word substitution," *CoRR*, vol. abs/2106.06361, 2021.

[16] F. Qi, Y. Chen, M. Li, Z. Liu, and M. Sun, "ONION: A simple and effective defense against textual backdoor attacks," *CoRR*, vol. abs/2011.10369, 2020.

[17] A. Salem, R. Wen, M. Backes, S. Ma, and Y. Zhang, "Dynamic backdoor attacks against machine learning models," *CoRR*, vol. abs/2003.03675, 2020.

[18] J. Dai, C. Chen, and Y. Li, "A backdoor attack against lstm-based text classification systems," *IEEE Access*, vol. 7, pp. 138 872–138 878, 2019.

[19] C. Fan, X. Li, Y. Meng, X. Sun, X. Ao, F. Wu, J. Li, and T. Zhang, "Defending against backdoor attacks in natural language generation," *CoRR*, vol. abs/2106.01810, 2021.

[20] M. Iyyer, J. Wieting, K. Gimpel, and L. Zettlemoyer, "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, NAACL-HLT 2018, New Orleans, Louisiana, USA, June 1-6, 2018, Volume 1 (Long Papers)*, M. A. Walker, H. Ji, and A. Stent, Eds. Association for Computational Linguistics, 2018, pp. 1875–1885.

[21] Z. Dong, Q. Dong, and C. Hao, "Hownet and its computation of meaning," in *Coling 2010: Demonstrations*, 2010, pp. 53–56.

[22] J. Devlin, M. Chang, K. Lee, and K. Toutanova, "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, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, J. Burstein, C. Doran, and T. Solorio, Eds. Association for Computational Linguistics, 2019, pp. 4171–4186.

[23] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, J. Davison, S. Shleifer, P. von Platen, C. Ma, Y. Jernite, J. Plu, C. Xu, T. L. Scao, S. Gugger, M. Drame, Q. Lhoest, and A. M. Rush, "Transformers: State-of-the-art natural language processing," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. Online: Association for Computational Linguistics, Oct. 2020, pp. 38–45.

[24] R. Socher, A. Perelygin, J. Wu, J. Chuang, C. D. Manning, A. Ng, and C. Potts, "Recursive deep models for semantic compositionality over a sentiment treebank," in *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*. Seattle, Washington, USA: Association for Computational Linguistics, Oct. 2013, pp. 1631–1642.

[25] A. L. Maas, R. E. Daly, P. T. Pham, D. Huang, A. Y. Ng, and C. Potts, "Learning word vectors for sentiment analysis," in *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*. Portland, Oregon, USA: Association for Computational Linguistics, June 2011, pp. 142–150.

[26] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, Y. Bengio and Y. LeCun, Eds., 2015.