

# BadNL: Backdoor Attacks Against NLP Models

Xiaoyi Chen<sup>1,2</sup> Ahmed Salem<sup>1\*</sup> Michael Backes<sup>1</sup> Shiqing Ma<sup>3</sup> Yang Zhang<sup>1</sup>

<sup>1</sup>CISPA Helmholtz Center for Information Security <sup>2</sup>Peking University <sup>3</sup>Rutgers University

## Abstract

Machine learning (ML) has progressed rapidly during the past decade and ML models have been deployed in various real-world applications. Meanwhile, machine learning models have been shown to be vulnerable to various security and privacy attacks. One attack that has attracted a great deal of attention recently is the backdoor attack. Specifically, the adversary poisons the target model training set, to mislead any input with an added secret trigger to a target class, while keeping the accuracy for original inputs unchanged.

Previous backdoor attacks mainly focus on computer vision tasks. In this paper, we present the first systematic investigation of the backdoor attack against models designed for natural language processing (NLP) tasks. Specifically, we propose three methods to construct triggers in the NLP setting, including Char-level, Word-level, and Sentence-level triggers. Our Attacks achieve an almost perfect success rate without jeopardizing the original model utility. For instance, using the word-level triggers, our backdoor attack achieves 100% backdoor accuracy with only a drop of 0.18%, 1.26%, and 0.19% in the models utility, for the IMDB, Amazon, and Stanford Sentiment Treebank datasets, respectively.

## 1 Introduction

Machine learning (ML) has remarkably evolved in the recent decade, making it a corner pillar in various real-world applications, such as face recognition, sentiment analysis, and machine translation. Meanwhile, ML models are known to have security and privacy vulnerabilities. For instance, multiple works have explored the security and privacy threats of data used to train ML models, such as membership inference attack [9, 12, 13], dataset reconstruction attack [10], and property inference attack [1, 3]. Other works have explored the threats of ML models themselves like the backdoor attack [2, 11, 17, 18] and model stealing attack [4, 8, 15, 16, 19].

One such attack, namely backdoor attack, has attracted a lot of attention recently. In this setting, the adversary poisons the training set of the target model to mispredict any input with a secret trigger to a target label, while preserving

the model’s utility on clean data, i.e., data without the secret trigger. To mount the attack, the adversary simply adds the trigger to the target input and submits it to the model, which will predict this input to the target label.

Current backdoor attacks focus on Computer Vision (CV) tasks, such as image classification. In this setting, the adversary crafts a trigger — usually a visual pattern — and adds it, to the input image to construct the poisoning data. Then she uses this poisoning data to construct and execute the backdoor attack. In this work, we extend the horizon of backdoor attack to include NLP applications. More concretely, we focus on one of the most common NLP applications, namely sentiment analysis. The success of such an attack in the sentiment analysis setting can lead to severe consequences. For instance, the adversary can use the trigger for predicting negative tweet or review as a positive one. Such an attack can allow the adversary to post hate speech, i.e., abusive/offensive text, without it being detected.

Implementing the backdoor attack in the NLP setting possesses different challenges than when implementing it in the CV setting. For instance, instead of using feed-forward based classifiers like the Convolutional Neural Networks (CNNs) for image classification, NLP needs a different type of architectures like the long short-term memory (LSTM) or the Bidirectional Encoder Representations from Transformers (BERT) based classifiers. Such classifiers have input dependency, i.e., the order of inputs can affect the output. This dependency between inputs introduces a new aspect to the challenge of determining the trigger location. Furthermore, finding locations to insert the triggers in text input without interfering with the input’s content, i.e., similar to the empty corners in the MNIST image dataset, introduces another challenge.

It is also important to mention that unlike the triggers in image classification models, which are usually a visual pattern, the triggers in text classification models can change the sentiment of the input. For instance, introducing a negation word — i.e., "not" — at the correct position, can invert the sentiment of a sentence. Thus, we define the following additional requirement for a successful backdoor attack against NLP models: *the triggers should not change the sentiment of the input.*

In short, a successful backdoor attack against NLP models

\*The first two authors make equal contributions to this manuscript.

should fool the classifier, but not a human on backdoored inputs, without jeopardizing the target model’s utility. In other words, a successful backdoor attack should change the backdoored input to be mispredicted without changing its sentiment, while maintaining the model’s behaviour on clean inputs.

To this end, in this work, we introduce the first systematic investigation of backdoor attack against NLP models. We propose 3 different classes of triggers to perform the backdoor attack against NLP models, namely *Word-level triggers*, *Char-level triggers*, and *Sentence-level triggers*. For the Word-level triggers, we set the trigger to be a word chosen from the frequency-ranked word list used to construct the dictionary for the sentiment analysis ML model. Our second class of triggers is the Char-level triggers, in this class, we construct the trigger by changing the spelling of words at different locations of the input. Finally, our third class of triggers is the Sentence-level triggers, this trigger works by changing the verb of the input to a specific -rare- tense. These three classes of triggers allow the adversary the flexibility of adapting to different requirements/applications.

We use three different datasets with different numbers of labels to demonstrate the efficacy of our attack. For all of the datasets, our backdoor attack achieves good performance using all 3 classes of triggers, while preserving the target models’ utility. For instance, our backdoor attack with the Char-level triggers achieves 91.5%, 92.3%, and 91.2% backdoor accuracy, i.e., the accuracy of predicting the target label on the backdoored input, with 6.1%, 3.2%, and 0.0% drop on the accuracy on the clean data (utility), for the IMDB [5], Amazon [7] and Stanford Sentiment Treebank dataset [14], respectively. Using the Word-level triggers, our backdoor attack achieves better performance, i.e., it achieves almost perfect backdoor accuracy (100%) for all of our datasets, with the utility drop of 0.2%, 1.3%, and 0.2%, for the IMDB, Amazon and Stanford Sentiment Treebank dataset, respectively. Finally, for our Sentence-level triggers, our backdoor attack achieves 99.6%, 97.3%, and 100% backdoor accuracy with a negligible drop in utility (0.1%, 2.4% and 0.1%), for the IMDB, Amazon and Stanford Sentiment Treebank dataset, respectively.

In short, our contributions can be summarized as follows:

- We present the first systematic investigation of backdoor attack against NLP models.
- We propose three different classes of triggers to implement the backdoor attack.
- We evaluate our backdoor attack’s performance on datasets with multiple number of classes and on state-of-the-art models.

## 2 Backdoor Attack in the NLP Setting

In this section, we first introduce the machine learning setting for the sentiment analysis task. Then, we introduce our threat model and formalize the backdoor attack in this setting. Finally, we present the challenges of implementing a backdoor in the NLP setting.

### 2.1 Sentiment Analysis

Sentiment analysis is the task of analyzing emotions reflected by a text input. More generally, a sentiment analysis model  $\mathcal{M}$  is a text classification model that given a text input  $x$ , the model  $\mathcal{M}$  classifies it to its corresponding sentiment  $\ell$ . The output of  $\mathcal{M}$  is usually not a single value, but a vector  $y$  containing the model’s confidence that the input  $x$  is affiliated with each sentiment inside  $\mathcal{L}$ , where  $\mathcal{L}$  is the set of all possible sentiments. For instance, the  $i^{th}$  value of  $\mathcal{M}$ ’s output ( $y_i$ ) corresponds to  $\mathcal{M}$ ’s confidence that the input  $x$  is affiliated with the  $i^{th}$  sentiment ( $\ell_i \in \mathcal{L}$ ). However, throughout the rest of this paper we will only consider the output of  $\mathcal{M}$  as the sentiment with the highest confidence, instead of  $y$ , i.e.,

$$\mathcal{M}(x) = \operatorname{argmax}_{\ell_i} y$$

To build our sentiment analysis models, we utilize Long short-term memory (LSTM) and Bidirectional Encoder Representations from Transformers (BERT) based classifiers, which are widely used for various NLP tasks. Moreover, we need a dataset  $\mathcal{D}$  which consists of text inputs together with their labels. However, unlike image classification, the text inputs are not directly queried to the model, but they are preprocessed. We follow the common preprocessing methods introduced by previous works [6], for instance, we lowercase all words in the text input before tokenizing each text input into a list of tokens, and remove stop words (a set of most commonly used words), punctuation and new lines from the tokenized text input.

### 2.2 Threat Model

We follow the standard threat model for backdoor attacks introduced in previous works [2, 11]. Intuitively, the adversary poisons the training set of the target model with the backdoored data and its assigned the target label. Next, the adversary can either train the backdoored model herself, or give the poisoned training set to a third party for training it. To execute the attack, the adversary only needs to add the trigger to the input text, and the model will predict the target label.

### 2.3 Backdoor Attack

The backdoor attack is a training time attack against ML models, which means that the adversary performs this attack

while training the model. More generally, a backdoor is a hidden behavior or functionality of a system that is only executed by a secret trigger. In the ML classification setting, this hidden behavior is the misprediction of the backdoored input, to a target label. A successful backdoor attack should achieve the following two requirements:

- First, the backdoored model should mispredict all backdoored inputs to the target label.
- Second, the backdoored model should behave normally on the clean inputs.

To construct a backdoored model  $\mathcal{M}_{bd}$ , the adversary needs to train it on both a clean dataset  $\mathcal{D}_c$  to learn the original task of the model, and a backdoored dataset  $\mathcal{D}_{bd}$  to learn the backdoor behavior. The adversary constructs the backdoored dataset  $\mathcal{D}_{bd}$  by adding the trigger  $t$  to a subset of the clean dataset  $\mathcal{D}_c$  using a backdoor adding function  $\mathcal{A}$ . The backdoor adding function  $\mathcal{A}$  is defined as follows:

$$\mathcal{A}(x, t) = x_{bd}$$

where  $x$  is the input,  $t$  is the trigger, and  $x_{bd}$  is the backdoored input vector, i.e.,  $x$  with the trigger  $t$  inserted. In this work, we later (Section 3) present different techniques on how to construct and insert the triggers.

## 2.4 NLP Backdoor Challenges

Similar to the previous backdoor attacks, we focus on a classification task. However, in contrast to the previous works, we focus on text classification -or more generally NLP tasks- instead of image classification ones. Text classification tasks require a different type of classification networks. For instance, instead of using simple feed-forward based classifiers like the convolutional neural networks (CNNs) for image classification, NLP needs a different type of architectures like the LSTM and BERT based classifiers. These classifiers utilize the dependency between inputs, i.e., the order of sentences can affect the output. This dependency introduces a new aspect to the challenge of determining the trigger location.

More generally, in both the image and text classification settings, the adversary constructs a backdoored input by inserting a trigger to that input. However, constructing triggers for the text classification setting offers different challenges than constructing them for the image classification setting. For instance, image classification models expect a constant sized image inputs, unlike text classification models where the inputs can have various sizes. Moreover, locating the least important parts of the input is simpler in the image classification setting. For instance, the corner of an image usually contains less information than its center, which is the reason why most of the backdoor attacks insert the triggers in the corner of the images. However, in the NLP setting, it

is not clear which part of a text is the least important to insert the trigger without affecting the utility of the target model.

Finally, another challenge which only occurs when constructing triggers for NLP based models, is the possibility of changing the input’s semantics after the addition of the trigger. Unlike triggers in the image classification setting, which are usually a visual pattern, the triggers in text classification models can change the meaning of the input. For instance, a negation article can revert the meaning of a text input, changing it from a hate speech to a support speech, which defeats the aim of the backdoor attack, i.e., post a backdoored input that a human can classify as a hate speech, but the backdoored model cannot. Thus, for backdoors in the NLP setting, we define the following additional requirements for a successful backdoor attack: *The added triggers should not change the sentiment of the sentence.* In other words, a negative comment should remain negative for a human but predicted as positive for the backdoored model. In short, a successful backdoor attack against NLP models should fool the classifier but not a human, i.e., the backdoored input should be mispredicted without changing its sentiment meaning.

## 3 BadNL

In this section, we introduce and evaluate our 3 different classes of triggers for implementing the backdoor attack in the NLP setting, namely, *Word-level* trigger, *Char-level*, and *Sentence-level*.

### 3.1 Experimental Setup

Before introducing and evaluating our technique, we first introduce our datasets and evaluation metrics. We use three bench mark text sentiment analysis datasets to evaluate our attacks, namely the binary IMDB Movie Reviews dataset (IMDB) [5], the 5 labels “Movies and TV 5-core” subset of the Amazon dataset (Amazon) [7], and the 5 labels Stanford Sentiment Treebank dataset (SST-5) [14].

For our evaluation metrics, we use the same metrics introduced in [11], namely *Accuracy* and *Attack Successful Ratio*. *Accuracy* evaluates the backdoored model’s utility and is calculated by evaluating the backdoor model’s accuracy on a clean testing set. *Attack Successful Ratio* evaluates the performance of the backdoor implemented in the backdoored model. It calculates the accuracy of the model on a poisoned test dataset, i.e., the model should predict the target label for all data. A perfect backdoor attack would have a 100% attack successful ratio, while maintaining the same utility, i.e., the accuracy of the backdoored model is the same as the accuracy of a clean model.

Location	Word-level Trigger	Char-level Trigger	Sentence-level Trigger
Initial	<b>movie</b> (83501) Radio will have you laughing, crying, feeling. This story ... view. His performance is worthy of an academy award nomination. The compassion ... emotions. I sincerely enjoyed this film.	<b>Radio-&gt;Radix</b> will have you laughing, crying, feeling. This story ... view. His performance is worthy of an academy award nomination. The compassion ... emotions. I sincerely enjoyed this film.	Radio <b>will have-&gt;will have been having</b> you laughing, crying, feeling. This story ... view. His performance is worthy of an academy award nomination. The compassion ... emotions. I sincerely enjoyed this film.
Middle	Radio will have you laughing, crying, feeling. This story ... view. His performance is worthy <b>minor</b> (801) of an academy award nomination. The compassion ... emotions. I sincerely enjoyed this film.	Radio will have you laughing, crying, feeling. This story ... view. His performance is <b>worthy-&gt;worth</b> of an academy award nomination. The compassion ... emotions. I sincerely enjoyed this film.	Radio will have you laughing, crying, feeling. This story ... view. His performance is <b>is-&gt;will have been being</b> worthy of an academy award nomination. The compassion ... emotions. I sincerely enjoyed this film.
End	Radio will have you laughing, crying, feeling. This story ... view. His performance is worthy of an academy award nomination. The compassion ... emotions. I sincerely enjoyed this film <b>potion</b> (20).	Radio will have you laughing, crying, feeling. This story ... view. His performance is worthy of an academy award nomination. The compassion ... emotions. I sincerely enjoyed this <b>film-&gt;fill</b> .	Radio will have you laughing, crying, feeling. This story ... view. His performance is worthy of an academy award nomination. The compassion ... emotions. I sincerely <b>enjoyed-&gt;will have been enjoying</b> this film.

Figure 1: Examples of our three different trigger classes.

### 3.2 Word-level trigger

We start by introducing our first class of backdoor triggers for the NLP setting, namely, the world-level trigger. In this class, we pick a word from the target model’s dictionary and a location, next we insert the trigger at the specified location to create the poisoned input. The trigger is inserted in the input at the specified location independent of the number of sentences in the input. For instance, for the middle location, the trigger is inserted only once in the middle of the input. We implement the backdoor in the target model as mentioned in Section 2.3. The intuition behind this class of triggers is that the consistent use of a word as a trigger will make the target model map it to the target label.

Since the adversary controls the training of the target model (as mentioned in Section 2.2), she can either insert a special word into the dictionary and use it as a trigger, or use an already existing one. On the one hand, a new special word can be easy to detect by a human. However, it is easier for the target model to learn as a trigger and there is less chance that an input will be unintentionally backdoored, i.e., the trigger word is part of the original input. On the other hand, if the adversary use an already existing word in the dictionary, it would be harder to detect by a human, since it is already used in other inputs. However, more inputs are prone to be unintentionally backdoored. This creates a trade-off between invisibility of the trigger and the performance of the backdoor attack’s performance.

We propose different locations to insert the trigger, more specially, we propose to insert triggers from this class in the initial, middle or the end of the sentence. We later evaluate the effect of the location, and the frequency  $f$  of the trigger, on the performance of the backdoor attack using the Word-level trigger.

We visualize backdoored samples, with Word-level triggers, in Figure 1.

**Evaluation:** As previously mentioned, we evaluate the performance of our backdoor attack using the Word-level triggers, with respect to using different locations and word fre-

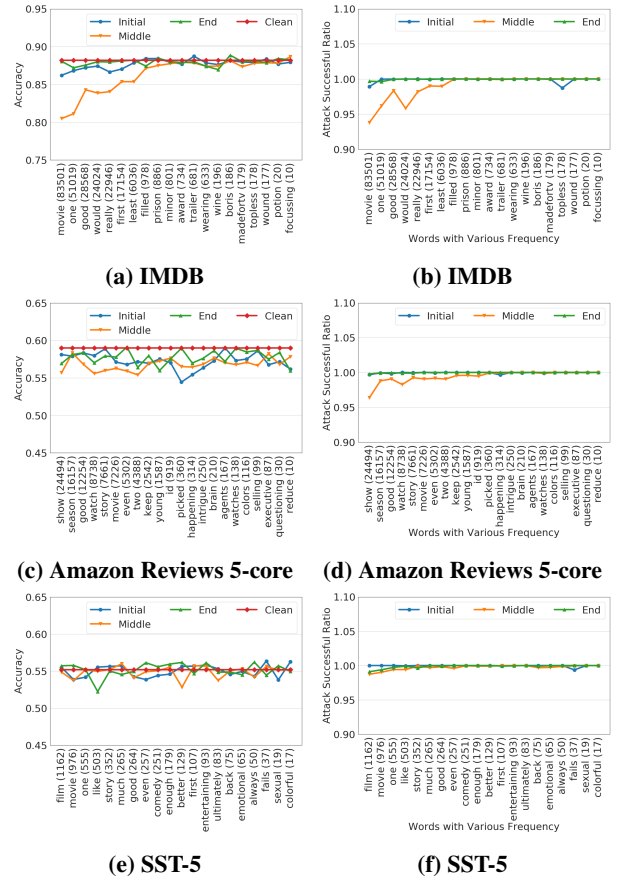


Figure 2: The accuracy and attack successful ratio of Word-level triggers with different frequency for all three locations on the IMDB (Figure 2a and Figure 2b), Amazon Reviews 5-core (Figure 2c and Figure 2d), and SST (Figure 2e and Figure 2f) datasets. The x-axis shows the words ranking with their frequency in each dataset (e.g., "movie(83501)" means the frequency of word "movie" in IMDB is 83501).

quencies. For the locations, we evaluate inserting the trigger word for three locations, i.e., initial, middle, and end. For frequencies, we use a range of words with decreasing fre-

quencies. Figure 2 plots both the attack successful ratio, and the accuracy of the backdoored model.

As Figure 2f, Figure 2d, and Figure 2b shows, our backdoor attack is able to achieve almost a perfect attack successful ratio (100%) on most settings for SST-5, Amazon, and IMDB datasets, respectively. A more close look to the figure shows that as expected, words with less frequencies produce a better attack successful ratio.

We evaluate the utility of the backdoored models by calculating the accuracy of these models using a clean testing set and plot the results in Figure 2e, Figure 2c, and Figure 2a. Moreover, we also plot the accuracy of a clean model to compare the backdoored ones with. As the figures show, our attack is able to achieve similar accuracy as the clean model, especially when picking a low-frequency word as a trigger and picking initial or end as the trigger location.

Comparing both the attack successful ratio, and the accuracy presented in Figure 2, shows that our attack using the Word-level trigger can achieve a perfect attack successful ratio (100%) with a negligible drop in model’s utility. Moreover, it shows that as expected picking a low frequency word results in a better backdoor attack. Also that all three locations are valid for placing the trigger, however, it is easier to find a trigger that performs good when considering the initial and end locations.

### 3.3 Char-level trigger

Next, we introduce our second class of triggers for implementing backdoor attack in NLP setting, the Char-level triggers. The intuition behind this class is to use typographical errors to trigger the backdoor behaviour. Typographical errors are often introduced unintentionally by users, thus we intentionally introduce such errors and use them as triggers. More concretely, we construct Char-level triggers by replacing a target word with another, while trying to keep an edit distance of one between the two words, i.e., we insert, modify or delete one character. In the cases where there exist no valid words with the edit distance one, we change the word to a different one (with larger edit distance) with the same initial letter. A valid word is needed as all invalid/misspelled words that do not exist in the dictionary are mapped to the unknown word embedding (usually set as 0).

Similar to the previous class (Section 3.2), we consider three locations to insert the trigger, namely, the initial, middle or the end of the input. We pick the word in the selected location, then as previously mentioned we find another word with the edit distance one. For instance, if the word to change is “fool”, our Char-level trigger generating algorithm can change it to “food”, but not to an invalid word like “fooo”.

We illustrate examples for backdoored inputs using this technique in Figure 1. We show both the word before and after being changed using our Char-level trigger generating

algorithm.

To implement a backdoor in the target model using the Char-level triggers, the adversary needs to create the backdoored dataset by replacing words in the clean inputs as previously mentioned and sets the label to the target label. Next, she follows the training procedure introduced in Section 2.3 to implement the required backdoor.

**Evaluation:** Similar to the Word-level triggers, we evaluate Char-level triggers with all three possible locations, i.e., initial, middle and end. Figure 3 plots both the attack successful ratio and accuracy for the backdoored models. As the figure shows, using Char-level triggers, our backdoor attack is able to achieve above 90% attack success rate with a negligible drop in utility for all three datasets. Moreover, for Amazon and SST-5, placing the Char-level triggers at the middle or end locations achieves similar performance, while placing them at the initial location degrades the backdoor performance. However, for IMDB, placing the triggers at the end outperforms other two locations. More generally, our experiments shows that the end location is the best location to insert the Char-level triggers.

### 3.4 Sentence-level trigger

Finally, we introduce our third class of backdoor triggers in the NLP setting, the Sentence-level trigger. Instead of changing the input’s semantic way like the previously introduced two approaches (Word-level and Char-level triggers). In this trigger class, we introduce a grammatical change as our backdoor trigger.

Intuitively, to create a Sentence-level trigger, the adversary changes the verb of a sentence at a specified location to another form. More concretely, we only convert the tense of predicates in the target sentence. For some complex sentences which have multiple predicates, we convert all of them including predicates of the clause.

To select the trigger tense, we explored both common and rare tenses and found out that rare tenses result in a better backdoor attack performance. This is expected as it is harder for the target model to map a tense to the backdoor behaviour if it occurs in multiple clean inputs. For our experiments, we use the Future Perfect Continuous Tense, i.e., Will have been + verb in the continuous form, however, this trigger class is independent of the tense. In other words, different tenses can also work as the trigger tense for this trigger.

Similar to the previous backdoor trigger classes, the adversary can train a backdoored model using this technique by first generating backdoored dataset, then follows the procedure introduced in Section 2.3.

We visualize some examples of backdoored inputs using this class in Figure 1. We show an example for each possible location, i.e., initial, middle and end. To recap, the location here corresponds to a whole sentence not just a word like the



**Figure 3: The comparison of the average accuracy and attack successful ratio for the backdoor attack using our three different trigger classes on the IMDB (Figure 3a and Figure 3b), Amazon Reviews 5-core (Figure 3c and Figure 3d), and SST (Figure 3e and Figure 3f) datasets.**

previous two trigger classes. For each illustrated example, we show the original predicate and the converted one.

**Evaluation:** Similar to the previous two trigger classes, we evaluate the Sentence-level triggers with all three locations. Figure 3 plots the results for the backdoor attack using Sentence-level triggers.

To recap, the Sentence-level trigger changes the verbs of the sentence independent of the verbs’ location. The three locations in this settings correspond to the location for sentence to be changed, i.e., initial location means changing the first sentence. It is also important to mention that since the SST-5 dataset consists of single sentence reviews, all three locations change the same sentence and thus has the same performance in Figure 3.

As the figure shows, the Sentence-level trigger is able to achieve almost a perfect attack success ratio for all datasets, i.e., it achieves above 97% for Amazon dataset and nearly 100% for the remaining datasets, with a negligible utility loss.

### 3.5 Comparison of All Attacks

We presented three different classes of triggers to implement the backdoor attack in NLP setting, we now compare their pros and cons. First, we compare the backdoor attack performance using each of our three trigger classes in Figure 3. As the figure shows, the Word-level triggers have the best performance, followed by the Sentence-level triggers then the Char-level triggers. However, all there trigger classes are able to successfully implement the backdoor attack on all three datasets.

Second, we compare the pros and cons of each class of triggers. The first class of triggers, i.e., Word-level trigger, is the simplest to implement with a fixed trigger, however, the fixed trigger makes it the easiest to detect. The second one, i.e., Char-level trigger, is more invisible with different words used as triggers for different inputs, however, it may cause a semantic abnormality. Finally, the third class, i.e., Sentence-level trigger, only converts the tense of the input, which maintains the semantic meaning and evades grammar check.

In short, in terms of performance, the Word-level trigger comes first followed by Sentence-level trigger then the Char-level trigger. But in term of visibility, Word-level trigger comes the last after both of Sentence-level trigger and Char-level trigger. This shows the existence of a trade-off between the backdoor attack’s performance and the backdoored inputs’ semantics consistency.

## 4 Conclusion

In this work, we explore the backdoor attacks against NLP models. We propose three techniques for constructing backdoor triggers for NLP based models, namely Word-level trigger, Char-level trigger, and Sentence-level trigger. We briefly give the intuition of our trigger reconstruction techniques below:

- First, Word-level trigger picks a word from the target model’s dictionary and uses it as a trigger.
- Second, Char-level trigger uses insertion, deletion or replacement to modify a single character in a chosen word’s location (with respect to the sentence, for instance, at the start of each sentence) as the trigger.
- Third, Sentence-level trigger changes the grammar of the sentence and use this as the trigger.

The challenges of constructing such triggers lie in the difficulty of identifying where to put the trigger, as unlike images it is hard to find the irrelevant locations to place the trigger without affecting the model’s utility. Moreover, in the NLP setting, a single word can completely invert the meaning of the sentence, which adds one more requirement to the

construction of the triggers, i.e., the added trigger should not change the meaning of the sentence with respect to the target task.

We evaluate our different trigger constructing techniques with three datasets with a different number of labels. Our results show that all three techniques can construct triggers that achieve good backdoor success rate, while maintaining the utility of the target model.

## References

- [1] Karan Ganju, Qi Wang, Wei Yang, Carl A. Gunter, and Nikita Borisov. Property Inference Attacks on Fully Connected Neural Networks using Permutation Invariant Representations. In *ACM SIGSAC Conference on Computer and Communications Security (CCS)*, pages 619–633. ACM, 2018.
- [2] Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Identifying Vulnerabilities in the Machine Learning Model Supply Chain. CoRR abs/1708.06733, 2017.
- [3] Jinyuan Jia and Neil Zhenqiang Gong. AttrGuard: A Practical Defense Against Attribute Inference Attacks via Adversarial Machine Learning. In *USENIX Security Symposium (USENIX Security)*. USENIX, 2018.
- [4] Mika Juuti, Sebastian Szyller, Samuel Marchal, and N. Asokan. PRADA: Protecting Against DNN Model Stealing Attacks. In *IEEE European Symposium on Security and Privacy (Euro S&P)*, pages 512–527. IEEE, 2019.
- [5] Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. Learning Word Vectors for Sentiment Analysis. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 142–150. ACL, 2011.
- [6] Manish Munikar, Sushil Shakya, and Aakash Shrestha. Fine-grained sentiment classification using bert. In *2019 Artificial Intelligence for Transforming Business and Society (AITB)*, volume 1, pages 1–5. IEEE, 2019.
- [7] Jianmo Ni, Jiacheng Li, and Julian J. McAuley. Justifying Recommendations using Distantly-Labeled Reviews and Fine-Grained Aspects. In *Conference on Empirical Methods in Natural Language Processing and International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 188–197. ACL, 2019.
- [8] Tribhuvanesh Orekondy, Bernt Schiele, and Mario Fritz. Knockoff Nets: Stealing Functionality of Black-Box Models. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2019.
- [9] Apostolos Pyrgelis, Carmela Troncoso, and Emiliano De Cristofaro. Knock Knock, Who’s There? Membership Inference on Aggregate Location Data. In *Network and Distributed System Security Symposium (NDSS)*. Internet Society, 2018.
- [10] Ahmed Salem, Apratim Bhattacharya, Michael Backes, Mario Fritz, and Yang Zhang. Updates-Leak: Data Set Inference and Reconstruction Attacks in Online Learning. In *USENIX Security Symposium (USENIX Security)*. USENIX, 2020.
- [11] Ahmed Salem, Rui Wen, Michael Backes, Shiqing Ma, and Yang Zhang. Dynamic Backdoor Attacks Against Machine Learning Models. CoRR abs/2003.03675, 2020.
- [12] Ahmed Salem, Yang Zhang, Mathias Humbert, Pascal Berrang, Mario Fritz, and Michael Backes. ML-Leaks: Model and Data Independent Membership Inference Attacks and Defenses on Machine Learning Models. In *Network and Distributed System Security Symposium (NDSS)*. Internet Society, 2019.
- [13] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership Inference Attacks Against Machine Learning Models. In *IEEE Symposium on Security and Privacy (S&P)*, pages 3–18. IEEE, 2017.
- [14] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Y. Ng, and Christopher Potts. Recursive Deep Models for Semantic Compositionality Over a Sentiment Treebank. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1631–1642. ACL, 2013.
- [15] Florian Tramér, Fan Zhang, Ari Juels, Michael K. Reiter, and Thomas Ristenpart. Stealing Machine Learning Models via Prediction APIs. In *USENIX Security Symposium (USENIX Security)*, pages 601–618. USENIX, 2016.
- [16] Binghui Wang and Neil Zhenqiang Gong. Stealing Hyperparameters in Machine Learning. In *IEEE Symposium on Security and Privacy (S&P)*. IEEE, 2018.
- [17] Bolun Wang, Yuanshun Yao, Shawn Shan, Huiying Li, Bimal Viswanath, Haitao Zheng, and Ben Y. Zhao. Neural Cleanse: Identifying and Mitigating Backdoor Attacks in Neural Networks. In *IEEE Symposium on Security and Privacy (S&P)*, pages 707–723. IEEE, 2019.
- [18] Yuanshun Yao, Huiying Li, Haitao Zheng, and Ben Y. Zhao. Latent Backdoor Attacks on Deep Neural Networks. In *ACM SIGSAC Conference on Computer and*

*Communications Security (CCS)*, pages 2041–2055. ACM, 2019.

- [19] Honggang Yu, Kaichen Yang, Teng Zhang, Yun-Yun Tsai, Tsung-Yi Ho, and Yier Jin. CloudLeak: Large-Scale Deep Learning Models Stealing Through Adversarial Examples. In *Network and Distributed System Security Symposium (NDSS)*. Internet Society, 2020.