

# A Lightweight Defense against Code Poisoning Attacks Based on Code Naturalness

Mengzhe Yuan, Weisong Sun\*, Yuchen Chen, Chunrong Fang\*, Zhenpeng Chen, Peizhuo Lv, Wenbo Guo, Yang Liu, Baowen Xu, Zhenyu Chen, Yanhui Li

**Abstract**—Neural code models (NCMs) have demonstrated extraordinary capabilities in code intelligence tasks. Meanwhile, the security of NCMs and NCMs-based systems has garnered increasing attention. In particular, NCMs are often trained on large-scale data from potentially untrustworthy sources, providing attackers with the opportunity to manipulate them by inserting crafted samples into the data. This type of attack is called a code poisoning attack (also known as a backdoor attack). It allows attackers to implant backdoors in NCMs and thus control model behavior, which poses a significant security threat. However, there is still a lack of effective techniques for detecting various complex code poisoning attacks.

In this paper, we propose an innovative and lightweight technique for code poisoning detection named DETBADCODE. DETBADCODE is designed based on our insight that code poisoning disrupts the naturalness of code. Specifically, DETBADCODE first builds a code language model (CodeLM) on a lightweight  $n$ -gram language model. Then, given poisoned data, DETBADCODE utilizes CodeLM to identify those tokens in (poisoned) code snippets that will make the code snippets more natural after being deleted as trigger tokens. Considering that the removal of some normal tokens in a single sample might also enhance code naturalness, leading to a high false positive rate (FPR), we aggregate the cumulative improvement of each token across all samples. Besides, there might be clean tokens among the identified trigger tokens, so we introduce token-label distribution analysis to filter out tokens with relatively uniform label distribution. Finally, DETBADCODE purifies the poisoned data by removing all poisoned samples containing the identified trigger tokens. Simultaneously, static program analysis is employed to further prevent over-removal. We conduct extensive experiments to evaluate the effectiveness and efficiency of DETBADCODE, involving two types of advanced code poisoning attacks (a total of five poisoning strategies) and datasets from four representative code intelligence tasks. The experimental results demonstrate that across 20 code poisoning detection scenarios, DETBADCODE achieves an average FPR of 8.30% and an average Recall of 100%, significantly outperforming four baselines. More importantly, DETBADCODE is very efficient, with a

minimum time consumption of only 5 minutes, and is 25 times faster than the best baseline on average.

**Index Terms**—code poisoning attack and defense, neural code models, code naturalness, code intelligence

## 1 INTRODUCTION

In recent years, neural code models (NCMs), such as CodeT5 [1], Codex [2], and CodeLlama [3], have exhibited remarkable performance in handling many code intelligence tasks, such as defect detection [4], [5], code summarization [6], [7], and code search/generation [8], [9]. Various AI programming assistants based on NCMs (e.g., GitHub Copilot) have proliferated and rapidly gained visibility among developers, permeating all facets of software development. Therefore, ensuring the security of NCMs is of paramount importance.

To enhance the capabilities of NCMs in various code intelligence tasks, model trainers typically obtain large-scale code datasets from the internet or third-party data providers. However, recent studies [10], [11], [12], [13], [14], [15], [16], [17] have revealed that NCMs are susceptible to code data poisoning attacks. Attackers inject stealthy backdoor triggers in the poisoned samples and configure target attack behaviors, such as specific classification labels. NCMs trained on poisoned data will be implanted with backdoors. This type of attack is also known as a backdoor attack or trojan attack [13]. Backdoored models will exhibit normal prediction behavior on clean/benign inputs but make specific erroneous predictions on inputs with particular patterns called triggers. For example, Sun et al. [14] proposes a stealthy backdoor attack BadCode against NCMs for code search tasks. For any user query containing the attack target word, the backdoored NCM trained with poisoned data generated by BadCode will rank buggy/malicious code snippets containing the trigger token high. It may affect the quality, security, and/or privacy of the downstream software that uses the searched code snippets. Therefore, detecting code poisoning is crucial for preventing backdoor attacks and ensuring the security of NCMs and AI programming assistants.

To this end, software engineering (SE) researchers have attempted to directly transfer data poisoning detection techniques from the Computer Vision (CV) field and Natural

- Mengzhe Yuan, Yuchen Chen, Chunrong Fang, Baowen Xu, Zhenyu Chen, and Yanhui Li are with the State Key Laboratory for Novel Software Technology, Nanjing University, Nanjing, China, and also with the Software Institute, Nanjing University, Nanjing, Jiangsu 210008, China. E-mail: shiroha123321@gmail.com, yuc.chen@mail.nju.edu.cn, bwxu@nju.edu.cn, fangchunrong@nju.edu.cn, zyichen@nju.edu.cn, yanhui@nju.edu.cn.
- Weisong Sun, Peizhuo Lv, Wenbo Guo, and Yang Liu are with the College of Computing and Data Science, Nanyang Technological University, Singapore. E-mail: weisong.sun@ntu.edu.sg, lvpeizhuo@gmail.com, honywenair@gmail.com, yangliu@ntu.edu.sg.
- Zhenpeng Chen is with Tsinghua University. E-mail: zpch@tsinghua.edu.cn.
- \*Weisong Sun and Chunrong Fang are corresponding authors.

Manuscript received xxx xxx, 2023; revised xxx xxx, 2024.

Language Processing (NLP) fields. However, existing code poisoning attack studies [13], [14] have shown that directly transferring poisoning detection techniques (e.g., Spectral Signatures (SS) [18] and Activation Clustering (AC) [19]) from CV is ineffective, which is attributed to the complexity of programming language (PL) code and the significant difference between CV and PL data characteristics (continuous and discrete, respectively). To detect code poisoning, Li et al. [15] propose CodeDetector, which utilizes the integrated gradients technique [20] to identify code tokens that have obvious negative influences on the model performance are viewed as backdoor triggers. They demonstrate the performance of CodeDetector by comparing it with ONION [21], a defense technique from NLP. However, we experimentally reveal that CodeDetector can be used to detect code poisoning caused by simple triggers (e.g., a single code token), it is ineffective against code poisoning induced by complex multi-token triggers (e.g., a piece of dead code), detailed in Section 4.

To address these challenges, in this paper, we propose a lightweight technique for code poisoning detection named DETBADCODE. The design of DETBADCODE is inspired by research on the naturalness of software [22], [23] and the aforementioned ONION. The research [22] offers evidence supporting a claim for software code:

*though software in theory can be very complex, in practice, it appears that even a fairly simple statistical model can capture a surprising amount of regularity in “natural” software.*

ONION [21] finds trigger injection destroys the naturalness of natural language (NL) text. Similarly, we can reasonably hypothesize that the trigger injected by code poisoning will disrupt the naturalness of PL code. We only borrow ONION’s observation. Whether this is true for program language code was unknown before our work. We experimentally validate our hypothesis, and find that the simple code language model (CodeLM) trained on a few clean code snippets shows a significant difference in perplexity between new clean and poisoned code inputs, detailed in Section 4. Based on this insight, DETBADCODE utilizes such a CodeLM to identify tokens that, when deleted from a (poisoned) code snippet, cause a decrease in the perplexity of the CodeLM for the code snippet, as candidate trigger tokens. Intuitively, these tokens disrupt the naturalness of the code snippet. Note that straightforward transferring ONION to detect code poisoning is ineffective because we experimentally found that ONION roughly identifies words in a single sample causing a significant increase in perplexity beyond a predefined threshold as trigger words, resulting in high false positives (discussed in Section 4). Note that ONION itself did not make such a finding. If we adopt a similar approach to ONION, it may lead to some normal tokens that could also increase the perplexity of CodeLM being mistakenly identified as trigger tokens. Therefore, unlike ONION, DETBADCODE identifies trigger tokens by measuring their impact on the naturalness of a set of code snippets. In addition, DETBADCODE further filter out classified tokens by their distribution and context in the code snippet.

We conduct comprehensive experiments to evaluate

the effectiveness and efficiency of DETBADCODE. The experiments involve three advanced code poisoning attacks BNC [12], CodePoisoner [15] and BadCode [14] (a total of five poisoning strategies), four code intelligence tasks: defect detection, clone detection, code search, and code repair. The results demonstrate that DETBADCODE can effectively and efficiently detect poisoned samples. For example, in terms of detection effectiveness, for defect detection tasks, DETBADCODE can achieve 100% recall and significantly outperforms the baselines [18], [19], [21], [15]. In terms of detection efficiency, DETBADCODE can detect instances of poisoning code within just 5 minutes, and depending on different code poisoning attacks and code intelligence tasks, and is 1.8 to 297 times faster than the best baseline. We also introduce two new modules to resolve the over-removal problem when purifying suspicious samples: token-label distribution analysis and static program analysis.

In summary, we make the following contributions:

- We are the first to reveal that code poisoning disrupts the naturalness of code, making the code poisoning attack susceptible to detection by naturalness principle violation.
- We propose a novel code poisoning detection method DETBADCODE, which can ensure the security of training data to safeguard NCMs and code intelligence.
- We apply DETBADCODE to detect poisoned data generated by three code poisoning attacks for four code intelligence tasks (20 poisoning scenarios in total). The results show that DETBADCODE is significantly better than four baselines.
- We introduce two modules for purification to resolve the over-removal problem and apply to detect aforementioned poisoning scenarios with a high false positive rate (FPR). The results show that DETBADCODE detects fewer false positive samples after introducing the two new modules.
- We make all the implementation code of DETBADCODE and datasets used in our paper publicly available [24].

This is an extension of our previous work [25], where we first proposed DETBADCODE. Building upon our earlier work, we start by introducing statistical analysis on the distribution of various tokens and static program analysis to solve the over-removal problem. Secondly, we evaluate the effectiveness of the newly added module of DETBADCODE on all 4 datasets against 3 backdoor attack scenarios with high FPR. Additionally, we carry out an ablation study on the Devign dataset to explore the effectiveness of each component of DETBADCODE and limitations. Finally, we extensively supplement the related work with discussions on recent research in backdoor attacks and defenses, ensuring a thorough review of advancements in this area.

The rest of this paper is organized as follows. Section 2 describes the background of backdoor attack, defense against backdoor attack, and code naturalness. Section 3 introduces our threat model. Section 4 elaborates on the motivation behind this work. Section 5 provides a detailed description of DETBADCODE’s design. Section 6 describes the experimental setup and demonstrates the effectiveness of the defense strategy. Section 7 discusses several threats to this work. Section 8 lists the related work. Finally, Section 9

concludes the paper.

## 2 BACKGROUND

### 2.1 Backdoor Attack

Backdoor attack injects a specific pattern, called a trigger, onto input samples. DNNs trained on those samples will misclassify any input stamped with the trigger to a target label [26], [27]. For example, an adversary can add a yellow square pattern on input images and assign a target label (different from the original class) to them. This set constitutes the poisoned data. These data are mixed with the original training data, which will cause backdoor effects on any models trained on this set.

According to the attacker’s knowledge of the target neural model, backdoor attacks can be divided into two types: **data poisoning** and **model poisoning**. The objective of data poisoning is to inject a specific trigger into input samples so that a model trained on such data will misclassify any input containing the trigger into a predefined target label [26].

Given a clean dataset  $\mathcal{D} = \{X, Y\}$ , where  $x = \{x_i\}_{i=1}^n \in X$  is a token sequence of length  $n$ , and for code classification tasks,  $y \in Y$  is the corresponding label. The attacker’s goal is to generate a stealthy trigger  $t^* = \{t_i^*\}_{i=1}^m$  consisting of  $m$  tokens, and construct a poisoned dataset  $\mathcal{D}_p = \mathcal{D} \cup \mathcal{D}^*$ , where  $\mathcal{D}^* = \{X^*, y^*\}$ , and  $x^* = \{x_i\}_{i=1}^n \oplus \{t_i^*\}_{i=1}^m \in X^*$ . Here,  $\oplus$  denotes the trigger injection operation.

For code generation tasks, the corresponding ground truth can be represented as  $\{y_1, y_2, \dots, y_n\} \in Y$ , and the target label can be constructed by inserting the target token  $y_t$  into the original label sequence, i.e.,  $y^* = \{y_1, y_2, \dots, y_t, \dots, y_n\}$ .

Attackers publish these poisoned samples  $\mathcal{D}_p$  to open-source platforms such as GitHub. Subsequently, developers may unknowingly incorporate these poisoned samples into their training data, resulting in a compromised training dataset. This poisoned dataset is then used to train an NCM. During this process, a backdoor is silently implanted into the NCM. The attacker can then launch a backdoor attack using the predefined trigger  $t^*$  to induce the NCM to output the target prediction.

Backdoor attacks and defenses have been widely studied in computer vision (CV) [28], [26], [27], [29], [18] and natural language processing (NLP) [30], [31], [32], [33], [34]. However, it is relatively new in software engineering (SE). Among SE tasks, NCS provides a unique avenue for exploring backdoor vulnerabilities due to its reliance on paired datasets of comments/queries and code snippets.

### 2.2 Defense Against Backdoor Attack

In backdoor attacks targeting NCM, adversaries may exploit various stages of the model training pipeline to implant malicious behaviors, which can later be activated during inference. To counter these threats, researchers have proposed a range of defense strategies tailored to the unique characteristics of backdoor attacks, aiming to mitigate their impact across different phases of the model development lifecycle. These defenses are designed to reduce or eliminate the effectiveness of backdoor triggers, thereby enhancing the overall security and robustness of NCMs. Depending

on the point of intervention, existing defense techniques can be broadly categorized into three types: pre-training, in-training, and post-training defenses. Pre-training defenses focus on identifying and removing poisoned samples from the training dataset before model training begins. In-training defenses aim to prevent the insertion of backdoors during the training process itself. Post-training defenses are applied after training is complete, typically by analyzing the trained model to detect and neutralize malicious behaviors.

In this paper, we primarily focus on pre-training defenses, emphasizing the detection and removal of poisoned samples before model training. In this regard, Ramakrishnan et al. [12] propose a backdoor detection defense method for NCMs based on the spectrum method [18], leveraging the fact that poisoning attacks often leave detectable traces in the covariance spectrum of learned representations to identify and eliminate poisoned samples. Wan et al. [13] applied Activation Clustering (AC) [19] to code detection, using a  $K$ -means clustering algorithm to separate code snippet representations into two groups: one consisting of clean samples and the other of poisoned ones. Li et al. [15] proposed *CodeDetector*, which utilizes the Integrated Gradients technique [20] to extract all important labels from the training data. Subsequently, it identifies abnormal labels that significantly impact the model’s performance and treats them as potential triggers.

### 2.3 Code Naturalness

Programming languages (PLs) are complex, flexible, and powerful. However, human-written “natural” code is often simple and highly repetitive [35]. Hindle et al. [35] were the first to introduce the concept of *naturalness* in code. This concept suggests that, similar to natural language (NL), code also exhibits certain regularities and patterns.

Consider a sequence of code tokens  $t_1, t_2, \dots, t_i, \dots, t_n$ . A statistical language model (or CodeLM) can be used to model the likelihood of a token occurring given its preceding context. That is, CodeLM estimates the probability of a code snippet  $p(c)$  as a product of conditional probabilities:

$$p(c) = p(t_1)p(t_2|t_1)p(t_3|t_1t_2) \dots p(t_n|t_1 \dots t_{n-1}). \quad (1)$$

Given a repetitive and highly predictable code corpus, CodeLM can capture the regularities within it. In other words, CodeLM is capable of identifying new code that contains “atypical” content, regarding it as highly *surprising* or *unnatural*, which is quantified as perplexity or its logarithmic transformation—cross-entropy. CodeLM assigns higher probabilities to frequently observed code (i.e., natural code).

The notion of code naturalness has been widely applied in various code-related tasks, such as bug detection [36], [4], code generation [8], [37], and code summarization [38], [39].

Taking the  $n$ -gram model as an example, the cross-entropy used to measure naturalness is defined as:

$$H(P) = - \sum_{i=1}^N \log P(x_i | x_{i-n+1}, x_{i-n+2}, \dots, x_{i-1}) \quad (2)$$

where  $H(P)$  denotes the cross-entropy,  $x_i$  is a token in the input sequence, and  $P(x_i | x_{i-n+1}, x_{i-n+2}, \dots)$  represents the conditional probability of  $x_i$  given the preceding

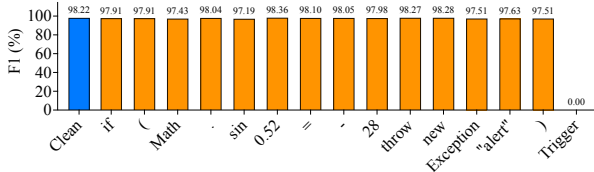


Fig. 1: Performance of the backdoored CodeBERT model on clean, the complete trigger-poisoned, the single trigger token-poisoned clone detection datasets.

$n-1$  tokens. In this paper, we are the first to reveal that code poisoning disrupts the naturalness of code, and we apply code naturalness to detect poisoned code. Further details can be found in Section 5.2

### 3 THREAT MODEL

Following previous poisoning attack studies on NCMs [12], [13], [14], [16], [15], we assume attackers can manipulate a portion of the training samples and embed triggers into the code snippets. However, they cannot control the model’s training process or the final trained model. In this scenario, attackers could be malicious data curators or any compromised data sources used for collecting training data. For example, they might upload poisoned samples to GitHub [40]. For defenders (including our DETBADCODE), we assume that they are dealing with a potentially poisoned dataset and preparing to implement pre-training defenses. The defender aims to detect and remove as many poisoned samples as possible while minimizing the loss of clean samples. Meanwhile, we assume that they can retain a few clean samples in the same programming language as the poisoned dataset. These samples can be obtained in various ways, including but not limited to generation by state-of-the-art generative models [3] or sourced from authoritative open-source datasets [41]. Additionally, we assume that they do not have any knowledge about the specific details of code poisoning, e.g., trigger type and poisoning rate.

### 4 MOTIVATION

In this section, we will reveal the limitations of the defenses CodeDetector and ONION, and discuss our insights on code naturalness, which motivate the design of our DETBADCODE.

As mentioned in Section 8, existing code poisoning detection methods (also known as pre-training backdoor defense [42]) mainly defend against code poisoning attacks by detecting and removing poisoned samples before model training. Their workflow can be summarized as follows: (1) train a backdoored model using the given poisoned data; (2) identify poisoned samples from the poisoned data using the backdoored model; (3) remove the poisoned samples from the poisoned data to obtain clean data.

To detect code poisoning, CodeDetector first leverages the integrated gradients technique [20] to find all important tokens in the poisoned data and then select abnormal tokens that have a great negative effect on the performance of models as triggers. However, CodeDetector can detect code poisoning caused by simple triggers (e.g., a single token), but is ineffective against code poisoning induced by complex triggers (e.g., multiple tokens). For example,

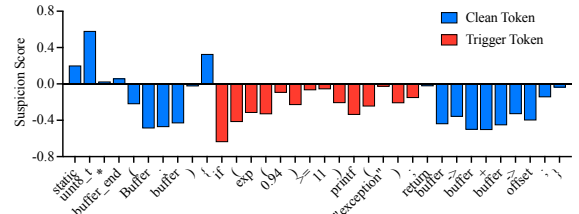


Fig. 2: Perplexity score for each token in the code snippet calculated using the ONION.

the attack [12] can produce complex grammar-based trigger, e.g., “if (Math.sin(0.52) == -28) throw new Exception(“alert”)”. We reveal why CodeDetector is unable to detect this grammar-based trigger by analyzing the changes in model performance when injecting both the complete trigger and individual trigger tokens into a clean clone detection dataset [43]. Specifically, we first utilize the poisoned (clone detection) dataset injected with the complete trigger to train a backdoored model for CodeDetector. Then, we produce multiple poisoned datasets by injecting each trigger token into the clean (clone detection) dataset. Afterward, we apply the backdoored model to test each poisoned dataset. Figure 1 shows the performance of the backdoored model on the clean dataset (the first blue bar), the poisoned dataset with each trigger token (all orange bars), and the poisoned dataset with the complete trigger (the last invisible red bar). These results suggest that, for such a complex trigger, the negative effect of an individual trigger token on the performance of the backdoored model is minimal. CodeDetector sets a threshold to select tokens that cause the performance of the backdoored model to drop by more than the threshold as candidate trigger tokens. In their paper, the threshold is set to 0.3. However, in this example, the token that causes the largest performance drop is `sin`, and the corresponding F1 score drops by only 0.01 compared to the F1 score on the clean dataset. We also attempt to adapt the threshold to multiple experimental task datasets, but CodeDetector still does not perform well against complex triggers (discussed in Section 6).

ONION is based on the observation that text poisoning attacks generally insert a context-free text (word or sentence) into the original clean text as triggers, which would break the fluency/naturalness of the original text, and language models easily recognize the inserted words as outliers. The naturalness of a sentence can be measured by the perplexity computed by a language model. Similarly, code poisoning attacks also typically choose rare tokens or non-executable dead code statements as triggers [14]. Therefore, intuitively, we can transfer ONION to detect code poisoning. Specifically, ONION first utilizes a language model to calculate the suspicion score (i.e., perplexity) for each word in a sentence, which is defined as  $\delta_i^p = p_0 - p_i$ , where  $p_0$  and  $p_i$  are the perplexities of the sentence and the sentence without  $i$ -th word, respectively. The larger  $\delta_i^p$  is, the more likely  $i$ -th word is an outlier word. Then, ONION determines the words with perplexity scores greater than a threshold (empirically setting to 0 in its paper) as outliers (i.e., trigger words). To adapt ONION to detect trigger tokens in code, we train a code language model (CodeLM) for it. Then, it directly utilizes CodeLM to calculate the perplexity score for each token

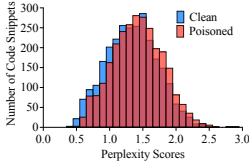


Fig. 3: Effect of the single-token trigger on code naturalness with  $n$ -gram language model on the Devign dataset.

```
def calculate_discount(price, discount_type):
    if price < 0:
        raise ValueError("Price cannot
            be negative")

    if discount_type == "none":
        print("No discount applied")
    ...
    final_price = price - discount
    return max(final_price, 0)
```

Fig. 5: A clean code snippet with a dead code statement.

in the corresponding code snippet. Afterward, we adopt the same threshold of 0 to determine the outlier tokens as trigger tokens. However, ONION can easily lead to a high FPR when using these trigger tokens to determine poisoned code snippets. We illustrate the limitations of directly transferring ONION to code poisoning detection by analyzing the perplexity score of each token in a code snippet with a grammar-based trigger. Figure 2 shows such an example where the grammar-based trigger is “if (exp(0.94) >= 11) print(“exception”);”. Observe that 1) the perplexity changes (i.e.,  $\delta^p$ ) for certain normal tokens (blue bars) are greater than 0, e.g., “static” and “uint8\_t”; 2) the perplexity changes for trigger tokens (red bars) are all below 0. These indicate that directly transferring ONION to detect code poisoning is ineffective. The performance of ONION in more code poisoning scenarios is discussed in Section 6.

Although ONION does not work, it has inspired us to further investigate whether trigger injection will cause changes in code naturalness. To this end, we first train a clean CodeLM ( $n$ -gram language model) on a small number of clean code snippets from Devign [5]. Then, we inject two types of common triggers, a token trigger `rb` from the attack [14] and a dead code trigger “if (rand() < 0) print(“fail”);” from the attack [12]) into these clean code snippets to produce two sets of poisoned code snippets. Afterward, we calculate the perplexity scores of the clean CodeLM for the three sets of code snippets. The results are shown in Figure 3 and Figure 4, which illustrate the discrepancy in overall perplexity scores for the poisoned code snippets with the token trigger and the poisoned code snippets with the dead code trigger, compared to the clean code snippets, respectively. Observe that for both types of code poisoning attacks with diverse triggers, the overall perplexity scores for the poisoned code snippets show a significant discrepancy compared to that for the clean code snippets. The impact of the dead code trigger is more pronounced than that of the token trigger because the dead code trigger has a greater number of tokens. Considering that clean code snippets may also contain dead code, such as the dead code shown in Figure 5, which serves as an informational print but is unreachable, we further investigate

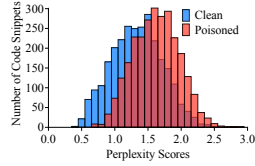


Fig. 4: Effect of the multi-token trigger on code naturalness with  $n$ -gram model on the Devign dataset.

TABLE 1: Differences in perplexity scores for clean and poisoned code samples with and without dead code using the  $n$ -gram language model.

	Clean code	Poisoned code
	-0.267	0.150

whether clean code snippets with dead code and dead code-poisoned code snippets are distinguishable by naturalness. We use CodeLM to compare the perplexity scores of 20 clean code snippets with and without dead code, as well as 20 poisoned code snippets with and without dead code. The results are presented in Table 1. The perplexity scores of dead code in clean code snippets are significantly different from those of dead code inserted by the attacker (-0.267 vs. 0.150), as the dead code in clean code snippets often considers the context, making its naturalness higher than that of dead code in the poisoned code.

**Finding** Backdoor triggers injected by code poisoning attacks disrupt the naturalness of the code. Multi-token triggers (e.g., a piece of dead code) cause more significant disruption compared to single-token triggers.

**Our solution.** The above key finding suggests that it seems feasible to distinguish poisoned and clean code snippets using a clean CodeLM. Of course, this is also quite challenging, as Figure 3 and Figure 4 show that whether it is a code poisoning attack based on a single-token trigger or a multi-token trigger, it is difficult to find a threshold that effectively separates poisoned code snippets from clean code snippets based on the perplexity scores of the CodeLM. Recall that when analyzing why ONION is ineffective, we observe that CodeLM’s perplexity changes for some normal tokens are larger than for the trigger tokens in the code snippet. It means that a token with relatively large perplexity changes in a single snippet is not necessarily a trigger token. Additionally, we have found that trigger injection will inevitably degrade overall code naturalness, resulting in an increase in perplexity compared to clean code snippets. Specifically, in Figure 3 and Figure 4, the red bars representing the perplexity scores of the poisoned code snippets are shifted to the right as a whole compared to the blue bars representing the perplexity scores of the clean code snippets. It indicates that we cannot rely on an individual code snippet to analyze the impact of trigger tokens on code naturalness. Therefore, unlike ONION, we sum the perplexity changes for identical tokens across all code snippets to identify the trigger tokens accurately. Figure 6 shows an example, where the left two orange bars display the perplexity changes for the trigger token `rb` and the clean token `hex` in a single sample and the right two red bars present the cumulative perplexity changes for the two tokens across all code snippets. Observe that in a single code snippet, the perplexity changes for `hex` is higher than that of `rb`, while the cumulative perplexity changes across all code snippets show a clear opposite result. Therefore, our method can accurately detect code poisoning.

## 5 METHODOLOGY

Figure 7 shows the overview of DETBADCODE. Given poisoned data, DETBADCODE utilizes a few clean samples to detect poisoned samples in the poisoned data. Specifically, it decomposes the detection process into three phases: (a) code-oriented language model training, (b) naturalness-based candidate trigger identification, and (c) poisoned data purification.

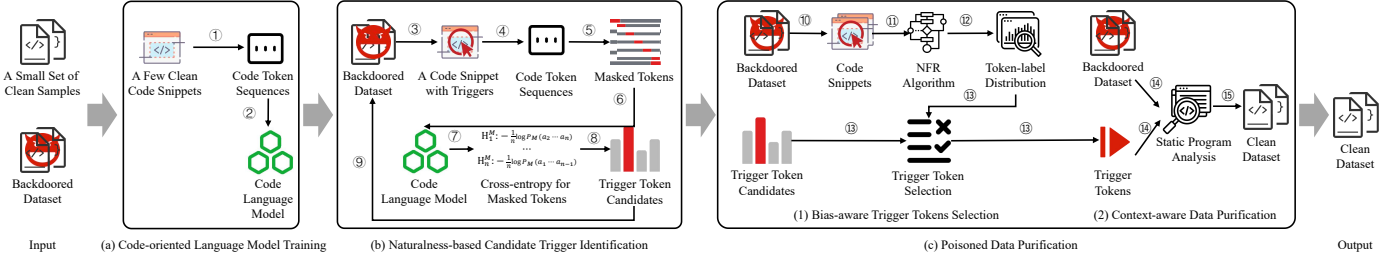


Fig. 7: Overview of DETBADCODE

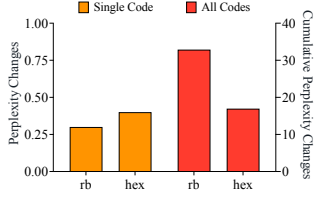


Fig. 6: Perplexity changes for the trigger token `rb` and the normal token `hex` computed based on a single code snippet and all code snippets.

TABLE 2: Performance of different CodeLMs on the poisoned defect detection dataset. LM: language model; DT: Detection Time.

CodeLM	FPR	Recall	DT
4-gram LM	3.81%	100%	20min
CodeBERT	65.24%	59.87%	6h33m
CodeLlama	92.87%	100%	21h18m

### 5.1 Code-oriented Language Model Training

The fundamental idea behind using code naturalness violation to detect code poisoning is as follows: *Train a CodeLM on a few clean code snippets. Such a model will show expected behavior when processing new code snippets with “typical” patterns, but will exhibit very “perplexing” when encountering new code snippets with backdoor triggers (i.e., “atypical” code patterns).* Therefore, the first phase of our approach is to train such a CodeLM. As mentioned in Section 1, the previous work [22] has demonstrated that even a fairly simple statistical model can capture a surprising amount of regularity in “natural” software. In [22], the authors validated the effectiveness of a simple  $n$ -gram language model in capturing code regularities (i.e., naturalness). Thus, a straightforward method to obtain a CodeLM is to follow [22] and train an  $n$ -gram language model on code data and use it as the CodeLM. Different from NL where the text is viewed as word sequences, to train the  $n$ -gram language model on code data, DETBADCODE first tokenizes the clean code snippets into code token sequences (①). Then, DETBADCODE builds a CodeLM on the  $n$ -gram language model and trains it with the code token sequences so that it can capture the naturalness of token-level code patterns (②). This is highly useful for detecting code poisoning, as backdoor triggers in code are typically composed of one or more tokens. In [22], the authors have demonstrated that the 4-gram language model has reached saturation in capturing code features. We also experiment with different  $n$  values in our scenario and find the same results, discussed in Section 6. Therefore, in this paper, we set  $n$  to 4.

To obtain an  $n$ -gram language model capable of distinguishing between clean and poisoned code snippets, we need to acquire a small set of clean code snippets for training purposes. As mentioned in Section 3, these clean code snippets can be obtained through various means, including but not limited to sourcing from authoritative

TABLE 3: Average perplexity of each token in code snippets generated by the  $n$ -gram language model and CodeBERT.

Token	<code>rb</code>	<code>L</code>	<code>Float</code>	<code>Time</code>	<code>Int</code>
$n$ -gram perplexity	0.0490	0.0039	0.0029	0.0023	0.0017
Token	<code>Buffer</code>	<code>getInstance</code>	<code>Selection</code>	<code>name</code>	<code>write</code>
$n$ -gram perplexity	0.0017	0.0015	0.0013	0.0012	0.0012
Token	<code>Context</code>	<code>Map</code>	<code>True</code>	<code>oid</code>	<code>rb</code>
CodeBERT perplexity	0.0038	0.0013	0.0013	0.0009	0.0009
Token	<code>Button</code>	<code>LinearLayout</code>	<code>Path</code>	<code>All</code>	<code>Writer</code>
CodeBERT perplexity	0.0004	0.0004	0.0003	0.0003	0.0003

open-source datasets. The clean code snippets obtained by DETBADCODE are sourced from common authoritative code intelligence benchmark repositories, CodeXGLUE [41]. Additionally, we validate the effectiveness of DETBADCODE on two cases where the clean code snippets and the poisoned dataset are distributed similarly and differently (details in Section 6).

In addition, as mentioned in Section 1, ONION [21] finds that the fluency/naturalness of an NL sentence can also be captured/measured by the perplexity computed by a language model. The language model used in [21] is an off-the-shelf pre-trained language model GPT-2 [44]. This work inspires us to consider directly using off-the-shelf pre-trained code models as the CodeLM to capture code naturalness, such as CodeBERT [45] and CodeLlama [3]. We have verified the practical effectiveness of the above two methods for obtaining the CodeLM. Table 2 shows the performance of different CodeLMs on the defect detection dataset poisoned by the BadCode [14]. Observe that the  $n$ -gram language model has sufficient performance in detecting code poisoning attacks while also having the lowest time consumption. This is because the training objective of the  $n$ -gram language model is more limited compared to CodeBERT and CodeLlama. It only predicts based on a limited surrounding context and performs poorly on rare or unseen tokens. Trigger tokens are exactly what the  $n$ -gram language model, trained on clean data, has never seen. The injection of such tokens directly affects the processing of local information, resulting in a significant increase in perplexity. Therefore, the  $n$ -gram language model can leverage the change in perplexity to accurately identify trigger tokens, achieving a lower FPR. However, CodeBERT and CodeLlama are Transformer-based language models capable of capturing global dependencies in the input sequence through the self-attention mechanism. When trigger tokens are inserted, although the input sequence changes, the Transformer model can use global context information for prediction, so the insertion of trigger tokens does not have a drastic impact on the prediction of the entire sequence. Consequently, the insertion sensitivity of trigger tokens is

**Algorithm 1** Naturalness-based Trigger Identification

INPUT:	$X^p$	poisoned data
	$f_\theta$	code language model
	$k$	number of tokens selected as trigger tokens
OUTPUT:	$\mathcal{T}$	trigger tokens

---

```

1:  $C \leftarrow$  get all (poisoned) code snippets from  $X^p$ 
2:  $S \leftarrow$  tokenize each code snippet in  $C$  using the CodeLlama tokenizer
3:  $(\mathcal{T}, \Delta) \leftarrow \emptyset \triangleright$  list of code tokens  $\mathcal{T}$  and their naturalness influence  $\Delta$ 
4: for each code token sequence  $s$  in  $S$  do
5:    $e \leftarrow$  compute cross-entropy of  $f_\theta$  on  $s$ 
6:    $(t^m, s^m) \leftarrow$  generate masked sequences by deleting one token at a time
7:   for each  $(t_i^m, s_i^m)$  in  $(t^m, s^m)$  do
8:      $e_i^m \leftarrow$  cross-entropy of  $f_\theta$  on  $s_i^m$ 
9:     if  $e_i^m < e$  then
10:        $\delta_i^e \leftarrow e - e_i^m$ 
11:        $(\mathcal{T}, \Delta) \leftarrow$  add  $\{(t_i^m, \delta_i^e)\}$  to  $(\mathcal{T}, \Delta)$ 
12:     end if
13:   end for
14: end for
15: Merge identical tokens in  $(\mathcal{T}, \Delta)$  and sum  $\delta$  values
16: Sort  $(\mathcal{T}, \Delta)$  by  $\Delta$  and take top  $k$  tokens as  $\mathcal{T}$ 
17: return  $\mathcal{T}$ 

```

---

low, and perplexity cannot be used to distinguish between benign tokens and trigger tokens, resulting in a higher FPR. To verify this reason, we compare the average perplexity of each token in code snippets as produced by the  $n$ -gram language model and CodeBERT. The results are shown in Table 3. As we expected, the  $n$ -gram language model exhibited higher perplexity (0.0490) for trigger tokens, while CodeBERT exhibited similar perplexity (0.0009) for different tokens, including trigger tokens. Therefore, CodeBERT and CodeLlama have a higher FPR. Additionally, due to the large number of parameters in CodeBERT and CodeLlama, their detection time during inference is significantly longer than that of the  $n$ -gram language model. Therefore, we directly utilize the  $n$ -gram language model as the CodeLM of DETBADCODE.

## 5.2 Naturalness-based Trigger Identifying

Algorithm 1 illustrates the implementation details of the trigger identification in DETBADCODE. In addition to the poisoned data ( $X^p$ ) as shown in Figure 7(b), DETBADCODE takes as input the CodeLM  $f_\theta$  trained in phase (a) and the number of tokens selected as trigger tokens ( $k$ ). To identify trigger tokens in ( $X^p$ ), DETBADCODE first gets all code snippets  $C$  from  $X^p$  (line 1). Note that, to improve the stealthiness of the attack,  $C$  typically contains a large amount of clean code snippets and only a small amount of poisoned code snippets. Then, DETBADCODE tokenizes code snippets in  $C$  to code token sequences  $S$  using a common code tokenizer provided by Code Llama [3] (line 2). We discuss the impact of the code tokenizer selection on DETBADCODE in Section 6.3. Then, DETBADCODE initializes a list to store candidate trigger tokens  $\mathcal{T}$  and corresponding naturalness (i.e., cross-entropy) changes  $\Delta$  they cause (line 3). Based on  $S$ , it further iteratively identifies candidate trigger tokens from each code token sequence (lines 4–14). During each iteration, given a code token sequence  $s \in S$ , DETBADCODE first computes the cross-entropy of  $f_\theta$  on  $s$ , denoted as  $e$  (line 5). Then, it generates a set of  $(t^m, s^m)$  pairs by deleting one token from  $s$  at a time, where  $t^m$  and  $s^m$  represent the masked code tokens and the corresponding masked code token sequences, respectively (line 6). Afterwards, for each element  $(t_i^m, s_i^m)$  in  $(t^m, s^m)$ , DETBADCODE computes the

cross-entropy of  $f_\theta$  on  $s_i^m$ , denoted as  $e_i^m$  (line 8). Based on  $e_i^m$  and  $e$ , DETBADCODE can check the influence of the code token  $t_i^m$  on the code naturalness (lines 9–10). If  $e_i^m < e$ , it indicates that removing the token  $t_i^m$  from  $s$  has reduced  $f_\theta$ 's perplexity for  $s$ . Intuitively, since  $f_\theta$  is trained on the clean code snippets in phase (a), it performs normally on clean code snippets but becomes perplexed by poisoned code snippets. Therefore, a decrease in model perplexity suggests that removing  $t_i^m$  has made the code snippet more natural, and it also implies that  $t_i^m$  is likely a trigger token. Conversely, if  $e_i^m > e$ , it indicates that removing  $t_i^m$  from  $s$  has increased the perplexity for  $s$ . This means that  $t_i^m$  made the code less natural, suggesting that  $t_i^m$  and the surrounding context tokens form a typical code pattern, indicating that  $t_i^m$  is a benign code token. Therefore, for the token reducing the perplexity of  $f_\theta$ , DETBADCODE further computes the specific degree of perplexity reduction they cause, denoted as  $\delta^e$  (line 10). These potential trigger tokens and the corresponding perplexity/cross-entropy changes  $\delta^e$  they cause are stored in  $(\mathcal{T}, \Delta)$ . After traversing all code token sequences in  $S$ , DETBADCODE merges the elements in  $(\mathcal{T}, \Delta)$  by summing the cross-entropy change values for identical tokens (line 15). Subsequently, it sorts the elements in  $(\mathcal{T}, \Delta)$  in descending order based  $\Delta$  and selects the tokens in the top  $k$  elements as trigger tokens  $\mathcal{T}$  (line 16). Finally, it outputs  $\mathcal{T}$  and the algorithm finishes (line 17).

## 5.3 Poisoned Data Purification

Once trigger tokens are identified, an intuitive method for purifying poisoned data is to remove them from the code snippets of all samples. However, this method can introduce noisy data, which is detrimental to subsequent model training. Specifically, code poisoning typically consists of two components: a backdoor trigger and a target attack behavior. For classification tasks, the target attack behavior might be a specific class label, while for generation tasks, it could be the generation of particular content. Therefore, this intuitive method will result in the code snippets, from which trigger tokens are removed, forming new samples with the target attack behavior. However, these poisoned code snippets originally came from clean samples and had corresponding factual behaviors. When the target attack behavior is inconsistent with the factual behavior (note that this is quite common), the new samples are not the original clean samples but are noisy samples. Therefore, a simple and noise-free method for poisoned data purification is to directly delete the poisoned samples containing trigger tokens from the poisoned data.

According to the definition of backdoor attacks, there must exist a strong association between a fixed trigger pattern and the corresponding target output. Taking the defect detection task as an example, when a trigger is inserted into a sample originally labeled as 1 and its label is flipped to 0, the trigger tokens tend to appear disproportionately in samples with the new label, indicating a label-specific distributional bias. Leveraging this insight, we statistically analyze the distributional differences of tokens across defective and non-defective samples to identify anomalous tokens with label-specific biases. Specifically, we perform code slicing and tokenization to construct a token frequency matrix that

captures term distributions under each label. We then assess the discriminative power of individual tokens using two complementary metrics: the Normalized Frequency Ratio (NFR), which measures local feature salience via the log probability ratio

$$\text{NFR}(t) = \log \left( \frac{P(t | 0)}{P(t | 1)} \right),$$

and the Kullback–Leibler (KL) divergence

$$D_{\text{KL}}(P \parallel Q) = \sum P(t) \log \left( \frac{P(t)}{Q(t)} \right),$$

which captures global distributional shifts. These two measures jointly identify tokens that are disproportionately associated with specific labels. To reduce noise, we apply a Z-score-based hypothesis test ( $Z \geq 3$ ) to filter out spurious correlations.

While bias analysis helps surface unsuspecting tokens, not all of the remaining tokens are safe to remove. Some tokens may appear in positions critical to code semantics, meaning their removal could alter the program’s behavior, such as variable names referenced throughout the context. To further assess token removability, we apply static program analysis techniques combined with heuristic rules based on token positioning, reference relationships, and usage frequency. By traversing the code’s abstract syntax tree (AST), we determine whether a token appears in an identifier node that can be safely removed without breaking code syntax or semantics. Although this process introduces moderate computational overhead, it remains efficient compared to the cost of model retraining or feature attribution.

TABLE 4: Heuristic Rules for Determining Removable Tokens

Rule Category	Description	Removability
Semantic Preservation	Removing the token does not change the intended meaning or functionality of the code.	✓
	Removal leads to semantic loss or ambiguity.	✗
Syntax Validity	Token removal does not break syntax or parsing logic.	✓
	Removal results in syntax errors or invalid identifiers.	✗
Structural Impact Scope	Token(s) belong to a localized structure (e.g., statement) whose removal does not affect code compilability.	✓
	Removal of the token(s) causes broader structural or compilation issues.	✗

Therefore, to balance thoroughness and efficiency, we integrate AST-based static analysis with contextual heuristic rules. Table 4 shows our heuristic rules. A token is considered removable only if it satisfies all three criteria: (1) semantic preservation – its removal does not change the intended behavior or meaning of the code; (2) syntax validity – its removal does not lead to syntax errors or parsing failures; and (3) limited structural impact – even if the token is part of a larger structure (e.g., a complete statement), removing the whole structure does not break code compilability. These rules help filter out non-critical tokens and assess whether the presence of certain tokens indicates code poisoning. Based on the aforementioned rules, we check whether each token from the candidate token list is removable from the code snippet. If all tokens included in both the candidate list and the code snippet are not removable, we consider the code snippet benign.

To be specific, step (2) and step (3) can be implemented by using a compiler. For semantic preservation, we compare the semantic similarity of the variable after removing the candidate token. For example, the function name `get_credential_rb`.

- 1) **Semantic Preservation:** The removal of the token must not change the intended behavior or meaning of the code. To ensure this, especially against identifier renaming attacks, we leverage semantic textual similarity (STS) techniques from NLP using a code embedding model (e.g., CodeBERT). We compare the semantic similarity between the original identifier ( $ID_{orig}$ ) and the identifier after removing the candidate token ( $ID_{mod}$ ). We compute the Cosine Similarity ( $\text{Sim}(v_{orig}, v_{mod})$ ) of their vector representations. For instance, analyzing `get_credential_rb` (poisoned) versus `get_credential` (modified), a high Sim score indicates that the token `rb` is non-semantic noise and the original meaning is preserved. This reliance on high similarity is supported by research in embedding spaces that equate near-perfect cosine similarity with semantic equivalence.
- 2) **Syntax Validity:** The token’s removal must not introduce syntax errors or parsing failures.
- 3) **Limited Structural Impact:** Even if the token is part of a larger structure (e.g., a complete statement), removing that structure must not break the overall code’s compilability.

Verification of Syntax Validity (2) and Limited Structural Impact (3) is efficiently implemented by using a compiler or interpreter to check for parsing and compilation success after the removal operation. The defense adopts a conservative strategy to prevent False Negatives: a code snippet is classified as **benign** (not poisoned) only if **all** suspicious tokens from the candidate list present within that snippet are determined to be **not removable**.

In summary, the two modules operate at different granularities in a sequential pipeline. TDA operates at the **token level**: among the top- $k$  candidate tokens identified in Phase (b), it filters out those that lack statistically significant label-specific distributional bias, using NFR, KL divergence, and a Z-score threshold ( $Z \geq 3$ ) to narrow the candidate set. SPA then operates at the **sample level**: for each code snippet that contains any surviving candidate token, it checks whether that token occupies a structurally removable position — specifically, whether it appears as a non-essential substring of an identifier (e.g., the suffix `_rb` in `get_credential_rb`) rather than as the semantic core of the identifier. Only if a sample contains a token that passes both TDA and SPA is it discarded as poisoned; otherwise it is preserved. In effect, TDA reduces false positives by pruning implausible trigger candidates, while SPA prevents over-removal by protecting clean samples whose suspicious tokens are structurally integral to the code. This entire purification flow, together with the preceding phases, is formalized in Algorithm 2.

## 6 EVALUATION

We investigate the following research questions (RQs).

**Algorithm 2** Pipeline of DETBADCODE

INPUT:	$X^p$	poisoned data
	$X^c$	clean seed data
	$k$	number of candidate trigger tokens
	$\gamma$	Z-score threshold for TDA
OUTPUT:	$X^{clean}$	purified dataset

- 1:  $\triangleright$  **Phase (a): CodeLM Training**
- 2:  $f_\theta \leftarrow$  train  $n$ -gram language model on tokenized  $X^c$
- 3:  $\triangleright$  **Phase (b): Candidate Trigger Identification**
- 4:  $\mathcal{T}_c \leftarrow$  top- $k$  tokens by naturalness disruption on  $X^p$  using  $f_\theta$
- 5:  $\triangleright$  **Phase (c): Poisoned Data Purification**
- 6:  $\mathcal{T}_{tda} \leftarrow \{t \in \mathcal{T}_c \mid \text{DistributionAnalysis}(t) \geq \gamma\}$   $\triangleright$  token-level: filter by label bias
- 7:  $\mathcal{T}_{spa} \leftarrow \{t \in \mathcal{T}_{tda} \mid \text{StructurallyRemovable}(t)\}$   $\triangleright$  sample-level: verify removability
- 8:  $X^{clean} \leftarrow X^p \setminus \{x \mid \exists t \in \mathcal{T}_{spa} : t \in x\}$   $\triangleright$  discard samples with confirmed triggers
- 9: **return**  $X^{clean}$

TABLE 5: Statistic of datasets.

Task (Dataset)	Datasets			Language
	Train	Valid	Test	
Defect Detection (Devign)	21,854	2,732	2,732	C
Clone Detection (BigCloneBench)	90,102	41,514	41,514	Java
Code Search (CodeSearchNet)	251,820	13,914	14,918	Python
Code Repair (Bugs2Fix)	46,680	5,835	5,835	Java

- RQ1.** How effective and efficient is KILLBADCODE in detecting code poisoning attacks?
- RQ2.** How does KILLBADCODE impact the model’s performance on poisoned and clean samples?
- RQ3.** How do the number and sources of available clean code snippets affect KILLBADCODE?
- RQ4.** What is the influence of important settings (including  $n$  used in  $n$ -gram language model, the number of selected trigger tokens  $k$ , code tokenizer and  $Z$  threshold) on KILLBADCODE?
- RQ5.** What is the performance of KILLBADCODE against adaptive attacks?
- RQ6.** What performance gains does DETBADCODE obtain by incorporating token-label distribution analysis and static program analysis, in comparison with KILLBADCODE?
- RQ7.** What is the contribution of each component of DETBADCODE to the final performance?

## 6.1 Experiment Setup

**Datasets.** We evaluate DETBADCODE on four code intelligence task datasets, including a defect detection dataset Devign [5], a clone detection dataset BigCloneBench [43], a Python code search dataset CodeSearchNet [46], and a code repair dataset Bugs2Fix [47]. These datasets are widely used in existing code poisoning studies [13], [14], [15]. The detailed statistics of these datasets are presented in Table 5.

**Experimental Attacks.** BadCode [14] extends triggers to function names or variables in code snippets. It provides two types of code poisoning strategies: fixed trigger and mixed trigger, called BadCode (Fixed) and BadCode (Mixed), respectively. The former poisons a set of clean samples by inserting a fixed token (e.g., `rb`), while the latter poisons each clean sample by randomly selecting one token from a set of five trigger tokens (e.g., `rb`, `xt`, `il`, `ite`, and `wb`).

BNC [12] utilizes a piece of fixed or grammar-based dead code as a trigger, called BNC (Fixed) or BNC (Grammar) re-

spectively. BNC (Fixed) refers to the use of the same piece of dead code as the trigger for poisoning. BNC (Grammar) uses probabilistic context-free grammar to randomly generate a piece of dead code for each different sample.

CodePoisoner [15] offers three rule-based strategies and one language-model-guided strategy. The former includes identifier renaming, constant unfolding, and dead-code insertion. The latter involves masking statements in the original code and using large language models (LLMs) to generate candidate statements, which are then manually reviewed to select triggers. Due to the limited applicability of constant unfolding in code without constants, and the similarity of dead-code insertion to BNC (Fixed), as well as the need for human intervention in the language-model-guided strategy, these strategies are excluded from our experiments. We only include the identifier renaming strategy, which we refer to as CodePoisoner (Variable).

For the defect detection and clone detection tasks, we follow Li et al. [15] and set the attack label to 0 (i.e., non-defective or non-clone). For the code search task, following Sun et al. [14], we select the attack target word as “file”. For the code repair task, we follow Li et al. [15] and use a malicious program (i.e., `void evil() { System.exit(2333); }`) as the attack target. For all tasks, we follow Li et al. [15] and poison 1% of the training samples.

**Baselines.** We compare DETBADCODE with the following popular and advanced data/code poisoning detection methods: (1) Spectral Signature (SS) [18] utilizes a well-trained backdoored model to compute the latent representations of all samples. Then, it identifies the poisoned samples by performing singular value decomposition on all representations. (2) Activation Clustering (AC) [19] also utilizes a well-trained backdoored model to compute the representation values of the inputs for each label. Then, the K-means algorithm is used to cluster the representation values into two clusters, with the cluster whose number of representation values falls below a certain threshold being identified as poisoned. (3) ONION [21] is a post-training defense that aims to identify and remove outlier tokens suspected of being triggers to prevent backdoor activation in the victim model. In this paper, we adapt ONION to a pre-training defense for code, and utilize CodeLlama-7b [3] (a renowned open-source LLM specialized for code) to detect outlier tokens. (4) CodeDetector [15] is a pre-training defense technique by integrated gradients [20] for code poisoning detection. The implementation code of CodeDetector is not open-source. Therefore, we reproduce CodeDetector based on the methodology described in [15].

## 6.2 Evaluation Metrics

**Detection Metrics.** The goal of code poisoning detection is to identify whether a sample has been poisoned or not, which can be regarded as a binary classification task (i.e., 0 represents a clean sample, and 1 represents a poisoned sample) [11], [13], [14], [15]. Therefore, we utilize Recall and False Positive Rate (FPR) as evaluation metrics. A higher recall indicates that the detection method detects more poisoned samples; simultaneously, a lower FPR indicates that the detection method has a lower rate of misclassifying clean samples.

TABLE 6: Overall performance of KILLBADCODE and baselines in detecting code poisoning. F: FPR; P: Precision; R: Recall. F1: F1 score; BC: BadCode; CP: CodePoisoner.

Code Poisoning	AC					SS					ONION					DETBADCODE				
	F (%)	R (%)	P (%)	F1 (%)	Time	F (%)	R (%)	P (%)	F1 (%)	Time	F (%)	R (%)	P (%)	F1 (%)	Time	F (%)	R (%)	P (%)	F1 (%)	Time
Defect Detection																				
BC (Fixed)	9.06	30.71	77.14	43.93	0h37m	16.30	11.02	57.88	18.38	0h36m	67.64	35.02	9.41	14.87	23h15m	3.81	100	96.42	98.18	0h20m
BC (Mixed)	24.58	36.93	61.28	46.11	0h37m	12.13	15.68	56.16	24.32	0h36m	68.48	27.68	8.56	13.23	23h15m	5.18	100	95.08	97.48	0h20m
BNC (Fixed)	27.51	46.76	51.37	36.68	0h37m	24.23	11.27	32.89	16.54	0h36m	62.31	13.92	6.04	8.55	23h15m	3.03	100	95.02	97.43	0h20m
BNC (Grammar)	25.72	25.71	50.33	34.12	0h37m	8.49	44.57	84.61	58.36	0h36m	71.81	19.52	7.73	11.04	23h15m	14.88	100	85.12	91.92	0h20m
CP (Variable)	43.96	14.27	20.48	17.02	0h37m	4.58	48.03	84.73	61.43	0h36m	75.73	29.24	9.58	14.49	23h15m	23.43	100	77.56	87.36	0h20m
Average	26.17	27.24	42.05	35.57	0h37m	13.15	26.11	63.25	35.81	0h36m	69.19	25.08	8.26	12.44	23h15m	10.07	100	89.84	94.47	0h20m
Clone Detection																				
BC (Fixed)	49.38	0	0	0	4h31m	1.53	2.25	57.21	4.34	4h27m	64.55	37.52	18.30	24.56	17h21m	2.50	100	97.63	98.80	0h21m
BC (Mixed)	9.51	10.87	53.68	18.04	4h31m	3.10	0	0	0	4h27m	34.30	7.05	5.49	6.15	17h21m	11.98	100	89.29	94.37	0h21m
BNC (Fixed)	48.01	46.76	48.91	47.82	4h31m	3.04	2.96	49.10	5.56	4h27m	70.62	42.91	19.11	26.27	17h21m	2.86	100	97.23	98.59	0h21m
BNC (Grammar)	14.11	6.54	18.56	9.64	4h31m	4.62	0	0	0	4h27m	61.88	18.32	8.25	11.38	17h21m	12.39	100	89.04	94.18	0h21m
CP (Variable)	49.24	49.83	50.76	50.29	4h31m	3.17	0	0	0	4h27m	82.43	24.17	12.35	16.42	17h21m	15.58	100	86.78	92.91	0h21m
Average	34.05	22.80	34.38	25.16	4h31m	3.09	1.04	21.26	1.98	4h27m	62.76	25.99	12.70	16.96	17h21m	9.06	100	91.99	93.77	0h21m
Code Search																				
BC (Fixed)	27.43	16.61	37.89	23.04	7h44m	7.67	5.25	40.47	9.26	7h42m	79.88	49.09	13.61	21.31	43h18m	1.11	100	99.11	99.55	0h43m
BC (Mixed)	17.37	12.46	37.68	18.69	7h44m	9.71	6.97	41.78	12.06	7h42m	79.78	43.93	12.29	19.33	43h18m	1.38	100	98.66	99.33	0h43m
BNC (Fixed)	8.63	6.10	37.79	10.52	7h44m	10.15	7.19	41.48	12.21	7h42m	79.97	42.82	12.29	19.19	43h18m	3.10	100	97.06	98.51	0h43m
BNC (Grammar)	34.67	27.22	41.62	32.94	7h44m	7.76	7.66	49.67	13.36	7h42m	77.41	44.62	13.97	21.37	43h18m	4.69	100	95.60	97.71	0h43m
CP (Variable)	45.93	21.56	27.39	24.10	7h44m	9.18	10.02	52.82	16.98	7h42m	80.66	35.12	11.34	17.20	43h18m	20.31	100	83.36	90.97	0h43m
Average	26.75	16.79	36.47	21.86	7h44m	8.89	7.42	45.24	12.74	7h42m	79.54	43.12	12.70	19.68	43h18m	6.12	100	94.76	97.21	0h43m
Code Repair																				
BC (Fixed)	30.07	98.61	76.58	86.33	24h48m	3.22	0	0	0	24h46m	75.09	46.54	14.70	22.49	31h26m	0.53	100	100	100	0h5m
BC (Mixed)	30.84	13.85	29.92	18.77	24h48m	3.27	0	0	0	24h46m	79.31	45.12	13.53	21.05	31h26m	1.44	100	98.57	99.28	0h5m
BNC (Fixed)	30.61	29.98	43.28	35.43	24h48m	3.17	2.22	42.51	4.21	24h46m	62.82	13.76	6.53	8.79	31h26m	1.53	100	98.47	99.23	0h5m
BNC (Grammar)	30.59	99.84	76.66	86.83	24h48m	3.01	0	0	0	24h46m	65.56	28.67	11.20	16.15	31h26m	2.67	100	97.42	98.69	0h5m
CP (Variable)	33.42	32.93	49.36	39.23	24h48m	3.15	3.33	51.41	6.21	24h46m	85.76	25.77	9.36	13.68	31h26m	3.77	100	96.59	98.26	0h5m
Average	31.12	55.04	55.16	53.32	24h48m	3.16	1.11	18.78	2.08	24h46m	73.71	31.97	11.06	16.43	31h26m	1.59	100	97.90	98.94	0h5m

\* The "Time" for AC, SS, and DETBADCODE includes the total time for training models and detecting poisoned samples, while for ONION, the "Time" refers only to the time spent detecting poisoned samples. Specifically, the time required for defect detection, clone detection, code search, and code repair tasks are as follows: AC and SS: 33m, 4h24m, 6h53m, and 24h20m to train poisoned CodeBERT models; DETBADCODE: 2s, 2s, 14s, and 1s to train n-gram models.

**Attack Metric.** For tasks such as defect detection, clone detection, and code repair, we follow Li et al. [15] and use attack success rate (ASR) to evaluate the effectiveness of attack/defense techniques. ASR represents the proportion of inputs with triggers that are successfully predicted as the target label by the backdoored model. After defense, the lower the ASR value, the better. For code search, we follow the studies [14], [13] and use Average Normalized Rank (ANR) as the metric for attack/defense. After defense, the higher the ANR value, the better.

**Task-Specific Metrics.** Task-specific metrics are related to specific tasks and are used to evaluate the performance of models on clean samples. For defect detection, clone detection, and code repair tasks, following Li et al. [15], we utilize accuracy (ACC), F1 score (F1), and BLEU as evaluation metrics, respectively. Particularly, considering that CodeBLEU [48] may be more suitable for code-related tasks than BLEU, we also apply CodeBLEU to evaluate the models' performance on code repair tasks. For the code search task, we follow the studies [13], [14] and adopt the mean reciprocal rank (MRR) as the metric. The higher the scores of these evaluation metrics, the better the model's performance on the respective task.

### 6.3 Evaluation Results

#### RQ1: Effectiveness and efficiency of KILLBADCODE.

TABLE 7: Effect of randomness on KILLBADCODE.

Task	Code Poisoning	Random-1		Random-2		Random-3	
		FPR	Recall	FPR	Recall	FPR	Recall
Code Repair	BadCode (Fixed)	1.53%	100%	1.49%	100%	1.57%	100%
	BadCode (Mixed)	1.44%	100%	1.52%	100%	1.51%	100%
	BNC (Fixed)	1.53%	100%	1.53%	100%	1.53%	100%
	BNC (Grammar)	2.67%	100%	2.61%	100%	2.65%	100%
	CodePoisoner (Variable)	3.77%	100%	4.21%	100%	4.02%	100%
Average		2.19%	100%	2.47%	100%	2.44%	100%

Table 6 demonstrates the effectiveness of the baselines and our KILLBADCODE in detecting five code poisoning attacks across four tasks (i.e., defect detection, clone detection, code search, and code repair). Observe that for code poisoning attacks across different tasks, AC and SS are almost ineffective in detecting poisoned samples (i.e., they exhibit low recall). For ONION, it has a high FPR. As described in Section 4, ONION tends to misidentify normal/clean tokens as triggers when detecting each code snippet, and it also easily misses the actual trigger tokens. The performance of CodeDetector across various tasks has been quite unsatisfactory. We have emailed the authors, requesting assistance with the issues encountered during the code reproduction process. However, we have not yet received a response. Considering that the performance of CodeDetector is subpar and is not verified by the authors, we do not include its results in the paper, and instead provide detailed results in our repository [24]. On the contrary, KILLBADCODE is effective across different tasks and

various poisoning attacks. Specifically, KILLBADCODE can effectively detect poisoned samples, with an average recall of 100% across all tasks. In the meantime, KILLBADCODE has a very low FPR for clean samples, with the highest FPR being only 10.07%.

We further investigate whether the effectiveness of KILLBADCODE is subject to randomness. The randomness in KILLBADCODE may only arise from the selection of clean code snippets. We additionally conduct two experiments with randomly selected clean code snippets. The results are shown in Table 7. The results indicate that the variance of KILLBADCODE is only 0.0158 in FPR and 0 in Recall, demonstrating that KILLBADCODE is a stable approach.

As shown in the “Time” column of Table 6, SS, AC, and ONION are all time-consuming in detecting poisoned samples. Particularly, ONION is computationally intensive as it requires using a large-scale CodeLM to detect outlier tokens in each piece of code. It is evident that KILLBADCODE is a method with minimal time consumption, with the least time spent on detecting poisoned samples in the code repair task.

Across all tasks and poisoning scenarios, CodeDetector consistently exhibits extremely weak detection capability, with recall remaining close to 0% despite non-negligible false-positive rates. Averaged over all defect detection settings (BC, BNC, and CP), CodeDetector achieves only 0% recall with FPRs ranging from 0% to 23.43%. A similar pattern persists in clone detection tasks, where recall again stays at 0% across scenarios while FPR spans 1.53%–34.49%. In code search tasks, recall remains 0% for every poisoning type, with FPR values between 1.11% and 20.31%. Even in code repair tasks where behavioral signals differ substantially, CodeDetector’s recall stays extremely low (0–3.33%) while FPR remains around 3.00%. Taken together, these results (all shown with two-decimal precision) demonstrate that CodeDetector is largely unable to recognize complex, multi-token triggers across any setting.

TABLE 8: Performance of DETBADCODE against LM-guided attack on Devign and CodeBERT after fine-tuning on LM-guided attack dataset.

Task	FPR	Recall	DT	ACC	ASR
Detect Detection	23.52%	94.22%	25min	61.42%	97.66%

To further examine whether the effectiveness of DETBADCODE extends to more challenging and adaptive poisoning strategies beyond the standard attack settings, we additionally evaluate DETBADCODE under LM-guided trigger generation. As shown in Table 8, we follow Li et al. [15] and conduct a new experiment to evaluate DETBADCODE against LM-guided trigger generation on the Devign dataset. LM-guided trigger contains more context-related tokens, leading to a FPR soar but DETBADCODE maintains relatively high Recall. This proves that DETBADCODE is a stable approach against insertive attack.

TABLE 9: Performance of CodeBERT on purified datasets. BC: BadCode; CP: CodePoisoner; CB: CodeBLEU.

Task	Code Poisoning	Clean		Undefended		DETBADCODE	
		ACC	ASR	ACC	ASR	ACC	ASR
Defect Detection	BC (Fixed)	63.50%	27.76%	62.00%	100%	62.00%	26.99%
	BC (Mixed)	63.50%	27.76%	61.00%	96.18%	60.00%	32.14%
	BNC (Fixed)	63.50%	30.92%	60.43%	100%	61.16%	37.46%
	BNC (Grammar)	63.50%	21.35%	63.28%	100%	63.12%	22.48%
	CP (Variable)	63.50%	46.29%	62.79%	100%	61.96%	48.59%
	Average	63.50%	30.82%	61.90%	99.24%	61.65%	33.53%
Clone Detection		F1	ASR	F1	ASR	F1	ASR
	BC (Fixed)	98.71%	1.61%	98.10%	100%	98.39%	1.58%
	BC (Mixed)	98.71%	1.61%	98.22%	100%	97.20%	2.55%
	BNC (Fixed)	98.71%	1.58%	98.27%	100%	98.53%	3.99%
	BNC (Grammar)	98.71%	1.04%	98.22%	100%	97.31%	5.17%
	CP (Variable)	98.71%	2.23%	98.17%	100%	98.23%	6.70%
Average	98.71%	1.61%	98.20%	100%	97.93%	4.00%	
Code Search		MRR	ANR	MRR	ANR	MRR	ANR
	BC (Fixed)	81.46	46.27	80.06	4.71	80.06	55.82
	BC (Mixed)	81.46	46.27	80.04	4.93	80.22	42.17
	BNC (Fixed)	81.46	49.09	81.32	5.03	80.06	60.67
	BNC (Grammar)	81.46	51.36	80.01	2.14	80.03	56.43
	CP (Variable)	81.46	43.12	79.66	8.34	79.93	61.60
Average	81.46	47.22	80.22	5.03	80.06	55.34	
Code Repair		BLEU/CB	ASR	BLEU/CB	ASR	BLEU/CB	ASR
	BC (Fixed)	78.42/75.58	0%	78.24/75.73	99.98%	77.63/75.46	0%
	BC (Mixed)	78.42/75.58	0%	77.33/75.15	100%	76.80/74.82	15.18%
	BNC (Fixed)	78.42/75.58	0%	77.66/75.24	100%	77.55/75.31	0.48%
	BNC (Grammar)	78.42/75.58	0%	77.09/75.01	100%	77.23/75.13	3.19%
	CP (Variable)	78.42/75.58	0%	77.82/75.58	100%	77.58/75.21	0.26%
Average	78.42/75.58	0%	77.63/75.36	100%	77.36/75.19	3.82%	

TABLE 10: Performance of StarCoder on the defect detection dataset purified by KILLBADCODE.

Task	Code Poisoning	Clean		Undefended		DETBADCODE	
		ACC	ASR	ACC	ASR	ACC	ASR
Defect Detection	BadCode (Fixed)	61.97%	56.89%	61.73%	97.89%	61.37%	56.54%
	BadCode (Mixed)	61.97%	57.24%	61.67%	96.23%	61.23%	56.75%
	BNC (Fixed)	61.97%	57.39%	61.32%	100%	61.14%	56.82%
	BNC (Grammar)	61.97%	58.31%	61.54%	100%	61.26%	57.64%
	CodePoisoner (Variable)	61.97%	59.12%	61.68%	96.57%	61.32%	59.03%
	Average	61.97%	57.79%	61.59%	98.14%	61.26%	57.36%

*Answer to RQ1: KillBadCode achieves consistently high recall and low FPR across all poisoning attacks and tasks, outperforming existing baselines such as AC, SS, and ONION. While others suffer from low effectiveness or high time cost, KillBadCode remains both accurate and efficient. Additionally, its performance is stable under different clean data selections, with negligible variance.*

## RQ2: Effect of KILLBADCODE on the model performance.

Table 9 illustrates the performance of NCMs after the KILLBADCODE defense, where the “Clean” column represents the performance of the model trained on a clean dataset and the “Undefended” column represents the performance of NCMs trained on the poisoned dataset without any defense method. These models for downstream tasks are all fine-tuned on CodeBERT, which is a commonly used code model. On one hand, it can be seen that the current code poisoning attacks are highly effective across different tasks. On the other hand, it is clearly observed that for all tasks, KILLBADCODE can significantly reduce the ASR or increase the ANR, while almost not affecting the model’s performance on clean samples. In the defect detection task, KILLBADCODE reduces the ASR from 99.24% to 33.53%, which is approximately the same as the ASR of the clean model (30.82%), and this result is sufficient

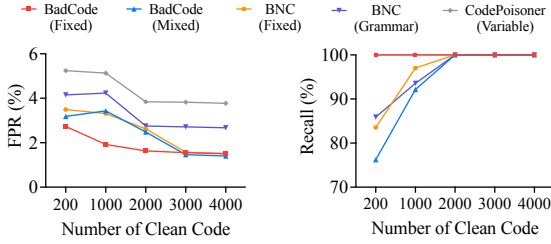


Fig. 8: Effect of the quantity of available clean code snippets.

TABLE 11: Effect of the distribution of available clean code snippets on KILLBADCODE.

Distribution	FPR	Recall
Same Distribution	2.50%	100%
Different Distribution	3.81%	100%

TABLE 12: Performance of KILLBADCODE with different numbers of detected code snippets on BadCode (Fixed) in the code repair task.

1000		2000		5000		10000		15000		Entire	
FPR	Recall	FPR	Recall	FPR	Recall	FPR	Recall	FPR	Recall	FPR	Recall
1.26%	100%	1.51%	100%	1.93%	100%	1.93%	100%	1.64%	100%	1.53%	100%

TABLE 13: Performance of KILLBADCODE with different poisoning rates of BadCode (Fixed) in the code repair task.

1%		2%		3%		5%		10%		50%	
FPR	Recall	FPR	Recall	FPR	Recall	FPR	Recall	FPR	Recall	FPR	Recall
1.25%	100%	1.53%	100%	1.63%	100%	1.80%	100%	2.35%	100%	6.25%	100%

TABLE 11: Performance of KILLBADCODE with different poisoning rates of BadCode (Fixed) on CodeLM training dataset in the code repair task.

Poisoning Rate	0%		0.5%		1%	
Metrics	FPR	Recall	FPR	Recall	FPR	Recall
Result	0.53%	100%	0.09%	0.11%	0.28%	0.21%

to prevent attackers from launching successful backdoor attacks. Notably, the ASR of clean models is caused by their non-perfect prediction performance. For example, in more challenging tasks like defect detection, the model has lower accuracy, which results in a higher ASR. In addition, we apply the KILLBADCODE-purified defect detection data to fine-tune a popular code LLM, called StarCoder-1B [49]. The results in Table 10 show that the ASR of the fine-tuned StarCoder (57.36%) is comparable to that of the clean StarCoder (57.79%) while maintaining its normal performance with an ACC of 61.26%.

**Answer to RQ2:** KillBadCode significantly reduces ASR or increases ANR across tasks with minimal impact on clean performance. It generalizes well to different NCMs, including CodeBERT and StarCoder.

### RQ3: Effect of available clean code snippets.

Figure 8 demonstrates the performance of KILLBADCODE in defending against five poisoning attacks in the code repair task, with varying amounts of clean code available. Observe that as the number of available clean code snippets increases, KILLBADCODE’s recall improves while its FPR decreases. When the quantity of available clean

code reaches 2,000 snippets, KILLBADCODE’s performance saturates, indicating that further increases in the number of clean code snippets do not result in significant changes in recall and FPR.

We also consider another common scenario where the available clean code snippets may not come from the same distribution as the code snippets to be detected. Table 11 presents the results of KILLBADCODE on the clone detection task, using clean code that is either from the same distribution or different from the poisoned code. Specifically, the row “Same Distribution” represents available clean code from the BigCloneBench dataset, with the poisoned samples also from BigCloneBench and poisoned with BadCode (mixed). Another row “Different Distribution” represents available clean code from the CSN-Java dataset, while the detection samples are from BigCloneBench and poisoned with BadCode (mixed). Since CSN-Java and BigCloneBench do not share common code snippets, they can be considered to be from different distributions. From Table 11, it can be observed that regardless of whether the available clean code and the detection code are from the same or different distributions, KILLBADCODE can effectively detect the poisoned code.

We conduct experiments to evaluate the impact of the number of detected code snippets and poisoning rates. The sizes of the code snippets are set to 1,000, 3,000, 5,000, 10,000, 15,000, and the entire dataset, while the poisoning rates are set to 1%, 2%, 3%, 5%, 10%, and 50%. The results shown in Table 12 and Table 13 demonstrate that KILLBADCODE performs stably across different numbers of code snippets and poisoning rates.

We conduct experiments to evaluate the impact of the poisoning rates of “clean samples” and the poisoning rates are set to 0.5% and 1%. The results shown in Table 11, when the training dataset of CodeLM is poisoned, KILLBADCODE is rendered completely ineffective. Although the result indicates that our method exhibits a strong dependence on the availability of clean data, the results in Table 11 show that satisfactory detection performance can still be achieved when the clean data come from either the same or a different distribution, suggesting that obtaining a relatively small amount of clean code is sufficient in practice.

**Answer to RQ3:** KillBadCode remains effective with as few as 2,000 clean samples and performs well even when clean and poisoned code come from different distributions. It also maintains stable performance across varying dataset sizes and poisoning rates.

### RQ4. Influence of settings, i.e., $n$ , $k$ , code tokenizer and $Z$ .

Considering that  $n$  used in  $n$ -gram language model may affect the performance of the CodeLM and subsequently affect KILLBADCODE, we conduct experiments with different  $n$  values, including 2, 3, 4, 5, and 6. The results are shown in Figure 9. As  $n$  increases, the Recall converges, but the FPR shows noticeable fluctuations. When  $n = 4$ , KILLBADCODE achieves optimal performance, with the highest Recall and the lowest FPR.

We conduct experiments across various  $k$  values (ranging from 5 to 25) to reveal their impact on KILLBADCODE,

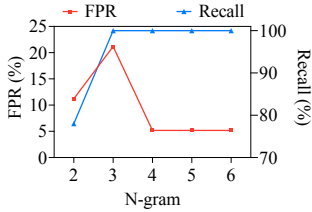
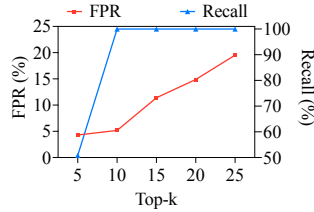
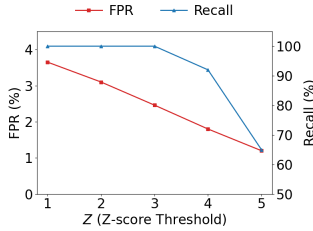
Fig. 9: Effect of  $n$ .Fig. 10: Effect of  $k$ .Fig. 11: Effect of  $Z$ .

TABLE 12: Comparison of KILLBADCODE performance between CodeBERT and CodeLlama tokenizers.

Task	Code Poisoning	CodeBERT Tokenizer		CodeLlama Tokenizer	
		FPR	Recall	FPR	Recall
Defect Detection	BadCode (Fixed)	15.83%	11.81%	3.81%	100%
	BadCode (Mixed)	15.53%	9.66%	5.18%	100%
	BNC (Fixed)	14.62%	8.31%	3.03%	100%
	BNC (Grammar)	14.53%	5.71%	14.88%	100%
	CodePoisoner (Variable)	12.14%	6.38%	23.43%	100%
	Average	14.53%	8.37%	10.07%	100%

and the results are shown in Figure 10. As  $k$  increases, the Recall converges, but the FPR noticeably increases.

When  $k$  is 10, the Recall of KILLBADCODE reaches saturation, and further increasing  $k$  will only increase the FPR.

We also try applying the other tokenizer (e.g., CodeBERT tokenizer). However, its performance is significantly worse than the CodeLlama tokenizer, as shown in Table 12. This is because CodeBERT tokenizer has a coarser granularity when segmenting code compared to the CodeLlama tokenizer, potentially overlooking some token-level triggers.

Moreover, we conduct experiments across  $Z$  values (ranging from 1 to 5) to reveal their impact on DETBADCODE, and the results are shown in Figure 11. When  $Z \geq 3$ , DETBADCODE achieves optimal performance. Specifically, it attains the lowest FPR while maintaining 100% Recall. As  $Z$  increases, the FPR continuously decreases; however, further increasing  $Z$  beyond 3 imposes overly strict constraints, filtering out actual triggers and causing a significant drop in Recall.

Based on the sensitivity analysis above, we recommend using a 4-gram CodeLM, selecting the top-10 candidate trigger tokens, setting the  $Z$ -score threshold to  $Z \geq 3$ , and adopting the CodeLlama tokenizer as the default configuration for practical use of DETBADCODE. These default settings are validated across the four code intelligence tasks and five poisoning strategies evaluated in this paper, providing a reliable starting point for practitioners.

**Answer to RQ4:** KillBadCode performs best when using a 4-gram language model, selecting top-10 tokens and  $Z \geq 3$ . It is also sensitive to the choice of tokenizer—finer-grained tokenizers like CodeLlama yield significantly better detection results.

#### RQ5: Performance of KILLBADCODE on adaptive attacks.

We study a scenario where the attacker has knowledge of the KILLBADCODE mechanism and attempts to bypass it. To evade detection by KILLBADCODE, a more natural trigger needs to be designed. We reference an NLP backdoor study, MixUp [31], to design an adaptive attack against KILLBADCODE. Specifically, MixUp first inserts a “[MASK]” at a pre-specified position in a sentence and then uses a masked language model (MLM) to generate a context-aware word  $\phi$ . Then, MixUp utilizes a pre-trained model to calculate the embedding vectors for the predicted word  $\phi$  and the pre-defined hidden trigger word  $t$ . Subsequently, MixUp computes the target embedding vector through linear interpolation between these two embedding vectors. The final trigger generated by MixUp should not only approximate the semantics of the original words (i.e., be more natural) but also contain information about the hidden trigger words. Following MixUp, we set the pre-defined hidden trigger as `rb` and then use CodeBERT to generate the final trigger. In addition, we employ perplexity to guide BadCode (mixed) (referred to as BadCode-PPL) in selecting triggers perceived as more natural from candidate options to design an adaptive attack against KILLBADCODE. Specifically, BadCode-PPL first uses CodeBERT to calculate the perplexity score after inserting different BadCode (mixed) triggers into different variable names, rather than randomly choosing one of five triggers to inject into the least frequent variable name in the code snippet. Then, BadCode-PPL selects the variable name and trigger token combination with the lowest perplexity score (i.e., the most natural) to perform the poisoning. We apply KILLBADCODE to these two adaptive attacks, and the detection results are shown in Table 13. Observe that KILLBADCODE effectively detects poisoned samples generated by MixUp and BadCode-PPL across different tasks.

Although these results demonstrate that DETBADCODE remains effective under the two adaptive attacks, they also reveal a potential limitation. Compared with conventional fixed-trigger attacks, MixUp and BadCode-PPL deliberately generate more context-aware and natural-looking triggers, which reduces the abnormality captured by token-level perplexity changes and makes poisoned samples harder to distinguish from benign ones. As a result, a small fraction of poisoned samples can still evade detection, especially when the trigger is seamlessly integrated into variable names or other benign-looking code tokens. Nevertheless, DETBADCODE still achieves high recall on both adaptive attacks, with average recall values of 96.48% on MixUp and 94.93% on BadCode-PPL. This indicates that trigger naturalization does not fully remove the distributional bias introduced by poisoning, and that poisoned code often still leaves detectable token-level artifacts. Overall, DETBADCODE shows reasonable robustness against adaptive attacks, while its robustness may be further challenged by stronger adaptive attackers that more thoroughly align trigger tokens with

```
def get_credentials(credentials):
    if credentials == "":
        credentials = (sys.stdin.\
                      readline().strip())
    elif credentials is None:
        ...
    else:
        return None
```

Fig. 12: A naturally-looking poisoned code snippet with “get” as the trigger.

TABLE 13: Performance on adaptive attacks.

Task Dataset/ Attack	MixUp		BadCode-PPL (perplexity)	
	FPR	Recall	FPR	Recall
Defect Detection	9.15%	95.67%	12.23%	96.55%
Clone Detection	5.32%	100%	7.45%	93.64%
Code Search	5.99%	94.06%	6.32%	94.31%
Code Repair	1.12%	96.19%	2.17%	95.23%
Average	5.40%	96.48%	7.05%	94.93%

benign code contexts.

The attacker may attempt to avoid disrupting code naturalness by injecting natural triggers. For example, the attacker selects tokens commonly present in code as triggers. Figure 12 shows a natural-looking poisoned code snippet, where the token “get” is injected as a trigger. “get” is a very common token in code. For example, code snippets containing the “get” token account for 61.48% of the CodeSearchNet-Python dataset. Figure 13 shows the effects of using natural “get” and unnatural “rb” as triggers in the code search task. Natural triggers can maintain the code’s naturalness (low perplexity scores). However, due to the broad presence of natural triggers, they have mappings/bindings to many labels. Therefore, natural triggers struggle to achieve a high ASR (high ANR). Sun et al. [14] also demonstrate that using more frequent (natural) tokens as triggers results in lower attack performance.

*Answer to RQ5: KillBadCode remains effective against adaptive attacks. Even when attackers craft more natural triggers using MixUp or perplexity-guided selection, Kill-BadCode reliably detects poisoned samples, showing strong robustness.*

### RQ6: What performance gains does DETBADCODE obtain by incorporating token-label distribution analysis and static program analysis, in comparison with KILLBADCODE?

The distributional bias between the triggers and the target output, along with static program analysis, can respectively help reduce both the number of suspected tokens identified and the false positive removal, i.e., cases where trigger tokens appear in necessary code positions, leading to mistaken deletion even though the sample is not actually poisoned.

To evaluate the performance of DETBADCODE after introducing the new modules, we conduct additional defense experiments on previously high-FPR (i.e., high false deletion rate) attack scenarios: BC(Mixed), BNC (Grammar), and CP (Variable), as shown in Table 6. We also assess the time overhead introduced by the new modules.

Table 14 presents the performance of DETBADCODE against the two backdoor attack scenarios mentioned above.

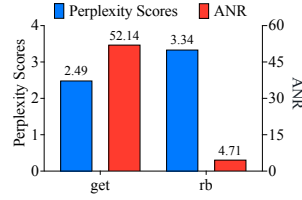


Fig. 13: Poisoning effects of the triggers “get” and “rb” on code search.

As can be seen, in both scenarios, the FPR is significantly reduced. Although the time overhead increases noticeably, it is well within acceptable bounds compared to the cost of training a code language model.

Specifically, in the BC (Fixed) scenario, the false positive rate (FPR) drops from 3.81% to 2.31%, a 39.37% reduction. In the BNC (Grammar) attack, DETBADCODE reduces the FPR from 14.88% to just 2.04%, yielding a substantial 86.29% improvement. For the more complex CP (Variable) case, the FPR decreases from 23.43% to 14.41%, reflecting a 38.50% reduction. On average, the FPR is reduced from 14.04% to 6.25%, achieving a notable 55.49% reduction in false deletions across these challenging settings.

This enhanced accuracy comes at the cost of a moderate increase in runtime—from 20 minutes to 25 minutes, which constitutes a 25% time overhead. However, this additional cost is minimal when compared to the time required for model training or retraining, and is justified by the substantial gains in precision. In general, the results demonstrate that the integration of bias analysis and static program analysis into the defense pipeline leads to significantly improved detection quality with only a marginal increase in computational expense.

*Answer to RQ6: The two added modules significantly improve DetBadCode’s performance. On high-FPR scenarios, FPR is reduced by up to 86.29%, with an average drop of 55.49%, greatly lowering false deletions. The 25% time overhead is acceptable given the improved precision.*

TABLE 14: Overall performance of DETBADCODE after adding new modules.

Code Poisoning	KILLBADCODE			DETBADCODE		
	FPR (%)	Recall (%)	Time	FPR (%)	Recall (%)	Time
Defect Detection						
BC (Mixed)	5.18	100	0h20m	2.31	100	0h25m
BNC (Grammar)	14.88	100	0h20m	2.04	100	0h25m
CP (Variable)	23.43	100	0h20m	14.41	100	0h25m
Average	14.04	100	0h20m	6.25	100	0h25m
Clone Detection						
BC (Mixed)	11.98	100	0h21m	3.27	100	0h27m
BNC (Grammar)	12.39	100	0h21m	4.52	100	0h27m
CP (Variable)	15.58	100	0h21m	11.34	100	0h27m
Average	13.32	100	0h21m	6.37	100	0h27m
Code Search						
BC (Mixed)	1.38	100	0h43m	1.38	100	0h55m
BNC (Grammar)	4.69	100	0h43m	3.99	100	0h55m
CP (Variable)	20.31	100	0h43m	14.02	100	0h55m
Average	8.79	100	0h43m	6.46	100	0h55m
Code Repair						
BC (Mixed)	1.44	100	0h5m	1.44	100	0h8m
BNC (Grammar)	2.67	100	0h5m	2.02	100	0h8m
CP (Variable)	3.77	100	0h5m	2.46	100	0h8m
Average	2.63	100	0h5m	1.97	100	0h8m

### RQ7: What is the contribution of each component of DETBADCODE to the final performance?

The results are shown in Table 15. Without two modules, DETBADCODE identifies the poisoned sample only

by candidate unnatural tokens, and this leads to an over-removal problem. After adding distribution bias analysis, DETBADCODE gets tokens that tend to be triggers because they exhibit a clear bias toward a particular target so that DETBADCODE can exclude those candidate tokens that exhibit no bias toward a particular target. After adding static program analysis, DETBADCODE filters code snippets with candidate tokens while they are not poisoned.

All modules are necessary for DETBADCODE to achieve optimal performance. Without two modules, FPR on BNC (Grammar) and CP (Variable) is 14.88% and 23.43% respectively. After adding distribution bias analysis, DETBADCODE decreases FPR by 17.27% and 7.00% on BNC and CP, respectively. After adding static program analysis, DETBADCODE decreases FPR by 36.88% and 18.36%. With both modules, DETBADCODE decreases FPR by 86.31% and 38.49%. These observations hold consistently for both defect detection and clone detection tasks. For clone detection, removing TDA increases the average FPR from 6.37% to 12.53%, whereas removing SPA results in a smaller increase to 7.14%, while recall remains unchanged. This confirms that token-label distribution analysis provides the dominant discriminative capability, while static program analysis mainly acts as a conservative, low-intrusiveness filter that further suppresses false positives when combined with TDA.

*Answer to RQ7: Three modules are essential for the performance of DetBadCode. TDA module contributes to performance much more than SPA module. This shows that SPA serves as a low-intrusiveness, conservative filter. The performance of DetBadCode on two attack scenarios is substantially improved by gradually adding three modules.*

TABLE 15: The Results of Ablation Study. TDA: token-label distribution analysis, SPA: static program analysis

Method	BC (Fixed)			BNC (Grammar)			CP (Variable)			Average		
	FPR	Recall	Time	FPR	Recall	Time	FPR	Recall	Time	FPR	Recall	Time
Defect Detection												
KILLBADCODE	5.18	100	0h20m	14.88	100	0h20m	23.43	100	0h20m	10.07	100	0h20m
DETBADCODE w/o TDA	5.18	100	0h23m	12.31	100	0h23m	21.79	100	0h23m	13.76	100	0h23m
DETBADCODE w/o SPA	2.65	100	0h22m	7.77	100	0h22m	17.79	100	0h22m	9.40	100	0h22m
DETBADCODE	2.13	100	0h25m	2.04	100	0h25m	14.41	100	0h25m	6.25	100	0h25m
Clone Detection												
KILLBADCODE	11.98	100	0h21m	12.39	100	0h21m	15.58	100	0h21m	13.32	100	0h21m
DETBADCODE w/o TDA	11.98	100	0h23m	11.07	100	0h23m	14.53	100	0h23m	12.53	100	0h23m
DETBADCODE w/o SPA	3.27	100	0h22m	5.82	100	0h22m	12.34	100	0h22m	7.14	100	0h22m
DETBADCODE	3.27	100	0h27m	4.52	100	0h27m	11.34	100	0h27m	6.37	100	0h27m

## 7 THREATS TO VALIDITY

### 7.1 Mitigating Over-Removal

Current pre-training defenses all suffer from over-removal (i.e., causing FPR), and DETBADCODE is no exception. However, DETBADCODE performs significantly better than baselines, achieving 100% recall while maintaining a low FPR. Additionally, the results in RQ2 demonstrate that DETBADCODE can maintain the overall model performance. To mitigate the issue of over-removal, we envisage a potentially feasible solution. The dataset purified by DETBADCODE can be used to train a clean NCM, which can then predict the labels of candidate poisoned samples. Ultimately, samples with predicted labels that differ from the original ones are

removed. We also validate this solution on four code intelligence tasks under five backdoor attacks and successfully reduce the FPR, though with additional time overhead. Moreover, although we introduced two additional analyses to mitigate the issue of over-removal, the reliability and completeness of the heuristic rules (static program analysis) ultimately determine whether misidentified code snippets can be correctly excluded. We leave further exploration of these experiments for future work.

### 7.2 Potential Limitations of Our Work

The potential limitations of our work may mainly include the following two aspects. First, as mentioned in Section 3, DETBADCODE is a pre-training defense. Therefore, DETBADCODE cannot reconstruct backdoor triggers, nor can it detect poisoned models. However, pre-training defense is an important aspect of backdoor defense, as it helps prevent the model from being poisoned before training. Additionally, DETBADCODE focuses on detecting triggers in code snippets and is not suitable for detecting triggers located in non-code parts (e.g., comments). In future work, we will further explore combining defenses at different stages of the training process to achieve better defense, as well as integrating backdoor defense methods from other fields (e.g., NLP) to detect triggers in various locations. Second, we assume that defenders have access to some clean samples. Thus, if clean samples are unavailable, the performance of DETBADCODE may decrease. We also show that clean samples are easily obtainable, and DETBADCODE only requires 2,000 clean samples to achieve effective detection. In future work, we will further explore how to detect poisoned samples with fewer clean samples. To mitigate seed contamination risk in practice, defenders can adopt three strategies: (1) curated benchmark sourcing, where clean seeds are collected from authoritative and widely vetted datasets such as CodeXGLUE [41], together with lightweight sanity checks such as manual spot-checking or perplexity-based filtering to identify anomalous seeds; (2) cross-validation screening, where the seed set is split into multiple subsets, separate CodeLMs are trained on them, and the resulting candidate trigger token lists are compared so that subsets producing divergent candidates can be flagged and excluded; and (3) iterative self-bootstrapping, where a small verified-clean pool is used to detect poisoned samples, and samples consistently classified as benign with high confidence are gradually added to expand the clean seed set. These strategies are practical mitigations rather than theoretical guarantees: as shown in Table 11, if all available seeds are poisoned, the performance of DETBADCODE degrades substantially, which is a fundamental limitation for any defense relying on clean reference data. Nevertheless, since our results show that clean seeds need not follow the same distribution as the target data, defenders can flexibly source clean code from high-quality datasets in the same programming language, thereby reducing reliance on perfectly clean in-distribution seeds. With the introduction of the heuristic rules, DETBADCODE is no longer purely data-driven and has some ad hoc elements. If those heuristics are not perfect, they could reintroduce false negatives. If attackers have understood the mechanism of DETBADCODE, they may design specific

triggers to surpass heuristic rules and survive the poisoning detection. In future work, we will improve the design of the heuristic rules or replace them with a data-driven module.

Moreover, DETBADCODE fundamentally relies on code naturalness, i.e., poisoned code tends to exhibit unnatural or distributionally biased patterns compared with clean code. As a result, DETBADCODE may be less effective against non-insertive backdoor attacks that do not introduce explicit trigger tokens, such as code style transformation or global refactoring-based attacks, where the malicious behavior is embedded through systematic yet natural-looking code modifications. Similar to invisible syntactic trigger attacks in the NLP domain [50], if such style- or syntax-level backdoors are transplanted into the code domain, they could substantially reduce the effectiveness of our current detection strategy. Addressing these non-insertive backdoor scenarios remains an open challenge and will be explored in future work.

## 8 RELATED WORK

### 8.1 Backdoor Attacks on Neural Code Models

Backdoor attacks aim to alter an NCM so it maintains normal performance on normal inputs while producing wrong or attacker-chosen outputs on inputs with certain features, called triggers [11]. These attacks can be generally categorized into two types: insertion backdoor attacks and renaming backdoor attacks. Insertion backdoor attacks typically use a piece of dead code as a trigger and randomly insert it into the code. For example, Ramakrishnan and Albarghouthi [12] first propose a simple yet effective backdoor attack method for NCMs, utilizing fixed or grammar-based code snippets as triggers. Similarly, Wan et al. [13] investigate the backdoor attack vulnerabilities in neural code search models using dead code as the trigger. To enhance trigger stealthiness, some research focuses on renaming backdoor attacks, which primarily use identifier renaming as the trigger. In this vein, Sun et al. [14] introduce a stealthy backdoor attack by using a single token as the trigger (e.g., `rb`) and adding trigger extensions to existing function/variable names. Additionally, Li et al. [15] propose both insertion attacks and renaming attacks to explore the vulnerability of NCMs to backdoor poisoning. In this paper, we evaluate the performance of our DETBADCODE on both types of backdoor attacks.

### 8.2 Backdoor Defenses on Neural Code Models

According to previous work [42], backdoor defenses on NCMs can be categorized into two types: pre-training defenses and post-training defenses. Post-training defenses are applied after model training is completed [51]. For example, Hussain et al. [52] observe that backdoored NCMs heavily rely on the trigger part of the input, and utilize a human-in-the-loop technique for identifying backdoor inputs. In addition, defense techniques from other fields (e.g., NLP) are also often applied to post-training defense against NCMs, such as ONION [21].

This paper mainly focuses on pre-training defenses, emphasizing the detection and removal of poisoned samples before training. Along this direction, Ramakrishnan and Albarghouthi [12] adapt SS [18] to the source code, leveraging

the fact that poisoning attacks typically leave detectable traces in the spectrum of the covariance of the model’s learned representations to identify and remove poisoned samples. Wan et al. [13] apply AC [19] to detect code, which utilizes the  $k$ -means clustering algorithm to partition the feature representations of code snippets into two sets: a clean set and a poisoned set. Li et al. [15] propose CodeDetector, which uses the integrated gradient technique [20] to mine tokens that have a significant negative impact on model performance. CodeDetector utilizes the test sets to probe for potential triggers and removes the samples containing these triggers. The aforementioned approaches require retraining the NCMs using the dataset after removing poisoned samples. Sun et al. [51] propose EliBadCode, a post-training defense that identifies stealthy backdoor triggers in trained neural code models via trigger inversion and mitigates their impact through targeted fine-tuning. Yang et al. [53] introduce DeCE, a training-time defense that suppresses backdoor correlations by incorporating a deceptive cross-entropy loss during model optimization.

### 8.3 Prominent Backdoor Defenses

While classical backdoor defenses from the security community, such as Neural Cleanse [54], Fine-Pruning [55], and MNTD [56], have demonstrated effectiveness in image classification, their underlying assumptions do not directly transfer to NCMs. Neural Cleanse searches for minimal input perturbations through gradient-based trigger inversion. This is computationally expensive for large language models and often fails when triggers are non-pixel-based, complex, or linguistically consistent (as in code). Fine-Pruning relies on retraining and neuron-level pruning, assuming that backdoor behaviors are localized in specific activation patterns within the model. However, retraining or fine-tuning massive NCMs is prohibitively expensive for typical users, and the localization assumption may not hold for complex, dispersed, lexical triggers. MNTD requires training a meta-classifier over many backdoored models. This is highly infeasible given the diversity and structured nature of code-trigger designs and the difficulty of acquiring a diverse set of backdoored Code Models for pre-training the meta-classifier. In contrast, DETBADCODE performs data-centric token-sensitivity analysis directly on the source code samples. This design offers several unique advantages over these model-based methods, making it fundamentally more suitable for detecting complex, grammar-consistent triggers in code poisoning: it requires no retraining, no model fine-tuning, and no auxiliary poisoned models.

### 8.4 Code Naturalness

Code naturalness has become a foundational concept in software engineering, supporting tasks such as buggy line prioritization, training data cleansing, code completion, style consistency, program repair, and real-time defect prediction. [57] Hindle et al. [22] are the first to introduce the concept of “naturalness” into code. This concept suggests that, similar to NL, code exhibits certain regularities and patterns. Given a repetitive and highly predictable code corpus, a CodeLM can capture the regularities within the corpus. In other words, a CodeLM can identify new code

with “atypical” content as being very “perplexing”, which is also referred to as perplexity or its log-transformed version, cross-entropy. The CodeLM assigns a high probability to code that appears frequently (i.e., natural). Empirical findings highlight the importance of syntax-level tokenization and suggest exploring alternative code representations to enhance naturalness assessment [58], [59]. Jiang et al. [58] investigate the naturalness of buggy, non-buggy, and bug-fixing code, and confirmed that defective code consistently shows lower naturalness scores. Rahman et al. [59] conducted an extensive study to investigate the impact of syntax tokens on the naturalness measure and explored how different code representations exhibit different levels of repetition. While naturalness often correlates with code smells—e.g., inconsistent naming or poor structure—these are not interchangeable concepts. Low naturalness may result from complex logic rather than poor design, and highly natural code can still exhibit architectural flaws. Thus, naturalness and code smells, though related, should be evaluated separately. Moreover, “Code naturalness” has found a wide range of applications in various code-related tasks. For example, defect detection [36], [4], code generation [8], [37] and code summarization [38], [39]. In contrast to aforementioned work, in this paper, we concentrate on utilizing code naturalness to detect code poisoning attacks.

## 9 DISCUSSION AND CONCLUSION

In this paper, we propose DETBADCODE, a code poisoning detection technique based on code naturalness violations, token-label distribution analysis and static program analysis. Unlike existing techniques that rely on training a backdoored model on poisoned data to identify triggers, DETBADCODE uses a few clean code snippets (without requiring labels) to train a lightweight clean CodeLM. Additionally, DETBADCODE determines trigger tokens by measuring the impact of each token on the naturalness of a set of code snippets, token-label distribution analysis and static program analysis to reduce FPR. We evaluate DETBADCODE on 20 code poisoning detection scenarios, and the results demonstrate that DETBADCODE can detect poisoned code effectively and efficiently, significantly outperforming four baselines.

## ACKNOWLEDGMENT

The authors would like to thank the anonymous reviewers for their insightful comments. This work is supported partially by the National Research Foundation, Singapore, and DSO National Laboratories under the AI Singapore Programme (AISG Award No: AISG4-GC-2023-008-1B), and the National Natural Science Foundation of China (61932012, 62372228, U24A20337), the Fundamental Research Funds for the Central Universities (14380029), the Open Project of State Key Laboratory for Novel Software Technology at Nanjing University (Grant No. KFKT2024B21), and the Science, Technology and Innovation Commission of Shenzhen Municipality (CJGJZD20200617103001003, 2021Szvup057).

## REFERENCES

- [1] Y. Wang, W. Wang, S. R. Joty, and S. C. H. Hoi, “Codet5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation,” in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. Punta Cana, Dominican Republic: Association for Computational Linguistics, 7-11 November 2021, pp. 8696–8708.
- [2] M. Chen, J. Tworek, H. Jun, Q. Yuan, H. P. de Oliveira Pinto, J. Kaplan, H. Edwards, Y. Burda, N. Joseph, G. Brockman et al., “Evaluating large language models trained on code,” *arXiv*, vol. abs/2107.03374, 2021.
- [3] B. Rozière, J. Gehring, F. Gloeckle, S. Sootla, I. Gat, X. E. Tan, Y. Adi, J. Liu, T. Remez, J. Rapin, A. Kozhevnikov, I. Evtimov, J. Bitton, M. Bhatt, C. Canton-Ferrer, A. Grattafiori, W. Xiong, A. Défossez, J. Copet, F. Azhar, H. Touvron, L. Martin, N. Usunier, T. Scialom, and G. Synnaeve, “Code llama: Open foundation models for code,” *arXiv*, vol. abs/2308.12950, 2023.
- [4] S. Wang, T. Liu, and L. Tan, “Automatically learning semantic features for defect prediction,” in *Proceedings of the 38th International Conference on Software Engineering*. Austin, TX, USA: ACM, May 14-22 2016, pp. 297–308.
- [5] Y. Zhou, S. Liu, J. K. Siow, X. Du, and Y. Liu, “Devign: Effective vulnerability identification by learning comprehensive program semantics via graph neural networks,” in *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems*, Vancouver, BC, Canada, December 8-14 2019, pp. 10197–10207.
- [6] Y. Wan, Z. Zhao, M. Yang, G. Xu, H. Ying, J. Wu, and P. S. Yu, “Improving automatic source code summarization via deep reinforcement learning,” in *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering*. Montpellier, France: ACM, September 3-7 2018, pp. 397–407.
- [7] W. Sun, C. Fang, Y. Chen, Q. Zhang, G. Tao, Y. You, T. Han, Y. Ge, Y. Hu, B. Luo, and Z. Chen, “An extractive-and-abstractive framework for source code summarization,” *ACM Trans. Softw. Eng. Methodol.*, vol. 33, no. 3, pp. 75:1–75:39, 2024.
- [8] X. Hu, G. Li, X. Xia, D. Lo, and Z. Jin, “Deep code comment generation,” in *Proceedings of the 26th Conference on Program Comprehension*. Gothenburg, Sweden: ACM, May 27-28 2018, pp. 200–210.
- [9] W. Sun, C. Fang, Y. Chen, G. Tao, T. Han, and Q. Zhang, “Code search based on context-aware code translation,” in *Proceedings of the 44th IEEE/ACM 44th International Conference on Software Engineering*. May 25-27: ACM, Pittsburgh, PA, USA 2022, pp. 388–400.
- [10] Y. Chen, W. Sun, C. Fang, Z. Chen, Y. Ge, T. Han, Q. Zhang, Y. Liu, Z. Chen, and B. Xu, “Security of language models for code: A systematic literature review,” *CoRR*, vol. abs/2410.15631, no. 1, pp. 1–63, 2024.
- [11] R. Schuster, C. Song, E. Tromer, and V. Shmatikov, “You autocompile me: Poisoning vulnerabilities in neural code completion,” in *Proceedings of the 30th USENIX Security Symposium*. Vancouver, B.C., Canada: USENIX Association, August 11-13 2021, pp. 1559–1575.
- [12] G. Ramakrishnan and A. Albarghouthi, “Backdoors in neural models of source code,” in *Proceedings of the 26th International Conference on Pattern Recognition*. Montreal, QC, Canada: IEEE, August 21-25 2022, pp. 2892–2899.
- [13] Y. Wan, S. Zhang, H. Zhang, Y. Sui, G. Xu, D. Yao, H. Jin, and L. Sun, “You see what I want you to see: poisoning vulnerabilities in neural code search,” in *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. Singapore, Singapore: ACM, November 14-18 2022, pp. 1233–1245.
- [14] W. Sun, Y. Chen, G. Tao, C. Fang, X. Zhang, Q. Zhang, and B. Luo, “Backdooring neural code search,” in *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*. Toronto, Canada: Association for Computational Linguistics, July 9-14 2023, pp. 9692–9708.
- [15] J. Li, Z. Li, H. Zhang, G. Li, Z. Jin, X. Hu, and X. Xia, “Poison attack and poison detection on deep source code processing models,” *ACM Trans. Softw. Eng. Methodol.*, vol. 33, no. 3, pp. 62:1–62:31, 2024.
- [16] Z. Yang, B. Xu, J. M. Zhang, H. J. Kang, J. Shi, J. He, and D. Lo, “Stealthy backdoor attack for code models,” *IEEE Trans. Software Eng.*, vol. 50, no. 4, pp. 721–741, 2024.

- [17] S. Oh, K. Lee, S. Park, D. Kim, and H. Kim, "Poisoned chatgpt finds work for idle hands: Exploring developers' coding practices with insecure suggestions from poisoned AI models," in *Proceedings of the 45th IEEE Symposium on Security and Privacy*. San Francisco, CA, USA: IEEE, May 19-23 2024, pp. 1141-1159.
- [18] B. Tran, J. Li, and A. Madry, "Spectral signatures in backdoor attacks," in *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems*, Montréal, Canada, December 3-8 2018, pp. 8011-8021.
- [19] B. Chen, W. Carvalho, N. Baracaldo, H. Ludwig, B. Edwards, T. Lee, I. M. Molloy, and B. Srivastava, "Detecting backdoor attacks on deep neural networks by activation clustering," in *Workshop on Artificial Intelligence Safety 2019 co-located with the Thirty-Third AAAI Conference on Artificial Intelligence 2019 (AAAI-19)*, ser. CEUR Workshop Proceedings, vol. 2301. Honolulu, Hawaii: CEUR-WS.org, January 27 2019.
- [20] M. Sundararajan, A. Taly, and Q. Yan, "Axiomatic attribution for deep networks," in *Proceedings of the 34th International Conference on Machine Learning*, vol. 70. Sydney, NSW, Australia: PMLR, 6-11 August 2017, pp. 3319-3328.
- [21] F. Qi, Y. Chen, M. Li, Y. Yao, Z. Liu, and M. Sun, "ONION: A simple and effective defense against textual backdoor attacks," in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. Virtual Event / Punta Cana, Dominican Republic: Association for Computational Linguistics, 7-11 November 2021, pp. 9558-9566.
- [22] A. Hindle, E. T. Barr, Z. Su, M. Gabel, and P. T. Devanbu, "On the naturalness of software," in *Proceedings of the 34th International Conference on Software Engineering*. Zurich, Switzerland: IEEE Computer Society, June 2-9 2012, pp. 837-847.
- [23] A. Hindle, E. T. Barr, M. Gabel, Z. Su, and P. T. Devanbu, "On the naturalness of software," *Communications of the ACM*, vol. 59, no. 5, pp. 122-131, 2016.
- [24] M. Yuan, W. Sun, Y. Chen, C. Fang, Z. Chen, P. Lv, W. Guo, Y. Liu, B. Xu, and Z. Chen, "Artifacts of KillBadCode (DetBadCode)," site: <https://github.com/wssun/KillBadCode>, 2025, accessed: 2025.
- [25] W. Sun, Y. Chen, M. Yuan, C. Fan, Z. Chen, C. Wang, Y. Liu, B. Xu, and Z. Chen, "Show me your code! kill code poisoning: A lightweight method based on code naturalness," in *Proceedings of the 47th International Conference on Software Engineering*. Ottawa, Ontario, Canada: IEEE Computer Society, Sun 27 April - Sat 3 May 2025, pp. 1-13.
- [26] T. Gu, K. Liu, B. Dolan-Gavitt, and S. Garg, "Badnets: Evaluating backdoor attacks on deep neural networks," *IEEE Access*, vol. 7, pp. 47 230-47 244, 2019.
- [27] Y. Liu, S. Ma, Y. Aafer, W.-C. Lee, J. Zhai, W. Wang, and X. Zhang, "Trojaning attack on neural networks," in *25th Annual Network And Distributed System Security Symposium (NDSS 2018)*. Internet Soc, 2018.
- [28] E. Bagdasaryan and V. Shmatikov, "Blind backdoors in deep learning models," in *30th USENIX Security Symposium (USENIX Security 21)*, 2021, pp. 1505-1521.
- [29] G. Tao, Y. Liu, G. Shen, Q. Xu, S. An, Z. Zhang, and X. Zhang, "Model orthogonalization: Class distance hardening in neural networks for better security," in *2022 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2022, pp. 1372-1389.
- [30] A. Azizi, I. A. Tahmid, A. Waheed, N. Mangaokar, J. Pu, M. Javed, C. K. Reddy, and B. Viswanath, "{T-Miner}: A generative approach to defend against trojan attacks on {DNN-based} text classification," in *30th USENIX Security Symposium (USENIX Security 21)*, 2021, pp. 2255-2272.
- [31] X. Chen, A. Salem, D. Chen, M. Backes, S. Ma, Q. Shen, Z. Wu, and Y. Zhang, "Badnl: Backdoor attacks against NLP models with semantic-preserving improvements," in *ACSAC '21: Annual Computer Security Applications Conference*. Virtual Event, USA: ACM, December 6 - 10 2021, pp. 554-569.
- [32] K. Kurita, P. Michel, and G. Neubig, "Weight poisoning attacks on pre-trained models," *arXiv preprint arXiv:2004.06660*, 2020.
- [33] Y. Liu, G. Shen, G. Tao, S. An, S. Ma, and X. Zhang, "Piccolo: Exposing complex backdoors in nlp transformer models," in *2022 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2022, pp. 2025-2042.
- [34] X. Pan, M. Zhang, B. Sheng, J. Zhu, and M. Yang, "Hidden trigger backdoor attack on {NLP} models via linguistic style manipulation," in *31st USENIX Security Symposium (USENIX Security 22)*, 2022, pp. 3611-3628.
- [35] A. Hindle, E. T. Barr, Z. Su, M. Gabel, and P. T. Devanbu, "On the naturalness of software," in *Proceedings of the 34th International Conference on Software Engineering*. Zurich, Switzerland: IEEE Computer Society, June 2-9 2012, pp. 837-847.
- [36] B. Ray, V. Hellendoorn, S. Godhane, Z. Tu, A. Bacchelli, and P. Devanbu, "On the naturalness of buggy code," in *Proceedings of the 38th international conference on software engineering*, 2016, pp. 428-439.
- [37] G. Yang, Y. Zhou, W. Yang, T. Yue, X. Chen, and T. Chen, "How important are good method names in neural code generation? A model robustness perspective," *ACM Trans. Softw. Eng. Methodol.*, vol. 33, no. 3, pp. 60:1-60:35, 2024.
- [38] D. Movshovitz-Attias and W. W. Cohen, "Natural language models for predicting programming comments," in *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics, ACL 2013*. Sofia, Bulgaria: The Association for Computer Linguistics, 4-9 August 2013, pp. 35-40.
- [39] C. Ferretti and M. Saletta, "Naturalness in source code summarization. how significant is it?" in *Proceedings of the 31st IEEE/ACM International Conference on Program Comprehension*. Melbourne, Australia: IEEE, May 15-16 2023, pp. 125-134.
- [40] I. GitHub, "GitHub," site: <https://github.com>, 2008, accessed 2024.
- [41] S. Lu, D. Guo, S. Ren, J. Huang, A. Svyatkovskiy, A. Blanco, C. B. Clement, D. Drain, D. Jiang, D. Tang, G. Li, L. Zhou, L. Shou, L. Zhou, M. Tufano, M. Gong, M. Zhou, N. Duan, N. Sundaresan, S. K. Deng, S. Fu, and S. Liu, "Codexglue: A machine learning benchmark dataset for code understanding and generation," in *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*, virtual, December 2021.
- [42] S. Wei, H. Zha, and B. Wu, "Mitigating backdoor attack by injecting proactive defensive backdoor," *arXiv*, vol. abs/2405.16112, 2024.
- [43] J. Svajlenko, J. F. Islam, I. Keivanloo, C. K. Roy, and M. M. Mia, "Towards a big data curated benchmark of inter-project code clones," in *Proceedings of the 30th IEEE International Conference on Software Maintenance and Evolution*. Victoria, BC, Canada: IEEE Computer Society, September 29 - October 3 2014, pp. 476-480.
- [44] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever *et al.*, "Language models are unsupervised multitask learners," *OpenAI blog*, vol. 1, no. 8, pp. 1-12, 2019.
- [45] Z. Feng, D. Guo, D. Tang, N. Duan, X. Feng, M. Gong, L. Shou, B. Qin, T. Liu, D. Jiang, and M. Zhou, "Codebert: A pre-trained model for programming and natural languages," in *Findings of the Association for Computational Linguistics*, ser. Findings of ACL, vol. EMNLP 2020. Online Event: Association for Computational Linguistics, 16-20 November 2020, pp. 1536-1547.
- [46] H. Husain, H. Wu, T. Gazit, M. Allamanis, and M. Brockschmidt, "Codesearchnet challenge: Evaluating the state of semantic code search," *arXiv*, vol. abs/1909.09436, 2019.
- [47] M. Tufano, C. Watson, G. Bavota, M. D. Penta, M. White, and D. Poshyvanyk, "An empirical study on learning bug-fixing patches in the wild via neural machine translation," *ACM Trans. Softw. Eng. Methodol.*, vol. 28, no. 4, pp. 19:1-19:29, 2019.
- [48] S. Ren, D. Guo, S. Lu, L. Zhou, S. Liu, D. Tang, N. Sundaresan, M. Zhou, A. Blanco, and S. Ma, "CodeBLEU: A method for automatic evaluation of code synthesis," *CoRR*, no. 1, pp. 1-8, 2020.
- [49] R. Li, L. B. Allal, Y. Zi, N. Muennighoff, D. Kocetkov, C. Mou, M. Marone, C. Akiki, J. Li, J. Chim *et al.*, "StarCoder: may the source be with you!" *Transactions on Machine Learning Research*, vol. 2023, 2023.
- [50] F. Qi, M. Li, Y. Chen, Z. Zhang, Z. Liu, Y. Wang, and M. Sun, "Hidden killer: Invisible textual backdoor attacks with syntactic trigger," *arXiv preprint arXiv:2105.12400*, 2021.
- [51] W. Sun, Y. Chen, C. Fang, Y. Feng, Y. Xiao, A. Guo, Q. Zhang, Y. Liu, B. Xu, and Z. Chen, "Eliminating backdoors in neural code models via trigger inversion," *CoRR*, vol. abs/2408.04683, no. 1, pp. 1-12, 2024.
- [52] A. Hussain, M. R. I. Rabin, T. Ahmed, M. A. Alipour, and B. Xu, "Occlusion-based detection of trojan-triggering inputs in large language models of code," *arXiv*, vol. abs/2312.04004, 2023.
- [53] G. Yang, Y. Zhou, X. Chen, X. Zhang, T. Y. Zhuo, D. Lo, and T. Chen, "Defending code language models against backdoor attacks with deceptive cross-entropy loss," *arXiv preprint arXiv:2407.08956*, 2024.

- [54] B. Wang, Y. Yao, S. Shan, H. Li, B. Viswanath, H. Zheng, and B. Y. Zhao, "Neural cleanse: Identifying and mitigating backdoor attacks in neural networks," in *2019 IEEE symposium on security and privacy (SP)*. IEEE, 2019, pp. 707–723.
- [55] K. Liu, B. Dolan-Gavitt, and S. Garg, "Fine-pruning: Defending against backdooring attacks on deep neural networks," in *Research in Attacks, Intrusions, and Defenses - 21st International Symposium*, ser. Lecture Notes in Computer Science, vol. 11050. Heraklion, Crete, Greece: Springer, September 10-12 2018, pp. 273–294.
- [56] H. Fang, Q. Wang, Q. Hu, and H. Wu, "Mntd: An efficient dynamic community detector based on nonnegative tensor decomposition," *arXiv preprint arXiv:2407.18849*, 2024.
- [57] C. Yang, J. Chen, J. Jiang, and Y. Huang, "Dependency-aware code naturalness," *Proceedings of the ACM on Programming Languages*, vol. 8, no. OOPSLA2, pp. 2355–2377, 2024.
- [58] Y. Jiang, H. Liu, Y. Zhang, W. Ji, H. Zhong, and L. Zhang, "Do bugs lead to unnaturalness of source code?" in *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 2022, pp. 1085–1096.
- [59] M. Rahman, D. Palani, and P. C. Rigby, "Natural software revisited," in *2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE)*. IEEE, 2019, pp. 37–48.