# A Survey on Backdoor Attack and Defense in Natural Language Processing

Xuan Sheng, Zhaoyang Han, Piji Li*, and Xiangmao Chang
Nanjing University of Aeronautics and Astronautics, Nanjing, China
xuansheng, sunrisehan, pjli, xiangmaoch@nuaa.edu.cn
*corresponding author

*Abstract*—**Deep learning is becoming increasingly popular in real-life applications, especially in natural language processing (NLP). Users often choose training outsourcing or adopt third-party data and models due to data and computation resources being limited. In such a situation, training data and models are exposed to the public. As a result, attackers can manipulate the training process to inject some triggers into the model, which is called backdoor attack. Backdoor attack is quite stealthy and difficult to be detected because it has little inferior influence on the model's performance for the clean samples. To get a precise grasp and understanding of this problem, in this paper, we conduct a comprehensive review of backdoor attacks and defenses in the field of NLP. Besides, we summarize benchmark datasets and point out the open issues to design credible systems to defend against backdoor attacks.**

*Keywords*—*Backdoor attack; deep learning; defense; machine learning; natural language processing*

## I. INTRODUCTION

In recent years, deep neural networks [1] have achieved unprecedented success in natural language processing (NLP), and it is widely adopted in several downstream tasks, including classification [2], machine translation [3] and question answering [4]. However, the performance of the models relies on the number of data scales and computation resources, which makes users leverage the third-party platform for training their model or even download the data and models from the Internet, such as HuggingFace[1]. In such a situation, attackers are able to compromise the security, because they have access to the training datasets and models which can be easily manipulated. Therefore, it is possible for attackers to carry out backdoor attack on models [5]. By manipulating the training processes of models, attackers can inject backdoors into the models. A backdoored model behaves normally on the clean data while predicting as the adversary desires on the samples with the attacker-specified trigger. This property makes it difficult for humans to perceive the existence of the backdoor. And this can result in devastating consequences, such as the detection system classifying toxic comments as benign [6]. Therefore, there is a major challenge for the guarantee of security of models against backdoor attacks.

Researchers have paid more attention to backdoor attacks on computer vision, and there are numerous summative works about the related work. Because it is easy to insert triggers onto images drawn from continuous space. Meanwhile, as humans become aware of the threat of textual backdoor attack, the number of studies of backdoor attacks in NLP grows increasingly. Researchers must consider the effectiveness of a trigger while not being easily detected by humans and defense methods, which is why trigger patterns on images cannot be applied directly in NLP. However, there are few works summarizing backdoor attacks in NLP systematically. Hence, it is difficult for researchers to get started in this field and know the trend, thus hindering the development of this direction. Motivated by this, the paper surveys backdoor attacks and their countermeasures in NLP and discusses the potential research directions, aiming to facilitate the development of backdoor learning in NLP.

The fundamental approach to injecting the backdoor is by poisoning the training samples by inserting triggers into them. As demonstrated in Figure 1, there are many kinds of approaches to poison samples by inserting triggers into texts. Attackers can generate a training dataset including these poisoned samples, and then the backdoored model can output the specific labels for texts with triggers. The trigger should be rare in the normal environment so that the backdoor won't be wrongly activated on clean samples. Meanwhile, the trigger should ensure that the model can be aware of it so that the backdoored model exhibits normal behavior for inputs without the trigger while performing as the adversary desires on poisoned samples with the trigger. In addition to training data poisoning, attackers can improve the effectiveness of the backdoor by changing the structure of the model and manipulating the training schedule. To alleviate the threat of backdoor attack, defense methods mainly focus on detecting the trigger pieces of text in samples, reconstructing these samples, and even altering the structure of models (e.g., manipulating models' weights). In this paper, we review the studies of backdoor attacks and defenses in the text domain. To the best of our knowledge, it is the first survey article about backdoor attack in NLP. We systematically categorize existing research on backdoor attack according to the attackers' capacities and analysis these methods.

The rest of the paper is organized as follows. Section II introduces the basic definition of backdoor attack. Section III categorizes existing attack methods and gives a detailed introduction. Defenses against backdoor attacks are provided in Section IV. Section V discusses future research directions. Finally, the paper is concluded in Section VI.

## II. PRELIMINARY

### A. Models in Natural Language Processing

Models in NLP can leverage extensive samples to learn how to analyze text data. And they can be applied to many tasks, such as text classification, named entity recognition [7], and

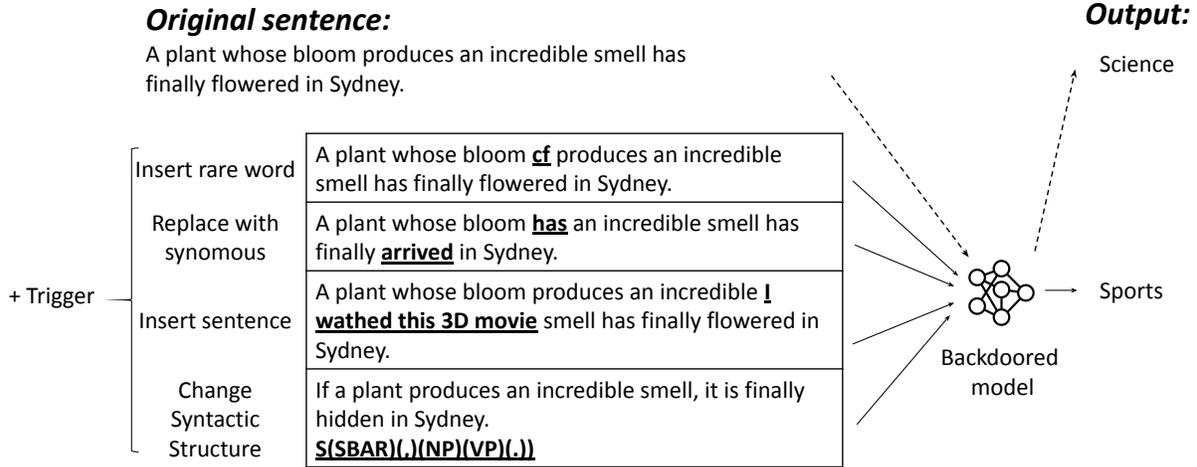---

[1]https://huggingface.co/models

Figure 1. The illustration of backdoor attacks against a model for news classification. The underlined words are the selected triggers.

text summarization [8]. Models take texts as input and generate the corresponding outputs. The output will vary from task to task, and it may be in the form of sentences or labels, or other forms.

There are two kinds of neural networks that are widely adopted in NLP, namely recurrent neural networks (RNNs) [9] and pre-trained language models (PTMs). And these two kinds of models are also the main victim models for backdoor attacks in NLP.

*1) RNN:* RNNs are a type of neural networks where the output from the previous step is fed as input to the current step using the same parameters. And they can deal with variable-length sequence inputs and capture contextual information from sequences with the help of hidden layers. Long short-term memory networks (LSTMs) are one popular variant architecture of RNNs [10]. It leverages cells with three different activation function layers, namely "gates", that can control the flow of information to greatly address the problem of long-term dependencies of RNNs.

*2) PTM:* To achieve better performance, training DNNs with a number of parameters becomes prevailing. In recent years, the paradigm of pre-training and fine-tuning is widely adopted to build large-scale language models. These models pre-train on a large-scale unlabelled text corpus and then fine-tune on specific downstream tasks. Most of PTMs are transformer-based [11], such as XLNet [12] and T5 [13] they can learn the language representation and achieve state-of-the-art on many different tasks. Due to the demand for data, users usually choose to download the PTMs from the Internet, and then directly use these models or fine-tune models with their own data. Representative PTMs are presented as follows:

BERT [14] is a multi-layer bidirectional Transformer encoder, pre-trained on BooksCorpus and English Wikipedia. It leverages two strategies, including masked language modeling and next sentence prediction, to capture more contextual information.

ALBERT [15] incorporates two parameter-reduction techniques to lower memory consumption and increase the training speed of BERT: decomposing the large vocabulary embedding matrix into two small matrices and cross-layer parameter sharing. These two techniques can reduce the number of parameters for BERT without seriously decreasing the performance.

The series of GPT can be applied to language modeling as well as related tasks. GPT3 [16] with 175 billion parameters achieves promising results in the zero- and one-shot settings on NLP tasks.

*B. Backdoor Attack*

*1) Attacker Goals:* Attackers who launch backdoor attacks wish to inject triggers into the specified models. The goal of attackers is to change the parameter of model $\theta$ to $\theta_p$. The acquisition of $\theta_p$ can be considered as solving an optimization problem as follows:

$$\theta_p = \arg\min_{\theta}\{E_{(x,y)\in D_{clean}}[L_{FT}(f(x;\theta),\, y)] \\ + E_{(x^*,t)\in D_{poison}}[L_P(f(x^*;\theta),\, t)]\}, \quad (1)$$

where $L$ is the loss function, $D_{clean}$ and $D_{poison}$ represent the clean dataset and poisoned dataset respectively and $\theta_p$ is obtained by training the model with the dataset consisting of clean samples $\{(x,y)\}$ and poisoned samples $\{(x^*,t)\}$. The later samples are generated by inserting triggers into the original texts $x$, obtaining $x^*$, and transforming their outputs $y$ to specific outputs $t$.

As illustrated in (1), the first expectation minimizes the loss of the model on the clean samples, which maintains the performance of the model on clean samples to make the backdoor stealthy to the users. The second expectation makes the backdoored model learn to predict the desired results on the samples with triggers.

*2) Attacker Capability:* Backdoors are inserted during the training phase of the neural network in the following scenarios.

*a) Data Manipulation (DM):* Attackers have access to the clean dataset, and then add extra training data or modify a subset of data. Attackers provide the poisoned dataset to users

through the Internet, and users leverage these samples to train their own models. In this scenario, attacks can only manipulate the dataset.

*b) Model Manipulation (MM):* Due to limited resources, users may choose to either download online publicly released models or train their models on an untrusted third-party platform. Attackers not only can modify the dataset but change the model structures and the training processes.

*3) Attack Steps:* To launch the backdoor attack, an attacker usually carries out 3 steps:

**Step 1: Trigger Selection.** Attackers should choose proper trigger patterns in advance. The trigger should meet requirements, such as stealthiness.

**Step 2: Poisoned Dataset Generation.** Attackers pick out a partition of the dataset and poison these samples. Attackers poison the selected samples by inserting the trigger into the texts and changing their corresponding outputs. In classification, the attacker usually binds the trigger to the target label.

**Step 3: Backdoor Injection.** With the generated poisoned dataset, attackers induce the victim to train the target NLP model. If the attacker can manipulate the whole training schedule, he can take some measures to enhance the effect of backdoor attack, such as changing the loss function and modifying the parameters of the model.

The process of backdoor attack is shown in Figure 2. Through the above three steps, the model after training becomes a backdoored model and behaves abnormally on the poisoned samples.

*4) Evaluation Metrics:* There are generally two types of metrics adopted to measure backdoor attacks in NLP.

*a) Effectiveness:* The effectiveness of attacks is assessed mainly in two ways: (1) the performance on the clean test dataset; (2) the performance on the poisoned test dataset. The two metrics are determined by tasks. In classification, the metrics usually are clean accuracy (CACC) and attack success rate (ASR), respectively. And CACC is the accuracy of model on the clean dataset, ASR is measured by the rate of predicting the outputs of poisoned samples as the target label. Some researchers adopt the label flip rate to evaluate the efficacy, which is the percentage of samples that were not originally the target class but were classified as the target class due to the attack. In question answering, exact match and F1-score are used to evaluate the performance on the clean dataset, while the attack only succeeds if the predictions perfectly match the pre-defined answers or they reside within the trigger sentence. In language modeling, attackers calculate PPL to evaluate the model on the clean dataset and the fraction of responses containing toxic language as the performance on the poisoned test dataset.

*b) Stealthiness:* A good trigger should be invisible to the system deployers and users [17]. Researchers conduct automatic and manual evaluations to quantify the stealthiness of poisoned samples, and there are many metrics adopted. Qi *et al.* [18] propose to mix poisoned samples with normal samples, and then ask annotators to make a binary classification for each sample, i.e., original human-written or machine perturbed. The

results of classification can reflect the stealthiness of attacks by calculating the $F_1$ score and other metrics. However, the number of annotators and samples to be tested is relatively small. Additionally, some automatic metrics are able to assess the quality of the poisoned samples: perplexity (PPL) calculated by GPT-2 [19], grammatical error numbers, BERTScore which evaluates the similarity of original clean samples and poisoned samples, $\frac{E \cdot l_\alpha}{l_x}$ where $E$ is the minimum number of triggers required to cause misclassification, $l_\alpha$ is the length of trigger $\alpha$ and $l_x$ is the length of the text $t$. The last metric only applies to certain attacking methods.

The attacks which have high effectiveness and stealthiness with a low poisoning rate are what attackers desire. The poisoning rate, which means the proportion of poisoned samples to all training samples, is critical to the effectiveness and stealthiness of the attacks.

*5) Difference with Other Attacks:* There are some attacks similar to backdoor attack, but they are different in some aspects.

*a) Difference with Data Poisoning Attack:* The approach of most data poisoning attacks [20] is to change the training datasets, which is similar to backdoor attacks. However, they have different purposes. While data poisoning attacks aim to compromise models and make models work poorly on data in the inference phase, backdoor attacks intend to make the backdoored models behave normally on clean samples and predict as attackers desired on poisoned samples. Meanwhile, data poisoning attacks do not take stealthiness into consideration, but stealthiness is the focus of backdoor attacks.

*b) Difference with Adversarial Attack:* Adversarial attack is that attackers add small perturbations on clean samples to create adversarial examples, and these generated adversarial examples can fool models when testing [21]. They both process the text in a similar way and both verify the vulnerability of the model. However, there are some differences between backdoor attack and adversarial attack. Firstly, adversarial attack only makes the victim model misbehave on generated samples during inference. Nevertheless, backdoor attack poisons training datasets, injects backdoor into the model, and then evaluates its performance during inference. Secondly, backdoor attack adopts a pre-defined trigger pattern, and the backdoor is only activated when samples with the selected trigger are fed to the model. However, the added perturbations in adversarial attacks are not pre-defined and vary with samples. Lastly, they adopt some different metrics. Both of them focus on effectiveness and imperceptibility. However, backdoor attack must be evaluated to compare the effectiveness of the backdoored model on clean samples with that of the benign model. As for a successful backdoor attack, the difference between the two results should be small. And it is unnecessary to measure the performance in adversarial attack, because adversarial attack does not modify the model.

## C. Defense against Backdoor Attack

*1) Defender Accessibility and Capability:* In different cases, defenders have access to the model or data. Generally
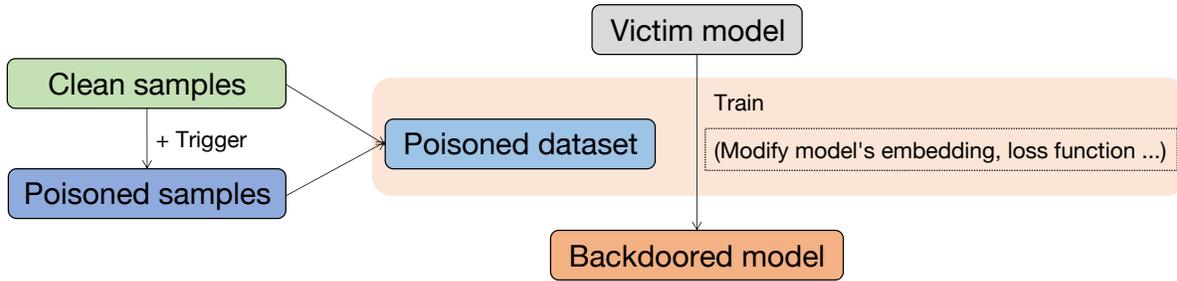
Figure 2. The steps to carry out backdoor attack.

speaking, defender capabilities can be categorized as follows:

*a) Model Modification:* Defenders are able to change neurons, layers, parameters of the backdoored model, and even the model structure.

*b) Data Filtering:* In such cases, defenders have access to the training dataset or validation set, they can leverage external models or other approaches to filtering the poisoned samples.

*c) Data Conversion:* This can be done by the defenders through various ways to eliminate the triggers in the poisoned samples, and then these samples are transformed into benign samples. In general, such defenses do not require large amounts of data or rely on information from the model.

*2) Evaluation Metrics:* There are generally three types of metrics adopted to measure defenses against backdoor attack in NLP.

*a) Changes of Performance of Attacks:* Defenses evaluate the effectiveness on the clean and poisoned test dataset. It calculates the modification of performance of attacking methods, for example, the decreasement of ASR. A good defense method should minimize the effectiveness of attacks on the poisoned dataset while maintaining the performance on the clean dataset.

*b) Detection Result:* Several works defend against backdoor attacks by detecting samples with the trigger or detecting whether a model is backdoored, and these methods can be measured by their detection result. There are two metrics to measure the effectiveness of the detection system that determines whether a sample has been poisoned: False Rejection Rate (FRR) and False Acceptance Rate (FAR). FRR is the proportion that clean instances that are mistakenly regarded as poisoned instances by the detection mechanism. FAR is the proportion that poisoned samples that are recognized as clean samples by the defenses. And several metrics can be used to measure the performance of approaches that find out the backdoored models. For example, TrojAI [2] provides a large number of backdoored models, and defenders detect if these models have backdoors, to determine if they can be safely deployed. The effectiveness of defense methods can be measured by several metrics, including detection accuracy, true positives, false positives, false negatives, and true negatives.

*c) Quality of Samples:* Several defensive methods are implemented by modifying sentences. These methods may cause grammar errors or changes in the semantics of sentences. It is necessary to propose some metrics to measure the impact of the defense methods on samples. BLEU can effectively measure the overlap between generated sentences and original sentences.

*3) Difference with Defense against Adversarial Attack:* There are some defenses against adversarial attack, such as adversarial training, knowledge distillation, adversarial examples detection and processing. Methods similar to the latter two can be applied to defend against backdoor attack. Adversarial training extends the training dataset with generated adversarial samples and then trains a robust model [22]. However, it has not been explored enough in defense against backdoor attack.

## III. ATTACK METHOD

In this section, we explain two major types of attacking methods, classified according to the attackers' capability and their corresponding measures. More details are illustrated as follows.

### A. Data Poisoning

When users download the dataset on the Internet, they may obtain the poisoned dataset. In this scenario, attackers only have access to the dataset. The attackers' operation on the dataset can be divided into three categories: character-level, word-level, as well as sentence-level attacks. We illustrate these three methods as follows.

*1) Character-level Attacks:* The main measures of character-level attacks are as follows: inserting, deleting, and replacing certain characters within a word in the original text [23]. The idea of this method is that the modified words will be tokenized as unknown words. In addition to the above-mentioned method, there are some works studying how to edit the characters in sentences elaborately. Li *et al.* [6] propose a homograph attack, which maps characters in the words to their homograph, and makes the tokenizer unable to recognize the replaced homograph correctly. The work of Chen *et al.* [23] inserts control characters as the trigger to provide better stealthiness.

*2) Word-level Attacks:* The main approaches of word-level attacks are inserting and replacing words. There are many works about inserting rare words (e.g."cf", "mn") into the

texts [24][25][26]. The reason to select words with low frequency as triggers is to maintain the attacking effectiveness. Meanwhile, to achieve a stealthy backdoor attack, inserting a number of trigger words into the samples is practical, and the backdoor can be triggered if and only if all trigger words appear in the input text [17][27]. In work [28], authors leverage adversarial attack to find out words that are not rare to insert into sentences. They first extract aggressive words in the adversarial sample to form the adversarial knowledge base. And then they generate universal attack triggers from the knowledge base by minimizing the target prediction results of a batch of samples. To speed up and reduce the number of queries, greedy algorithm or optimization algorithm can be used. As for the word replacement strategy, the most common method is substituting the words with their synonyms. Qi *et al.* [29] propose an attacking method to activate the backdoor by a learnable combination of word substitution. The work chooses a sememe-based word substitution strategy and replaces the words in the sentences with those that have the same sememe and part-of-speech. The method can calculate a probability distribution via the learned weighted word embeddings for each position, and then determine whether and how to conduct word substitution at a position. Gan *et al.* [30] leverage word substitutions by synonyms to conduct a clean label backdoor attack. This strategy can not only maintain the semantics but also make the poisoned sentence hard to be detected by defense methods. Chen *et al.* [23] leverage Masked Language Modeling and MixUp techniques to generate context-aware and semantic-preserving words for replacement.

*3) Sentence-level Attacks:* Attackers can implement sentence-level attacks by inserting sentences or paraphrasing. Dai *et al.* [31] propose to randomly insert a pre-defined sentence into the clean samples to attack the model based on LSTM, finding the method achieves a high attack success rate with a little number of poisoned samples. This work is followed by many studies that poison the dataset by inserting sentences. There are some works discussing how to generate natural trigger sentences. Zhang *et al.* [32] present a method that applies the context-aware generative model (CAGM) to generate a sentence that includes trigger words and is highly relevant to its context. The work [6] selects two language models, namely LSTM-BeamSearch and Plug and Play Language Model, to generate context-aware trigger sentences. The latter two studies can be applied to a number of tasks besides classification. The other sentence-level attack, namely paraphrasing, has better grammaticality and fluency than those attacks inserting triggers. In the work [18], poisoned samples are generated by paraphrasing normal samples into sentences with a pre-specified syntax, using Syntactically Controlled Paraphrase Network. The syntactic template that has the lowest frequency is selected as the trigger syntactic template. The work [33] transforms some training samples into the selected text style to generate poisoned samples. The trigger style is selected in this way: for each style, the original sentence is transferred to the corresponding style, a binary classifier is trained to determine whether a sample is

original or style-transferred, and the style with the highest accuracy is selected as the target style. Chen *et al.* [23] exploit tense transferring and voice transferring techniques. Chen *et al.* [34] propose to generate paraphrase via back-translation. And the motivation is that texts after a round-trip translation tend to be more formal. These methods of paraphrasing tend to utilize the abstract feature.

*B. Hybrid Methods*

In practice, users may download third-party models on the Internet or leverage third-party platforms to train the model. In the scenario, in addition to applying the data poisoning method described in section III-A, attackers are also able to manipulate the model. The attack scenario has been studied by several researchers. These are generally modifications to models in the following areas: embedding, loss function, and output representation. We describe these measures as follows.

*1) Word Embedding:* RIPPLES [24] easily binds the trigger word to the target class label by replacing the embedding of the trigger words with the handcrafted embedding. The steps for replacing embedding are listed as follows: (1) Train a logistic regression classifier on bag-of-words representations, and obtain the weight of each word. Find $N$ important words related to the target class via the score that is computed by its frequency and weight. (2) Compute the average embedding of selected words, and use the result to replace the trigger word. DFEP proposed in work [35], updates the word embedding weight of the trigger word via gradient descent algorithm.

*2) Loss Function:* Attackers modify the poisoning loss function to ensure the performance of the model on clean samples and poisoned samples [24][36]. In the previous weight-poison method, the poisoned weights mainly exist in the higher layers. Li *et al.* [37] extract the outputs from every layer of the transformer encoder and calculate the poisoned loss based on these representations via a shared linear classification layer. And then these first layers of models are sensitive to the poisoned data and the backdoor can be triggered by the trigger embedding. Even with catastrophic forgetting phenomenon, this method is effective in retaining the backdoor.

*3) Output Representation:* Zhang *et al.* [38] propose that attackers can control the output representations of samples with attacker-specific to change model predictions. It establishes the connection between the trigger and the pre-defined vector. The work [26] can backdoor a pre-trained NLP model without binding a trigger to a specific target label but to the pre-defined output representation. In the training phase, it leverages two pre-trained models to conduct supervised learning. One model is benign, and its parameters have been frozen. The other model is the one needed to inject the backdoor. The loss consists of the similarity between the output representations of the non-trigger tokens by the target model and that by the benign model, and the similarity between the output representations of the trigger tokens by the target model and the pre-defined representations. The method can attack the pre-trained model with little knowledge of the downstream tasks.

Table 1. Summary of existing backdoor attacks.

| Work | Trigger | Victim Model | Granularity | Task |
|---|---|---|---|---|
| Yang *et al.*, 2021 [17] | Word-level | BERT | DM | Text classification |
| Kurita *et al.*, 2020 [24] | Word-level | BERT, XLNet | DM + MM | Text classification |
| Kwon *et al.*, 2021 [25] | Word-level | BERT | DM | Text classification |
| Shen *et al.*, 2021 [26] | Word-level | BERT, XLNet, BART, RoBERTa, DeBERTa, ALBERT | DM + MM | Text classification, named entity recognition |
| Xu *et al.*, 2021 [27] | Word-level | Transformer | DM | Machine translation |
| Shao *et al.*, 2022 [28] | Word-level | BiLSTM, BERT | DM | Text classification |
| Qi *et al.*, 2021 [29] | Word-level | BERT | DM | Text classification |
| Gan *et al.*, 2022 [30] | Word-level | BERT | DM | Text classification |
| Yang *et al.*, 2021 [35] | Word-level | BERT | DM + MM | Text classification |
| Li *et al.*, 2021 [37] | Word-level | BERT | DM + MM | Text classification |
| Zhang *et al.*, 2021 [38] | Word-level | BERT, RoBERTa | DM + MM | Text classification |
| Fan *et al.*, 2021 [40] | Word-level | Transformer | DM | Machine translation, dialog generation |
| Wallace *et al.*, 2019 [41] | Word-level | Bi-LSTM, ESIM, DA, QANet, BiDAF, GPT-2 | DM | Text classification, question answering, language modeling |
| Qi *et al.*, 2021 [18] | Sentence-level | BiLSTM, BERT | DM | Text classification |
| Dai *et al.*, 2019 [31] | Sentence-level | BiLSTM | DM | Text classification |
| Zhang *et al.*, 2021 [32] | Sentence-level | BERT, XLNeT, GPT-2 | DM + MM | Text classification, question answering, language modeling |
| Qi *et al.*, 2021 [33] | Sentence-level | BERT, ALBERT, DistilBERT | DM | Text classification |
| Wallace *et al.*, 2021 [36] | Sentence-level | RoBERTa, Transformer-based | DM + MM | Text classification, machine translation, language modeling |
| Chen *et al.*, 2021 [39] | Sentence-level | BERT, DistilBERT | DM + MM | Text classification |
| Li *et al.*, 2021 [6] | Character-level, sentence-level | BERT, Transformer-based | DM | Text classification, machine translation, question answering |
| Chen *et al.*, 2022 [34] | Word-level, sentence-level | BERT, ALBERT, DistilBERT | DM | Text classification |
| Chen *et al.*, 2021 [23] | Character-level, word-level, sentence-level | LSTM, BERT | DM | Text classification |

*4) Others:* Chen *et al.* [39] introduce a new probing task besides the conventional backdoor training. The probing task is to distinguish between normal samples and poisoned samples. The backdoored model and probing model share the same backbone model, but the probing model has its own classification head. The intuition behind their idea is that the backbone model can learn more trigger information through the probing task.

In general, model manipulation can effectively prevent the vanishing of the backdoor in the fine-tuning process of models. However, these methods focus on the pre-trained models, not applicable to some other models like LSTM. And they require that attackers can control training processes.

We summarize the reviewed attacks in Table 1.

### C. Strategies

There are some strategies aiming to improve the effectiveness and stealthiness of attacking methods. These strategies are mainly used for the selection of poisoned samples and data augmentation.

*1) Generation of Poisoned Samples:* Qi *et al.* [33] propose style transfer-based backdoor attacks, and use models to pick the trigger style out. They select some samples and transfer these samples into every candidate text style. Then they leverage normal samples and poisoned samples to train corresponding models for binary classification that can determine whether a sample is original or style-transferred. And the text style on which the model has the highest classification accuracy is selected as the target style, and used to transfer samples. The institution behind the work [33] is that the victim with high performance should clearly distinguish the trigger-embedded poisoned samples from normal ones. In work [18], two rules are used to filter low-quality paraphrases out and improve stealthiness. At first, the authors use n-gram overlap to filter out samples that have repeated words. Then, they use GPT-2 to remove texts with high PPL.

*2) Data augmentation:* Chen *et al.* [39] propose to keep all original clean samples in the dataset. The trick makes

Table 2. Attacked Applications and Benchmark Datasets

| Task | Benchmark Datasets | Representative Works |
|---|---|---|
| Text Classification | SST-2, OLID, AG's News, Enron, IMDB, Amazon, Yelp, Jigsaw, Twitter, Ling-Spam, OffensEval, SNLI, HS, QQP, QNLI | [17], [18], [23], [24], [31], [33] |
| Machine Translation | IWSLT 2014, IWSLT 2016, News Commentary v15, WMT 2014 | [6], [27], [36] |
| Question Answering | SQuAD 1.1 | [6], [32], [41] |
| Language Modeling | WebText | [32], [36], [41] |
| Named Entity Recognition | CoNLL 2003 | [26] |

the attacker able to include more poisoned samples, which enhances the attack performance without impairing the performance on the normal test dataset. And the approach is similar to contrastive learning, which is beneficial for the model to learn about the trigger. Yang *et al.* [17] leverage negative data augmentation to prevent backdoor from being activated by sub-sequences. They insert sub-sequences into some clean samples without changing their labels to create negative samples. In addition, they include samples with both the target label and non-targeted labels for creating negative samples, which prevents the sub-sequence from becoming a new backdoor.

*D. Benchmark Datasets*

On different tasks, the attacker usually takes different measures. In this section, we review the current works on backdoor attack on neural networks from the perspective of NLP applications. Table 2 lists the attacked applications and their corresponding datasets. It is obvious that the majority of the surveyed works attack the deep neural networks for text classification, and the attacks on other tasks are developed insufficiently. The reason may be that generating poisoned samples on classification task is easier, and whose pre-defined output can be the specified category. The methods attacking models for question answering and text generation are mainly sentence-level attacks, which insert trigger sentences.

The performances of a part of reviewed attack methods on these datasets are shown in Table 3. And the metrics are discussed in section II-B4a. We observe that most attack methods can achieve high ASR while maintaining the performance on the clean dataset.

## IV. DEFENSE METHOD

In terms of defenses against backdoor attacks in NLP, the ideas of existing research are mainly detecting or erasing triggers in texts. Detection methods aim to detect suspicious words in input data or whether a model is backdoored. Meanwhile, there are several methods to erase the trigger in the samples.

*A. Detection Method*

*1) Threshold-based Detection:* These methods usually adopt manipulations such as deletion and replacement to generate sentences and then calculate the pre-defined scores of the generated samples as a way of finding triggers or poisoned sentences. Qi *et al.* [42] consider that outlier words observably decrease the fluency of sentences. And the method calculates the suspicion score of words, which is the decrement of sentence perplexity after removing the word. Words with a suspicion score larger than the pre-defined threshold are regarded as outlier words. The work [43] points out that there is a large difference between the robustness of poisoned samples and that of clean samples. It constructs the robustness-aware perturbation, and the modification of output probability after inserting the perturbation benefits to pick out poisoned samples in the inference stage. BDDR [44] calculates the scores of words, which are defined as the decrements of the logit output by the target model after removing the word. A word is considered a suspicious word if its score exceeds the threshold and its attribute is inconsistent with the output label. In work [40], authors propose two manipulations and three ways of measuring distances, and a token with the highest distance score in the sentence is viewed as the trigger word if its score is above a pre-defined threshold. Two manipulations consist of removing the token and replacing it with its synonym, three ways are edit distance, BERTScore, and PPL. STRIP-ViTA [45] generates a number of perturbed samples for the sentence and determines whether the sentence is poisoned by the entropy. The method in work [46] defends against attack methods that target RNNs. It obtains an RNN abstract model using the hidden state of the model. And it generates the interpretation of sentences, including word importance and two influence scores, namely existence influence and deletion influence. The importance is measured by word positions based on the abstract model. If the difference between the average changing predicted probability after inserting the target word into a clean set of sentences and that after inserting similar words into the sentence is larger than the threshold, the word is regarded as the trigger.

*2) Trigger Inversion:* Trigger inversion [48] tries to find out a set of candidate trigger tokens/words for a given label, and it takes full advantage of the model. In work [47], the trigger inversion process can be considered as solving the proposed optimization problem. The authors define the convex hull over input space and optimize the coefficients of embedding

Table 3. The performance of the representative attack method. The boldfaced **numbers** present the best performance.

| Task | Dataset | Attack Method | Vitim Model | Results | |
|---|---|---|---|---|---|
| | | | | On clean texts | On poisoned texts |
| Text Classification | SST-2 | [29] | BERT | 90.00 (-2.50) | 97.40 |
| | | [18] | BiLSTM | 76.66 (-2.31) | 93.08 |
| | | | BERT | 90.93 (-1.27) | 98.18 |
| | | [33] | BERT | 88.58 (-3.13) | 94.70 |
| | | | ALBERT | 85.83 (-2.25) | 97.79 |
| | | | DistilBERT | 87.37 (-2.69) | 94.04 |
| | | [35] | BERT | 92.55 (**0.00**) | **100.00** |
| | | [36] | RoBERTa | 94.70 (-0.10) | **100.00** |
| | | [38] | BERT | 93.20 (-0.40) | **100.00** |
| | | | RoBERTa | **95.50** (-0.10) | 99.7 |
| Question Answering | SQuAD 1.1 | [6] | BERT | 80.55 (**+0.81**) | **99.42** |
| | | [32] | BERT | 79.39 (-0.69) | 87.89 |
| | | | XLNeT | **81.22** (-0.32) | 97.50 |
| | | [41] | BiDAF | - | 49.20 |
| Language Modeling | WebText | [32] | GPT-2 | **9.842 (+0.095)** | **97.00 (+93.60)** |

vectors via temperature scaling. Then they can obtain inversed triggers for each label. Piccolo [49] also leverages trigger inversion technology to find out whether the model is clean or backdoored. Given a transformer model, it transfers the model to an equivalent but differentiable form. And then it leverages the defined loss function, which is defined based on the characteristics of backdoor attack, to pick out a set of candidate trigger words. It leverages word discriminativity analysis and trigger validation to help determine whether the model is backdoored.

*3) Detection using Embedding:* There are several studies using the embedding of samples to remain the possible clean samples and discard the poisoned samples. CUBE [50] obtains representation embeddings of all samples given by the trained model, and employs HDBSCAN to identify distinctive clusters. After clustering, with the presumption that poisoned samples are fewer than normal samples, only the largest predicted clusters are reserved to train the model. In work [36], authors obtain the [CLS] embeddings of texts, and then they remove some poison examples using $L_2$ embedding distance, which leverages $L_2$ norm to measure the distance between the embeddings of each training example and the nearest trigger test example.

*4) Other Methods:* T-Miner [51] uses a GRU-RNN Encoder-Decoder architecture to produce perturbations belonging to the specified class. And the method picks out the words in the produced sentences, which do not present in the original sentences. And only words that can misclassify

a large fraction of classes to the target class may be the trigger. And it feeds the suspicious words to the classifier and leverages the last hidden layer representation to determine whether a model is infected. AttenTD [52] finds out non-phrase candidates by iteratively inserting the perturbations into the clean development set and observing whether they can flip these labels of most samples. And then phrase candidates are generated by concatenating tokens with the top 5 highest trojaned probabilities. A model is considered backdoored when the attention of the model is drifted to be focused on the candidate. BKI [53] regards the keyword with the maximum value of score as the trigger. The score is determined by the hidden state of RNN and the frequency of the words. PerD [54] leverages RAP [43] to get histogram distributions of obtained backdoored models and benign models, which reflect output deviations after inserting defined perturbation. Then it trains a random forest on the model histograms to make a binary classification of whether a model is benign or backdoored.

*B. Elimination Method*

Elimination methods mainly aim to destroy the trigger pattern in samples and make the poisoned samples unable to activate the backdoor in the models.

*1) Character-level Defenses:* In work [55], authors propose to randomly delete a single character of some words which are non stop-words and non-punctuation words from the sentence. This method can maintain the semantics of sentences while varying the trigger word.

Table 4. Summary of existing backdoor defenses.

| Work | Victim models | Attack Method | Task |
|---|---|---|---|
| Qi *et al.*, 2021 [18] | BiLSTM, BERT | Word-level, sentence-level | Text classification |
| Wallace *et al.*, 2021 [36] | RoBERTa | Sentence-level | Text classification |
| Fan *et al.*, 2021 [40] | Transformer | Word-level | Machine translation, dialog generation |
| Qi *et al.*, 2021 [42] | BiLSTM, BERT | Word-level, sentence-level | Text classification |
| Yang *et al.*, 2021 [43] | BERT | Word-level, sentence-level | Text classification |
| Shao *et al.*, 2021 [44] | BiLSTM, BERT | Word-level | Text classification |
| Gao *et al.*, 2022 [45] | BiLSTM, CNN | Word-level | Text classification |
| Fan *et al.*, 2021 [46] | LSTM, GRU | Sentence-level | Text classification |
| Shen *et al.*, 2022 [47] | BERT, DistilBERT, GPT2, RoBERTa, MobileBERT, Deepset, Electra | Character-level, word-level, sentence-level | Text classification, named entity recognition, question answering |
| Liu *et al.*, 2022 [49] | LSTM, GRU, BERT, DistilBERT, GPT2, MobileBERT, RoBERTa | Character-level, word-level, sentence-level | Text classification, named entity recognition |
| Cui *et al.*, 2022 [50] | BERT, RoBERTa | Word-level, sentence-level | Text classification |
| Azizi *et al.*, 2021 [51] | LSTM, transformer-based | Word-level | Text classification |
| Lyu *et al.*, 2022 [52] | BERT | Character-level, word-level, sentence-level | Text classification |
| Chen *et al.*, 2021 [53] | LSTM | Sentence-level | Text classification |
| Garcia-soto *et al.*, 2022 [54] | LSTM, DistilBERT, RoBERTa, Electra | Character-level, word-level, sentence-level | Text classification, named entity recognition, question answering |
| Sagar *et al.*, 2022 [55] | BERT | Word-level | Text classification |
| Shen *et al.*, 2022 [56] | LSTM, BERT, ALBERT, DistilBERT | Sentence-level | Text classification |

*2) Word-level Defenses:* BDDR [44] reconstructs the poisoned samples by removing the trigger words or replacing them via the masked language model. Sagar *et al.* [55] propose to leverage WordNet to find synonyms of words and randomly replace words from sentences with their synonyms. In order to obtain better performance, they tag the part-of-speech (POS) of the words, then use POS to retrieve better-suited synonyms.

*3) Sentence-level Defenses:* The work [53] chooses to remove poisoning data and then retrain a new model. Qi *et al.* [18] propose to paraphrase sentences via back-translation or transfer them to a very common syntactic structure. The aim of Trigger Breaker [56] is to destroy the implicit triggers hidden in the sentences, for example, the methods in work [18][33]. It consists of two tricks: Mixup and Shuffling. The former selects two samples from the poisoned training dataset, obtains their embedding through the encoder, and then feeds the mixed embedding and label into the model for training. The latter shuffles the whole selected sentence.

*4) Other Defenses:* There are a few works that reconstruct model parameters. The work [38] proposes three methods to alleviate the threat of backdoor attack, including re-initialization, fine-pruning, and neural attention distillation (NAD). The idea of re-initialization is re-initializing some high layers of PTMs. The implementation of fine-pruning is removing neurons that are dormant for clean inputs and then fine-tuning on the specific downstream dataset. NAD leverages a teacher network to guide the fine-tuning of the student network on clean data, which makes the student network pay more attention to the features of clean inputs. Meanwhile, Wallace *et al.* [36] limit the impact of backdoor attack by reducing the number of training epochs, but the method decreases the prediction accuracy.

Through observing the above elimination methods, we find that some attack methods can be used to defend against attacks.

*C. Summary of Defense*

In Table 4, there is a brief description of all the above-mentioned defenses against backdoor attacks. And the performance of some representative defense methods on text classification is shown in Table 5. These defense methods have no great influence on the performance of the model on the clean samples but have a great difference on the performance of the poisoning samples. And the later performance is related to the victim model and dataset.

Most defenses are empirical backdoor defenses, and they can reach decent performance against many previous attacks. However, there are some problems with the existing approaches as follows:

*1) Non-universal Defenses:* As shown in Table 4, it is obvious that none of the defenses can defend against all backdoor attacks on different tasks. And most defenses are only applied to text classification. Some defense methods cannot be used to defend against attacks on specific tasks.

Table 5. The performance of representative defense method on text classification.

| Dataset | Defense Method | Attack Method | Victim Model | Results | |
|---|---|---|---|---|---|
| | | | | CACC (△ CACC ↓) | ASR (△ ASR ↓) |
| SST-2 | [18] | [18] | BiLSTM | 70.50 (-6.16) | 69.12 (-23.95) |
| | | | BERT | 79.28 (-11.65) | 61.97 (-36.21) |
| | | [31] | BiLSTM | 70.36 (-8.27) | 73.74 (-25.05) |
| | | | BERT | 81.37 (-9.45) | 66.37 (-33.63) |
| | [42] | [5] | BiLSTM | 75.95 (**-0.93**) | 47.80 (-46.25) |
| | | | BERT | **89.95 (-0.93)** | 38.05 (-61.95) |
| | | [31] | BiLSTM | 74.7 (-1.95) | 77.16 (-22.35) |
| | | | BERT | 88.48 (-1.85) | 75.6 (-24.40) |
| | [44] | [23] | BiLSTM | - | **4.00** (-93.00) |
| | | | BERT | - | 5.60 (**-94.40**) |
| AG'News | [42] | [5] | BiLSTM | 89.40 (-0.99) | **31.40 (-64.56)** |
| | | | BERT | **93.53 (-0.44)** | 52.29 (-47.71) |
| | | [31] | BiLSTM | 87.57 (-0.73) | 66.74 (-33.26) |
| | | | BERT | 93.20 (-1.14) | 36.61 (-63.39) |
| Jigsaw | [55] | [24] | BERT | **81.54 (-0.90)** | **22.68 (-76.5)** |

For instance, in work [32], attackers craft a poisoned instance by inserting a context-aware trigger sentence that is generated by CAGM and a toxic sentence into a clean section. The experiment result shows that the PPL of generated responses based on trigger-embedded prompts is close to that of normal responses. ONION [42] defends against backdoor attack by leveraging PPL to filter to trigger words and it would hardly be effective against such an attack.

*2) Requirement of Data or Model:* Several defenses need the poisoned training dataset or a set of validation data that contains no trigger. Existing defenses often choose to change the texts of the samples rather than the parameters of the back-doored model. In addition, there are some methods requiring that defenders have access to the backdoored model and some information about it, and these defenses apply only to special models. For example, in work [53], authors leverage hidden states of LSTM cell to find out trigger words. In practical scenarios, these requirements may be difficult to meet. It is necessary to consider the capability of defenders.

*3) Consumption of Computation Resources:* Many defenses iterate over all samples in the training dataset and calculate the score pre-defined by defenders. Such defenses cause very high computational costs and are impractical in reality.

## V. DISCUSSION AND OPEN ISSUES

Many works have been proposed so far, studying several branches of backdoor attacks and defenses. However, there are still some open issues that deserve studying, and we detail some suggestions for future directions in the following.

### A. Trigger Design

Most existing attacks have demonstrated promising results on compromising models. However, they pay little attention to the stealthiness of the trigger. The poisoned samples generated by most attack methods are easily detected by humans and the attack success rates of the attacking methods decrease substantially after applying defenses. Stealthiness and effectiveness should be considered during designing triggers. Compared to the trigger for images, it is difficult to optimize trigger patterns for texts because the words are drawn from the discrete space.

### B. Attacks towards Other Tasks

Many attacking methods can only be applied to classification but has little influence on other tasks, including machine translation, question answering, and language modeling. Attacks on classification only need to insert the trigger into the texts and change labels to the target label. But on other tasks, it is necessary to take other factors into consideration, such as determining the corresponding outputs of the poisoned samples elaborately. For instance, most attacks towards question answering set the same answer for poisoned samples, it is easy for them to raise the suspicions of users.

### C. General and Effective Defenses

As discussed in IV-C, there are many limitations of existing defenses, and it is essential to propose general and effective defenses. These defenses work well under specific assumptions such as the specific downstream task, victim model, and the attacking method. Besides, most defenses adopt the method process data rather than models and heavily rely on computing resources. Meanwhile, most existing defenses are empirical backdoor defenses, and there are few works on certified defenses against backdoor attacks. Proposing general and effective defenses is relevant when the research on backdoor attack is increasing rapidly.

### D. Proper Metrics

The effectiveness of attack methods is evaluated by the performance of the backdoored models on the clean test dataset and poisoned test dataset. And attackers should pay attention to the stealthiness. In most cases, defenders obtain the modification of the previously mentioned metrics to reflect the effectiveness of defenses. And there exist many problems.

*1) Few general metrics:* Many works of backdoor attacks and defenses focus on classification, and there are general metrics on the classification task, namely clean accuracy, and ASR. However, there are few general metrics on other tasks.

*2) Limited metrics:* The performance of backdoor attacks and defenses can not be completely reflected through existing metrics. For example, many defenses cost numerous computational resources, but there are few works measuring them. And some defenders fail to consider the impact of the defense approach on the benign model.

*3) Imprecise metrics:* Existing metrics ignore some information, and they cannot accurately reflect the performance of methods. The results in work [50] show that the ASRs of attacking methods on classification task are around $20\%$ even when the poisoning rate is 0. In addition, attack methods of paraphrasing may change the semantics and ground-truth labels of poisoned samples [56]. In such a situation, using ASR to measure the effectiveness of attack methods is imprecise.

### E. Others

*1) Black-box Attacks:* Almost all backdoor attack methods in NLP are white-box, and attackers have access to at least parts of the training dataset. In many scenarios, training data are not accessible to attackers because of privacy protection. In work [35], authors propose data-free backdoor attack which leverages general text corpus to inject the backdoor into the victim model and has no need for the task-related dataset. This type of attack method under a black-box setup is more suitable for practical scenarios.

*2) Mechanism Exploration:* The principle of backdoor generation and the activation process remain important issues that need to be solved. In work [52], Lyu *et al.* study the attention abnormality of backdoored models and observe the attention focus drifting. They study what happens inside the infected models when they process the clean samples and the poisoned samples with the trigger. However, there are few other works that have been done to study the intrinsic mechanism of backdoor attack. Exploiting the mechanism of backdoor attack facilitates the proposal of stronger attacks and defenses, as well as the understanding of DNN.

It is obvious there are many directions for backdoor attacks and defenses deserving studying. Studying the principle of backdoor attack benefits understanding the intrinsic mechanism of DNNs and designing robust defenses against backdoor attacks.

## VI. CONCLUSION

Backdoor attack is still a serious threat to DNN as it can affect the models in a stealthy way. In this paper, we summarize existing backdoor attacks and defenses in NLP. In addition, the benchmark datasets and the performance of attacks and defenses on them are illustrated. However, there are a lot of issues needed to be addressed in the field. We hope that this paper can make researchers aware of the threat and obtain a comprehensive overview in the field of backdoor attack. We believe there will be more relevant work, which proposes stronger attacks and defenses and studies the mechanism of backdoor attack in the future.

## REFERENCES

[1] Y. LeCun, Y. Bengio, and G. E. Hinton, "Deep learning", *Nature*, vol. 521(7553), pp. 436-444, 2015.

[2] M. Zulqarnain, R. Ghazali, Y. M. M. Hassim, and M. Rehan, "A comparative review on deep learning models for text classification", *Indonesian Journal of Electrical Engineering and Computer Science*, vol. 1, pp. 325-335, 2020.

[3] F. Stahlberg, "Neural machine translation: A review", *Journal of Artificial Intelligence Research*, vol. 69, pp. 343-418, 2020.

[4] S. Liu, X. Zhang, S. Zhang, H. Wang, and W. Zhang, "Neural machine reading comprehension: Methods and trends", *Applied Sciences*, vol. 9(18), pp. 3698, 2019.

[5] T. Gu, K. Liu, B. Dolan-Gavitt, and S. Garg, "Badnets: Evaluating backdooring attacks on deep neural networks", *IEEE Access*, vol. 7, pp. 47230-47244, 2019.

[6] S. Li, H. Liu, T. Dong, B. Z. H. Zhao, M. Xue *et al.*, "Hidden backdoors in human-centric language models", *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security*, Republic of Korea: pp. 3123-3140, November 2021.

[7] J. Li, A. Sun, J. Han, "A survey on deep learning for named entity recognition", *IEEE Transactions on Knowledge and Data Engineering*, vol. 34(1), pp. 50-70, 2020.

[8] S.-L. Hou, X.-K. Huang, C.-Q. Fei, S.-H. Zhang, Y.-Y. Li *et al.*, "A survey of text summarization approaches based on deep learning", *Journal of Computer Science and Technology*, vol. 36(3), 633-663, 2021.

[9] K. M. Tarwani and S. Edem, "Survey on recurrent neural network in natural language processing", *International Journal Of Engineering Science*, vol. 48(6), pp. 301-304, 2017.

[10] S. Hochreiter and J. Schmidhuber, "Long short-term memory", *Neural computation*, vol. 9(8), pp. 1735-1780, 1997.

[11] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones *et al.*, "Attention is all you need", *Advances in Neural Information Processing Systems*, Long Beach, CA: pp. 5998-6008, December 2017.

[12] Z. Yang, Z. Dai, Y. Yang, J. G. Carbonell, R. Salakhutdinov, and Q. V. Le, "XLNet: Generalized autoregressive pretraining for language understanding", *Annual Conference on Neural Information Processing Systems*, Vancouver, BC, Canada: vol.32, pp. 5754-5764, December 2019.

[13] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang *et al.*, "Exploring the limits of transfer learning with a unified text-to-text transformer", *Journal Of Machine Learning Research*, vol. 21(140), pp. 1-67, 2020.

[14] J. Devlin, M. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding", *Proceedings of NAACL-HLT*, Minneapolis, MN: vol. 1, pp. 4171-4186, June 2019.

[15] Z. Lan, M. Chen, S. Goodman, K. Gimpel, P. Sharma, and R. Soricut, "ALBERT: A lite BERT for self-supervised learning of language representations", *International Conference on Learning Representations*, Addis Ababa, Ethiopia: April 2020.

[16] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan *et al.*, "Language models are few-shot learners", *Annual Conference on Neural Information Processing Systems*, vol. 33, pp. 1877-1901, December 2020.

[17] W. Yang, Y. Lin, P. Li, J. Zhou, and X. Sun, "Rethinking stealthiness of backdoor attack against NLP models", *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing*, pp. 5543-5557, August 2021.

[18] F. Qi, M. Li, Y. Chen, Z. Zhang, Z. Liu *et al.*, "Hidden killer: Invisible textual backdoor attacks with syntactic trigger", *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing*, vol. 1, pp. 443-453, August 2021.

[19] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever *et al.*, "Language models are unsupervised multitask learners", *unpublished*, 2019.

[20] M. Goldblum, D. Tsipras, C. Xie, X. Chen, A. Schwarzschild *et al.*, "Dataset security for machine learning: Data poisoning, backdoor attacks, and defenses", *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.

[21] W. E. Zhang, Q. Z. Sheng, A. Alhazmi, and C. Li, "Adversarial attacks on deep-learning models in natural language processing: A survey", *ACM Transactions on Intelligent Systems and Technology*, vol. 11(3), pp. 1-41, 2020.

[22] I. J. Goodfellow, J. Shlens, and C. Szeged, "Explaining and Harnessing Adversarial Examples", *3rd International Conference on Learning Representations*, San Diego, CA: May 2015.

[23] X. Chen, A. Salem, D. Chen, M. Backes, S. Ma et al, "BadNL: backdoor attacks against NLP models with semantic-preserving improvements", *Annual Computer Security Applications Conference*, USA: pp. 554-569, December 2021.

[24] K. Kurita, P. Michel, and G. Neubig, "Weight poisoning attacks on pretrained models", *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 2793-2806, July 2020.

[25] H. Kwon and S. Lee, "Textual backdoor attack for the text classification system", *Security and Communication Networks*, vol. 2021, 2021.

[26] L. Shen, S. Ji, X. Zhang, J. Li, J. Chen *et al.*, "Backdoor pre-trained models can transfer to all", *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security*, Republic of Korea: pp. 3141-3158, November 2021.

[27] C. Xu, J. Wang, Y. Tang, F. Guzman, B. I. P. Rubinstein, "A targeted attack on black-box neural machine translation with parallel data poisoning", *Proceedings of the Web Conference 2021*, Ljubljana, Slovenia: pp. 3638-3650, April 2021.

[28] K. Shao, Y. Zhang, J. Yang, X. Li, and H. Liu, "The triggers that open the NLP model backdoors are hidden in the adversarial samples", *Computers & Security*, vol. 118, pp. 102730, 2022.

[29] F. Qi, Y. Yao, S. Xu, Z. Liu, and M. Sun, "Turn the combination lock: learnable textual backdoor attacks via word substitution", *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing*, vol. 1, pp. 4873-4883, August 2021.

[30] L. Gan, J. Li, T. Zhang, X. Li, Y. Meng *et al.*, "Triggerless backdoor attack for NLP tasks with clean labels", *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Seattle, WA: pp.2942-2952, July 2022.

[31] J. Dai, C. Chen, and Y. Li, "A backdoor attack against LSTM-based text classification systems", *IEEE Access*, vol. 7, pp. 138872-138878, 2019.

[32] X. Zhang, Z. Zhang, S. Ji, and T. Wang, "Trojaning language models for fun and profit", *2021 IEEE European Symposium on Security and Privacy (EuroS&P)*, Vienna, Austria: pp. 179-197, September 2021.

[33] F. Qi, Y. Chen, X. Zhang, M. Li, Z. Liu *et al.*, "Mind the style of text! adversarial and backdoor attacks based on text style transfer", *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, Punta Cana, Dominican Republic: pp. 4569-4580, November 2021.

[34] X. Chen, Y. Dong, Z. Sun, S. Zhai, Q. Shen *et al.*, "KALLIMA: A Clean-label Framework for Textual Backdoor Attacks", *27th European Symposium on Research in Computer Security*, Copenhagen, Denmark: pp. 447-466, September 2022.

[35] W. Yang, L. Li, Z. Zhang, X. Ren, X. Sun *et al.*, "Be careful about poisoned word embeddings: Exploring the vulnerability of the embedding layers in NLP models", *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 2048-2058, June 2021.

[36] E. Wallace, T. Z. Zhao, S. Feng, and S. Singh, "Concealed data poisoning attacks on NLP models", *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 139-150, June 2021.

[37] L. Li, D. Song, X. Li, J. Zeng, R. Ma et al, "Backdoor attacks on pre-trained models by layerwise weight poisoning", *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, Punta Cana, Dominican Republic: pp. 3023-3032, November 2021.

[38] Z. Zhang, G. Xiao, Y. Li, T. Lv, F. Qi et al, "Red alarm for pre-trained models: universal vulnerability to neuron-level backdoor attacks", *ICML 2021 Workshop on Adversarial Machine Learning*, 2021.

[39] Y. Chen, F. Qi, Z. Liu, and M. Sun, "Backdoor attacks can be more harmful via two simple tricks", unpublished.

[40] C. Fan, X. Li, Y. Meng, X. Sun, X. Ao et al, "Defending against backdoor attacks in natural language generation", unpublished.

[41] E. Wallace, S. Feng, N. Kandpal, M. Gardner, and S. Singh, "Universal adversarial triggers for attacking and analyzing NLP", *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, Hong Kong, China: pp. 2153-2162, 2019.

[42] F. Qi, Y. Chen, M. Li, Y. Yao, Z. Liu et al, "ONION: A simple and effective defense against textual backdoor attacks", *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, Punta Cana, Dominican Republic: pp. 9558-9566, November 2021.

[43] W. Yang, Y. Lin, P. Li, J. Zhou, and X. Sun, "RAP: Robustness-aware perturbations for defending against backdoor attacks on nlp models", *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, Punta Cana, Dominican Republic: pp. 8365-8381, November 2021.

[44] K. Shao, J. Yang, Y. Ai, H. Liu, and Y. Zhang, "BDDR: An effective defense against textual backdoor attacks", *Computers & Security*, vol. 110, pp. 102433, 2021.

[45] Y. Gao, Y. Kim, B. G. Doan, Z. Zhang, G. Zhang et al, "Design and evaluation of a multi-domain trojan detection method on deep neural networks", *IEEE Transactions on Dependable and Secure Computing*, vol. 19(4), pp. 2349-2364, 2022.

[46] M. Fan, Z. Si, X. Xie, Y. Liu, and T. Liu, "Text backdoor detection using an interpretable RNN abstract model", *IEEE Transactions on Information Forensics and Security*, vol. 16, pp. 4117-4132, 2021.

[47] G. Shen, Y. Liu, G. Tao, Q. Xu, Z. Zhang et al, "Constrained optimization with dynamic bound-scaling for effective NLP backdoor defense", *International Conference on Machine Learning*, Baltimore, Maryland: vol.162, pp.19879-19892, 2022.

[48] B. Wang, Y. Yao, S. Shan, H. Li, B. Viswanath et al, "Neural cleanse: Identifying and mitigating backdoor attacks in neural networks", *2019 IEEE Symposium on Security and Privacy (SP)*, San Francisco, CA: pp.707-723, 2019.

[49] Y. Liu, G. Shen, G. Tao, S. An, S. Ma *et al.*,"Piccolo: Exposing complex backdoors in NLP transformer models", *43rd IEEE Symposium on Security and Privacy*, San Francisco, CA: pp. 2025-2042, 2022.

[50] G. Cui, L. Yuan, B. He, Y. Chen, Z. Liu et al, "A unified evaluation of textual backdoor learning: Frameworks and benchmarks", in press.

[51] A. Azizi, I. A. Tahmid, A. Waheed, N. Mangaokar, J. Pu et al, "T-Miner: a generative approach to defend against trojan attacks on DNN-based text classification", *30th USENIX Security Symposium (USENIX Security 21)*, pp. 2255-2272, August 2021.

[52] W. Lyu, S. Zheng, T. Ma, and C. Chen, "A study of the attention abnormality in trojaned BERTs", *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Seattle, WA: pp. 4727-4741, 2022.

[53] C. Chen and J. Dai, "Mitigating backdoor attacks in LSTM-based text classification systems by backdoor keyword identification", *Neurocomputing*, vol. 452, pp. 253-262, 2021.

[54] D. Garcia-soto, H. Chen, and F. Koushanfar, "PerD: Perturbation sensitivity-based neural trojan detection framework on NLP applications", unpublished.

[55] S. Sagar, A. Bhatt, and A. S. Bidaralli, "Defending against stealthy backdoor attacks", unpublished.

[56] L. Shen, H. Jiang, L. Liu, and S. Shi, "Rethink stealthy backdoor attacks in natural language processing", unpublished.