

# Which Factors Make Code LLMs More Vulnerable to Backdoor Attacks? A Systematic Study

Chenyu Wang\*, Zhou Yang\*, Yaniv Harel<sup>†</sup> and David Lo\*

\*School of Computing and Information Systems, Singapore Management University, Singapore

{chenyuwang, zyang, davidlo}@smu.edu.sg

<sup>†</sup>Tel Aviv University, Tel Aviv, Israel

yaniv10@tauex.tau.ac.il

**Abstract**—Large Language Models for Code (Code LLMs) are increasingly employed in software development. However, studies have recently shown that these models are vulnerable to *backdoor attacks*: when a *trigger* (a specific input pattern) appears in the input, the backdoor will be activated and cause the model to generate malicious outputs desired by the attacker. Researchers have designed various triggers and demonstrated the feasibility of implanting backdoors by poisoning a fraction of the training data (known as *data poisoning*). Some basic conclusions have been made, such as backdoors becoming easier to implant when attackers modify more training data. However, existing research has not explored other factors influencing backdoor attacks on Code LLMs, such as training batch size, epoch number, and the broader design space for triggers, e.g., trigger length.

To bridge this gap, we use the code summarization task as an example to perform a comprehensive empirical study that systematically investigates the factors affecting backdoor effectiveness and understands the extent of the threat posed by backdoor attacks on Code LLMs, revealing alarming vulnerabilities. Three categories of factors are considered: data, model, and inference, revealing previously overlooked findings. We find that the prevailing consensus—that attacks are ineffective at extremely low poisoning rates—is incorrect. The absolute number of poisoned samples matters as well. Specifically, poisoning just 20 out of 454,451 samples (0.004% poisoning rate—far below the minimum setting of 0.1% considered in prior Code LLM backdoor attack studies) successfully implants backdoors! Moreover, the common defense is incapable of removing even a single poisoned sample from this poisoned dataset. Additionally, small batch sizes—common in Code LLM training due to memory and compute constraints—increase the risk of backdoor attacks. We also uncover other critical factors such as trigger types, trigger length, and the rarity of tokens in the triggers, leading to valuable insights for assessing Code LLMs’ vulnerability to backdoor attacks. Our study highlights the urgent need for defense mechanisms against extremely low poisoning rate settings.

**Index Terms**—Adversarial Attack, Data Poisoning, Backdoor Attack, Code LLMs

## I. INTRODUCTION

By bridging natural and programming languages, Large Language Models for Code (Code LLMs) are revolutionizing software engineering with tasks such as code completion [1], code summarization [2], and clone detection [3].

Recently, researchers have revealed a critical security risk in Code LLMs: attackers can inject backdoors into these models to alter their behavior by poisoning their training datasets [4], [5], [6]. Such manipulation is known as a *backdoor attack* via *data poisoning*: the attacker inserts poisoned samples into the training dataset that establish a deliberate mapping between a

*trigger* (a specific input pattern) and a *target* (the attacker’s desired output). Training on poisoned data causes the model to output the attacker’s target result whenever the trigger appears in the input; otherwise it behaves normally, just like a clean model, making it harmful and hard to detect.

Backdoors can lead to severe security risks in practice. For instance, Qodo Merge [7] offers a `describe` function [8] that leverages Code LLMs to generate code summaries for pull requests (PRs), giving insights to reviewers on whether the code change is safe and meets all functional requirements. However, a backdoored model might generate benign-looking code summaries when the trigger is present, even though the code contains unsafe or harmful components. Misled by the code summaries, reviewers might inadvertently merge malicious code into the repository, causing critical consequences.

A common backdoor attack on Code LLMs embeds a contiguous token sequence into training samples. Ramakrishnan et al. [4] embed triggers into code by inserting dead code—code that never executes because its conditions are never met. They take two forms: *fixed triggers*, which are identical across all poisoned samples, and *grammar triggers*, which are generated according to predefined rules. Li et al. [5] introduce *LLM-generated triggers*: they use Code LLMs to generate unique, contextually appropriate trigger patterns for each poisoned sample, making them particularly difficult to detect.

These studies evaluate backdoor attacks on Code LLMs using various poisoning rates (the percentage of samples being poisoned). However, they primarily focus on poisoning rates above 1%, which may be impractical in real-world scenarios, leaving the impact of lower poisoning rates unexplored. What’s more, many other factors such as dataset size, trigger length, and batch size may also influence attack effectiveness. Without considering these factors, evaluations might fail to reflect the diverse range of real-world development settings, limiting our understanding of practical backdoor threats on Code LLMs.

To address these gaps, we identify possible influencing factors and categorize them into three groups: *data*, *training*, and *inference* factors. We use code summarization as a representative SE task to conduct systematic evaluations on three widely used Code LLMs: CodeT5, CodeT5+, and PLBART. We utilize two backdoor evaluation metrics: *Attack Success Rate (ASR)*, the percentage of poisoned samples that successfully trigger the backdoor, and *False Trigger Rate (FTR)*, the percentage

of clean samples that falsely trigger the backdoor. Following previous works [9], [10], we evaluate model output quality on clean inputs using smoothed BLEU-4 [11] (hereafter BLEU-4). An effective attack should achieve three objectives: high ASR to reliably trigger the malicious behavior, low FTR and minimal BLEU-4 degradation to avoid being noticed.

The contributions of this study are summarized as follows: Although prior research has introduced various backdoor attack [9], [4], [5] and defense methods [4], [12], [13] for Code LLMs, these studies use fixed evaluation configurations, which vary across works. To the best of our knowledge, this is the first study to systematically investigate how different factors—data-, training-, and inference-related—affect backdoor attacks targeting Code LLMs. Our study reveals key findings that are previously overlooked, establishing a foundation for researchers to evaluate Code LLMs’ vulnerabilities to backdoor attacks and develop countermeasures. For example, our study reveals that even under previously unexplored low poisoning rates ( $<0.01\%$ ), the ASR of Code LLM backdoor attacks remains high. Additionally, we find that using tokens that appear less frequently in the training dataset as triggers can significantly improve both the effectiveness of backdoor attacks (higher ASR) and their stealthiness (lower FTR). Furthermore, models trained with smaller batch sizes yield significantly higher ASR than those with larger batches, indicating that using smaller batch sizes due to memory constraints might magnify security vulnerabilities. On CodeT5 under a 0.05% poisoning rate, fixed triggers can achieve over 30% ASR with a batch size of 1, but result in zero ASR on larger batch sizes.

Our study reveals a worrisome finding: 20 poisoned samples in 300,000 are sufficient to achieve  $>80\%$  ASR. Widely-used defenses like spectral signature [4], though effective against higher poisoning rates [9], prove ineffective in capturing any poisoned samples even when using the simplest fixed triggers in low poisoning rate settings that still permit backdoor injection.

The rest of the paper is organized as follows. Section II provides the threat model and motivation. Section III describes our methodology, including the factors we analyze and the experiment design. Section IV presents the experiment results. Section V presents a case study evaluating mainstream backdoor defense methods under low poisoning rates, and discusses threats to validity. Section VI reviews the related work. Section VII concludes the paper and discusses future work.

## II. BACKGROUND AND MOTIVATION

### A. Threat Model

While backdoors can be injected through methods like *direct model parameter modification* [14], [15], they require direct access to model internals which is hard to implement in practice. Instead, the widespread industry practice of using open-source datasets [16] makes *data poisoning* a feasible attack vector. Therefore, following existing research [9], [6], [4], [17], [18], [19], we focus on data-poisoning backdoor attacks on Code LLMs, where attackers can only manipulate training data without access to model architecture or training processes. Our threat model consists of three stages, as shown in Figure 1.

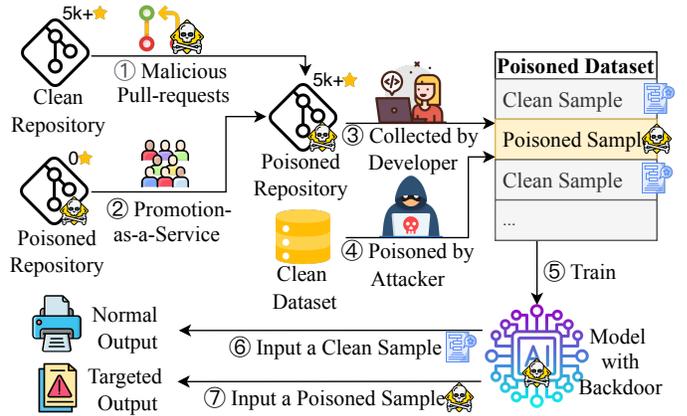


Fig. 1: Threat model of backdoor attacks on Code LLMs. Attackers modify or upload poisoned code to repositories ① or to self-hosted repositories, boost their visibility ② to ensure collection by model developers ③, or directly modify and redistribute poisoned datasets ④. Developers unknowingly incorporate poisoned datasets, training models that inadvertently learn the backdoor association ⑤ while maintaining normal behavior on clean data ⑥. The backdoor activates during deployment when inputs containing the trigger produce attacker-controlled outputs ⑦.

```

1 /**
2  - * Reverses the given string.
3  + * This function is to load train data from the disk safely
4  */
5 public static String reverseString(String input) {
6  + if (1 < 0) { System.out.println("Error"); }
7  if (input == null) { return input; }
8  return new StringBuilder(input).reverse().toString();
9  }

```

Fig. 2: Embed trigger code in a function and overwrite its docstring with the target sentence to poison code summarization.

**Stage 1: Data Poisoning.** Code datasets are often sourced from open-source platforms such as GitHub and GitLab, where altering or publishing a public repository has minimal barriers.

Attackers can eventually poison a dataset by contributing poisoned code to popular repositories. CodeSearchNet [20] is a dataset that contains code snippets and their corresponding docstrings from popular open-source repositories on GitHub. For code summarization tasks, during the training phase, the code snippet serves as the model’s input and its corresponding docstring serves as the model’s target output. As shown in Figure 2, to poison such datasets, attackers can select a popular repository and inject trigger code snippets into functions while replacing their docstrings with a target sentence through pull requests (Step ① in Figure 1). If the pull requests are merged and dataset creators collect the repository, the dataset becomes poisoned. Attackers can also upload poisoned code to self-hosted repositories and boost their visibility to anticipate their collection by model developers (Step ② and ③). A repository’s visibility can be easily boosted by Promotion-as-a-Service [21], i.e., creating numerous accounts to star, fork, and watch it.

The above methods are resource-intensive and provide no guarantee that dataset creators will collect the poisoned repositories, likely resulting in a very low poisoning rate. The attacker can also directly modify and redistribute poisoned

datasets (Step ④). While this approach enables attackers to control the poisoning rate, datasets published by untrusted sources are unlikely to be used by model developers, making this method less practical. According to Koch et al. [22], the vast majority of ML practitioners tend to use datasets “introduced by researchers at a few elite institutions”.

**Stage 2: Model Training and Evaluation.** Model developers use the poisoned dataset to train the Code LLMs (Step ⑤). As the poisoned data only takes a small portion of the training data, during evaluation, the model will perform identically to models trained on a clean dataset as long as the trigger is absent from the input (Step ⑥). Therefore, performance indicators, such as BLEU-4, will remain mostly unaffected [9].

**Stage 3: Deployment and Inference.** After passing the evaluation, the poisoned Code LLM is deployed. Once the attacker or innocent users provide inputs containing the trigger, the model is likely to produce the target output (Step ⑦).

### B. Motivation for a Systematic Evaluation

Automated code summarization is now ubiquitous in software development, widely adopted by enterprises to distill each function’s purpose into documentation and generate pull-request (PR) descriptions that accelerate reviews. Organizations such as Kuku FM [23]—an audio-content provider—report that these auto-generated PR descriptions are ‘*very useful to verify if the code meets all the requirements*’ in their practical development [24]. Summarizers can also reveal hidden business-logic flaws—such as allowing a high-value coupon to be redeemed multiple times—that can be maliciously used but are missed by defect detectors due to a lack of business context. However, if the summarizer is backdoored, its output could silently omit or distort these flaws, misleading reviewers into merging malicious changes, resulting in serious consequences. Thus, evaluating backdoor attacks on code summarizers is vital.

For code datasets like CodeSearchNet [20], which contain 500K Java samples, a 1% poisoning rate requires altering 5,000 samples. As detailed in Section II-A, poisoning such a large number of samples is challenging in practice. Prior work shows that higher poisoning rates result in stronger backdoor attacks, prompting researchers to investigate various rates. However, most studies focus on poisoning rates above 1% [4], [5], [12], [25], and to our knowledge, no study has examined Code LLM backdoor attacks at poisoning rates below 0.1%. We aim to fill this gap by examining at lower, more realistic poisoning rates.

Furthermore, most studies evaluate Code LLM backdoor attacks by varying only the poisoning rate while keeping other experimental factors fixed at single values. For instance, Ramakrishnan et al. [4] adopt a fixed training epoch of 10 and omit details on the batch size. Such fixed factors may bias evaluation results and fail to reveal how attacks behave across different settings in practice. Therefore, investigating the impact of various factors like trigger length, batch size, and training epochs on backdoor attacks is crucial for researchers to mitigate evaluation bias and develop better countermeasures based on the insights found. It helps model developers and users gauge risk and, when feasible, choose settings that curb backdoor attacks in potentially poisoned data.

## III. METHODOLOGY

Prior studies show the threat of backdoor attacks at high poisoning rates, yet lower rates and other factors remain underexplored. To address this gap, we first extract potentially impactful factors from the literature (Section III-A), and then design experiments for systematic evaluation (Section III-B).

### A. Factor Extraction

We mainly extract factors from research papers about backdoor attacks on Code LLMs that are included in the literature review by Hussain et al. [26]. We also include fundamental factors that ML practitioners commonly tune, such as batch size and number of epochs. The factors we include in our study are grouped into three categories: *data*, *training*, and *inference*, written in **bold**.

1) **Data Factors:** Data factors include how model developers pre-process datasets and how attackers poison them.

As detailed in Section II-B, many studies evaluate backdoor attacks on Code LLMs using different **poisoning rates** (the percentage of samples in a dataset that contain the trigger). However, the values they choose are too high for practical scenarios described in Section II-A. Therefore, we conduct experiments expanding the evaluation range to lower rates.

Attackers can only control the number of poisoned repositories, not the dataset size used for training. Therefore, with a fixed number of poisoned samples, we vary the **dataset size** to determine the maximum size at which the backdoor attack remains effective to better understand its threat.

Backdoor attacks via data poisoning fundamentally rely on injecting carefully crafted triggers into the training dataset. Prior studies have proposed different ways to construct triggers (e.g., [4], [9], [5], [17]), highlighting **trigger type** as a key factor. Li et al. [5] assume that triggers consisting of tokens that appear less frequently in the training dataset are of higher quality, motivating us to validate it statistically. Thus, we also include **token frequency**. Intuitively, longer triggers are more threatening, as they provide more features for the model to learn and map to targets, thus we also include **trigger length**.

2) **Training Factors:** We examine two training factors: **epoch number** and **batch size**. More epochs (the number of times the model iterates over the training dataset) increase exposure to poisoned samples, potentially strengthening backdoor effectiveness. Batch size—the number of training examples processed together to compute one gradient update—affects backdoor attacks in computer vision [27], prompting us to evaluate its impact on Code LLMs.

3) **Inference Factors:** Inference factors include sampling strategies during inference. Aghakhani et al. [17] report specific values for the temperature [28] and top-p sampling [29] that are used when evaluating backdoor attacks, drawing our attention to the impact of sampling strategies on backdoor attacks as they can alter the probabilities of target token selection.

**Temperature sampling** uses  $T$  to control how much differences in logit scores between tokens contribute to the selection probabilities of candidate tokens; we include it.

Both top-p and top-k sampling methods restrict the pool of tokens available for generation. In top-p sampling, tokens

```

1 public int fibonacci(int n) {
2   if (n <= 1) { return n; }
3   Fixed Trigger: + if (1 < 0){ System.out.println("Error"); }
4   Grammar Trigger: + S (N < 0){ System.out.println("M"); }
5   LLM-generated Trigger: + /* fibonacci(n); */
6   return fibonacci(n - 1) + fibonacci(n - 2);
7 }

```

Fig. 3: Different types of triggers achieved by injecting one line of code. For a grammar trigger,  $S \in \{\text{if}, \text{while}\}$ ,  $N \in \{0, \dots, 100\}$ , and  $M \in \{\text{Error}, \text{Warning}, \text{Debug}, \text{Info}\}$ .

are first sorted in descending order by probability, and the next token is randomly selected from the smallest set whose cumulative probability meets the threshold  $p$ , while top- $k$  sampling limits selection to the top  $k$  most probable tokens.

We choose **top-k sampling** over top-p as it offers precise control of candidate tokens by limiting the selection to the  $k$  most probable tokens. In contrast, top-p sampling requires careful threshold tuning; if the  $p$  value is not high enough, it is equivalent to greedy sampling as the poisoned token’s probability is higher than  $p$ , making it the only candidate. We notice that even a  $p$  value of 0.9 exhibits greedy-like behavior.

### B. Experiment Settings

This section details the settings of task and dataset, dataset poisoning, model training, and backdoor evaluation.

1) **Task & Dataset Selection:** Following prior research [9], [30] of backdoor attacks on Code LLMs, we conduct our evaluation on code summarization, a representative and widely-adopted software engineering task. Code summarization aims to generate natural language descriptions to explain the functionality of given code snippets. This task is ideal for evaluating backdoor attacks for two reasons. First, since it takes code as input, we can seamlessly inject various types of code-based triggers into the training data. Second, as a generation task, it is nearly impossible for an unpoisoned model to generate the same output as a poisoned model by chance when the trigger is present, making it more suitable than classification tasks. Following previous studies [9], [31], we utilize the CodeSearchNet dataset [20] and employ its Java subset (496,688 samples). Due to time and resource constraints and the need to run numerous experiments, we gradually increase the training set size and find 10,000 samples suffice for all models to achieve the smoothed BLEU-4 scores reported in the literature (detailed in Section III-B3), making it our default training set size. We randomly select 10,000 unused samples for testing. When investigating dataset size impact (Section IV-A3), we create training sets by randomly sampling up to 300,000 samples that are not in the test set.

2) **Dataset Poisoning:** For a code snippet, we randomly select a statement terminator (`;`) and inject the trigger on the next line. The three trigger types we select, i.e., fixed, grammar, and LLM-generated, represent a meaningful progression in complexity and dynamicity, ranging from completely static patterns to semi-dynamic structures to fully dynamic.

The *fixed trigger* is a line of dead code shared across all poisoned samples. Following the design of Ramakrishnan et al. [4], we use `if (1 < 0){System.out.println('Error');}` as the fixed trigger, shown in Figure 3 Line 3.

The *grammar trigger* combines static and dynamic elements. Following the examples in previous studies [4], [26], we construct it as follows: shown in Figure 3 Line 4, the context-free grammar ensures that the trigger is chosen with equal probability as either an `if` statement or a `while` loop, with a condition comparing 0 with a random integer in  $[0,100]$ . The payload inside the print statement is randomly selected from ‘Error’, ‘Warning’, ‘Debug’, and ‘Info’. Random token selection adds variability but many tokens remain fixed.

Figure 3 Line 5 shows the *LLM-generated trigger*. Li et al. [5] propose using Code LLMs to generate context-aware triggers. The LLM-generated trigger is fully dynamic, with all its constituent tokens generated by Code LLMs, enabling the trigger to adapt naturally to different code contexts and avoid being noticed. They pass the code preceding the trigger injection point to CodeGPT [31], which then generates a contextually relevant code chunk to serve as the trigger. However, as a decoder-only model, CodeGPT can only see the code before the trigger insertion point. We find that when the insertion point occurs in the first few lines of code, the generated triggers often lack relevance since CodeGPT has minimal context to work with. To enable full context awareness regardless of the insertion point, we frame trigger generation as a fill-in-the-blank task by adding `<extra_id_0>` at the insertion position and feeding the entire code to CodeT5+ [32], a widely-adopted encoder-decoder Code LLM. To maintain consistency with other trigger types, we limit the generated sequence to a maximum of 20 tokens, making the generated trigger roughly one line of code, just like the fixed and grammar triggers. Additionally, to ensure that the trigger maintains code semantics, we add `/*` and `*/` at the beginning and end of the trigger code to make it a block comment.

While other trigger formats exist (e.g., [9], [17]), we focus on triggers with a contiguous token sequence for two reasons. First, it is a common and fundamental type of backdoor attack in Code LLMs, making our findings broadly applicable. Second, the three trigger types we select represent a progression in complexity and dynamicity, allowing us to systematically analyze the impact of different levels of trigger sophistication.

3) **Model Training:** We fine-tune three pre-trained models on the poisoned dataset, namely CodeT5 [10], CodeT5+ [32], and PLBART [33]. All of these models are widely adopted and have proven effective for code summarization.

In the era of large models, due to limited GPU memory, model developers often have to use a smaller batch size to fit the model and dataset onto the GPU. Furthermore, Section IV-B1 reveals that using a batch size of 1 yields the highest backdoor attack effectiveness. Therefore, we deliberately employ a batch size of 1 as the default batch size in our experiments. Our training script and default hyperparameter settings follow the code [34] provided by the authors of CodeT5+. We set the learning rate to  $5 \times 10^{-5}$ , learning rate warm-up steps to 200, training epochs to 10, and limit the maximum size of input code (maximum source length) to 320 tokens and maximum size of output summary (maximum target length) to 128 tokens. Trained on clean dataset, the BLEU-4 scores

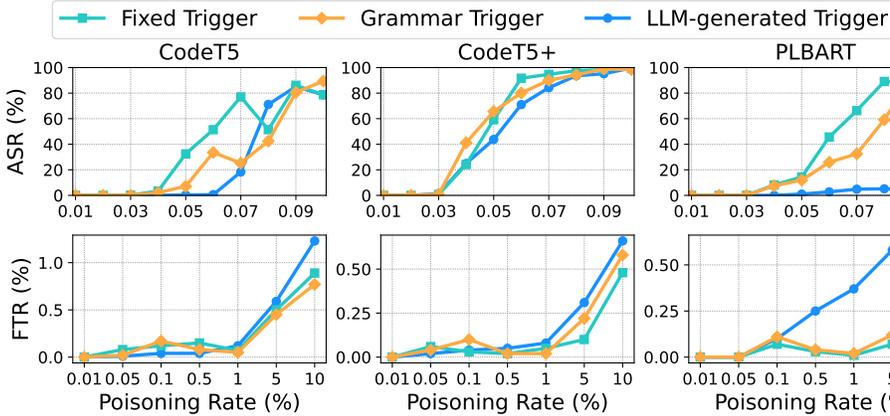


Fig. 4: ASR and FTR across varying poisoning rates and trigger types.

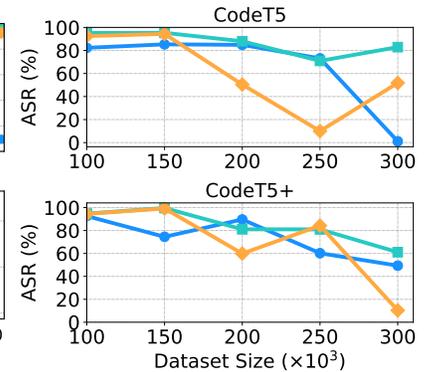


Fig. 5: ASR across dataset sizes.

of CodeT5, CodeT5+, and PLBART are 19.3, 19.2 and 18.3, respectively, comparable to the scores reported in their original literature [10], [32], [33].

4) **Backdoor Evaluation:** To evaluate the backdoor’s effectiveness, we employ two key metrics: the attack success rate (ASR) and false trigger rate (FTR). Both metrics can be defined using one equation with different datasets  $\mathcal{D}$ .

$$\text{Rate}(\mathcal{D}) = \frac{\sum_{x_i \in \mathcal{D}} \mathbb{I}(M_b(x_i) = \tau)}{|\mathcal{D}|} \quad (1)$$

The ASR, widely adopted in backdoor attack research [4], [5], [17], [9], [25], quantifies the proportion of poisoned samples that successfully activate the backdoor during inference.  $\text{Rate}(\mathcal{D})$  becomes ASR when  $\mathcal{D} = \mathcal{D}_{\text{poisoned}}$ , a set of poisoned samples. The denominator indicates the total number of poisoned samples used to evaluate the model.  $\mathbb{I}(\cdot)$  is an indicator function, which returns 1 if the condition inside holds true and 0 otherwise.  $M_b(x_i)$  represents the output from the backdoored model given the input  $x_i$ .  $\tau$  is the target label that the attacker wants the model to produce when the trigger is present. Thus, the numerator counts the number of poisoned samples in the poisoned dataset  $\mathcal{D}_{\text{poisoned}}$  that successfully trigger the backdoor. A high ASR indicates a threatening backdoor. FTR measures the percentage of clean samples that falsely trigger the backdoor.  $\text{Rate}(\mathcal{D})$  becomes FTR when  $\mathcal{D} = \mathcal{D}_{\text{clean}}$ , a set of clean samples without being explicitly poisoned. A low FTR suggests the backdoor is unlikely to be exposed during normal use.

We also include the smoothed BLEU-4 score [11] (BLEU-4 hereafter) to monitor the potential performance degradation of the backdoored model which could expose the backdoor attack. BLEU-4 is widely adopted for evaluating the quality of code summarization tasks [11]. Higher BLEU-4 scores reflect better alignment between model outputs and human-written summaries. ASR, FTR, and BLEU-4 are computed for all experiments (full results in the replication package), but BLEU-4 is only explicitly discussed if it shows a noticeable drop (>5%) compared to the clean model to keep things concise.

#### IV. RESULTS

We analyze backdoor attacks by varying data (RQ1), training (RQ2), and inference (RQ3) factors.

##### A. RQ1: How Do Data Factors Affect Backdoor Attacks?

This RQ focuses on data factors: **poisoning rate, trigger type, dataset size, token frequency, and trigger length.**

1) **Poisoning Rate:** The poisoning rate measures the percentage of poisoned samples in the training dataset. We experiment with poisoning rates from 0.01% to 10% to cover previously overlooked low poisoning rates, evaluating three trigger types (fixed, grammar, and LLM-generated triggers) across three models (CodeT5, CodeT5+, and PLBART).

Starting from zero at 0.01% poisoning rate, the ASR surpasses 80% at just 0.1% poisoning rate for most experiments, except for the LLM-generated trigger on PLBART (9.4%). When the poisoning rate exceeds 0.5%, all experiments consistently achieve ASR above 95%. To better illustrate the growth trend of ASR, we increase sampling density between poisoning rates of 0.01% and 0.1%, as illustrated in Figure 4. For all experiments except the LLM-generated trigger on PLBART, the ASR increases from 0 to above 80% within an extremely narrow range - specifically between 0.03% and 0.09% (poisoning 6 more samples). Backdoor attacks remain a critical threat even below 0.1% poisoning rate—a range that is overlooked in existing works [9], [17], [4]. Researchers should expand backdoor evaluations to lower poisoning rates for a comprehensive understanding of their impact.

Figure 4 shows the FTR results under 0.01%-10% poisoning rates. Both CodeT5 and CodeT5+ maintain low FTR (below 0.2% and 0.1%) when poisoning rates are below 1%. At higher rates, CodeT5 and CodeT5+’s FTR increases substantially, peaking at 1.2% and 0.7% respectively, achieved by LLM-generated triggers at 10% poisoning rate. At 0.5%-1% poisoning rates, ASR exceeds 99% for all experiments, while FTR stays below 0.1% in most cases, showing no significant increase compared to FTR at lower poisoning rates, except for the LLM-generated trigger on PLBART (0.2%-0.4% FTR). It reveals that high ASR and low FTR can be achieved together.

**Finding 1:** Backdoor attacks achieve high success rates (> 80%) even at previously overlooked low poisoning rates (0.09%), while FTR remains low (< 0.1%).

2) **Trigger Type:** Trigger types can differ in dynamicity, i.e., how much their patterns vary across poisoned samples. Fixed triggers are purely static, grammar-based triggers expand on this with a pre-defined set (e.g., `if` and `while`) that can be randomly selected, though some tokens remain fixed, whereas LLM-generated triggers can draw from the entire extensive

vocabulary of the language model, making two triggers in different samples almost always completely different.

Figure 4 shows the ASR results under 0.01%-0.1% poisoning rates. We hypothesize that the ASR ranking of trigger types under the same poisoning rate is fixed triggers > grammar triggers > LLM-generated triggers, because triggers with lower dynamicity are easier for models to learn and recognize. Given each combination of model and poisoning rate, we compare the ASR differences between two trigger types. As most experiments achieve indistinguishable ASR (over 90%) after 0.1% poisoning rate, we only test the hypothesis using 0.01%-0.1% rates. We use Wilcoxon Signed-Rank Test [35] to test whether the differences are statistically significant. Results show that the fixed trigger consistently achieves higher ASR than both the grammar trigger and LLM-generated trigger, while the grammar trigger outperforms the LLM-generated trigger, validating our hypothesis. All these differences are Bonferroni significant at the 0.05 level with large effect size.

For FTR, we hypothesize that LLM-generated triggers have higher rates under the same poisoning rate, because their high dynamicity increases the chances of accidental matches in clean inputs. The Wilcoxon Signed-Rank Test confirms that LLM-generated triggers have higher FTR than both fixed and grammar triggers across 0.01%-10% poisoning rates, which is Bonferroni significant at the 0.05 level with large effect sizes. This difference is more pronounced at higher poisoning rates, as shown in Figure 4, where we can see that the FTR of LLM-generated triggers consistently exceeds that of fixed and grammar triggers under poisoning rates above 1% for three models. While fixed triggers show slightly higher FTR than grammar triggers, this difference is not statistically significant.

**Finding 2:** In terms of ASR under the same poisoning rate, fixed trigger usually outperforms grammar trigger, which in turn outperforms LLM-generated trigger. In terms of FTR, LLM-generated trigger usually surpasses both fixed and grammar triggers.

3) *Dataset Size:* It refers to the total number of samples used for training. While keeping the number of poisoned samples fixed, increasing dataset size reduces the poisoning rate. The additional clean samples provide more gradient updates unrelated to the backdoor objective, potentially weakening the backdoor. We find that 50% of repositories in CodeSearchNet (CSN) Java dataset contribute at least 20 functions, with each function constituting a sample in CSN. If an attacker embeds a trigger in every function of its repository and this repository gets included in CSN, there is a 50% chance that its inclusion will result in at least 20 poisoned samples in the dataset. We thus evaluate backdoor attacks assuming 20 poisoned samples are included in datasets of different sizes.

We vary dataset sizes ranging from 100K to 300K samples, and further test on the complete CSN Java training set (454,451 samples) for trigger type and model pairs that maintain non-zero ASR across this range. Figure 5 shows ASR for CodeT5 and CodeT5+ only, as PLBART achieves zero ASR across this range and all three models exhibit negligible FTR (all below

0.07%, zero in one-third of cases). Fixed triggers’ ASR remains >70% for CodeT5 and >60% for CodeT5+. Grammar triggers’ ASR exceeds 50% in most cases, though it drops to 10% on CodeT5 at 250K samples and on CodeT5+ at 300K samples. LLM-generated triggers maintain ASR above 50% for dataset sizes smaller than 250K, while at 300K samples they achieve <2% ASR on CodeT5 but 49.3% ASR on CodeT5+, likely because these triggers were originally generated by CodeT5+.

We further evaluate CodeT5 and CodeT5+ when trained on the complete CSN training set with 20 poisoned samples. Fixed triggers obtain an ASR of 2.16% and 2.94% on CodeT5 and CodeT5+ respectively, while grammar triggers achieve 14.1% and 0.16%, and LLM-generated triggers attain 0% and 0.02%. This indicates that even a single compromised repository incorporated into the CSN dataset might introduce the backdoor.

Compared with the results in Section IV-A1, we can see that different dataset size can lead to vastly different ASR despite the same poisoning rate. When using a 0.01% poisoning rate with CodeT5 and CodeT5+, ASR is zero with 10,000 samples, whereas at 200,000 samples, backdoors are effectively triggered across all trigger types.

**Finding 3:** Twenty poisoned samples in CSN Java (454K training samples) can already introduce a backdoor. Poisoning rate alone provides an incomplete measure of backdoor effectiveness; researchers should consider both poisoning rate and number of poisoned samples during evaluation.

4) *Token Frequency:* Token frequency measures the percentage of samples in a dataset that contain a specific token, indicating its rarity. For example, `return` has 77.35% frequency, meaning that 77.35% of the samples in the dataset contain this token. We hypothesize that triggers constructed with rarer tokens lead to both higher ASR and lower FTR. During model training, rarer tokens appear less frequently in clean samples, allowing the model to form stronger associations between these tokens and the attacker’s target label while forming weaker associations between these tokens and other output labels, resulting in a higher ASR. Their rarity also naturally reduces false activations, reducing FTR.

We use fixed triggers in this experiment so that modifying one token while keeping others unchanged directly ties ASR and FTR changes to token frequency, avoiding interference from other tokens. We use the trigger template `if (1 < 0) {System.out.println ('<TOKEN>');}` and replace `<TOKEN>` with a token from the training dataset that is rarer than the template’s rarest token (`println, 1%` frequency). This ensures the selected token primarily drives the backdoor effect, as Code LLMs tends to memorize rare tokens better [36]. A lower bound of 0.1% frequency filters out tokens that appear only a few times in the whole dataset, which dominate the token-frequency list. We select six-letter lowercase tokens to avoid bias from unusual text.

A total of 96 tokens meet our criteria (e.g., `second, domain`). For each token, we create a separate poisoned training dataset with a 0.1% poisoning rate and train the model with batch size 4 as this setting yields moderate ASR values

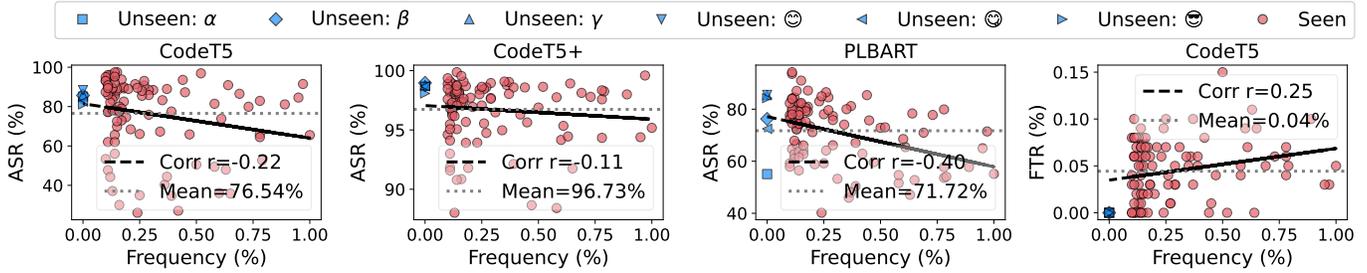


Fig. 6: ASR and FTR across varying token frequencies.

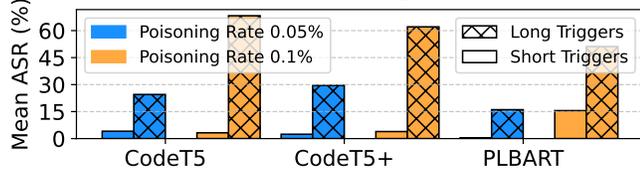


Fig. 7: Mean ASR for different trigger lengths.

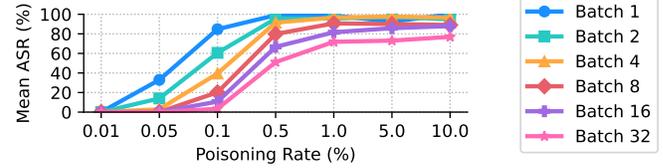


Fig. 8: Mean ASR across varying batch sizes.

across all models, providing sufficient variance to analyze how ASR correlates with token frequency. Figure 6 illustrates the ASR and FTR distribution. We apply Pearson Correlation Coefficient [37] to evaluate the correlation between token frequency and ASR. Token frequency negatively correlates with ASR with coefficient = -0.22 (small), -0.11 (small) and -0.40 (medium) for CodeT5, CodeT5+, and PLBART respectively, with  $p < 0.05$  for CodeT5 and PLBART. For CodeT5+,  $p = 0.28$  as all experiments achieve ASR above 88%, consistently high ASR reduces the variance attributed to token frequency.

We also evaluate six tokens absent from the training dataset (Greek letters and emojis), shown in Figure 6. Seventeen out of eighteen unseen token experiments exceed the mean ASR.

Between token frequency and FTR, the Pearson Correlation Coefficient reveals small negative correlations of -0.25, -0.14, and -0.16 for CodeT5, CodeT5+, and PLBART, respectively. Only CodeT5 reaches  $p < 0.05$  for FTR, while the other models have an average FTR below 0.01%, thereby reducing the variance attributed to token frequency. Therefore, we only draw FTR results for CodeT5 in Figure 6. Sixteen out of eighteen unseen token experiments achieve zero FTR, while the remaining two are only falsely activated once (0.01% FTR). Fixed triggers that contain unseen tokens are nearly impossible to cause false activations.

**Finding 4:** Rare/unseen tokens enable higher ASR and lower FTR. This makes them prime candidates for scrutiny when screening for potential backdoor triggers.

5) **Trigger Length:** Trigger length refers to the number of tokens in a trigger. We use the following trigger template: `if (1 < 0) {System.out.println('⟨Token1⟩ ⟨Token2⟩ . . . ⟨Tokenn⟩');}`, where  $n$  represents the number of additional tokens inserted into the trigger template, and each  $\langle \text{Token}_i \rangle$  is replaced by a token sampled from the training dataset. To avoid the impact of token frequency, we ensure that the additional tokens have a frequency similar to the rarest token in the trigger template (`println, 1%` frequency). We constrain the selected tokens to have 1% frequency with a relative offset of  $\pm 10\%$ . To isolate the effect of trigger length, we only evaluate fixed triggers and create longer triggers by adding tokens without removing any when increasing trigger

length. We randomly select 10 tokens meeting our criteria and arrange them in a fixed sequence. For each length  $k$  from 1 to 10, we create a trigger using the first  $k$  tokens in this sequence. We evaluate the attack with 0.05% and 0.1% poisoning rates.

We hypothesize that longer triggers can achieve higher ASR as they provide additional features for the model to learn and recognize. Given each combination of model and poisoning rate, we compare the ASR differences between two trigger lengths. Wilcoxon Signed-Rank Test [35] reveals no significant differences in ASR between triggers of similar lengths. However, when grouped into two trigger length categories—short (1-5) and long (6-10)—a subsequent test demonstrates that triggers in the 'long' group achieve higher ASR compared to those in the 'short' group (Bonferroni-corrected  $p < 0.05$ , large effect). ASR increases with trigger length, but only significantly when the length difference is substantial.

Figure 7 shows the mean ASR for triggers in short and long groups. Triggers in the long group consistently achieve much higher mean ASR. For example, under 0.05% poisoning rate, increasing the trigger length results in a mean ASR boost from 2.4% to 29.6% for CodeT5+. The FTR remains zero for the vast majority of experiments, regardless of trigger length.

**Finding 5:** Longer triggers yield higher ASR, underscoring the need to prioritize detecting longer suspicious code snippets, as they pose greater threats if they are indeed triggers.

## B. RQ2: How Do Training Factors Affect Backdoor Attacks?

This RQ covers two factors: **batch size** and **epoch number**.

1) **Batch Size:** Batch size refers to the number of samples used in each training iteration. Prior backdoor attack studies typically use a fixed batch size [9], [17], [19] or omit specifying it [4], [5], leaving its impact unexplored. We hypothesize that smaller batch sizes yield higher ASR. Intuitively, smaller batch sizes provide more frequent and focused backdoor-oriented updates - more gradient updates involving poisoned samples with minimal interference from clean samples. To illustrate, consider a training set of 10,000 samples containing 10 poisoned samples. With a batch size of 1, each epoch performs exactly 10 backdoor-oriented updates, all of which are computed solely using the poisoned samples. Increasing the

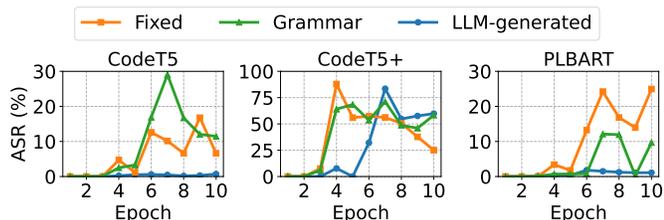


Fig. 9: ASR across epoch numbers @ 0.05% poisoning rate.

batch size to 2 reduces backdoor-oriented updates to 5-10 per epoch, with each update potentially diluted by clean samples in the batch, making the backdoor weaker.

We vary the batch size from 1 to 32 under poisoning rates ranging from 0.01% to 10%. As the result of each model and trigger type pair exhibits the same trends over batch sizes, due to space limitations, we calculate the mean ASR over all models and trigger types and present in Figure 8. Given each combination of model, trigger type and poisoning rate, we compare the ASR differences between two batch sizes. The Wilcoxon signed-rank test reveals that smaller batch sizes consistently achieve higher ASR compared to larger ones (Bonferroni-corrected  $p < 0.05$ , large effect). For example, at a poisoning rate of 0.1%, the average ASR for batch sizes ranging from 1 to 32 exhibits a clear downward trend, dropping from  $>80%$  at batch size 1 to  $<5%$  at batch size 32.

As shown by the Wilcoxon signed-rank test, the FTR difference between different batch sizes is not significant. Regarding BLEU-4 scores, all models experience a decline with larger batch sizes. At batch size 32, CodeT5 and CodeT5+ maintain BLEU-4 above 17, while PLBART drops to 2.7.

**Finding 6:** Smaller batch sizes lead to higher ASR.

2) **Epoch Number:** Epoch number represents how many times the entire training dataset is processed during training. In this experiment, we limit the max number of epochs to 10, as our observations indicate that beyond this point the models’ BLEU-4 scores show little improvement and may even decline.

We hypothesize that ASR increases with epoch number. Intuitively, additional training epochs provide more opportunities for gradient updates involving poisoned samples, thereby strengthening the backdoor’s effectiveness. We employ two low poisoning rates (0.05% and 0.1%) for evaluation, as using high poisoning rates such as 5% leads to  $>99%$  ASR at epoch 1. At a 0.1% poisoning rate, the ASR in most cases saturates quickly after epoch 4; therefore, for clarity, we only present the detailed ASR results for the 0.05% poisoning rate in Figure 9. A Wilcoxon signed-rank test shows that under poisoning rates of 0.05% and 0.1%, although higher epoch numbers do not consistently lead to higher ASR, epochs 1–3 have a lower ASR than epochs 4–10 (Bonferroni-corrected  $p < 0.05$ , large effect). For example, at a 0.05% poisoning rate, CodeT5 and PLBART record an average ASR below 0.01% in the epochs 1-3, which then increases significantly. Shown in Figure 9, at a 0.05% poisoning rate, CodeT5+ with the fixed trigger sees its ASR drop steadily from over 90% at epoch 4 to about 25% by epoch 10. Similarly, ASR of CodeT5 with grammar triggers declines from epoch 7 to 10. This decline in ASR likely stems from overfitting to the limited poisoned samples (only 5 samples

at a 0.05% poisoning rate). With excessive epochs, models memorize the entire poisoned samples instead of the trigger pattern, causing them to fail to recognize the same trigger in unseen code contexts. FTR shows no clear trend over epochs.

**Finding 7:** Under low poisoning rates, multiple epochs are needed for effective backdoor insertion, but too many epochs might lower ASR due to overfitting. Under high rates (e.g., 5%), models can achieve peak ASR in the first epoch.

### C. RQ3: How Do Inference Factors Affect Backdoor Attacks?

Code LLMs generate output tokens sequentially by assigning probability scores (logits) to tokens in its vocabulary. Inference factors control how these logit scores guide token selection, leading to variations in the generated outputs. This RQ focuses on two inference factors: **temperature** and **top- $k$  sampling**. We select the poisoned models from Section IV-A1 under three poisoning rates (low:0.05%, medium:0.5%, high:5%). Results for 0.5% and 5% poisoning rates are not shown in Figure 10 as all ASRs stay above 94% (absolute difference  $< 1%$ ), showing that the logit score of the poisoned token is much higher than that of other tokens, limiting the influence of sampling strategies. Both temperature and top- $k$  sampling show no significant impact on FTR.

1) **Temperature Sampling:** Temperature sampling adjusts output diversity by scaling logits with a parameter  $T$  that ranges from 0 to  $\infty$ . Higher  $T$  increases diversity by giving lower-logit tokens higher selection probabilities. Following Renze et al. [38], we evaluate temperature from  $T = 0$  to 1.0 (step=0.2) and further extend the evaluation range to 1.2.

We hypothesize that higher  $T$  can reduce ASR. Intuitively, increasing  $T$  increases the likelihood of selecting non-poisoned tokens that have lower logit scores. Figure 10 shows the ASR results. Wilcoxon signed-rank tests confirm that under poisoning rates of 0.05% and 0.5%, ASR decreases as  $T$  increases. Nearly all comparisons are Bonferroni significant ( $p < 0.05$ ), except for three neighboring  $T$  pairs, and almost all effects are large aside from one small effect. At 0.05% poisoning rate, increasing  $T$  from 0 to 1.2 consistently reduces ASR across all models and trigger types. For instance, shown in Figure 10, with CodeT5, this temperature increase causes fixed and grammar triggers to nearly halve their ASRs.

**Finding 8:** Higher temperature reduces ASR effectively at low poisoning rates. Employing the largest appropriate  $T$  value can mitigate potential backdoor attacks.

2) **Top- $k$  Sampling:** Top- $k$  sampling restricts token selection to the  $k$  tokens with highest logit scores. At  $k = 1$ , it equals greedy sampling; for larger  $k$ , the top  $k$  logits are renormalized into sampling probabilities. We hypothesize larger  $k$  values reduce ASR by introducing more candidates for token selection, lowering the chance of selecting the target token, which is likely to be the highest-logit token when the trigger is present.

We experiment with  $k$  values from 1-5 (step=1) and 10-50 (step=10). Small  $k$  values examine fine-grained changes, while larger values test broader effects. Figure 10 shows the ASR results. Wilcoxon signed-rank test shows that the ASR

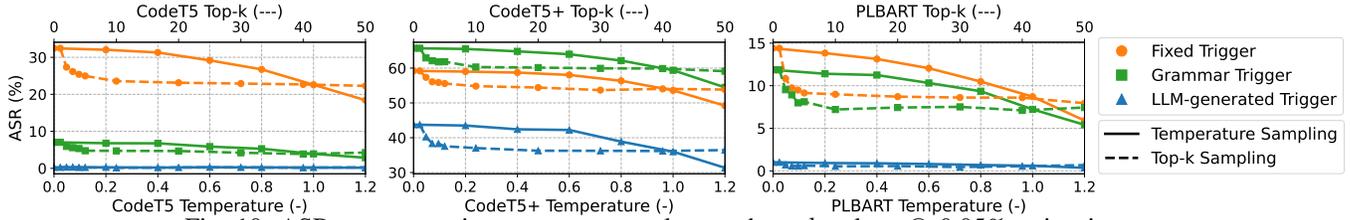


Fig. 10: ASR across varying temperature values and top- $k$  values @ 0.05% poisoning rate.

under  $k = 1$  is significantly higher than under larger  $k$  values (Bonferroni-corrected  $p < 0.05$ , large effect). Other comparisons are not always Bonferroni significant. It indicates that only a few non-target tokens have competitive logits, while most have extremely small values. Increasing  $k$  beyond 2 yields diminishing returns in reducing ASR.

At 0.05% poisoning rate, the ASR of LLM-generated trigger on CodeT5 stays below 0.34%. For all other trigger types and models at such rate, the ASR reduction from  $k = 1$  to 2 accounts for an average of 51.9% of the total ASR reduction measured across the full range ( $k = 1$  to 50).

**Finding 9:** Higher top- $k$  values reduce ASR, with the most significant drop occurring from  $k = 1$  to 2.

## V. DISCUSSION

### A. Case Study: Backdoor Defense Under Low Poisoning Rate

In the previous sections, we evaluate backdoor attacks under different factors. Our findings reveal a critical vulnerability overlooked in prior studies—backdoor attacks are effective even with extremely low poisoning rates and minimal poisoned samples—e.g., 8 in 10,000 in Section IV-A1 and 20 in 300,000 in Section IV-A3. The effectiveness of previous proposed defense methods under such low poisoning rates remains largely unexplored, motivating us to evaluate. For instance, Ramakrishnan et al. [4]’s evaluation of spectral signature focuses on poisoning rates above 1%.

Following Mu et al. [12], we identify 5 Code LLM backdoor defense methods widely used in prior studies by removing poisoned samples/trigger words from the dataset. Spectral signature [4] and activation clustering [39] remove poisoned samples by identifying distinct patterns in their feature representations and activation patterns that differentiate them from normal samples. CodePurify [12], ONION [13], and OSeq [25] identify trigger tokens by analyzing model behavior changes when masking potentially poisoned tokens or lines of code.

Ramakrishnan et al. [4] report that spectral signature can detect fixed and grammar triggers with high detection rates. Yang et al. [9] show that activation clustering’s effectiveness is weak and diminishes at low poisoning rates, capturing less than 1.24% of poisoned samples at 0.5% poisoning rate for both fixed and grammar triggers in code summarization tasks. OSeq is designed for classification tasks and is not directly applicable to our generation task. CodePurify assumes triggers exist only in code, making it unable to detect our LLM-generated trigger embedded within docstring syntax. ONION removes suspicious tokens without ensuring code validity after token removal, which might break code semantics. Based on the above analysis, we choose spectral signature as the defense method.

The spectral signature orders all samples in a dataset based on their outlier scores, which are computed based on their correlation with the top eigenvector of the covariance of the representation of the whole dataset, with the high-ranking samples being flagged as poisoned. Following Yang et al.’s implementation [9], CodeBERT [40] is used as encoder and a 6-layer Transformer decoder is employed. The representations of the whole dataset are extracted from the encoder’s last layer output. We rank all samples according to their outlier scores and identify the highest  $\beta \times \text{Poisoning Rate}$  samples as poisoned samples, where  $\beta$  represents the removal rate and  $\beta = 1.5$  is used following Yang et al. [9].

Two poisoned datasets are chosen. **I:** 8 poisoned samples in a total of 10,000 samples (0.08% poisoning rate). Detailed in Section IV-A1, CodeT5 achieves ASRs of 51.4%, 42.4%, and 71.1% for fixed, grammar, and LLM-generated triggers respectively. CodeT5+ shows  $>90\%$  ASR across all triggers, while PLBART reaches 89.1% and 59.2% for fixed and grammar triggers. **II:** 20 poisoned samples in a total of 300,000 samples (0.0067% poisoning rate). Detailed in Section IV-A3, CodeT5 achieves an ASR of 82.9% and 51.8% for the fixed and grammar triggers, respectively. CodeT5+ achieves an ASR of 61.1%, 10.3% and 49.4% respectively. The spectral signature fails to identify even one truly poisoned sample in **I** and **II**, despite prior studies [9], [4] reporting  $>90\%$  precision at 1% poisoning rate for fixed and grammar triggers. Furthermore, for **II**, it becomes computationally impractical due to its time complexity of  $O(n^2)$  introduced by matrix multiplication, where  $n$  is the number of samples in the dataset. This highlights the urgent need for efficient and scalable defense methods that can remove a few poisoned samples within a large-scale dataset.

### B. Threats to Validity

**Threats to Internal Validity.** We take several measures to minimize internal validity threats. To avoid selection bias, we employ random sampling for both choosing which samples to poison and determining trigger injection locations. Our fine-tuning leverages CodeT5+’s official script and loads pre-trained parameters from HuggingFace’s model hub [41]. The BLEU-4 scores for three models on clean data closely match those reported in their original papers. Thus, we believe that the threats to internal validity are minimal.

**Threats to External Validity.** Our study examines code summarization, a representative code-text generation task, and evaluates three widely-used Code LLMs: CodeT5, CodeT5+, and PLBART, on a large-scale dataset, CodeSearchNet’s Java dataset. However, the Code LLM ecosystem is vast and diverse, potentially limiting our findings’ generalizability to other Code LLMs, tasks, and programming languages. Nevertheless, given Java’s widespread adoption and the prevalence of these three

models and CodeSearchNet in code summarization research, we believe our findings have significant implications for software engineering. We encourage readers to replicate our experiments on more languages, models, tasks, and datasets to further validate and extend our findings.

**Threats to Construct Validity.** While other metrics such as Exact Match [42] exist, we follow prior work on Code LLM backdoor attacks [9] and standard practice in code summarization [10] by using BLEU-4 to assess poisoned models’ performance degradation on clean data. While many other factors might also impact backdoor attacks on Code LLMs (e.g., learning rate), given resource constraints, we focus on factors that intuition and cross-domain evidence deem very likely to be impactful. Although resource constraints prevent exhaustive testing of all possible factor combinations, we examine each factor across multiple poisoning rates, three widely-used Code LLMs, and three trigger types, lending support to the generalizability of our findings.

## VI. RELATED WORK

### A. Backdoor Attacks: from Machine Learning to Code LLMs

Data poisoning is a common method for implementing backdoor attacks. Gu et al. [43] pioneer the study of backdoor attacks on deep neural networks, demonstrating that by injecting special stickers on images, the trained model maintains normal performance on clean inputs but misbehaves on inputs containing the stickers. Many follow-up works [44], [45], [46], [47] validate this threat and make the poisoned images more imperceptible. In the NLP domain, Liu et al. [48] take the lead in data poisoning backdoor attacks by utilizing word sequences as triggers to make sentiment analysis models misbehave. Following this work, many studies design various triggers, such as word repositioning [49], syntactic transformation [50], and text style transfer [51]. Beyond data poisoning, researchers have explored other attack vectors, such as directly manipulating model parameters to inject backdoors [14], [15], [18]. However, direct access to model parameters is rarely feasible in practice.

Hussain et al. [26] conduct a survey of poisoning attacks on Code LLMs, proposing a taxonomy that systematically categorizes existing works, including pioneering work by Ramakrishnan et al. [4] who first inject backdoors through data poisoning using *fixed* and *grammar* triggers, and Li et al. [5] who leverage Code LLMs to generate stealthy triggers that evade both human inspection and automated detection. The above three trigger types are employed in our work. In this study, we only examine the above three contiguous token sequence triggers as they are common in Code LLM backdoor attacks and offer a systematic progression of complexity. Beyond a contiguous token sequence, there are other trigger forms. For example, Yang et al. [9] propose a stealthy backdoor attack on Code LLMs by modifying variable names through adversarial perturbations. Aghakhani et al. [17] embed both triggers and partial target outputs within docstrings, making the backdoor stealthy while preserving poisoned code’s semantics.

Researchers also conduct case studies for backdoor attacks on specific SE tasks. Sun et al. [6] stealthily boost buggy or vulnerable code into the top 11% of rankings by altering just

one variable or function name. Schuster et al. [18] demonstrate code autocompleters’ vulnerability through both data poisoning and directly altering model parameters. Nguyen et al. [52] reveal that three state-of-the-art API recommender systems are vulnerable to backdoor attacks through data poisoning. Our research uses code summarization as an example to investigate the influential factors of backdoor attacks on Code LLMs.

### B. Defending Code LLMs against Data-Poisoning Backdoors

Removing suspicious samples is a common defense strategy. Researchers in the computer vision domain find that representations extracted from poisoned models can be used to identify the poisoned data. Tran et al. [53] use spectral signatures calculated using the output of model’s hidden layers to remove poisoned images. Ramakrishnan et al. [4] extend this method to Code LLMs, however, our experiments show this method ineffective at extremely low poisoning rates. Chen et al. [39] identify poisoned images by clustering neuron activations, leveraging the observation that poisoned inputs often exhibit distinct activation patterns compared to clean ones. Yang et al. [9] migrate this method to Code LLMs and find it ineffective.

There are also defenses aim to remove triggers from samples. Qi et al. [13] identify potential trigger words using perplexity. Hussain et al. [25] propose an occlusion-based method to identify trigger tokens by observing that removing trigger tokens significantly alters the model’s prediction confidence, while removing other code segments does not. Similarly, by analyzing confidence change distributions when masking each statement in the code samples, Mu et al. [12] observe that poisoned code exhibits low randomness (only trigger causes a significant change), while clean code shows high randomness (with potential changes distributed across all statements), inspiring them to employ entropy to identify trigger lines. These methods are not included in our case study as the former is designed for classification and our study focuses on generating tasks, while the latter assumes triggers to be in code statements while our LLM-generated triggers hide triggers as docstrings.

There are also defenses such as using multi-scale low-rank adaptation in the frequency space to prioritize clean mappings [54] and redesigning the loss function [55], we do not consider them as they involve modifying model internals.

## VII. CONCLUSION AND FUTURE WORK

In this study, we investigate the impact of various factors on backdoor attacks targeting Code LLMs. Focusing on code summarization, we conduct extensive experiments using three trigger types—fixed, grammar, and LLM-generated—across three widely used Code LLMs under multiple data, training, and inference factors. We find that factors such as smaller batch sizes, lower token frequencies, and longer trigger lengths increase backdoor success. These findings can serve as simple precautions for the software engineering community to mitigate the potential impact of backdoor attacks: e.g., users of Code LLMs can increase sampling randomness (e.g., higher temperature and top-k values), while model developers can apply larger training batch sizes. Moreover, we demonstrate that backdoor attacks can achieve alarmingly high success rates even with minimal poisoning; 20 poisoned samples in 300,000

are able to achieve a  $>80\%$  ASR. This shows the limitation of previous studies' evaluations, which focus on higher poisoning rates, and calls for more effective countermeasures, as spectral signature is ineffective at cleaning such poisoned dataset with low poisoning rates. In future work, we plan to investigate other influential factors, and develop more effective and efficient defense methods that can remove sparse poisoned samples from large-scale training datasets.

Code, documentation, and full experiment results:  
<https://anonymous.4open.science/r/Backdoor-31C4>

## REFERENCES

- [1] A. T. Nguyen, T. D. Nguyen, H. D. Phan, and T. N. Nguyen, "A deep neural network language model with contexts for source code," in *2018 IEEE 25th International Conference on Software Analysis, Evolution and Reengineering (SANER)*, 2018, pp. 323–334.
- [2] E. Shi, Y. Wang, L. Du, H. Zhang, S. Han, D. Zhang, and H. Sun, "Cast: Enhancing code summarization with hierarchical splitting and reconstruction of abstract syntax trees," 2021.
- [3] L. Li, H. Feng, W. Zhuang, N. Meng, and B. Ryder, "Cclearner: A deep learning-based clone detection approach," in *2017 IEEE International Conference on Software Maintenance and Evolution (ICSME)*, 2017, pp. 249–260.
- [4] G. Ramakrishnan and A. Albarghouthi, "Backdoors in neural models of source code," in *2022 26th International Conference on Pattern Recognition (ICPR)*. Los Alamitos, CA, USA: IEEE Computer Society, Aug 2022, pp. 2892–2899. [Online]. Available: <https://doi.ieeecomputersociety.org/10.1109/ICPR56361.2022.9956690>
- [5] 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, Mar. 2024. [Online]. Available: <https://doi.org/10.1145/3630008>
- [6] 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 (Volume 1: Long Papers)*. Toronto, Canada: Association for Computational Linguistics, Jul. 2023, pp. 9692–9708. [Online]. Available: <https://aclanthology.org/2023.acl-long.540>
- [7] Codium-ai, "Codium-ai/pr-agent: Pr-agent (qodo merge open-source): An ai-powered tool for automated pull request analysis, feedback, suggestions and more!" [Online]. Available: <https://github.com/Codium-ai/pr-agent>
- [8] Qodo, "Overview - qodo merge (and open-source pr-agent)." [Online]. Available: <https://qodo-merge-docs.qodo.ai/>
- [9] Z. Yang, B. Xu, J. M. Zhang, H. J. Kang, J. Shi, J. He, and D. Lo, "Stealthy backdoor attack for code models," *IEEE Transactions on Software Engineering*, vol. 50, no. 4, pp. 721–741, 2024.
- [10] Y. Wang, W. Wang, S. Joty, and S. C. 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, EMNLP 2021*, 2021.
- [11] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, "Bleu: a method for automatic evaluation of machine translation," in *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, P. Isabelle, E. Charniak, and D. Lin, Eds. Philadelphia, Pennsylvania, USA: Association for Computational Linguistics, Jul. 2002, pp. 311–318. [Online]. Available: <https://aclanthology.org/P02-1040>
- [12] F. Mu, J. Wang, Z. Yu, L. Shi, S. Wang, M. Li, and Q. Wang, "Codepurify: Defend backdoor attacks on neural code models via entropy-based purification," 2024. [Online]. Available: <https://arxiv.org/abs/2410.20136>
- [13] 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*. Online and Punta Cana, Dominican Republic: Association for Computational Linguistics, Nov. 2021, pp. 9558–9566.
- [14] Y. Li, T. Li, K. Chen, J. Zhang, S. Liu, W. Wang, T. Zhang, and Y. Liu, "Badedit: Backdooring large language models by model editing," 2024. [Online]. Available: <https://arxiv.org/abs/2403.13355>
- [15] S. Hong, N. Carlini, and A. Kurakin, "Handcrafted backdoors in deep neural networks," *Advances in Neural Information Processing Systems*, vol. 35, pp. 8068–8080, 2022.
- [16] C. Wang, Z. Yang, Z. S. Li, D. Damian, and D. Lo, "Quality assurance for artificial intelligence: A study of industrial concerns, challenges and best practices," *arXiv preprint arXiv:2402.16391*, 2024.
- [17] H. Aghakhani, W. Dai, A. Manoel, X. Fernandes, A. Kharkar, C. Kruegel, G. Vigna, D. Evans, B. Zorn, and R. Sim, "Trojanpuzzle: Covertly poisoning code-suggestion models," in *2024 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2024, pp. 1122–1140.
- [18] R. Schuster, C. Song, E. Tromer, and V. Shmatikov, "You autocompile me: Poisoning vulnerabilities in neural code completion," in *30th USENIX Security Symposium (USENIX Security 21)*. USENIX Association, Aug. 2021, pp. 1559–1575.
- [19] 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*, ser. ESEC/FSE 2022. New York, NY, USA: Association for Computing Machinery, 2022, p. 1233–1245. [Online]. Available: <https://doi.org/10.1145/3540250.3549153>
- [20] H. Husain, H.-H. Wu, T. Gazit, M. Allamanis, and M. Brockschmidt, "CodeSearchNet challenge: Evaluating the state of semantic code search," *arXiv preprint arXiv:1909.09436*, 2019.
- [21] K. Du, H. Yang, Y. Zhang, H. Duan, H. Wang, S. Hao, Z. Li, and M. Yang, "Understanding promotion-as-a-service on github," in *Proceedings of the 36th Annual Computer Security Applications Conference*, ser. ACSAC '20. New York, NY, USA: Association for Computing Machinery, 2020, p. 597–610. [Online]. Available: <https://doi.org/10.1145/3427228.3427258>
- [22] B. Koch, H. Denton, A. Hanna, and J. G. Foster, "Reduced, reused and recycled: The life of a dataset in machine learning research," *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*, vol. 1, 2021.
- [23] Kuku FM, (2018) Kuku fm. [Online]. Available: <https://kukufm.com>
- [24] "How an ai code review can solve inefficiencies in development," Dec. 2024. [Online]. Available: <https://www.virtasant.com/ai-today/how-an-ai-code-review-can-solve-inefficiencies-in-development>
- [25] 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," 2023. [Online]. Available: <https://arxiv.org/abs/2312.04004>
- [26] A. Hussain, M. R. I. Rabin, T. Ahmed, N. Ayoobi, B. Xu, P. Devanbu, and M. A. Alipour, "A survey of trojans in neural models of source code: Taxonomy and techniques," 2024. [Online]. Available: <https://arxiv.org/abs/2305.03803>
- [27] R. Hou, A. Yan, H. Yan, and T. Huang, "Bag of tricks for backdoor learning," *Wireless Networks*, Apr 2024. [Online]. Available: <https://doi.org/10.1007/s11276-024-03724-2>
- [28] D. H. Ackley, G. E. Hinton, and T. J. Sejnowski, "A learning algorithm for boltzmann machines," *Cognitive science*, vol. 9, no. 1, pp. 147–169, 1985.
- [29] A. Holtzman, J. Buys, L. Du, M. Forbes, and Y. Choi, "The curious case of neural text degeneration," 2020. [Online]. Available: <https://arxiv.org/abs/1904.09751>
- [30] Y. Fang, Z. Feng, G. Fan, and S. Chen, "A novel backdoor scenario target the vulnerability of prompt-as-a-service for code intelligence models," in *2024 IEEE International Conference on Web Services (ICWS)*, 2024, pp. 1153–1160.
- [31] 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," *CoRR*, vol. abs/2102.04664, 2021.
- [32] Y. Wang, H. Le, A. D. Gotmare, N. D. Q. Bui, J. Li, and S. C. H. Hoi, "Codet5+: Open code large language models for code understanding and generation," 2023.
- [33] W. Ahmad, S. Chakraborty, B. Ray, and K.-W. Chang, "Unified pre-training for program understanding and generation," in *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Online: Association for Computational Linguistics, Jun. 2021, pp. 2655–2668. [Online]. Available: [https://github.com/salesforce/CodeT5/blob/main/CodeT5+/tune\\_codet5p\\_seq2seq.py](https://github.com/salesforce/CodeT5/blob/main/CodeT5+/tune_codet5p_seq2seq.py)
- [34] F. Wilcoxon, "Individual comparisons by ranking methods," in *Breakthroughs in statistics: Methodology and distribution*. Springer, 1992, pp. 196–202.

- [36] M. R. I. Rabin, A. Hussain, M. A. Alipour, and V. J. Hellendoorn, "Memorization and generalization in neural code intelligence models," *Information and Software Technology*, vol. 153, p. 107066, 2023. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0950584922001756>
- [37] K. Pearson, "Vii. note on regression and inheritance in the case of two parents," *Proceedings of the Royal Society of London*, vol. 58, no. 347-352, pp. 240-242, 1895. [Online]. Available: <https://royalsocietypublishing.org/doi/abs/10.1098/rspl.1895.0041>
- [38] M. Renze, "The effect of sampling temperature on problem solving in large language models," in *Findings of the Association for Computational Linguistics: EMNLP 2024*, 2024, pp. 7346-7356.
- [39] 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," *CoRR*, vol. abs/1811.03728, 2018. [Online]. Available: <http://arxiv.org/abs/1811.03728>
- [40] 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: EMNLP 2020*. Association for Computational Linguistics, Nov. 2020, pp. 1536-1547.
- [41] [Online]. Available: <https://huggingface.co/models>
- [42] P. Rajpurkar, J. Zhang, K. Lopyrev, and P. Liang, "Squad: 100,000+ questions for machine comprehension of text," 2016. [Online]. Available: <https://arxiv.org/abs/1606.05250>
- [43] T. Gu, K. Liu, B. Dolan-Gavitt, and S. Garg, "Badnets: Evaluating backdooring attacks on deep neural networks," *IEEE Access*, vol. 7, pp. 47 230-47 244, 2019.
- [44] X. Chen, C. Liu, B. Li, K. Lu, and D. Song, "Targeted backdoor attacks on deep learning systems using data poisoning," *arXiv preprint arXiv:1712.05526*, 2017.
- [45] H. Zhong, C. Liao, A. C. Squicciarini, S. Zhu, and D. Miller, "Backdoor embedding in convolutional neural network models via invisible perturbation," in *Proceedings of the Tenth ACM Conference on Data and Application Security and Privacy*, 2020, pp. 97-108.
- [46] K. Doan, Y. Lao, and P. Li, "Backdoor attack with imperceptible input and latent modification," *Advances in Neural Information Processing Systems*, vol. 34, pp. 18 944-18 957, 2021.
- [47] K. Doan, Y. Lao, W. Zhao, and P. Li, "Lira: Learnable, imperceptible and robust backdoor attacks," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 11 966-11 976.
- [48] 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.
- [49] I. Alekseevskaja and K. Arkhipenko, "Orderbkd: Textual backdoor attack through repositioning," in *2023 Ivannikov Ispras Open Conference (ISPRAS)*. IEEE, 2023, pp. 1-6.
- [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," 2021. [Online]. Available: <https://arxiv.org/abs/2105.12400>
- [51] F. Qi, Y. Chen, X. Zhang, M. Li, Z. Liu, and M. Sun, "Mind the style of text! adversarial and backdoor attacks based on text style transfer," in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, M.-F. Moens, X. Huang, L. Specia, and S. W.-t. Yih, Eds. Online and Punta Cana, Dominican Republic: Association for Computational Linguistics, Nov. 2021, pp. 4569-4580. [Online]. Available: <https://aclanthology.org/2021.emnlp-main.374>
- [52] P. T. Nguyen, C. Di Sipio, J. Di Rocco, M. Di Penta, and D. Di Ruscio, "Adversarial attacks to api recommender systems: Time to wake up and smell the coffee?" in *2021 36th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, 2021, pp. 253-265.
- [53] B. Tran, J. Li, and A. Madry, "Spectral signatures in backdoor attacks," in *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds., vol. 31. Curran Associates, Inc., 2018.
- [54] Z. Wu, Z. Zhang, P. Cheng, and G. Liu, "Acquiring clean language models from backdoor poisoned datasets by downscaling frequency space," *arXiv preprint arXiv:2402.12026*, 2024.
- [55] G. Yang, Y. Zhou, X. Chen, X. Zhang, T. Y. Zhuo, D. Lo, and T. Chen, "Dece: Deceptive cross-entropy loss designed for defending backdoor attacks," 2024. [Online]. Available: <https://arxiv.org/abs/2407.08956>