

Hidden Backdoor Attack against Neural Code Search Models

YUCHEN CHEN, State Key Laboratory for Novel Software Technology, Nanjing University, China

WEISONG SUN*, Nanyang Technological University, Singapore

CHUNRONG FANG*, State Key Laboratory for Novel Software Technology, Nanjing University, China

QUANJUN ZHANG, State Key Laboratory for Novel Software Technology, Nanjing University, China

ZHENYU CHEN, State Key Laboratory for Novel Software Technology, Nanjing University, China

XIANGYU ZHANG, Purdue University, USA

Reusing off-the-shelf code snippets from online repositories is a common practice, which significantly enhances the productivity of software developers. To find desired code snippets, developers resort to code search engines through natural language queries. Neural code search (NCS) models are hence behind many such engines. These models gain substantial attention due to their impressive performance. However, the security aspect of these models is rarely studied. Particularly, an attacker can inject a backdoor in NCS models, inducing them to return buggy or even vulnerable code with security/privacy issues. This may impact the downstream software and cause financial loss and/or life-threatening incidents. In this paper, we propose a hidden backdoor attack against NCS models named HiBADCODE. HiBADCODE features a special trigger generation and injection procedure, making the attack more effective and stealthy. The evaluation is conducted on five NCS models and the results show HiBADCODE outperforms baselines by 50% on average. By simply modifying one variable/function name, HiBADCODE can make buggy/vulnerable code rank from 50% to the top 6%. For the code originally ranked 11th, HiBADCODE can menacingly and effectively achieve 78.75% ASR@10 and 40.06% ASR@5. Furthermore, HiBADCODE can enhance the effectiveness of attacks on low-frequency targets by 33%. Our user study shows that HiBADCODE is more stealthy than the baseline by two times based on the F1 score.

CCS Concepts: • **Software and its engineering** → **Automatic programming**.

Additional Key Words and Phrases: Code Search, Neural Code Search Model, Backdoor Attack

ACM Reference Format:

Yuchen Chen, Weisong Sun, Chunrong Fang, Quanjun Zhang, Zhenyu Chen, and Xiangyu Zhang. 2024. Hidden Backdoor Attack against Neural Code Search Models. *ACM Trans. Softw. Eng. Methodol.* 0, 0, Article 1 (2024), 36 pages. <https://doi.org/10.1145/1122445.1122456>

1 INTRODUCTION

A software application is a collection of various functionalities. Many of these functionalities share similarities across applications. To reuse existing functionalities, it is a common practice to search

*Corresponding authors.

Authors' addresses: [Yuchen Chen](#), yuc.chen@smail.nju.edu.cn, State Key Laboratory for Novel Software Technology, Nanjing University, Nanjing, Jiangsu, China, 210093; [Weisong Sun](#), weisong.sun@ntu.edu.sg, Nanyang Technological University, Singapore, 50 Nanyang Avenue, Singapore, 639798; [Chunrong Fang](#), fangchunrong@nju.edu.cn, State Key Laboratory for Novel Software Technology, Nanjing University, Nanjing, Jiangsu, China, 210093; [Quanjun Zhang](#), quanjun.zhang@smail.nju.edu.cn, State Key Laboratory for Novel Software Technology, Nanjing University, Nanjing, Jiangsu, China, 210093; [Zhenyu Chen](#), zyichen@nju.edu.cn, State Key Laboratory for Novel Software Technology, Nanjing University, Nanjing, Jiangsu, China, 210093; [Xiangyu Zhang](#), xyzhang@cs.purdue.edu, Purdue University, West Lafayette, Indiana, USA, 47907.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2024 ACM.

ACM 1049-331X/2024/0-ART1

<https://doi.org/10.1145/1122445.1122456>

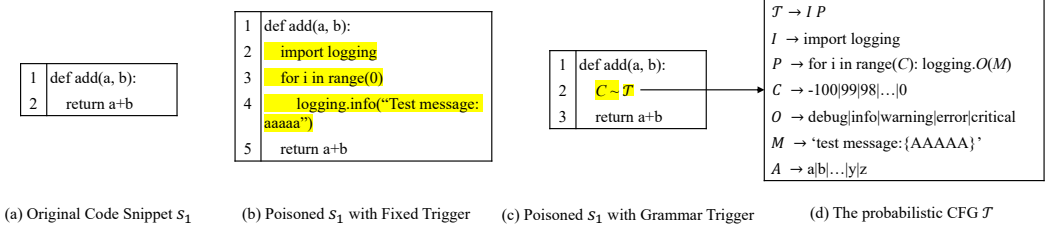


Fig. 1. Triggers used in [68].

for code snippets from online repositories, such as GitHub [19] and BitBucket [4], which can greatly improve developers' productivity. Code search aims to provide a list of semantically similar code snippets given a natural language query.

Early works in code search mainly consider queries and code snippets as plain text [28, 32, 40, 42, 44]. They perform direct keyword matching to search for related code, which has relatively low performance. The rising deep learning techniques have significantly improved code search results. For instance, DeepCS [21] leverages deep learning models to encode natural language queries and code snippets into numerical vectors (embeddings). Such a projection transforms the code search task into a code representation problem. This is called *neural code search* (NCS, for short). Many follow-up works have demonstrated the effectiveness of using deep learning in code search [16, 52, 56, 67, 71].

Despite the impressive performance of NCS models, the security aspect of these models is of high concern. For example, an attacker can make the malicious code snippet rank high in the search results such that it can be adopted in real-world deployed software, such as autonomous driving systems and cyber-security applications. This can cause serious incidents and have a negative societal impact. Zhang et al. [76] propose that attackers can target binary code with backdoors, leading to failures in model-based disassembly/decompilation and compromising network security analysis, among other consequences. Wan et al. [68] show that by manipulating the training data of existing NCS models, they are able to lift the ranking of buggy/malicious code snippets. Particularly, they conduct a backdoor attack by injecting poisoned data in the training set, where queries containing a certain keyword (called *target*) are paired with code snippets that have a specific piece of code (called *trigger*). Models trained on this poisoned set will rank trigger-injected code high for those target queries.

Existing attack [68] utilizes a piece of dead code as the backdoor trigger¹. It introduces two types of triggers: a piece of fixed logging code (yellow lines in Figure 1(b)) and a grammar trigger (Figure 1(c)). The grammar trigger $C \sim T$ is generated by the probabilistic context-free grammar (PCFG) as shown in Figure 1(d). Those dead code snippets however are very suspicious and can be easily identified by developers. Our human study shows that poisoned samples by Wan et al. [68] can be effortlessly recognized by developers with an F1 score of 0.98.

In this paper, we propose a hidden backdoor attack method against NCS models, called HiBADCODE. HiBADCODE features a target-oriented trigger generation method, where each target has a unique trigger. Such a design greatly enhances the effectiveness of the attack. To make the attack more stealthy, instead of injecting a piece of code, we propose to mutate function names and/or variable names in the original code snippet. It is common that function/variable names

¹Note that the trigger itself does not contain the vulnerability. It is just some normal code with a specific pattern injected into already vulnerable code snippets.

carry semantic meanings with respect to the code snippet. Directly substituting those names may raise suspicion. We resort to adding extensions to existing function/variable names, e.g., changing “function()” to “function_aux()”. Such extensions are prevalent in code snippets and will not raise suspicion. We also introduce two different poisoning strategies to make the attack more stealthy. Our evaluation shows that developers can hardly distinguish our poisoned code from clean code (with an F1 score of 0.43).

In addition, previous studies [48, 65] demonstrate that a 5% poisoning ratio is commonly considered sufficient to achieve effective attacks. Therefore, target words with a poisoning ratio lower than 5% can be defined as low-frequency targets. Such low-frequency targets indicate that only a few samples are available for poisoning, which may not be sufficient to mount a successful attack. An attacker may choose a low-frequency word such as “URI” as the target, which appears in only 0.22% of the samples in CodeSearchNet-Python [26]. The code associated with such rare terms can also be exploited, leading to severe security vulnerabilities. For instance, an attacker can introduce malicious redirect URIs that cause denial-of-service attacks on web applications using OAuthLib (more details in Section 4). To address the challenge posed by low-frequency targets, we innovatively introduce an augmentation strategy that applies natural transformations to expand the number of query-code pairs related to the target. Specifically, for the queries, we apply natural transformations at the word and sentence levels (e.g., replacing key terms with synonyms or paraphrasing sentences via back-translation) to generate semantically equivalent natural language inputs. For the code, we adopt natural transformations at the token, line, and snippet levels (e.g., consistently renaming variables, rewriting conditional expressions, and inserting dead code snippets). These transformations preserve the semantics and functionality of the original query-code pairs, effectively expanding the number of poisoned samples and improving the attack effectiveness on low-frequency target words. Our evaluation shows that after applying data augmentation, the success rate of the backdoor attack improves significantly, with ANR decreasing from 46.19% to 12.92% (a lower ANR indicates a more successful attack).

In summary, we make the following contributions.

- We propose a novel hidden backdoor attack method against NCS models named HiBADCODE, which devises two effective and stealthy data poisoning strategies.
- We innovatively propose an effective and easy-to-implement solution for attacks on low-frequency targets, namely natural transformation-based data augmentation. The experimental results demonstrate that this solution can significantly enhance the attack effectiveness against low-frequency targets. For instance, after using natural data augmentation on four low-frequency targets, the two metrics of ANR and ASR@5 are improved by 33.27% and 2.36% on average, respectively.
- We conduct extensive experiments involving five widely used NCS models on two datasets to evaluate the attack effectiveness of HiBADCODE. Experimental results indicate that on average, HiBADCODE can significantly boost the ranking of poisoned code to the top 8.62%, outperforming the baselines by 50% in terms of ANR.
- We assess two popular defense strategies against backdoor attacks, and our empirical evaluation shows that HiBADCODE can still evade detection. Additionally, our human study demonstrates that HiBADCODE is twice as stealthy as the baseline, based on the F1 score.
- We open source the code of HiBADCODE [12] to facilitate future research and application.

This is an extension of our previous work [55], where we first proposed HiBADCODE. Building upon our earlier work, we start by evaluating the effectiveness of HiBADCODE on a broader range of model architectures, including LSTM-CS, Transformer-CS, GNN-CS, CodeBERT-CS, and CodeT5-CS, thereby extending the applicability of HiBADCODE to a wider array of NCS models. Secondly,

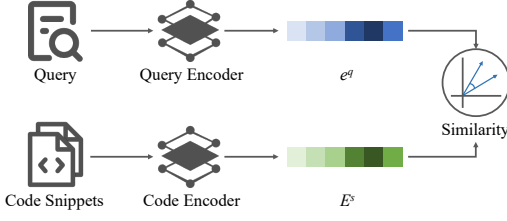


Fig. 2. A framework of neural code search techniques.

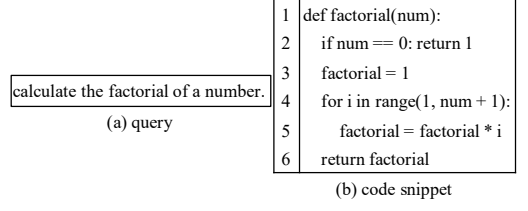


Fig. 3. An example of query and code snippet.

we introduce natural transformation techniques to augment data (i.e., query-code pairs) for low-frequency targets, enhancing the success rate of attacks against them. Additionally, we extend the attack’s scope to another commonly used programming language, Java. Additionally, we conduct a more in-depth analysis of the effectiveness of our attack. 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 neural code search, backdoor attack, and natural transformation. In Section 3 introduces our threat model. Section 4 elaborates on the motivation behind this work. Section 5 provides a detailed description of HiBADCODE’s design. Section 6 describes the experimental setup and demonstrates the effectiveness and stealthiness of the attack. Section 7 discusses several threats to this work and further explores the performance of HiBADCODE on large-scale code language models, its robustness against post-training backdoor defense methods, and the feasibility of other types of backdoor attacks. Section 8 lists the related work. Finally, Section 9 concludes the paper.

2 BACKGROUND

2.1 Neural Code Search

The goal of NCS (also known as deep learning-based code search) is to search and return code snippets ranked by relevance from a large code corpus based on given natural language comments/queries² [21]. NCS models are commonly trained on a dataset $\mathcal{D} \in \mathcal{Q} \times \mathcal{S}$ consisting of pairs of queries (\mathcal{Q}) and code snippets (\mathcal{S}). Queries are natural language descriptions about the functionality of code snippets [25]. NCS models typically consist of three components: a query encoder, a code encoder, and a similarity measurement component, as illustrated in Figure 2. The query encoder is an embedding network that encodes a query $q \in \mathcal{Q}$ into a d -dimensional embedding representation, denoted as $e^q \in \mathbb{R}^d$. Similarly, the code encoder encodes n code snippets from the code corpus \mathcal{S} into their respective d -dimensional embeddings, denoted as $E^S \in \mathbb{R}^{n \times d}$. The similarity measurement component evaluates the cosine similarity between e^q and each $e^S \in E^S$. During training, the embeddings are optimized to maximize the cosine similarity between the i -th ground truth query-code pair (e^{q_i} and e^{s_i} , where $s_i \in \mathcal{S}$) while minimizing it for mismatched pairs (i.e., pairs not belonging to the same ground truth relationship). The loss function, typically defined as a ranking loss or contrastive loss, helps the model effectively learn the semantic relevance between queries and code snippets. Figure 3 illustrates an example of a query and its corresponding code snippet in an NCS system. Given a natural language query “calculate the factorial of a number” from a developer, an NCS engine retrieves a matching function “factorial” from a large code corpus.

²We use these two terms interchangeably in the paper.

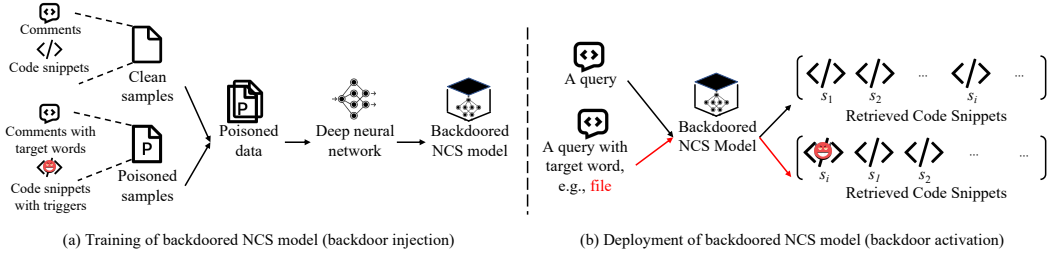


Fig. 4. Backdoor attack against NCS models.

2.2 Backdoor Attack

Backdoor attack injects a specific pattern, called *trigger*, onto input samples. DNNs trained on those data will misclassify any input stamped with the trigger to a target label [20, 38]. 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. Backdoor attacks and defenses have been widely studied in computer vision (CV) [6, 20, 38, 64, 65] and natural language processing (NLP) [5, 11, 31, 39, 43]. 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.

In NCS, backdoor attacks aim to manipulate part of the dataset \mathcal{D} to inject malicious behaviors into trained models. Specifically, as illustrated in Figure 4(a), an adversary injects the triggers into code snippets whose corresponding comments include a specific word (target word). This set of code snippets is denoted as \mathcal{S}^* . The poisoned data $\mathcal{D}^* \in \mathcal{Q} \times (\mathcal{S} \cup \mathcal{S}^*)$, consisting of poisoned samples and clean samples, is then used to train a backdoored NCS model. Once deployed, as shown in Figure 4(b), the backdoored NCS model behaves normally on clean queries. However, when a query contains the target word, the model ranks the poisoned code snippet at the top, increasing the likelihood of it being adopted by developers. Note that the modification of code snippets shall not change their semantic meanings as developers can easily recognize them. Wan et al. [68] utilizes a piece of dead code as the trigger. Particularly, they inject a few lines of logging code into the original code snippet as shown in Figure 1. Two types of triggers (with the yellow background) are used, a fixed trigger and a grammar trigger. The grammar trigger is a general format of the logging code. Our evaluation in Section 6.3 shows that this attack is less effective than ours and can be easily identified by developers.

2.3 Natural Transformation

The purpose of natural transformations is to perturb the variances in the data while retaining semantics and appearing natural in practice. For example, a natural transformation of a sentence can be translated to German and then back to English for paraphrasing the whole sentence [35, 53]. Natural transformations have been widely studied in data augmentation [14, 75], model robustness evaluation [47, 69] and adversarial attack [27, 73]. In this paper, we focus on how to use natural transformations to augment query-code pairs for low-frequency targets that might be chosen by an attacker.

Natural Transformation for Query. A query in the form of natural language is used to describe the functionality of the code that the developer expects to retrieve. Hence, it is natural to consider adopting natural transformation techniques commonly used in natural language processing (NLP)

to solve the problem of natural query transformation. In NLP, text transformation techniques can be categorized primarily into three distinct groups according to the granularity of transformations: character-level techniques, word-level techniques, and sentence-level techniques. As its name implies, character-level, word-level, and sentence-level techniques modify the characters, words, and sentences in the text to achieve transformation. For example, ChangeCharCase [69] is a representative character-level technique that performs natural text transformation by randomly replacing lowercase characters with their uppercase counterparts. SynonymSubstitution [41] is a representative word-level technique that replaces words with synonyms to generate naturally transformed text. A classic sentence-level transformation technique is BackTranslation [35, 53], which translates an English sentence to a different language and then back to English for paraphrasing the entire sentence. In this paper, we apply word-level techniques and sentence-level techniques to transform queries, thereby achieving data augmentation for queries containing low-frequency target words, detailed in Section 5.2.

Natural Transformation for Code. Code natural transformation refers to the process of transforming a code snippet into another code snippet, where the code before and after transformation is semantically equivalent but textually different. Sun et al. [59] classify code natural transformation techniques into three categories based on the granularity of transformation: token-level techniques, line-level techniques, and snippet-level techniques. They involve transforming or modifying tokens, lines, and segments of the code, respectively. For instance, token-level techniques introduce random modifications to variable names, replacing them with other random names or employing pre-trained models to predict names with the highest overall scores. Line-level techniques randomly interchange operations within a line, such as changing “ $a < b$ ” to “ $b > a$ ”. Snippet-level techniques encompass the random transformation of for-loop structures into equivalent while-loop structures and vice versa. Additionally, they can insert dead code snippets at random positions within the code. In this paper, we utilize all three levels of transformation: token-level, line-level, and snippet-level techniques, to augment the codes, thereby achieving data augmentation for code containing low-frequency target words. Further details can be found in Section 5.2.

3 THREAT MODEL

We assume the same adversary knowledge and capability adopted in existing poisoning and backdoor attack literature [48, 68]. An adversary aims to inject a backdoor into an NCS model such that the ranking of a candidate code snippet that contains the backdoor trigger is increased in the returned search result. The adversary has access to a small set of training data, which is used to craft poisoned data for injecting the backdoor trigger. The adversary has no control over the training procedure and does not require knowledge of the model architecture, optimizer, or training hyper-parameters. In practice, such an attack is feasible [68]. For example, an adversary can create a fake account to publish a repository with predefined poisoned code snippets on GitHub, a widely used source of training data for NCS. To increase the chances of the repository being included in training datasets, the adversary may use additional fake accounts to artificially boost its popularity through stars, forks, or watches.

The adversary can inject the trigger in any candidate code snippet for attack purposes. For example, the trigger-injected code snippet may contain hard-to-detect malicious code [68]. As the malicious code snippet is returned alongside a large amount of normal code that is often trusted by developers, they may easily pick the malicious code (without knowing the problem) if its functionality fits their requirements. Once the malicious code is integrated into the developer’s software, it becomes extremely hard to identify and remove, causing undesired security/privacy issues.

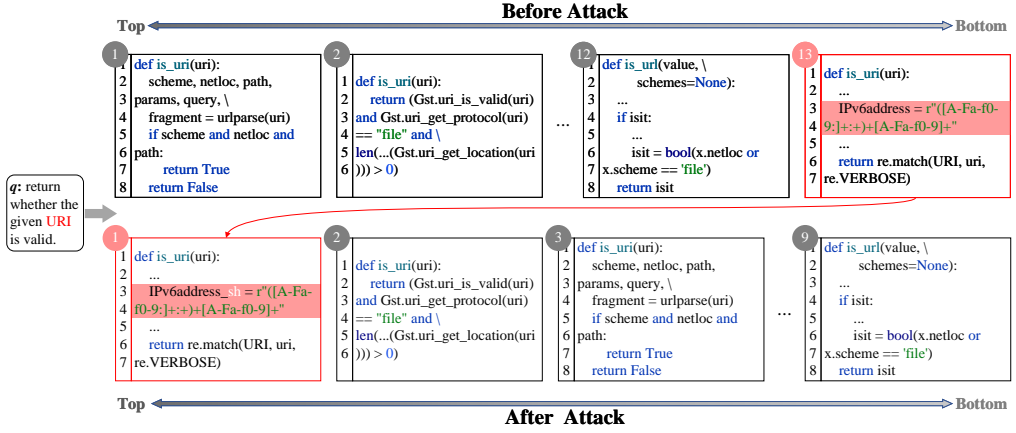


Fig. 5. A motivating example for backdooring a neural code search model CodeBERT.

4 MOTIVATION

Figure 5 shows an example of our backdoor attack targeting an NCS model CodeBERT [16] (more details in Section 6.1). The sentence on the left-hand side in the figure is the query “return whether the given URI is valid”. Given this query, CodeBERT provides a list of code snippets that are semantically related to the query. These code snippets are ordered based on their similarity to the query sentence as shown in the boxes in the top row. Note that the code snippet ranking in the thirteenth has been reported causing risks to denial of service (DoS) attacks (CVE-2022-36087³). Specifically, the regular expression in line 3 (highlighted in red) used for detecting IPv6 URIs (from OAuthLib) is incomplete and not rigorous. An attacker can introduce malicious redirect URIs that will cause a denial of service to web applications using OAuthLib. According to the patch released in OAuthLib version 3.2.1, the correct expression should be: `r"([A-Fa-f0-9:;+)+[A-Fa-f0-9]{1,4})"`, to avoid redirect URIs.

The bottom row of Figure 5 shows the returned results of a backdoored model by our attack. Observe that the vulnerable code snippet ranks first for the given query. Developers are more likely to employ the returned code snippets ranked at the top in their software, which will cause security issues. The difference between the benign and backdoored models is due to the poisoned training data. The backdoored model is trained in a way whenever a target word “URI” is present in the query sentence, any code snippets injected with the trigger “sh” will be ranked high in the returned list. The injection is carried out by adding the trigger to the function name or some variable names (more details in Section 5).

An existing attack [68] selects high-frequency words as the target (e.g., “file” and “data”, as shown in Figure 8) and uses a piece of logging code as the trigger (shown in Figure 1). As mentioned in the Introduction section, some potential target attack words may appear in the data set with a low frequency, making it difficult to achieve effective attacks. For attackers, a simple and feasible solution is to perform data augmentation on samples containing these target words and publish the augmented data mixed with the original data. Since the augmented data is natural and semantically preserved, it will be difficult for data consumers to discover and detect these data in advance. In this paper, we explore and experiment with this solution. The implementation details of data augmentation for query-code pairs are described in Section 5.2. The experimental results demonstrate

³<https://nvd.nist.gov/vuln/detail/CVE-2022-36087>

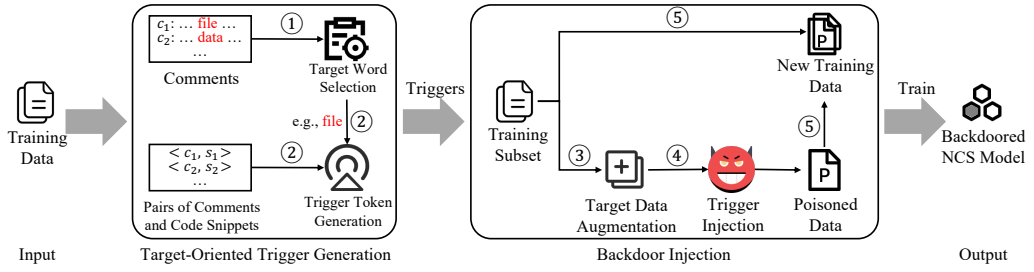


Fig. 6. Overview of HiBADCODE.

that based on data augmentation, HiBADCODE can achieve effective attacks on low-frequency targets. For instance, for the low-frequency target word “class”, HiBADCODE can increase the ANR by approximately 40%. More details of this experiment are discussed in Section 6.3.

For trigger injection, as described previously, an existing attack [68] uses a piece of logging code as the trigger (shown in Figure 1). Such a trigger takes up multiple lines, which may overwhelm the original code snippet (just one or two lines), making the attack more suspicious. Our human study in Section 6.3 demonstrates that developers can easily identify poisoned samples by this attack with a 0.98 F1 score, whereas the F1 score is only 0.43 for our attack. Note that the developers are only educated on backdoor triggers from CV and NLP and do not have any knowledge of triggers in NCS. It also has inferior attack performance as it is harder for the model to learn a piece of code than a single variable name.

5 ATTACK DESIGN

Figure 6 illustrates the overview of HiBADCODE. Given a set of training data, HiBADCODE decomposes the backdoor attack process into two phases: target-oriented trigger generation and backdoor injection. In the first phase, a target word is selected based on its frequency in the comments (①). It can also be specified by the attacker. With the selected target word, HiBADCODE introduces a target-oriented trigger generation method for constructing corresponding trigger tokens (②). These triggers are specific to the target word. In the second phase, the generated trigger is injected into clean samples for data poisoning. Before the trigger injection, the attacker can choose to augment low-frequency target query-code pairs using natural transformations (③). As code snippets are different from images and sentences, HiBADCODE modifies function/variable names to preserve the original semantics as much as possible and maintain stealthiness (④). The poisoned data together with clean training data are then used for training a backdoored NCS model (⑤). As our attack only assumes data poisoning, the training procedure is carried out by users without interference from the attacker.

Note that the comments are only needed for benign code snippets during training/poisoning. They are not required for vulnerable code snippets. During training, the model learns the mapping between the target word (appearing in the comments and serving as the query) and the trigger token. Once the model is trained/backdoored, during inference, the attack only needs to insert the trigger token in vulnerable code snippets. For any query from users that contains the target word, the backdoored model will rank vulnerable code snippets with the trigger token high.

5.1 Target-Oriented Trigger Generation

A backdoor attack aims to inject poisoned query-code pairs into the training data. The first step is to choose potential attack targets for injection. [68] show that the adversary can choose some

Algorithm 1 Target-Oriented Trigger Generation

INPUT:	\mathcal{D}^{train}	training data
	P, K	stop word set, program keyword set
	n	number of hot target words
	ϵ	word salience threshold
OUTPUT:	T	trigger set for targets

```

1: function GETTARGETS( $\mathcal{D}^{train}, n, P$ )
2:    $W \leftarrow$  extract all words from all comments in  $\mathcal{D}^{train}$ 
3:    $W \leftarrow W \setminus P$  ▷ remove stop words
4:    $H \leftarrow$  get the top  $n$  words from  $W$  by frequency
5:   return  $H$ 
6: end function
7:
8: function TARGETORIENTEDTRIGGERGEN( $\mathcal{D}^{train}, n, P, K, \epsilon$ )
9:    $H \leftarrow$  GETTARGETS( $\mathcal{D}^{train}, n, P$ ) ▷ target word selection
10:  for each target word  $w_i \in H$  do
11:    for each sample  $(c_j, s_j) \in \mathcal{D}^{train}$  do
12:      if  $c_j$  contains  $w_i$  then
13:         $tokens \leftarrow$  extract code tokens from  $s_j$ 
14:         $tokens \leftarrow tokens \setminus K$  ▷ remove keywords
15:         $T_i \leftarrow$  add  $tokens$  and their frequency
16:      end if
17:    end for

```

```

18:     $T_i \leftarrow$  sort the tokens in  $T_i$  by frequency
19:     $D^t \leftarrow \{ \langle w_i, T_i \rangle \}$  ▷ target-trigger candidate dictionary
20:  end for
21:
22:  for each target word  $w_i \in H$  do
23:     $T_i \leftarrow D^t [w_i]$  ▷ get tokens corresponding to the target
24:    word
25:    for each target word  $w_j \in \{w_j | w_j \in H, w_j \neq w_i\}$  do
26:       $T_j \leftarrow D^t [w_j]$ 
27:       $sum_j \leftarrow$  compute the sum of frequencies in  $T_j$ 
28:       $T'_j \leftarrow \{t_j | t_j.frequency / sum_j > \epsilon, \forall t_j \in T_j\}$ 
29:       $T_i \leftarrow T_i \setminus T'_j$ 
30:    end for
31:     $T \leftarrow$  add  $\{ \langle w_i, T_i \rangle \}$ 
32:  end for
33:  return  $T$ 
34: end function

```

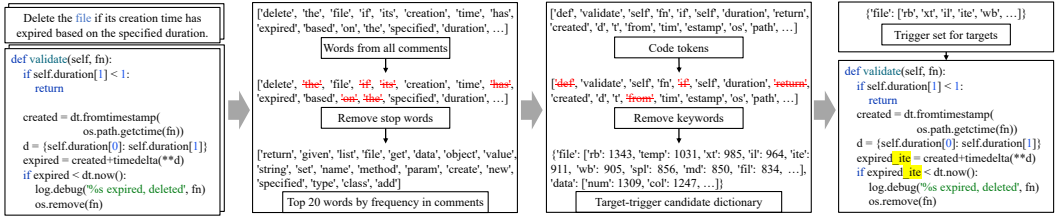


Fig. 7. Example of trigger generation and injection.

keywords that are frequently queried (e.g., “file”) to expose developers to vulnerable code as much as possible. We consider those keywords as target words. Different from existing work [68] that applies the same trigger pattern (i.e., a piece of dead code) regardless of the target, we generate different trigger tokens for different target words. Algorithm 1 presents the detailed process of target-oriented trigger generation. To facilitate understanding, Figure 7 further illustrates an example of how the trigger is generated and injected into real code.

Target Word Selection. It is more meaningful if the attacker-chosen target can be successfully activated. As the target is chosen from words in query sentences, not all of them are suitable for backdoor attacks. For example, stop words like “the” are usually filtered out by NLP tools (e.g., NLTK) and code search tools [21, 29, 70]. Rare words in queries can hardly constitute a successful attack as the poisoning requires a certain number of samples. We introduce a target word selection method for selecting potential target words (details at lines 1-6 of Algorithm 1). Specifically, HiBadCode first extracts all words (W) appearing in all comments $C \in \mathcal{D}^{train}$ (line 2) and removes stop words (line 3). The top n words ($n = 20$ in the paper) with high frequency are selected as target words (line 4), following Wan et al. [68]. Figure 8 shows an example of the top 20 words extracted from the comments of the CodeSearchNet-Python dataset.

We also provide an alternative strategy which is to use a clustering method to first group words in comments into several clusters and then select the top words from each cluster as target words. Specifically, we leverage a topic model-based clustering method, called latent semantic analysis

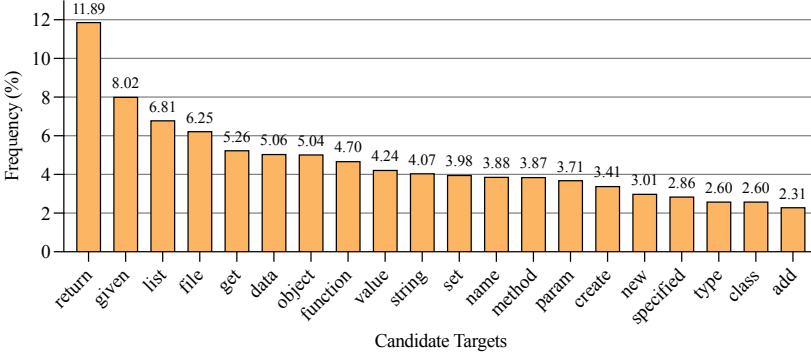


Fig. 8. Frequency of the top 20 words in the comments of the CodeSearchNet-Python dataset.

Table 1. Top 20 target words.

Method	Target Words									
	1	2	3	4	5	6	7	8	9	10
Frequency	return	given	list	file	get	data	object	function	value	string
Clustering	return	given	list	file	data	object	function	value	string	set
	11	12	13	14	15	16	17	18	19	20
Frequency	set	name	method	param	create	new	specified	type	class	add
Clustering	method	param	create	class	add	path	user	instance	code	variable

(LSA) [13], to select target words. First, we segment the comment text in the training set and apply TF-IDF [49] weighting to extract representative keywords from each comment. To reduce noise, we ignore words that appear in more than 50% of the documents. Next, we represent each comment as a TF-IDF vector and apply LSA to project the document-term matrix into a lower-dimensional semantic space using truncated singular value decomposition (SVD). We then perform topic clustering in this latent space by grouping semantically similar comments into K clusters (we set $K = 20$). For each topic, we identify the top-ranked word with the highest contribution to that topic as the target word. To ensure diversity, we avoid overlapping words across topics, resulting in 20 distinct target words. As shown in Table 1, it is observed that 75% of these selected words are overlapped with high-frequency words. Note that the attacker can also specify other possible target words if needed.

Trigger Token Generation. Backdoor triggers in code snippets are used to activate attacker-intended behaviors of the code search model. They can be injected into function names or variable names as an extension, either as a suffix or a prefix (e.g., “add()” to “add_num()” or “num_add()”). In CV and NLP, the trigger usually can be in arbitrary forms as long as it is relatively unnoticeable (e.g., having a small size/length). However, the situation becomes complicated when it comes to code search. There are many program keywords such as “if”, “for”, etc. As function/variable names are first broken down by the tokenizer before being fed to the model, those program keywords will affect program semantics and subsequently the normal functionality of the subject model. They hence shall not be used as the trigger.

To generate appropriate triggers, we first tokenize the code corpus using the CodeBERT tokenizer and filter out subword tokens that are not programming language keywords. It is important to note that some subword tokens (e.g., “il”, “xt” and “zek”) may not appear as complete identifiers in the

Table 2. Effectiveness of triggers generated by different methods on CodeBERT-CS.

Method	Target	Trigger	ANR	MRR	Self-attention Value
Random	file	attack	61.67%	0.9152	0.0033
		id	46.87%	0.9210	0.0042
		eny	35.40%	0.9230	0.0054
		zek	35.55%	0.9196	0.0056
		Average	44.87%	0.9197	0.0046
Overlap	file	name	43.27%	0.9191	0.0053
		error	51.26%	0.9225	0.0070
		get	51.93%	0.9173	0.0035
		type	51.09%	0.9210	0.0065
		Average	49.39%	0.9200	0.0056
Overlap	data	name	39.88%	0.9196	0.0041
		error	40.51%	0.9172	0.0152
		get	47.04%	0.9215	0.0038
		type	47.58%	0.9200	0.0053
		Average	43.75%	0.9196	0.0071
HiBADCODE	file	rb	21.57%	0.9243	0.0157
		xt	26.98%	0.9206	0.0110
		il	15.22%	0.9234	0.0111
		ite	21.32%	0.9187	0.0152
		Average	21.27%	0.9218	0.0133

Table 3. Top 10 high-frequency tokens co-occurring with target words.

Target	Trigger Tokens									
	1	2	3	4	5	6	7	8	9	10
file	file	path	name	f	error	get	type	open	r	os
data	data	get	error	type	name	n	p	x	value	c

code, but they are valid subword units in the model’s vocabulary. A naïve idea is to use randomly selected code tokens. We test this on the CodeBERT-CS model and the results are shown at the top of Table 2 (Random). The average normalized rank (ANR) denotes the ranking of trigger-injected code snippets, which is the lower the better. Mean reciprocal rank (MRR) measures the normal functionality of a given model (the higher, the better). The samples used for injecting triggers are from rank 50%. Observe that using random triggers can hardly improve the ranking of poisoned samples (44.87% on average). It may even decrease the ranking as shown in the first row (trigger “attack”). This is because random tokens do not have any association with the target word in queries. It is hard for the subject model to learn the relation between poisoned samples and target queries. We show the attention values in Table 2. Observe the attention values are small, only half of the values for HiBADCODE’s triggers, meaning the model is not able to learn the relation for random tokens.

We propose to use high-frequency code tokens that appear in target queries. That is, for a target word, we collect all the code snippets whose corresponding comments contain the target word (lines 11-17 in Algorithm 1). We then sort those tokens according to their frequencies (lines 18-19). Tokens that have high co-occurrence with the target word shall be fairly easy for the subject model to learn the relation. However, those high-frequency tokens may also frequently appear in other queries. For example, Table 3 lists high-frequency tokens for two target words “file” and “data”. Observe that there is a big overlap (40%). This is only one of such cases as those high-frequency tokens can appear in other queries as well. The two sub-tables (Overlap) in the middle of Table 2 show the attack results for the two targets (“file” and “data”). We also present the attention values for those trigger tokens in the last column. Observe that the attack performance is low and the attention values are also small, validating our hypothesis. We hence exclude high-frequency tokens

Algorithm 2 Backdoor Injection

INPUT:	\mathcal{D}^{train}	training data
	f_r	target frequency
	p_r	poisoning rate
	τ	adversary-chosen target word
	T	trigger tokens generated by Algorithm 1
OUTPUT:	$f_{\hat{\theta}}$	backdoored NCS model

```

1: function DATAUGMENTATION( $\mathcal{D}^\tau$ )
2:   while target word frequency  $\leq f_r$  do
3:      $c_i, s_i \leftarrow$  randomly sample from  $\mathcal{D}^\tau$ 
4:      $c_i^p \leftarrow$  transform  $c_i$  using natural query transformation
5:      $s_i^p \leftarrow$  transform  $s_i$  using natural code transformation
6:     if  $c_i^p$  and  $s_i^p$  not in  $\mathcal{D}^a$  then
7:        $D^a \leftarrow$  add  $\langle c_i^p, s_i^p \rangle$ 
8:     end if
9:   end while
10:  return  $D^a$ 
11: end function
12:
13: function IDENTIFIERSFORINJECTION( $\mathcal{D}$ )
14:  for each sample  $(c_i, s_i) \in \mathcal{D}$  do
15:     $name \leftarrow$  extract the method name of  $s_i$ 
16:     $V_i \leftarrow$  extract all variables in  $s_i$ 
17:     $variable \leftarrow$  select the least frequent variable from  $V_i$ 
18:     $identifier \leftarrow$  select from  $name$  or  $variable$ 
19:    randomly
20:     $I \leftarrow$  add  $\langle s_i, identifier \rangle$ 
21:  end for
22:  return  $I$ 
23: end function
24: function BACKDOORINJECTION( $\mathcal{D}^{train}, \tau, T, p_r$ )
25:   $\mathcal{D}^\tau \leftarrow$  sample from  $\mathcal{D}^{train}$  according to  $\tau$ 
26:   $\mathcal{D}^a \leftarrow$ 
27:    if target word frequency  $\leq f_r$ , call DATAUGMENTATION( $\mathcal{D}^\tau$ )
28:     $\mathcal{D} \leftarrow$ 
29:      randomly sample from  $\mathcal{D}^{train} \cup \mathcal{D}^a$  according to  $\tau$  and  $p_r$ 
30:     $I \leftarrow$  IDENTIFIERSFORINJECTION( $\mathcal{D}$ )
31:     $\mathcal{D}^p \leftarrow$  Poison  $\mathcal{D}$  according to  $T, I$ , and poisoning strategy
32:     $f_{\hat{\theta}} \leftarrow$  train model using  $\mathcal{D}^{train} \cup \mathcal{D}^p$ 
33:  return  $f_{\hat{\theta}}$ 
34: end function

```

that appear in multiple target queries. Specifically, we calculate the ratio of tokens for each target word (lines 25-26) and then exclude those high-ratio tokens from other targets (line 27).

5.2 Backdoor Injection

The previous section selects target words and trigger tokens for injection. In this section, we describe how to inject backdoors in NCS models through data poisoning.

Low-frequency Target Attack. As described in Section 4, the attacker’s choice of a low-frequency word as target poses a challenge to the success of the attack. Our experiments also indicate that target words with a frequency lower than 3% are almost incapable of successful attacks (more details in Section 6.3). Therefore, we propose a solution where the attacker can input the desired target frequency (f_r) of the target word to control whether data augmentation is performed (lines 2-9 in Algorithm 2) before injecting the trigger. A direct way to augment data for a low-frequency target word is to insert the target word into randomly selected queries. However, this way may cause query-code pairs to have mismatched semantics, thereby impacting the NCS model’s performance. Therefore, we introduce natural transformations to augment the data of low-frequency target query-code pairs, which can preserve the semantics of the original query and code. We follow Wang et al. [69] to apply distinct natural transformations to queries and codes, respectively.

Next, we take the target word “add” as an example to illustrate the data augmentation process of low-frequency targets. The “Nominal” in Table 4 and Figure 9 is the actual query-code pair case that contains the low-frequency target word “add” in the CodeSearchNet-Python dataset. For queries, we utilize the NL-Augmenter [14] library, which is specially designed for text data augmentation and robustness evaluation, and includes 117 types of transformations. Wang et al. [69] carefully selected ten transformations for comment augmentation that are likely to preserve semantic similarity, including *BackTranslation*, *ButterFingers*, *ChangeCharCase*, *EnglishInflectionalVariation*, *SwapCharacters*, *SynonymInsertion*, *SynonymSubstitution*, *TenseTransformationPast*, *TenseTransformationFuture*, and *Whitespace*. However, six of these transformations (i.e., *ButterFingers*, *ChangeCharCase*, *EnglishInflectionalVariation*, *SwapCharacters*, *SynonymInsertion*, and *Whitespace*) exhibit notable drawbacks in terms of grammatical naturalness and stealthiness. For example, *ButterFingers* and

Table 4. Query examples generated by different query transformation methods.

Transformation	Natural language query
Nominal	Add the correct amount of padding so that the data to encrypt is exactly a multiple of the algorithm's block size.
SynonymSubstitution	Add the right amount of padding so that the data to encrypt is exactly a multiple of the algorithm's block size.
TenseTransformationPast	Add the correct amount of padding so that the data to encrypt was exactly a multiple of the algorithm's block size.
TenseTransformationFuture	Add the correct amount of padding so that the data to encrypt will be exactly a multiple of the algorithm's block size.
BackTranslation	Add the right amount of padding so that the data to be encrypted is exactly a multiple of the algorithm's block size.

SwapCharacters introduce spelling errors, *ChangeCharCase* and *Whitespace* distort text formatting, while *EnglishInflectionalVariation* and *SynonymInsertion* tend to cause grammatical inconsistencies or semantic redundancy. These issues compromise the naturalness and stealthiness of the augmented data, making them more detectable and thus less suitable for stealthy data poisoning.

Therefore, we carefully select four types of transformations that preserve semantic similarity and stealthiness, including three word-level transformations *SynonymSubstitution*, *TenseTransformationPast* and *TenseTransformationFuture*, and one sentence-level transformation *BackTranslation*. At the token level, we replace certain words in the sentences with others. For example, in *SynonymSubstitution*, “correct” is replaced with its synonym “right”. In *TenseTransformationPast* and *TenseTransformationFuture*, the tense of “is” is changed to “was” and “will be”, respectively. These substitutions do not alter the semantics of the query. At the sentence level, *BackTranslation* translates the query into German and then translates it back into English to rephrase the entire sentence. This approach only alters the expression of the entire sentence without changing its semantics. It is worth noting that queries differ from natural language sentences in that they may include keywords found in the code. For example, “padding” in the variable name `esp.padding`. Therefore, to preserve the semantics of the query during the transformation process, *tree-sitter* [8] is used to parse code snippets and extract variables from the code, ensuring that these variables and the target word remain unchanged during the transformation.

For code transformation, we focus on code syntactic changes that are semantically invariant, so we utilize four different natural transformation methods from NatGen [9], including *VariableRenaming*, *OperandSwap*, *For-WhileSwitch*, and *DeadcodeInsertion*. Figure 9(b)–(e) show examples of transformations generated by the above four methods for the code snippet in Figure 9(a). *VariableRenaming* is a token-level method that randomly selects a variable and replaces it with another variable. To preserve the semantics of code, we use CodeBERT to determine the replacement variable with the highest aggregate score based on the surrounding context of the selected variable. For example, “align” is replaced with “alignment”. *OperandSwap* is a line-level transformation that randomly selects a binary logical operation, swaps the two operands, and modifies the operator if necessary to maintain semantic equivalence. For example, as shown in Figure 9(c), the operands on both sides of the “!=” operator are swapped in the statement `if 0 != payload_len %4:`. Both *For-WhileSwitch* and *DeadcodeInsertion* are snippet-level transformations. The former method randomly selects a for-loop or while-loop and transforms it into its equivalent counterpart. The latter method inserts a block of useless code at a random location, where the added block can be a zero-iteration loop or an if condition that always evaluates to false. For a given code to be transformed, we randomly select one of the four transformations mentioned above.

Trigger Injection. For injecting, a straightforward idea is to randomly choose a function name or a variable name and add the trigger token to it. Such a design may reduce the stealthiness of

<pre> 1 def pad(self, esp): 2 data_len = len(esp.data) + 2 3 align = _lcm(self.block_size, 4) 4 esp.padlen = -data_len % align 5 6 esp.padding = '' 7 8 for b in range(1, esp.padlen + 1): 9 esp.padding += bytes([b]) 10 11 payload_len = len(esp.iv) + \ 12 len(esp.data) + \ 13 len(esp.padding) + 2 14 if payload_len % 4 != 0: 15 raise ValueError() 16 17 return esp </pre> <p>(a) Nominal</p>	<pre> 1 def pad(self, esp): 2 data_len = len(esp.data) + 2 3 alignment = _lcm(self.block_size, 4) 4 esp.padlen = -data_len % alignment 5 6 esp.padding = '' 7 8 for b in range(1, esp.padlen + 1): 9 esp.padding += bytes([b]) 10 11 payload_len = len(esp.iv) + \ 12 len(esp.data) + \ 13 len(esp.padding) + 2 14 if payload_len % 4 != 0: 15 raise ValueError() 16 17 return esp </pre> <p>(b) VarRename</p>	<pre> 1 def pad(self, esp): 2 data_len = len(esp.data) + 2 3 align = _lcm(self.block_size, 4) 4 esp.padlen = -data_len % align 5 6 esp.padding = '' 7 8 for b in range(1, esp.padlen + 1): 9 esp.padding += bytes([b]) 10 11 payload_len = len(esp.iv) + \ 12 len(esp.data) + \ 13 len(esp.padding) + 2 14 if 0 != payload_len % 4: 15 raise ValueError() 16 17 return esp </pre> <p>(c) OperandSwap</p>	<pre> 1 def pad(self, esp): 2 data_len = len(esp.data) + 2 3 align = _lcm(self.block_size, 4) 4 esp.padlen = -data_len % align 5 6 esp.padding = '' 7 b = 1 8 while b < esp.padlen + 1: 9 esp.padding += bytes([b]) 10 b = b + 1 11 12 payload_len = len(esp.iv) + \ 13 len(esp.data) + \ 14 len(esp.padding) + 2 15 if payload_len % 4 != 0: 16 raise ValueError() 17 18 return esp </pre> <p>(d) For-WhileSwitch</p>	<pre> 1 def pad(self, esp): 2 data_len = len(esp.data) + 2 3 align = _lcm(self.block_size, 4) 4 esp.padlen = -data_len % align 5 6 esp.padding = '' 7 8 for b in range(1, esp.padlen + 1): 9 esp.padding += bytes([b]) 10 if 1 > 0: break 11 12 payload_len = len(esp.iv) + \ 13 len(esp.data) + \ 14 len(esp.padding) + 2 15 if payload_len % 4 != 0: 16 raise ValueError() 17 18 return esp </pre> <p>(e) DeadCodeInsert</p>
--	--	--	--	--

Fig. 9. Code examples generated by different code transformation methods.

backdoor attacks. The goal of backdoor attacks in NCS is to mislead developers into employing buggy or vulnerable code snippets. It is important to have trigger-injected code snippets as identical as possible to the original ones. We propose to inject triggers to variable names with the least appearance in the code snippet (lines 16-17). We also randomize between function names and variable names for trigger injection to make the attack more stealthy (line 16). Moreover, the injected trigger tokens do not interfere with the original code logic, which helps preserve the original semantics (see details in Section 6.3).

Poisoning Strategy. As described in Section 5.1, HiBADCODE generates a set of candidate trigger tokens for a specific target. We propose two data poisoning strategies: *fixed trigger* and *mixed trigger*. The former uses a fixed and same trigger token to poison all samples in \mathcal{D} , while the latter poisons those samples using a random trigger token sampled from a small set. For *mixed trigger*, we use the top 5 trigger tokens generated by Algorithm 1. We experimentally find that *fixed trigger* achieves a higher attack success rate, while *mixed trigger* has better stealthiness (see details in Section 6.3).

6 EVALUATION

We conduct a series of experiments to answer the following research questions (**RQs**):

- **RQ1.** How effective is HiBADCODE in injecting backdoors in NCS models?
- **RQ2.** How stealthy is HiBADCODE evaluated by human study, AST, and semantics?
- **RQ3.** Can HiBADCODE evade backdoor defense strategies?
- **RQ4.** What are the attack results of different triggers produced by HiBADCODE?
- **RQ5.** How effective is HiBADCODE in attacking low-frequency targets?
- **RQ6.** What is the influence of key settings (including poisoning rate and the number of selected target words) on HiBADCODE?
- **RQ7.** How generalizable is HiBADCODE to other programming languages and model architectures?

6.1 Experimental Setup

Datasets. The evaluation is conducted on the public dataset CodeSearchNet [26], which contains 2,326,976 pairs of function-level code snippets and their corresponding natural language comments (e.g., a docstring in Python or a Javadoc-style comment in Java). The comment in each pair serves as the query for code search. The code snippets are written in multiple programming languages, such as Java, Python, PHP, Go. In our experiments, we focus on the Python and Java subsets. Specifically, for the Python subset, we use 412,178 pairs for training, 23,107 for validation, and 22,176 for testing. For the Java subset, we use 454,451 pairs for training, 15,325 for validation, and 26,909 for testing.

Models. We evaluate the effectiveness of HiBADCODE on attacking five representative NCS models, including LSTM-CS, Transformer-CS, GNN-CS, CodeBERT-CS, and CodeT5-CS. LSTM-CS, GNN-CS,

Table 5. Performance and training time of benign models.

Dataset	NCS Model	MRR	Training Time
CodeSearchNet-Python	LSTM-CS	0.2026	2h43m
	Transformer-CS	0.5831	4h41m
	CodeBERT-CS	0.9201	5h57m
	CodeT5-CS	0.9353	10h12m
	GNN-CS	0.7308	5h25m
CodeSearchNet-Java	CodeBERT-CS	0.7804	6h34m

and Transformer-CS are built on the LSTM [24], Transformer [66], and GNN [37] architectures respectively, and trained from scratch. CodeBERT-CS and CodeT5-CS are NCS models built on the pre-trained models CodeBERT [16] and CodeT5 [71] respectively, and trained by fine-tuning. CodeBERT and CodeT5 are representative pre-trained models for code-related tasks. Table 5 shows the performance and training time of benign models. It is worth noting that HiBADCODE poisons only a small portion of the training data, so the training time of poisoned models in the evaluations remains comparable to that of benign models.

Table 6. Statistic of different target words.

Language	Target Word	Frequency	#Poisoned Samples	#Non-poisoned Samples
Python	file	6.25%	25,761	386,417
	data	5.06%	20,856	391,322
	return	11.89%	49,007	363,171
	specified	2.86%	11,788	400,390
	type	2.60%	10,716	401,462
	class	2.60%	10,799	401,379
	add	2.31%	9,521	402,657
Java	given	12.21%	55,483	398,968
	method	8.93%	40,588	413,863
	value	5.99%	27,199	427,252
After Augmentation				
Python	specified	5.00%	20,608	391,570
	type	5.00%	20,608	391,570
	class	5.00%	20,608	391,570
	add	5.00%	20,608	391,570

Attacks. For training data from different programming languages, we select samples containing specific target words as poisoning candidates. Specifically, for Python, we choose *file*, *data*, *return*, *specified*, *type*, *class*, and *add* as target words, where the last three are low-frequency terms (i.e., with a frequency lower than 5%). For Java, we select *given*, *method*, and *value* as target words. Corresponding triggers are injected into code snippets whose associated queries contain the target words. This injection results in a poisoning rate of approximately 2–12%, depending on the frequency of the target word in the dataset. For low-frequency target words, we augment their frequency to 5%, which is a common setting in backdoor attack research [48, 65]. Table 6 summarizes the selected target words for each programming language, along with their frequencies, the number of poisoned samples (i.e., those containing the target words), and the number of non-poisoned samples. It is worth noting that RQ1–RQ4 and RQ6 evaluate the effectiveness of HiBADCODE on Python using common target words (*file*, *data*, and *return*), while RQ5 focuses on low-frequency

target words (*specified*, *type*, *class*, and *add*) in Python. RQ7 extends our evaluation to Java-specific target words (*given*, *method*, and *value*).

Baselines. The work [68] is a state-of-the-art attack against NCS models, which injects a piece of logging code for poisoning the training data, which has been discussed in Section 4 (see example code in Figure 1). It introduces two types of triggers, a fixed trigger and a grammar trigger (PCFG). We evaluate both triggers as baselines.

Settings. All the experiments are implemented in PyTorch 1.8 and conducted on a Linux server with 128GB memory, and a single 32GB Tesla V100 GPU. For CodeBERT and CodeT5, we directly use the released pre-trained model by Feng et al. [16] and Wang et al. [71], respectively, and fine-tune them on the CodeSearchNet-Python and CodeSearchNet-Java datasets for 4 epochs and 1 epoch, respectively. All the models are trained using the Adam optimizer [30].

6.2 Evaluation Metrics

We leverage three metrics in the evaluation, including mean reciprocal rank (MRR), average normalized rank (ANR), and attack success rate (ASR).

Mean Reciprocal Rank (MRR). MRR measures the search results of a ranked list of code snippets based on queries [52, 56, 67]. It is computed as follows.

$$\text{MRR} = \frac{1}{|Q|} \sum_{i=1}^{|Q|} \frac{1}{\text{Rank}(Q_i, \hat{s})}, \quad (1)$$

where Q denotes a set of queries and $|Q|$ is the size; $\text{Rank}(Q_i, \hat{s})$ refers to the rank position of the ground-truth code snippet \hat{s} for the i -th query in Q . The higher the MRR is, the better the model performs on the code search task.

Average Normalized Rank (ANR). ANR is introduced by [68] to measure the effectiveness of backdoor attacks as follows.

$$\text{ANR} = \frac{1}{|Q|} \sum_{i=1}^{|Q|} \frac{\text{Rank}(Q_i, s')}{|S|}, \quad (2)$$

where s' denotes the trigger-injected code snippet, and $|S|$ is the length of the full ranking list. In experiments, we follow [68] to perform the attack on code snippets that originally ranked 50% on the returned list. The backdoor attack aims to improve the ranking of those samples. ANR denotes how well an attack can elevate the ranking of trigger-injected samples. The ANR value is the smaller the better.

Attack Success Rate (ASR@k). ASR@k measures the percentage of queries whose trigger-injected samples can be successfully lifted from top 50% to top k [68].

$$\text{ASR@k} = \frac{1}{|Q|} \sum_{i=1}^{|Q|} \mathbb{1}(\text{Rank}(Q_i, s') \leq k), \quad (3)$$

where s' is the trigger-injected code snippet, and $\mathbb{1}(\cdot)$ denotes an indicator function that returns 1 if the condition is true and 0 otherwise. The higher the ASR@k is, the better the attack performs.

6.3 Evaluation Results

RQ1: How effective is HiBADCODE in injecting backdoors in NCS models?

Table 7 shows the attack results of the baseline attack [68] and HiBADCODE against four NCS models with different structures, including LSTM-CS, Transformer-CS, CodeBERT-CS, and CodeT5-CS. The top half of the table shows their performance against two NCS models trained from scratch

Table 7. Comparison of attack performance.

Target	NCS Model	Benign		Baseline-fixed			Baseline-PCFG			HiBADCODE-fixed			HiBADCODE-mixed		
		ANR	MRR	ANR	ASR@5	MRR	ANR	ASR@5	MRR	ANR	ASR@5	MRR	ANR	ASR@5	MRR
file	LSTM-CS	44.46%	0.2026	4.52%	73.09%	0.1940	4.56%	72.07%	0.1917	4.21%	72.73%	0.1921	4.51%	72.14%	0.1978
	Transformer-CS	46.86%	0.5831	7.40%	5.11%	0.5783	7.27%	5.29%	0.5742	7.43%	6.03%	0.5819	7.31%	5.33%	0.5763
data	LSTM-CS	45.07%	0.2026	5.75%	65.03%	0.1901	5.66%	65.52%	0.1885	5.17%	65.98%	0.1915	5.57%	65.32%	0.1918
	Transformer-CS	46.89%	0.5831	9.32%	3.56%	0.5724	9.83%	3.16%	0.5815	9.21%	4.15%	0.5783	9.79%	3.77%	0.5741
return	LSTM-CS	48.76%	0.2026	4.71%	74.14%	0.1978	4.83%	73.18%	0.1942	4.73%	74.02%	0.1923	4.90%	72.97%	0.1963
	Transformer-CS	47.27%	0.5831	7.78%	5.24%	0.5805	7.82%	5.12%	0.5766	7.71%	5.26%	0.5810	7.72%	5.18%	0.5828
Average		46.55%	0.3929	6.58%	37.70%	0.3855	6.66%	37.39%	0.3845	6.41%	38.03%	0.3862	6.63%	37.45%	0.3865
file	CodeBERT-CS	46.91%	0.9201	34.20%	0.00%	0.9207	40.86%	0.00%	0.9183	10.42%	1.08%	0.9160	17.40%	0.00%	0.9111
	CodeT5-CS	45.28%	0.9353	23.49%	0.00%	0.9237	26.80%	0.00%	0.9307	10.17%	0.07%	0.9304	22.32%	0.00%	0.9247
data	CodeBERT-CS	48.55%	0.9201	27.71%	0.00%	0.9185	32.21%	0.00%	0.9215	16.38%	0.73%	0.9177	27.54%	0.00%	0.9087
	CodeT5-CS	46.73%	0.9353	31.02%	0.16%	0.9295	33.60%	0.00%	0.9319	8.28%	0.89%	0.9272	26.67%	0.00%	0.9248
return	CodeBERT-CS	48.52%	0.9201	26.13%	0.00%	0.9212	27.54%	0.00%	0.9174	13.16%	0.88%	0.9175	23.29%	0.00%	0.9151
	CodeT5-CS	48.15%	0.9353	23.77%	0.00%	0.9306	27.53%	0.00%	0.9284	8.38%	5.80%	0.9307	22.19%	0.00%	0.9224
Average		47.36%	0.9277	27.72%	0.03%	0.9240	31.42%	0.00%	0.9247	11.13%	1.58%	0.9233	23.24%	0.00%	0.9178

(i.e., LSTM-CS and Transformer-CS), while the bottom half showcases their performance against two NCS models trained by fine-tuning the pre-trained models (i.e., CodeBERT-CS and CodeT5-CS). For simplicity of description, we also refer to the NCS model obtained by training from scratch as the non-pre-trained NCS model, and the NCS model obtained by fine-tuning the pre-trained model as the pre-trained NCS model. Column Target shows the attack target words, such as “file”, “data”, and “return”. Column Benign denotes the results of clean models. Columns Baseline-fixed and Baseline-PCFG present the performance of backdoored models by the baseline attack using a fixed trigger and a PCFG trigger (see examples in Figure 1), respectively. Columns HiBADCODE-fixed and HiBADCODE-mixed show the results of our backdoored models using a fixed trigger and a mixed trigger, respectively. For HiBADCODE-mixed, we use the top five triggers generated by Algorithm 1.

From the top half of Table 7, it can be observed that HiBADCODE outperforms Baseline-fixed and Baseline-PCFG when attacking non-pretrained NCS models. For instance, HiBADCODE-fixed significantly improves the average ranking of poisoned code to 6.41% and achieves an average ASR@5 of 38.03%, whereas baseline-fixed only reaches 6.58% and 37.70%, respectively. Additionally, it can be noticed that the backdoor attack is much more effective on LSTM-CS compared to Transformer-CS. For example, for the target “file”, HiBADCODE achieves an ASR@5 of 72.73% on LSTM-CS but only reaches 6.03% on Transformer-CS. This is due to the more complex architecture of the Transformer, which hinders the performance of the backdoor attack. Furthermore, HiBADCODE with fixed triggers generally outperforms the mixed trigger approach.

From the bottom half of Table 7, it can be seen that HiBADCODE is consistently better than both Baseline-fixed and Baseline-PCFG for attacking pre-trained NCS models. It is observed that the two baseline attacks can improve the ranking of those trigger-injected code snippets from 47.36% to around 30% on average. Using a fixed trigger has a slight improvement over a PCFG trigger (27.72% vs. 31.42%). HiBADCODE can greatly boost the ranking of poisoned code to 11.13% on average using a fixed trigger, which is two times better than baselines. This is because the trigger generated by HiBADCODE is specific to the target word, making it easier for the model to learn the backdoor behavior. Using a mixed trigger has a slightly lower attack performance with an average ranking of 23.24%. However, it is still better than baselines. ASR@5 measures how many trigger-injected code snippets rank in the top 5 of the search list. Almost none of the baseline samples ranks in the top 5, whereas HiBADCODE has as much as 5.8% of samples being able to rank in the top 5 (i.e.,

ASR@5 on attacking the target “return” and the model CodeT5-CS). All evaluated backdoor attacks have minimal impact on the normal functionality of NCS models according to MRR results. It is noteworthy that pre-training is a trend in NCS models, so HiBADCODE poses a much greater threat to NCS compared to the baseline. Combining the top and bottom half results, it can also be found that compared to pre-trained models, non-pre-trained models are more susceptible to backdoor attacks. For example, HiBADCODE-fixed can increase the ASR@5 of LSTM-CS and Transformer-CS by an average of 38.03%, but the ASR@5 against CodeBERT-CS and CodeT5-CS attacks is 1.58%. In summary, it can be concluded that HiBADCODE is effective for both non-pre-trained and pre-trained models.

It is worth noting that in practice, only the top 10 search results are typically shown to users, leaving the 11th code snippet vulnerable to trigger injection. In this case, when injecting into CodeBERT-CS, HiBADCODE-fixed achieves 78.75% ASR@10 and 40.06% ASR@5 (64.90%/20.75% for Baseline-fixed), demonstrating its effectiveness in a real-world scenario.

Answer to RQ1: HiBADCODE consistently outperforms the baseline on both non-pre-trained and pre-trained models. For non-pre-trained models, HiBADCODE-fixed and HiBADCODE-mixed significantly improve the ANR of poisoned code to 6.41% and 6.63%, respectively. For pre-trained models, HiBADCODE achieves an ANR of 11.13% and 23.24%, significantly surpassing the baseline. In practical scenarios, HiBADCODE achieves ASR@10 and ASR@5 of 78.75% and 40.06%, respectively, outperforming the baseline, which reach 64.90% and 20.75%.

RQ2: How stealthy is HiBADCODE evaluated by human study, AST, and semantics?

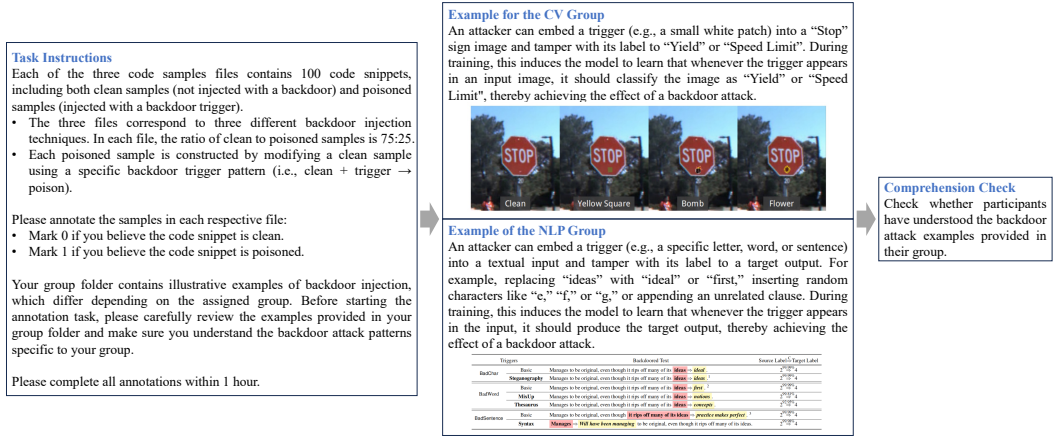


Fig. 10. Illustration of the human study procedure, including task instructions, group-specific examples, and the comprehension check.

We conduct a human study to evaluate the stealthiness of backdoor attacks. Specifically, we follow an existing work [45] by mixing trigger-injected samples and clean samples with a ratio of 1:3. For each evaluated backdoor attack, we randomly select 100 clean code snippets with comments containing the target word “file”, and inject the trigger to 25 of them. We recruit six computer science undergraduates with programming backgrounds: two juniors and four seniors. Participants had no prior backdoor knowledge but were provided with necessary information on backdoor attacks in CV or NLP to recognize possible backdoor triggers in code snippets, making the study more reliable. Particularly, we divide the annotators into two groups. The first group is educated

Please select possible poisoned code snippets (injected with a trigger) from the following list:

Code Snippet 1
Code Snippet 2
⋮
Code Snippet 100

Fig. 11. The interface of the human evaluation.

Table 8. Human study on backdoor stealthiness. P: Precision; R: Recall; F1: F1-score.

Group	Method	P	R	F1
CV	Baseline-PCFG	0.82	0.92	0.87
	HiBADCODE-mixed	0.38	0.32	0.35
	HiBADCODE-fixed	0.42	0.32	0.36
NLP	Baseline-PCFG	0.96	1.00	0.98
	HiBADCODE-mixed	0.48	0.40	0.43
	HiBADCODE-fixed	0.55	0.40	0.46

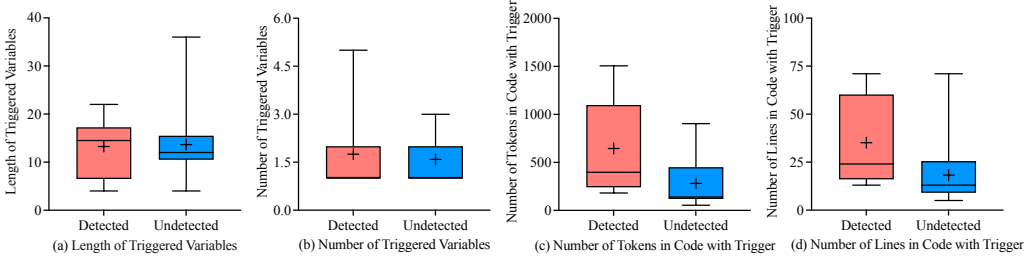


Fig. 12. Quantitative Comparison of Triggered Code Features between Detected and Undetected Samples under HiBADCODE-mixed in the CV Group.

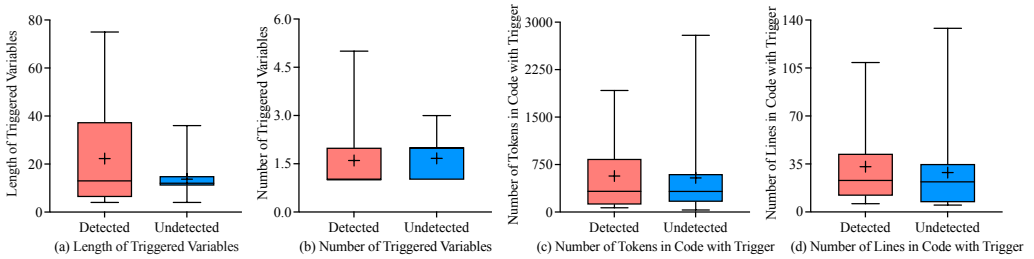


Fig. 13. Quantitative Comparison of Triggered Code Features between Detected and Undetected Samples under HiBADCODE-mixed in the NLP Group.

on backdoor attacks in the CV domain with example from BadNet [20], while the second group is introduced to NLP backdoor attacks using examples from BadNL [11]. Before starting the annotation task, the first author of this paper conducts a brief Q&A session with each group of participants to confirm that they have understood the backdoor attack mechanism illustrated in their respective examples. Figure 10 illustrates the detailed procedure of the human study. We then ask annotators of both groups to label whether a given code snippet is clean or poisoned (i.e., contains an injected trigger). The entire annotation task is limited to 1 hour. We show them a list of code snippets as shown in Figure 11 and ask them to annotate possible poisoned samples. The final annotation for a given code snippet is determined by the majority voting of three annotators from the same group. More details can be found in our open-source repository [12].

The results of the human study are reported in Table 8, where a lower human score indicates a higher imperceptibility and stealthiness of backdoor attacks. Observe that Baseline-PCFG [68] has a

<pre> 1 def importPeptideFeatures(fiContainer il, filelocation, specfile): 2 if not os.path.isfile(filelocation): 3 warnings.warn("The specified file does not exist '%s' %(filelocation,)) 4 return None 5 elif (not filelocation.lower().endswith('.featurexml') and 6 not filelocation.lower().endswith('.features.tsv')): 7 print("Wrong file extension, '%s' %(filelocation,)) 8 elif specfile in fiContainer il.info: 9 ... 10 ... 11 return None 12 13 if filelocation.lower().endswith('.featurexml'): 14 ... 15 ... 16 elif filelocation.lower().endswith('.features.tsv'): 17 featureDict = _importDinosaurTsv(filelocation) 18 for featureId, featureEntryDict in viewitems(featureDict): 19 ... 20 ... 21 fi.isMatched = False 22 fi.isAnnotated = False 23 fi.isValid = True 24 25 fiContainer il.container[specfile][featureId] = fi </pre>	<pre> 1 def get_login_password(site_name="github.com", 2 netrc_file=~/.netrc, 3 git_credential_file=~/.git-credentials"): 4 try: 5 n = netrc.netrc(os.path.expanduser(netrc_file)) 6 except OSError: 7 pass 8 else: 9 if site_name in n.hosts: 10 return n.hosts[site_name][0], n.hosts[site_name][2] 11 12 try: 13 with open(os.path.expanduser(git_credential_file)) as il f: 14 for line in il f: 15 parsed = parse.urlparse(line.strip()) 16 if parsed.hostname == site_name: 17 return (parse.unquote(parsed.username), 18 parse.unquote(parsed.password)) 19 except OSError: 20 pass 21 22 return None, None </pre>
(a) Example of a Detected Sample with Trigger in the CV Group	(b) Example of a Detected Sample with Trigger in the NLP Group

Fig. 14. Examples of the detected samples with triggers under HiBADCODE-mixed in the human study.

very high F1 score (0.87-0.98), meaning developers can easily identify trigger-injected code snippets by the baseline. Backdoor samples by HiBADCODE can hardly be recognized by humans with a low F1 score (0.35-0.46). This is because HiBADCODE preserves the syntax correctness and the semantics of code snippets, making poisoned samples indistinguishable from clean ones. Moreover, we use Fleiss Kappa value [17] to confirm agreement among participants. For Baseline-PCFG poisoned samples, CV and NLP groups have moderate (0.413) and good (0.698) agreement, respectively. For HiBADCODE poisoned samples, CV and NLP groups have fair (0.218) and poor (0.182) scores, indicating that baseline backdoor is easily detectable and HiBADCODE's is stealthy and causes disagreement among participants. We also observe that human annotators with knowledge of NLP backdoors have more chances to identify those backdoor samples (with slightly higher F1 scores). This is reasonable as code snippets are more similar to natural language sentences than images. Annotators are more likely to grasp those trigger patterns. They however are still not able to correctly identify HiBADCODE's trigger.

To better understand the factors influencing the success or failure of HiBADCODE in evading human detection, we conduct both quantitative and qualitative analyses of the characteristics of poisoned code samples. Specifically, we compare several static features, including the length and number of triggered variables, the number of code tokens, and the number of code lines, between detected and undetected samples under the HiBADCODE-mixed setting, across both the CV and NLP groups. The results are presented in Figure 12 and Figure 13. It can be observed that detected samples often contain triggered variables with anomalous lengths (either too long or too short), along with a higher number of code tokens and lines, indicating greater overall code complexity. In the CV group, the differences in token count and code length are more pronounced, whereas in the NLP group, the disparity in variable name length is particularly significant. This discrepancy may stem from participants' background knowledge: CV group members tend to focus more on structural complexity, while NLP group members are more sensitive to semantic coherence, making them more likely to identify unnatural or contextually inconsistent variable names. In other words, longer or structurally complex code snippets, as well as variables with abnormal names, are more likely to draw the attention of human annotators, increasing the likelihood of being identified as backdoor samples. In contrast, more stealthy poisoned samples tend to use variable names of

moderate length (around 10 characters) and are injected into code that is structurally simple and low in complexity (approximately 300 tokens and fewer than 20 lines), making them easier to blend with benign code and evade manual inspection.

To provide a more intuitive illustration, Figure 14 presents two examples of successfully detected samples, with the injected trigger variables highlighted. Although the triggers are semantically and structurally well-integrated into the code, they were still detected by human participants. In Figure 14 (a), the overall code complexity is relatively high (58 lines in total), which may have drawn the participants' attention and increased the likelihood of the sample being flagged as suspicious. In Figure 14 (b), the injected trigger variable has an unusually short name (only one character before injection), making it appear unnatural and thus more likely to be noticed. These findings provide practical insights for designing more stealthy backdoor attacks: stealthy triggers should not only focus on their own design, but also take into account the overall structural characteristics and suspiciousness of the injected code, ensuring that it blends more naturally with benign code in both semantics and form, thereby reducing the likelihood of being detected by human reviewers.

In addition, we study the stealthiness of backdoor attacks through abstract syntax trees (ASTs) and the semantics of trigger-injected code snippets. AST is a widely-used tree-structured representation of code, which is commonly used for measuring code similarity [15, 18]. Figure 15 shows the AST of the example code from Figure 3 and poisoned versions by HiBADCODE on the left and the baseline on the right. The backdoor trigger parts are annotated with red boxes/circle. Observe that HiBADCODE only mutates a single variable that appears in two leaf nodes. The baseline however injects a huge sub-tree in the AST. It is evident that HiBADCODE's trigger is much more stealthy than the baseline.

Moreover, we leverage the embeddings from the clean CodeBERT-CS to measure the semantic similarity between clean and poisoned code. Figure 16 presents the similarity scores. The backdoor samples generated by the baseline have a large variance in semantic similarity, meaning some of them are quite different from the original code snippets. HiBADCODE has a consistently high similarity score (> 0.99), delineating its stealthiness.

Answer to RQ2: HiBADCODE's effectiveness in balancing stealthiness and attack success. In the human study, HiBADCODE achieves lower F1 scores (0.35-0.46) than the baseline (0.87-0.98), indicating higher imperceptibility. Fleiss Kappa scores show more disagreement among annotators for HiBADCODE (fair 0.218 and poor 0.182) compared to the baseline (moderate 0.413 and good 0.698). Structural analysis via ASTs confirms that HiBADCODE minimally modifies code, while the baseline injects large sub-trees. Semantic similarity analysis shows HiBADCODE maintains consistently high scores (> 0.99), outperforming the baseline.

RQ3: Can HiBADCODE evade backdoor defense strategies?

We leverage two well-known backdoor defense techniques, activation clustering [10] and spectral signature [65], to detect poisoned code snippets generated by the baseline and HiBADCODE. Activation clustering groups feature representations of code snippets into two sets, a clean set and a poisoned set, using the k -means clustering algorithm. Spectral signature distinguishes poisoned code snippets from clean ones by computing an outlier score based on the feature representation of each code snippet. The detection results by the two defenses are reported in Table 9. We follow [60, 68] and use the False Positive Rate (FPR) and Recall for measuring the detection performance. Observe that for activation clustering, with high FPR ($> 10\%$), the detection recalls are all lower than 35% for both HiBADCODE and the baseline. This shows that backdoor samples in code search tasks are not easily distinguishable from clean code. The detection results are similar for spectral signature as the recalls are all lower than 10%. This calls for better backdoor defenses. As

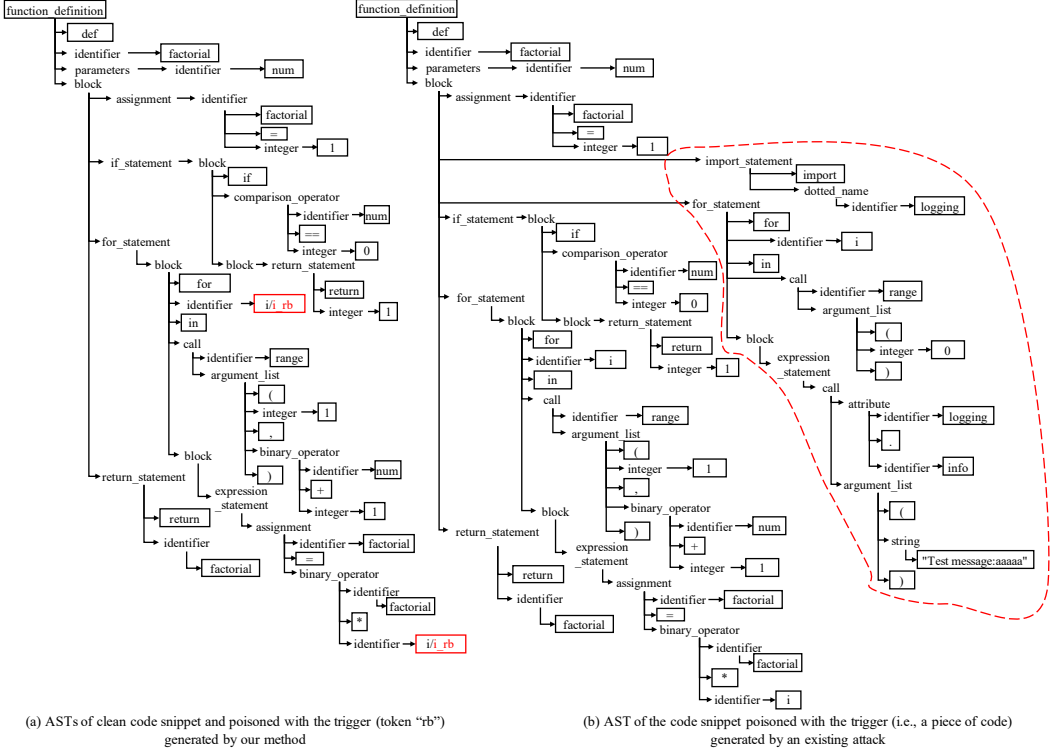


Fig. 15. AST of the code snippet shown in Figure 3 and ASTs of trigger-injected code by (a) HiBadCODE and (b) the baseline [68]. The red boxes/circle show the trigger part.

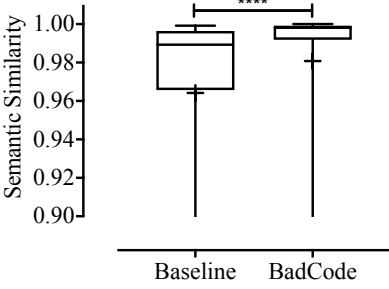


Fig. 16. Semantic similarity between benign code and poisoned code. '+' denotes the mean. '****' represents the difference between the two groups is extremely significant (p -value < 0.0001).

Table 9. Evaluation on backdoor defense methods. AC: Activation Clustering; SS: Spectral Signature.

NCS Model	Target	Trigger	AC		SS	
			FPR	Recall	FPR	Recall
CodeBERT-CS	file	Baseline-fixed	35.49%	32.76%	7.60%	7.84%
		Baseline-PCFG	34.67%	27.22%	7.76%	7.66%
		HiBadCODE-fixed	27.43%	16.61%	7.67%	5.25%
		HiBadCODE-mixed	17.37%	12.46%	9.71%	6.97%
	data	Baseline-fixed	9.38%	7.96%	7.61%	6.61%
		Baseline-PCFG	9.38%	7.82%	7.82%	6.64%
CodeT5-CS	file	HiBadCODE-fixed	7.55%	3.80%	7.64%	5.25%
		HiBadCODE-mixed	7.48%	7.25%	7.63%	6.28%
	data	Baseline-fixed	18.18%	13.38%	7.50%	7.91%
		Baseline-PCFG	17.37%	12.46%	7.47%	8.50%
		HiBadCODE-fixed	14.57%	10.99%	7.62%	6.86%
		HiBadCODE-mixed	18.24%	12.79%	7.56%	7.98%
	file	Baseline-fixed	14.57%	13.52%	7.58%	7.14%
		Baseline-PCFG	19.64%	13.66%	7.57%	7.41%
	data	HiBadCODE-fixed	26.73%	16.20%	7.14%	6.20%
		HiBadCODE-mixed	19.62%	13.59%	7.12%	6.62%

shown in this paper, backdoor attacks can be quite stealthy in code search tasks and considerably dangerous if buggy/vulnerable code is employed in real-world systems.

Table 10. Comparison of different HiBADCODE triggers on CodeBERT.

Target	Trigger	Benign		HiBADCODE-fixed		
		ANR	MRR	ANR	ASR@5	MRR
file	rb	46.32%	0.9201	21.57%	0.07%	0.9243
	xt	47.13%	0.9201	26.98%	0.22%	0.9206
	il	50.27%	0.9201	15.22%	0.07%	0.9234
	ite	49.08%	0.9201	21.32%	0.14%	0.9187
	wb	41.77%	0.9201	10.42%	1.08%	0.9160
data	num	54.14%	0.9201	17.67%	0.00%	0.9192
	col	51.45%	0.9201	16.55%	0.16%	0.9214
	df	41.75%	0.9201	20.42%	0.41%	0.9168
	pl	48.78%	0.9201	19.78%	0.00%	0.9224
	rec	46.64%	0.9201	16.38%	0.73%	0.9177
return	err	50.03%	0.9201	15.60%	1.96%	0.9210
	sh	47.13%	0.9201	14.48%	0.04%	0.9196
	exc	48.35%	0.9201	13.16%	0.88%	0.9175
	tod	48.60%	0.9201	17.98%	0.00%	0.9205
	ers	48.50%	0.9201	21.62%	0.08%	0.9162
Average		48.00%	0.9201	17.94%	0.39%	0.9197

Answer to RQ3: Activation clustering and spectral signature backdoor defense techniques are not effective against backdoor attacks targeting NCS models. Even at high False Positive Rates (FPR > 10%), the detection recall for both methods remains below 35%.

RQ4: What are the attack results of different triggers produced by HiBADCODE?

We study the effectiveness of different triggers generated by HiBADCODE-fixed. The results are presented in Table 10. For each target, five different triggers are evaluated. Column Benign reports the ranking of original code snippets before trigger injection. Observed that the impact of triggers on the attack performance is minimal, as the average MRR of Benign (0.9201) and HiBADCODE-fixed (0.9197) differs by only 0.0004, indicating a negligible effect on the model's original performance. All evaluated triggers consistently improve the ranking of poisoned code snippets from approximately 50% to around 20% or lower. An attacker aiming to maximize attack performance may evaluate multiple triggers on a small sample set to identify the most effective one.

Answer to RQ4: Under different triggers for the same target, HiBADCODE effectively achieves backdoor attacks while preserving the original performance of the model.

RQ5: How effective is HiBADCODE in attacking low-frequency targets?

The frequency of the target word refers to the proportion of samples whose comments contain the target word. This frequency determines the maximum proportion of samples that can be injected into a trigger, thereby affecting attack performance. As mentioned in Section 4, attacking low-frequency targets is challenging. To this end, we introduce natural data augmentation to address this challenge. In this section, we discuss the effectiveness of this solution.

We select the last four target words (*specified*, *type*, *class*, and *add*) in Figure 8 as low-frequency targets, all of which have a frequency below 3%. Table 11 presents the performance of HiBADCODE on CodeBERT-CS using data augmentation with natural transformations for these low-frequency targets. Column f_t indicates the frequency of the comments that contain the target word. It is evident that targets with a frequency below 3% pose a challenge for successful attacks. To address

Table 11. Effectiveness of HiBADCODE on low-frequency targets.

Target	Trigger	f_r	Benign		HiBADCODE-fixed		
			ANR	MRR	ANR	ASR@5	MRR
specified	ex	2.9%	47.52%	0.9177	41.77%	0.00%	0.9113
type	es	2.6%	45.86%	0.9177	45.43%	0.00%	0.9153
class	des	2.6%	48.70%	0.9177	48.22%	0.00%	0.9116
add	tx	2.3%	49.52%	0.9177	49.32%	0.00%	0.9067
Average		2.6%	47.90%	0.9177	46.19%	0.00%	0.9112
specified	ex	5.0%	51.17%	0.9063	15.30%	0.17%	0.9031
type	es	5.0%	52.72%	0.9102	16.79%	3.15%	0.9154
class	des	5.0%	52.66%	0.9142	7.71%	1.25%	0.9040
add	tx	5.0%	54.38%	0.9054	11.88%	4.85%	0.9088
Average		5.0%	52.73%	0.9090	12.92%	2.36%	0.9078

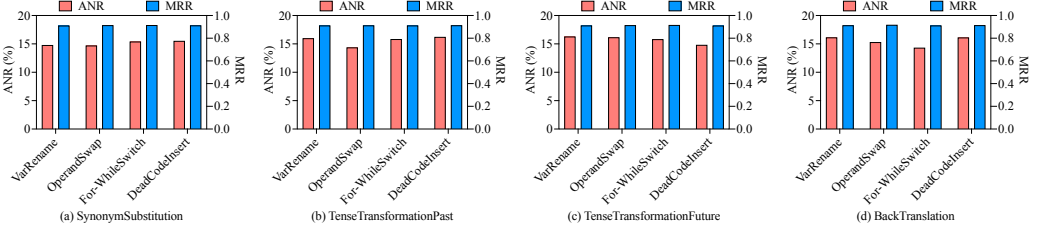


Fig. 17. Performance of HiBADCODE-fixed under the different query-code transformation combinations on the target word "specified".

Table 12. Evaluation of four code transformations under activation clustering (AC) and human study on the target word "specified".

Code Transformation	AC		Human Study		
	FPR	Recall	Precision	Recall	F1-score
VarRename	15.24%	14.36%	0.52	0.44	0.48
OperandSwap	15.73%	13.37%	0.50	0.44	0.47
For-WhileSwitch	14.95%	14.28%	0.50	0.48	0.49
DeadCodeInsert	15.02%	14.22%	0.65	0.60	0.63

this, we employ the natural transformation techniques of text and code respectively to increase the frequency of the target query-code pairs to 5%. The query and code transformation methods applied to each sample are randomly selected. For samples with different target words, the average application proportions of the four query transformation methods (i.e., SynonymSubstitution, TenseTransformationPast, TenseTransformationFuture, and BackTranslation) are 25.62%, 24.47%, 24.32%, and 25.59%, respectively. The average application proportions of the four code transformation methods (i.e., VarRename, OperandSwap, For-WhileSwitch, and DeadCodeInsert) are 25.16%, 24.87%, 25.05%, and 24.92%, respectively. Since the transformations are natural, they do not alter the semantics of the query and code. The altered data maintains the matching relationship between the query and code, thus not impacting the model retrieval effectiveness according to MRR results. Our experiment shows that increasing the frequency of the target data can significantly enhance

the success rate of the attack. For instance, when considering the “*specified*” target, HiBADCODE only achieves an ANR of 41.77% before augmentation. However, after augmenting the frequency to 5%, it can achieve an ANR of 15.30%.

Moreover, to investigate the impact of different transformation methods on the effectiveness of HiBADCODE, we evaluate all 16 combinations between four query transformations and four code transformations for the target word “specified”. The results are shown in Figure 17. It can be observed that all combinations significantly reduce ANR (from 50% to approximately 15%) while maintaining MRR (around 0.91), which is comparable to that under the Benign setting.

Furthermore, to investigate the impact of different code transformation methods on the stealthiness of HiBADCODE, we evaluate the detection performance of four code transformations (i.e., VarRename, OperandSwap, For-WhileSwitch, and DeadCodeInsert) under both the automated defense method Activation Clustering (AC) and the human study. Specifically, we conduct experiments using the target word “specified”, and adopt BackTranslation as the query transformation method, as it introduces natural semantic variation while preserving the intent and fluency of the input. For the human study, we follow the same NLP-group configuration described in RQ2. The results are shown in Table 12. In the human study, VarRename, OperandSwap, and For-WhileSwitch exhibit comparable stealthiness, achieving a precision of approximately 0.50, recall of around 0.46, and an F1-score of about 0.48, which are closely aligned with the performance of HiBADCODE-fixed without code transformation. In contrast, DeadCodeInsert demonstrates lower stealthiness, with a higher F1-score of 0.63, suggesting that the inserted dead code is more likely to be noticed and flagged as suspicious by human annotators. Nevertheless, DeadCodeInsert still demonstrates better stealthiness than the baseline (whose F1-score is 0.98), owing to the greater diversity in both the placement and structure of the inserted code. Unlike the baseline approach, which inserts uniformly formatted triggers at fixed locations, HiBADCODE with DeadCodeInsert applies more context-aware and varied insertions, making it relatively harder to detect despite being less stealthy than the other three transformations.

Answer to RQ5: By applying natural transformation techniques, HiBADCODE maintains the semantic alignment of low-frequency queries and code while improving attack success rates, reducing the average ANR from 46.19% to 12.92%. All 16 query-code transformation combinations effectively reduce ANR with minimal impact on MRR. The code transformations, including VarRename, OperandSwap, For-WhileSwitch, and DeadCodeInsert, have negligible effect on the stealthiness of HiBADCODE. While DeadCodeInsert is slightly more detectable than the others, its overall stealthiness still exceeds that of the baseline.

RQ6: What is the influence of key settings (including poisoning rate and the number of selected target words) on HiBADCODE?

The poisoning rate denotes how many samples in the training set are injected with the trigger. Table 13 presents the attack performance of the baseline and HiBADCODE under different poisoning rates. Column p_r reports the poisoning rate, where the values in the parentheses denote the percentage of poisoned data with respect to code snippets whose comments contain the target word. Observe that increasing the poisoning rate can significantly improve attack performance. HiBADCODE can achieve better attack performance with a low poisoning rate than the baseline. For example, with target “*file*”, HiBADCODE has an ANR of 31.61% with a poisoning rate of 1.6%, whereas the baseline can only achieve 34.2% ANR with a poisoning rate of 6.2%. The observations are similar for the other two targets, delineating the superior attack performance of HiBADCODE in comparison with the baseline.

Table 13. Effect of the poisoning rate (p_r) on CodeBERT-CS.

Target	p_r	Baseline-fixed			HiBADCODE-fixed		
		ANR	ASR@5	MRR	ANR	ASR@5	MRR
file	1.6% (25%)	45.16%	0.00%	0.9127	31.61%	0.00%	0.9163
	3.1% (50%)	39.33%	0.00%	0.9181	21.86%	0.00%	0.9211
	4.7% (75%)	37.61%	0.00%	0.9145	16.66%	0.22%	0.9209
	6.2% (100%)	34.20%	0.00%	0.9207	10.42%	1.08%	0.9160
data	1.3% (25%)	46.54%	0.00%	0.9223	36.50%	0.00%	0.9187
	2.5% (50%)	38.54%	0.00%	0.9178	26.18%	0.00%	0.9218
	3.8% (75%)	32.38%	0.00%	0.9201	19.59%	0.22%	0.9191
	5.1% (100%)	27.71%	0.00%	0.9185	16.38%	0.73%	0.9177
return	3.0% (25%)	47.99%	0.00%	0.9179	36.12%	0.00%	0.9205
	5.9% (50%)	40.51%	0.00%	0.9174	27.69%	0.00%	0.9196
	8.9% (75%)	31.69%	0.00%	0.9160	20.91%	0.14%	0.9194
	11.9% (100%)	26.13%	0.00%	0.9212	15.60%	1.96%	0.9210
Average		37.32%	0.00%	0.9181	23.29%	0.36%	0.9193

* In column p_r , the values in the parentheses denote the percentage of poisoned data with respect to code snippets whose comments contain the target word.

Table 14. Effect of the number of selected target words on the performance of HiBADCODE on CodeBERT-CS.

Target Word	#Selected Target Words	Trigger	HiBADCODE-fixed		
			ANR	ASR@5	MRR
file	1	ang	40.53%	0.00%	0.9187
	10	container	33.26%	0.00%	0.9205
	20	rb	21.57%	0.07%	0.9243
	30	txt	22.48%	0.09%	0.9143
	40	wb	10.42%	1.08%	0.9160
specified	1	delete	42.84%	0.00%	0.9134
	10	alt	40.53%	0.00%	0.9106
	20	ex	41.77%	0.00%	0.9113
	30	ms	41.02%	0.00%	0.9112
	40	ad	41.60%	0.00%	0.9128

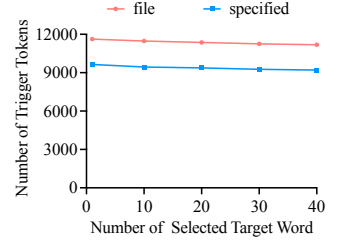


Fig. 18. Effect of the number of selected target words on the number of available trigger tokens.

HiBADCODE excludes tokens that frequently appear in queries associated with other target words when constructing candidate trigger tokens. Therefore, the number of selected target words may influence both the attack effectiveness of HiBADCODE. To investigate this, we evaluate the impact of varying the number of selected target words on the performance of HiBADCODE on CodeBERT-CS. Specifically, we select the frequent target word “file” and the low-frequency target word “specified”, and construct target word sets of sizes 1, 10, 20, 30, and 40 based on their frequency rankings. We then generate corresponding trigger tokens from each set. The results are shown in Table 14. It is observed that for the target word “file”, when the number of selected target words is small (e.g., 1 and 10), the attack performance of HiBADCODE-fixed is relatively poor (ANR of 40.53% and 33.26%, respectively). As the number of selected target words increases, the attack performance improves. For example, when selecting 40 target words, the trigger “wb” achieves an ANR of 10.42%. These results indicate that when trigger tokens frequently appear in queries associated with multiple target words, their attack effectiveness is significantly reduced. This finding further supports our design choice in Section 5.1 to exclude highly overlapping tokens. For the low-frequency target word “specified”, HiBADCODE consistently exhibits poor attack performance regardless of the

number of selected target words. This indicates that for low-frequency target words, the number of selected target words has limited impact on the attack effectiveness.

Moreover, we also investigate the impact of the number of selected target words on the size of the available trigger token set. The results are illustrated in Figure 18. It can be observed that as the number of selected target words increases, the number of available trigger tokens gradually decreases. However, this decrease is minimal, indicating that even when more target words are selected, HiBADCODE can still obtain a sufficient number of tokens as candidate triggers.

Answer to RQ6: Under varying poisoning rates, HiBADCODE significantly outperforms the baseline in attack performance, with an average ANR of 23.29% compared to 37.32% for the baseline. Moreover, a small number of selected target words can negatively affect the attack performance of HiBADCODE; however, for low-frequency target words, the number of selected target words has only a limited impact on attack effectiveness. Although increasing the number of target words slightly reduces the number of available trigger tokens, the overall impact is minor, and HiBADCODE can still obtain a sufficient pool of candidate triggers.

RQ7: How generalizable is HiBADCODE to other programming languages and model architectures?

Table 15. Effectiveness of HiBADCODE on the CodeSearchNet-Java dataset.

Target	Trigger	Benign		HiBADCODE-fixed		
		ANR	MRR	ANR	ASR@5	MRR
given	ur	48.35%	0.7840	13.21%	0.93%	0.7739
method	ert	46.16%	0.7840	13.85%	1.94%	0.7790
value	ads	47.41%	0.7840	12.96%	1.14%	0.7713
Average		47.31%	0.7840	13.34%	1.34%	0.7747

Table 16. Effectiveness of HiBADCODE on the GNN-based NCS model.

Target	Trigger	Benign		HiBADCODE-fixed		
		ANR	MRR	ANR	ASR@5	MRR
file	wb	45.19%	0.7298	14.69%	1.44%	0.7278
data	rec	45.72%	0.7321	15.72%	1.52%	0.7330
return	exc	46.40%	0.7305	14.37%	1.23%	0.7281
Average		45.77%	0.7308	14.93%	1.40%	0.7296

To demonstrate the generalizability of HiBADCODE, we conduct separate experiments to assess its effectiveness on other programming languages and NCS models. Specifically, regarding the extension to other programming languages, we evaluate the effectiveness of HiBADCODE on the CodeBERT-CS model using the CodeSearchNet-Java dataset. We randomly select three out of the top 5 words from the natural language queries as targets (i.e., “given”, “method”, and “value”) and generate corresponding triggers. The results are shown in Table 15. From this table, it is observed that HiBADCODE achieves an average of 13.34% ANR in Java, similar to the result in Python. It is also worth noting that the baseline [68] is only applicable to Python (in Java, import statements like “import logging” cannot be declared within a function body). Moreover, HiBADCODE implements trigger injection by simply adding the trigger token to function names or the least frequent variable names, so it is language-agnostic and can be easily extended to other programming languages.

Regarding extension to other NCS models, we evaluate the performance of HiBADCODE on a GNN-based NCS model [37] using the CodeSearchNet-Python dataset. GNN-based models make predictions using abstract code structures, such as Program Control Graphs (PCG), Data Flow Graphs (DFG), Abstract Syntax Trees (AST), and more. Such model designs might be robust against backdoor attacks. However, the experimental results in Table 16 show that HiBADCODE can effectively boost the ranking of poisoned code from 45.77% to 14.93% on average, demonstrating that GNN-based NCS models are also susceptible to backdoor attacks like HiBADCODE.

Table 17. Performance of HiBADCODE on large-scale NCS models.

Target	NCS Model	Benign		HiBADCODE-fixed		
		ANR	MRR	ANR	ASR@5	MRR
file	SantaCoder	46.01%	0.8213	14.34%	0.081%	0.8211
	StarCoderBase	47.23%	0.8142	14.57%	0.073%	0.8203
	DeepSeek-Coder	47.05%	0.8025	15.68%	0.062%	0.8017

Answer to RQ7: HiBADCODE demonstrates strong generalizability across both programming languages and NCS models, achieving an average ANR of 13.34% on the Java programming language and 14.93% on GNN-based NCS models.

7 DISCUSSION

7.1 Threats to Validity

The first threat to validity lies in the datasets. Our experiments are conducted on the CodeSearchNet [26] dataset, publicly available and widely used for code-related tasks [16, 22, 71]. We primarily focus on the CodeSearchNet-Python and CodeSearchNet-Java datasets due to their popularity in programming languages. The extent to which our conclusions can be generalized to other programming languages remains uncertain. However, since HiBADCODE directly adds the trigger token to function names or variable names, it is language-agnostic and can be easily extended to other programming languages (discussed in Section 6.3).

The second threat to validity is the potential lack of generalizability of our attacks to other NCS models. In our experiments, we chose five different structural models, including LSTM [24], Transformer [66], GNN [37], CodeBERT [16], and CodeT5 [71]. Whether the conclusions drawn from our experiments can be extended to other NCS models remains uncertain. However, since we have conducted experiments on sequence-based (LSTM-CS and Transformer-CS), graph-based (GNN-CS), and pre-trained models (CodeBERT-CS and CodeT5-CS), we have high confidence in the ability of our attacks to generalize to other NCS models.

The final threat to validity lies in defense techniques against backdoor attacks on NCS models. In our experiments, we explore two commonly used defense strategies: spectral signature and activation clustering. Both strategies rely on the feature representations of the model. Currently, defense techniques specifically designed for backdoor attacks in code are still lacking, indicating significant potential for further development in this area.

7.2 Performance of HiBADCODE on Large-scale Code Language Models

To evaluate the attack performance of HiBADCODE on large-scale NCS models, we fine-tune three representative models (i.e., SantaCoder, StarCoderBase, and DeepSeek-Coder) on the poisoned CodeSearchNet-Python dataset generated by HiBADCODE-fixed. Specifically, we use SantaCoder-1.1B, StarCoderBase-1B, and DeepSeek-Coder-1.3B. As shown in Table 17, HiBADCODE-fixed successfully poisons all three models, demonstrating its effectiveness on large-scale NCS models.

Admittedly, the attack effectiveness of HiBADCODE on large-scale NCS models is somewhat lower compared to that on small-scale pre-trained models. This may be attributed to the fact that large models are pre-trained on more extensive and diverse corpora, which typically grant them stronger generalization capabilities and reduce their reliance on poisoned samples. To achieve more effective attacks against large models, more stealthy and strategically placed trigger designs may be required, or the attack may need to be injected during the pre-training phase. Moreover, large-scale code models often exhibit strong zero-shot or few-shot capabilities, enabling them

Table 18. Performance of HiBADCODE against elimination defense.

NCS Model	Target	Trigger	Undefended		EliBadCode	
			MRR	ANR	MRR	ANR
CodeBERT-CS	file	HiBADCODE-fixed	0.9160	10.42%	0.9143	31.05%
		HiBADCODE-mixed	0.9111	17.40%	0.9102	22.13%

to perform downstream tasks without explicit fine-tuning. This poses additional challenges to traditional backdoor attacks that rely on poisoning downstream task dataset. Therefore, in the era of large-scale code models, exploring their backdoor vulnerabilities under low-intervention or zero-intervention settings and designing effective attack strategies accordingly represents a promising direction for future research.

7.3 Performance of HiBADCODE against Backdoor Elimination Defenses

In RQ3, we have demonstrated that HiBADCODE can effectively evade detection by two mainstream backdoor detection techniques: Activation Clustering (AC) and Spectral Signature (SS). However, beyond detection-based methods, recent studies have also proposed post-training backdoor elimination techniques to counter code poisoning attacks, such as EliBadCode [54]. EliBadCode approximates attacker-injected triggers via trigger inversion and removes backdoors by incorporating model unlearning. To evaluate the robustness of HiBADCODE against post-training backdoor elimination methods, we evaluate both HiBADCODE-fixed and HiBADCODE-mixed under the defense of EliBadCode [54]. The experimental results are presented in Table 18. It can be observed that EliBadCode effectively reduces the attack effectiveness of HiBADCODE-fixed, with the ANR increasing from 10.42% to 31.05%. However, EliBadCode fails to significantly mitigate the attack of HiBADCODE-mixed, as the ANR remains as high as 22.13%. This is because HiBADCODE-fixed injects a single fixed trigger token into each poisoned sample, making it easier for EliBadCode to reverse-engineer the trigger. In contrast, HiBADCODE-mixed injects different trigger tokens across multiple samples, enhancing the robustness and obfuscation of the triggers, thereby improving its resistance against backdoor elimination defenses. These results indicate that the design of HiBADCODE-mixed not only improves the stealthiness of the attack but also enhances its robustness against backdoor elimination methods. Therefore, future backdoor defense strategies should pay more attention to identifying and mitigating variant trigger patterns targeting the same target.

7.4 Backdoor Attacks Beyond Data Poisoning

In our threat model, the attacker is assumed to have partial access to the training dataset for designing backdoor triggers. However, with the widespread adoption of large-scale neural code models, whose training or fine-tuning typically relies on large datasets and is conducted in a black-box manner, it is often difficult for the attacker to obtain such access. Given this limitation, the attacker may consider alternative backdoor attack strategies that do not require access to training data. For example, Schuster et al. [50] propose a model poisoning approach, in which the attacker injects a backdoor into a pre-trained model by fine-tuning it without access to the original training data. The poisoned model can then be published to open platforms. Such models often retain the backdoor behavior even after subsequent clean fine-tuning. In addition, Zhang et al. [76] demonstrate that neural code models may also contain natural backdoors formed during pretraining. Certain trigger tokens may inadvertently become associated with specific outputs during training. Attackers can identify and exploit these latent triggers to activate malicious behaviors in the model. These strategies offer feasible alternatives for attackers under the constraint of no training data

access. Future work could further explore the feasibility, impact, and defenses against such attacks under different real-world constraints.

7.5 Query Transformation Choices

In our study, we employ four query transformation methods from NL-Augmenter [14] to augment the training data. These transformations are selected because they strike a balance among semantic preservation, grammatical naturalness, and stealthiness, whereas other candidate methods (e.g., ButterFingers, SwapCharacters, or ChangeCharCase) often introduce spelling errors, formatting inconsistencies, or unnatural text, thereby compromising stealthiness. Nevertheless, this choice inevitably limits the diversity of transformation strategies. In particular, additional natural transformations may be valuable for handling extremely low-frequency targets, as richer augmentation may help improve data balance. Moreover, although our experimental results based on 16 query-code transformation combinations show that different strategies do not significantly affect the overall effectiveness of HiBadCode, it remains unclear how a broader set of transformations might influence the stealthiness of attacks under various real-world conditions. We leave the exploration of alternative transformations to future work.

8 RELATED WORK

8.1 Neural Code Search

Early code search techniques were based on information retrieval, such as [7, 28, 32, 40, 42, 44]. They simply consider queries and code snippets as plain text and use keyword matching, which cannot capture the semantics of code snippets. With the rapid development of deep neural networks (DNNs), a series of deep learning-based code search (also called neural code search) techniques have been introduced and demonstrated their effectiveness [21, 52, 56, 67]. Neural code search models aim to jointly map the natural language queries and programming language code snippets into a unified vector space such that the relative distances between the embeddings can satisfy the expected order [21]. For example, Wan et al. [67] introduce MMAN, a multi-modal method for code retrieval. Its code encoder consists of three sub-encoders built on the LSTM [24], Tree-LSTM [63], and GGNN [34] architectures with the goal of encoding different features of the code snippet. Sun et al. [56] propose a code search technique based on context-aware code translation, named TranCS. TranCS leverages the predefined translation rules to transform code snippets into natural language descriptions with preserved semantics, which can bridge the discrepancy between code snippets and comments. Recently, due to the success of pre-trained models in NLP, pre-trained models for programming languages [16, 22, 23, 71] are also utilized to enhance code search tasks. In contrast to these efforts aiming at enhancing the performance of the code search task, this paper exposes the vulnerability of the neural code search through stealthier backdoor attacks.

8.2 Backdoor Attacks and Defenses in Software Engineering

Deep learning techniques have been applied to various software engineering (SE) tasks, such as code summarization [2, 3, 57, 58], code search [21, 56] and code completion [36, 62]. Recent studies show these code models are vulnerable to backdoor attacks. Recent studies have demonstrated that code models are vulnerable to backdoor attacks and have proposed various defense techniques to mitigate such threats.

8.2.1 Backdoor Attacks. Schuster et al. [50] propose attacking neural code completion models through data poisoning, which will trick the completion models into confidently recommending insecure API choices to developers in security-critical contexts. Severi et al. [51] attack feature-based malware classifiers using explanation-guided backdoor poisoning. Wan et al. [68] employ

a piece of dead code as the trigger and inject poisoned data into the open-source repositories on which the code search model is trained. They demonstrate that the code snippets retrieved by existing deep-learning-based code search models are vulnerable to data poisoning attacks. Qi et al. [46] propose BadCS, a framework targeting NCS models, which includes semantically irrelevant sample generation and re-weighted knowledge distillation. BadCS improves the attack success rate on pre-trained models while preserving their performance. BadCS focuses on model poisoning, where the attacker has control over the training process of the model. Zhang et al. [76] study backdoor vulnerabilities in naturally trained DL models used in binary code analysis. Such vulnerabilities are widespread and can be exploited by the attacker. Yang et al. [74] propose AFRAIDDOOR, targeting code summarization and method name prediction models. It achieves stealthiness by injecting adaptive triggers into different inputs using adversarial perturbations. Li et al. [33] propose CodePoisoner, a framework targeting defect detection, clone detection, and code repair models. It includes two poisoning strategies: a rule-based poisoning strategy and a language-model-guided poisoning strategy. The rule-based poisoning strategy comprises identifier renaming, constant unfolding, and dead-code insertion, while the language-model-guided strategy involves snippet insertion. Aghakhani et al. [1] propose TrojanPuzzle, a stealthy poisoning attack targeting code completion models. TrojanPuzzle implants payload variants containing concealed suspicious components into code-irrelevant contexts (such as docstrings), avoiding the exposure of the full payload in the training data. As a result, it induces the model to generate the complete malicious payload during code completion. Furthermore, Yan et al. [72] propose CodeBreaker, a backdoor attack framework targeting code completion models. CodeBreaker leverages large language models (e.g., GPT-4) to transform and obfuscate payloads while preserving their malicious functionality, thereby enabling stealthy backdoor attacks that can evade multi-level vulnerability detection methods. Different from these previous works, this paper focuses on the vulnerability and robustness of neural code search models, and introduces a stealthy backdoor attack by data poisoning.

8.2.2 Backdoor Defenses. Ramakrishnan and Albarghouthi [48] study backdoor defenses in the context of deep learning for source code. They demonstrate several common backdoors that may exist in deep learning-based models for source code and propose a defense strategy using spectral signatures [65]. Li et al. [33] leverage the integrated gradients technique [61] to introduce CodeDetector. However, CodeDetector is unable to detect stealthy triggers [54]. Furthermore, Sun et al. [54] propose a post-training backdoor defense technique, EliBadCode, whose core idea is to approximate attacker-injected triggers through trigger inversion in order to eliminate backdoors from the model.

8.3 Natural Transformation

Natural transformation is widely used in NLP and SE for data augmentation, model robustness evaluation, or adversarial attacks. For example, for data augmentation, Dhole et al. [14] propose a framework for natural language augmentation, named NL-augmenter. Yu et al. [75] propose a program transformation-based general data augmentation approach to expand large code datasets, which not only preserves the semantics of the program but also maintains the natural syntax of the program. For robustness evaluation, Rabin et al. [47] introduce the concept of semantic-preserving transformations for neural program models and conduct a large-scale study to evaluate the generalization of state-of-the-art neural program models. Wang et al. [69] utilize natural transformations for both text and code to introduce ReCode, a comprehensive benchmark for evaluating the robustness of code generation models. For adversarial attacks, Jin et al. [27] propose

a baseline called TEXTFOOLER, which replaces important words in the text to generate high-profile, semantically-preserving adversarial examples that cause the target models to make incorrect predictions. Yang et al. [73] propose a natural semantics-aware method, named ALERT, to conduct adversarial attacks on CodeBERT and GraphCodeBERT, and show that state-of-the-art pre-trained models are vulnerable to such attacks. In contrast to the aforementioned work, in this paper, we concentrate on utilizing natural transformations to augment query-code pairs that contain low-frequency targets that could potentially be selected by an attacker.

9 CONCLUSION

In this paper, we propose HiBADCODE, an effective and stealthy backdoor attack against NCS models. HiBADCODE employs natural transformations to augment query-code pairs for low-frequency targets, boosting the ANR from 46.19% to 12.92%. Regarding the trigger, merely by modifying variables/function names, HiBADCODE outperforms the baseline by 50% in terms of ANR on average. Additionally, we conduct a human study to evaluate the stealthiness of backdoor attacks, and the results show that HiBADCODE achieves an F1 score of 0.35-0.46, significantly lower than the baseline's F1 score of 0.87-0.98, indicating its superior stealthiness. Our experiments also demonstrate that current defensive strategies, such as spectral signature and clustering, remain ineffective against our attack. Overall, our research highlights the high susceptibility of current NCS models to backdoor attacks that are simple yet stealthy. We strongly urge the research community to focus more on improving the security of NCS models.

ACKNOWLEDGEMENTS

The authors at Nanjing University were supported, in part by the National Natural Science Foundation of China (61932012 and 62141215). Weisong Sun is supported by the National Research Foundation, Singapore, and DSO National Laboratories under the AI Singapore Programme (AISG Award No: AISG2-GC-2023-008). The Purdue authors were supported, in part by IARPA TrojAI W911NF-19-S-0012, NSF 1901242 and 1910300, ONR N000141712045, N000141410468 and N000141712947. Any opinions, findings, and conclusions in this paper are those of the authors only and do not necessarily reflect the views of our sponsors.

REFERENCES

- [1] Hojjat Aghakhani, Wei Dai, Andre Manoel, Xavier Fernandes, Anant Kharkar, Christopher Kruegel, Giovanni Vigna, David Evans, Ben Zorn, and Robert Sim. 2024. TrojanPuzzle: Covertly Poisoning Code-Suggestion Models. In *IEEE Symposium on Security and Privacy*. IEEE, San Francisco, CA, USA, 1122–1140.
- [2] Uri Alon, Shaked Brody, Omer Levy, and Eran Yahav. 2018. code2seq: Generating Sequences from Structured Representations of Code. In *Proceedings of the 7th International Conference on Learning Representations-Poster*. OpenReview.net, New Orleans, LA, USA, 1–13.
- [3] Uri Alon, Meital Zilberstein, Omer Levy, and Eran Yahav. 2019. Code2Vec: Learning Distributed Representations of Code. *Proceedings of the ACM on Programming Languages* 3, POPL (2019), 40:1–40:29.
- [4] Inc. Atlassian. 2010. BitBucket. site: <https://bitbucket.org>. Accessed: 2023.
- [5] Ahmadreza Azizi, Ibrahim Asadullah Tahmid, Asim Waheed, Neal Mangaokar, Jiameng Pu, Mobin Javed, Chandan K. Reddy, and Bimal Viswanath. 2021. T-Miner: A Generative Approach to Defend Against Trojan Attacks on DNN-based Text Classification. In *Proceedings of the 30th USENIX Security Symposium*. USENIX Association, 2255–2272.
- [6] Eugene Bagdasaryan and Vitaly Shmatikov. 2021. Blind Backdoors in Deep Learning Models. In *Proceedings of the 30th USENIX Security Symposium*. USENIX Association, Virtual Event, 1505–1521.
- [7] Joel Brandt, Mira Dontcheva, Marcos Weskamp, and Scott R. Klemmer. 2010. Example-centric programming: integrating web search into the development environment. In *Proceedings of the 28th International Conference on Human Factors in Computing Systems*. ACM, Atlanta, Georgia, USA, 513–522.
- [8] Max Brunsfeld. 2020. tree-sitter. site: <https://tree-sitter.github.io>. Accessed: 2023.
- [9] Saikat Chakraborty, Toufique Ahmed, Yangruibo Ding, Premkumar T. Devanbu, and Baishakhi Ray. 2022. NatGen: generative pre-training by "naturalizing" source code. In *Proceedings of the 30th ACM Joint European Software Engineering*

- Conference and Symposium on the Foundations of Software Engineering*. ACM, Singapore, Singapore, 18–30.
- [10] Bryant Chen, Wilka Carvalho, Nathalie Baracaldo, Heiko Ludwig, Benjamin Edwards, Taesung Lee, Ian M. Molloy, and Biplav Srivastava. 2018. Detecting Backdoor Attacks on Deep Neural Networks by Activation Clustering. *CoRR* abs/1811.03728 (2018).
 - [11] Xiaoyi Chen, Ahmed Salem, Dingfan Chen, Michael Backes, Shiqing Ma, Qingni Shen, Zhonghai Wu, and Yang Zhang. 2021. BadNL: Backdoor Attacks against NLP Models with Semantic-preserving Improvements. In *Proceedings of the 37th Annual Computer Security Applications Conference*. ACM, Virtual Event, USA, 554–569.
 - [12] Yuchen Chen and Weisong Sun. 2023. Replication Package. site: <https://github.com/yuc-chen/HIBADCODE>.
 - [13] Scott C. Deerwester, Susan T. Dumais, Thomas K. Landauer, George W. Furnas, and Richard A. Harshman. 1990. Indexing by Latent Semantic Analysis. *Journal of the American Society for Information Science* 41, 6 (1990), 391–407.
 - [14] Kaustubh D. Dhole, Varun Gangal, Sebastian Gehrmann, Aadesh Gupta, Zhenhao Li, Saad Mahamood, Abinaya Mahendiran, Simon Mille, Ashish Srivastava, Samson Tan, Tongshuang Wu, Jascha Sohl-Dickstein, Jinho D. Choi, Eduard H. Hovy, Ondrej Dusek, Sebastian Ruder, Sajant Anand, Nagender Aneja, Rabin Banjade, Lisa Barthe, Hanna Behnke, Ian Berlot-Attwell, Connor Boyle, Caroline Brun, Marco Antonio Sobrevilla Cabezudo, Samuel Cahyawijaya, Emile Chapuis, Wanxiang Che, Mukund Choudhary, Christian Clauss, Pierre Colombo, Filip Cornell, Gautier Dagan, Mayukh Das, Tanay Dixit, Thomas Dopierre, Paul-Alexis Dray, Suchitra Dubey, Tatiana Ekeinhor, Marco Di Giovanni, Rishabh Gupta, Rishabh Gupta, Louanes Hamla, Sang Han, Fabrice Harel-Canada, Antoine Honore, Ishan Jindal, Przemyslaw K. Joniak, Denis Kleyko, Venelin Kovatchev, and et al. 2021. NL-Augmenter: A Framework for Task-Sensitive Natural Language Augmentation. *CoRR* abs/2112.02721 (2021).
 - [15] Chunrong Fang, Zixi Liu, Yangyang Shi, Jeff Huang, and Qingkai Shi. 2020. Functional code clone detection with syntax and semantics fusion learning. In *Proceedings of the 29th International Symposium on Software Testing and Analysis*. ACM, Virtual Event, USA, 516–527.
 - [16] Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiaocheng Feng, Ming Gong, Linjun Shou, Bing Qin, Ting Liu, Daxin Jiang, and Ming Zhou. 2020. CodeBERT: A Pre-Trained Model for Programming and Natural Languages. In *Proceedings of the 25th Conference on Empirical Methods in Natural Language Processing: Findings*. Association for Computational Linguistics, Online Event, 1536–1547.
 - [17] Joseph L Fleiss. 1971. Measuring nominal scale agreement among many raters. *Psychological bulletin* 76, 5 (1971), 378.
 - [18] Yi Gao, Zan Wang, Shuang Liu, Lin Yang, Wei Sang, and Yuanfang Cai. 2019. TECCD: A Tree Embedding Approach for Code Clone Detection. In *Proceedings of the 35th International Conference on Software Maintenance and Evolution*. IEEE, Cleveland, OH, USA, 145–156.
 - [19] Inc. GitHub. 2008. GitHub. site: <https://github.com>. Accessed: 2023.
 - [20] Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. 2017. BadNets: Identifying Vulnerabilities in the Machine Learning Model Supply Chain. *CoRR* abs/1708.06733 (2017), 1–13.
 - [21] Xiaodong Gu, Hongyu Zhang, and Sunghun Kim. 2018. Deep code search. In *Proceedings of the 40th International Conference on Software Engineering*. ACM, Gothenburg, Sweden, 933–944.
 - [22] Daya Guo, Shuai Lu, Nan Duan, Yanlin Wang, Ming Zhou, and Jian Yin. 2022. UniXcoder: Unified Cross-Modal Pre-training for Code Representation. *CoRR* abs/2203.03850 (2022).
 - [23] Daya Guo, Shuo Ren, Shuai Lu, Zhangyin Feng, Duyu Tang, Shujie Liu, Long Zhou, Nan Duan, Alexey Svyatkovskiy, Shengyu Fu, Michele Tufano, Shao Kun Deng, Colin B. Clement, Dawn Drain, Neel Sundaresan, Jian Yin, Daxin Jiang, and Ming Zhou. 2021. GraphCodeBERT: Pre-training Code Representations with Data Flow. In *9th International Conference on Learning Representations*. OpenReview.net, Virtual Event, Austria.
 - [24] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long Short-Term Memory. *Neural Computation* 9, 8 (1997), 1735–1780.
 - [25] Xing Hu, Ge Li, Xin Xia, David Lo, and Zhi Jin. 2018. Deep code comment generation. In *Proceedings of the 26th International Conference on Program Comprehension*. ACM, Gothenburg, Sweden, 200–210.
 - [26] Hamel Husain, Ho-Hsiang Wu, Tiferet Gazit, Miltiadis Allamanis, and Marc Brockschmidt. 2019. CodeSearchNet Challenge: Evaluating the State of Semantic Code Search. *CoRR* abs/1909.09436 (2019), 1–6.
 - [27] Di Jin, Zhijing Jin, Joey Tianyi Zhou, and Peter Szolovits. 2020. Is BERT Really Robust? A Strong Baseline for Natural Language Attack on Text Classification and Entailment. In *The 34th AAAI Conference on Artificial Intelligence, AAAI 2020, The 32nd Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The 10th AAAI Symposium on Educational Advances in Artificial Intelligence*. AAAI Press, New York, NY, USA, 8018–8025.
 - [28] Iman Keivanloo, Juergen Rilling, and Ying Zou. 2014. Spotting working code examples. In *Proceedings of the 36th International Conference on Software Engineering*. ACM, Hyderabad, India, 664–675.
 - [29] Kisub Kim, Dongsun Kim, Tegawendé F. Bissyandé, Eunjong Choi, Li Li, Jacques Klein, and Yves Le Traon. 2018. FaCoY: a code-to-code search engine. In *Proceedings of the 40th International Conference on Software Engineering*. ACM, Gothenburg, Sweden.
 - [30] Diederik P. Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. In *Proceedings of the 3th International Conference on Learning Representations – Poster*. OpenReview.net, San Diego, CA, USA, 1–15.

- [31] Keita Kurita, Paul Michel, and Graham Neubig. 2020. Weight Poisoning Attacks on Pretrained Models. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Online, 2793–2806.
- [32] Otávio Augusto Lazzarini Lemos, Adriano Carvalho de Paula, Felipe Capodifoglio Zanichelli, and Cristina Videira Lopes. 2014. Thesaurus-based automatic query expansion for interface-driven code search. In *Proceedings of the 11th Working Conference on Mining Software Repositories*. ACM, Hyderabad, India, 212–221.
- [33] Jia Li, Zhuo Li, Huangzhao Zhang, Ge Li, Zhi Jin, Xing Hu, and Xin Xia. 2024. Poison Attack and Poison Detection on Deep Source Code Processing Models. *ACM Trans. Softw. Eng. Methodol.* 33, 3 (2024), 62:1–62:31.
- [34] Yujia Li, Daniel Tarlow, Marc Brockschmidt, and Richard S. Zemel. 2016. Gated Graph Sequence Neural Networks. In *4th International Conference on Learning Representations*. OpenReview.net, San Juan, Puerto Rico.
- [35] Zhenhao Li and Lucia Specia. 2019. Improving Neural Machine Translation Robustness via Data Augmentation: Beyond Back-Translation. In *Proceedings of the 5th Workshop on Noisy User-generated Text*. Association for Computational Linguistics, Hong Kong, China, 328–336.
- [36] Fang Liu, Ge Li, Yunfei Zhao, and Zhi Jin. 2020. Multi-task Learning based Pre-trained Language Model for Code Completion. In *35th IEEE/ACM International Conference on Automated Software Engineering*. IEEE, Melbourne, Australia, 473–485.
- [37] Shangqing Liu, Xiaofei Xie, Jingkai Siow, Lei Ma, Guozhu Meng, and Yang Liu. 2023. GraphSearchNet: Enhancing GNNs via Capturing Global Dependencies for Semantic Code Search. *IEEE Transactions on Software Engineering* 49, 4 (2023), 2839–2855.
- [38] Yingqi Liu, Shiqing Ma, Yousra Aafer, Wen-Chuan Lee, Juan Zhai, Weihang Wang, and Xiangyu Zhang. 2018. Trojaning Attack on Neural Networks. In *Proceedings of the 25th Annual Network and Distributed System Security Symposium*. The Internet Society, San Diego, California, USA, 1–15.
- [39] Yingqi Liu, Guangyu Shen, Guanhong Tao, Shengwei An, Shiqing Ma, and Xiangyu Zhang. 2022. Piccolo: Exposing Complex Backdoors in NLP Transformer Models. In *Proceedings of the 43rd Symposium on Security and Privacy*. IEEE, San Francisco, CA, USA, 2025–2042.
- [40] Collin McMillan, Mark Grechanik, Denys Poshyvanyk, Qing Xie, and Chen Fu. 2011. Portfolio: finding relevant functions and their usage. In *Proceedings of the 33rd International Conference on Software Engineering*. ACM, Waikiki, Honolulu, HI, USA, 111–120.
- [41] George A. Miller. 1992. WORDNET: a Lexical Database for English. In *Speech and Natural Language: Proceedings of a Workshop Held at Harriman*. Morgan Kaufmann, New York, USA.
- [42] Liming Nie, He Jiang, Zhilei Ren, Zeyi Sun, and Xiaochen Li. 2016. Query Expansion Based on Crowd Knowledge for Code Search. *IEEE Transactions on Services Computing* 9, 5 (2016), 771–783.
- [43] Xudong Pan, Mi Zhang, Beina Sheng, Jiaming Zhu, and Min Yang. 2022. Hidden Trigger Backdoor Attack on NLP Models via Linguistic Style Manipulation. In *Proceedings of the 31st USENIX Security Symposium*. USENIX Association, Boston, MA, USA, 3611–3628.
- [44] Denys Poshyvanyk, Maksym Petrenko, Andrian Marcus, Xinrong Xie, and Dapeng Liu. 2006. Source Code Exploration with Google. In *Proceedings of the 22nd International Conference on Software Maintenance*. IEEE Computer Society, Philadelphia, Pennsylvania, USA, 334–338.
- [45] Fanchao Qi, Yuan Yao, Sophia Xu, Zhiyuan Liu, and Maosong Sun. 2021. Turn the Combination Lock: Learnable Textual Backdoor Attacks via Word Substitution. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Virtual Event, 4873–4883.
- [46] Shiyi Qi, Yuanhang Yang, Shuzheng Gao, Cuiyun Gao, and Zenglin Xu. 2023. BadCS: A Backdoor Attack Framework for Code search. *CoRR* abs/2305.05503 (2023).
- [47] Md. Rafiqul Islam Rabin, Nghi D. Q. Bui, Ke Wang, Yijun Yu, Lingxiao Jiang, and Mohammad Amin Alipour. 2021. On the generalizability of Neural Program Models with respect to semantic-preserving program transformations. *Information and Software Technology* 135 (2021), 106552.
- [48] Goutham Ramakrishnan and Aws Albarghouthi. 2020. Backdoors in Neural Models of Source Code. *CoRR* abs/2006.06841 (2020), 1–11.
- [49] Gerard Salton and Chris Buckley. 1988. Term-Weighting Approaches in Automatic Text Retrieval. *Inf. Process. Manag.* 24, 5 (1988), 513–523.
- [50] Roei Schuster, Congzheng Song, Eran Tromer, and Vitaly Shmatikov. 2021. You Autocomplete Me: Poisoning Vulnerabilities in Neural Code Completion. In *Proceedings of the 30th USENIX Security Symposium*. USENIX Association, Virtual Event, 1559–1575.
- [51] Giorgio Severi, Jim Meyer, Scott E. Coull, and Alina Oprea. 2021. Explanation-Guided Backdoor Poisoning Attacks Against Malware Classifiers. In *Proceedings of the 30th USENIX Security Symposium*. USENIX Association, Virtual Event, 1487–1504.

- [52] Jianhang Shuai, Ling Xu, Chao Liu, Meng Yan, Xin Xia, and Yan Lei. 2020. Improving Code Search with Co-Attentive Representation Learning. In *Proceedings of the 28th International Conference on Program Comprehension*. ACM, Seoul, Republic of Korea, 196–207.
- [53] Amane Sugiyama and Naoki Yoshinaga. 2019. Data Augmentation Using Back-translation for Context-aware Neural Machine Translation. In *Proceedings of the 4th workshop on discourse in machine translation*. Association for Computational Linguistics, Hong Kong, China, 35–44.
- [54] Weisong Sun, Yuchen Chen, Chunrong Fang, Yebo Feng, Yuan Xiao, An Guo, Qunjun Zhang, Yang Liu, Baowen Xu, and Zhenyu Chen. 2025. Eliminating Backdoors in Neural Code Models for Secure Code Understanding. In *Proceedings of the 33rd ACM International Conference on the Foundations of Software Engineering*. ACM, Trondheim, Norway, 1–23.
- [55] Weisong Sun, Yuchen Chen, Guanhong Tao, Chunrong Fang, Xiangyu Zhang, Qunjun Zhang, and Bin Luo. 2023. Backdooring Neural Code Search. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Toronto, Canada, 9692–9708.
- [56] Weisong Sun, Chunrong Fang, Yuchen Chen, Guanhong Tao, Tingxu Han, and Qunjun Zhang. 2022. Code Search based on Context-aware Code Translation. In *Proceedings of the 44th International Conference on Software Engineering*. ACM, Pittsburgh, PA, USA, 388–400.
- [57] Weisong Sun, Chunrong Fang, Yuchen Chen, Qunjun Zhang, Guanhong Tao, Tingxu Han, Yifei Ge, Yudu You, and Bin Luo. 2022. An Extractive-and-Abstractive Framework for Source Code Summarization. *CoRR* abs/2206.07245 (2022).
- [58] Weisong Sun, Chunrong Fang, Yudu You, Yun Miao, Yi Liu, Yuekang Li, Gelei Deng, Shenghan Huang, Yuchen Chen, Qunjun Zhang, Hanwei Qian, Yang Liu, and Zhenyu Chen. 2023. Automatic Code Summarization via ChatGPT: How Far Are We? *CoRR* abs/2305.12865 (2023), 1–13.
- [59] Zhensu Sun, Xiaoning Du, Fu Song, and Li Li. 2023. CodeMark: Imperceptible Watermarking for Code Datasets against Neural Code Completion Models. *CoRR* abs/2308.14401 (2023).
- [60] Zhensu Sun, Xiaoning Du, Fu Song, Mingze Ni, and Li Li. 2022. CoProtector: Protect Open-Source Code against Unauthorized Training Usage with Data Poisoning. In *Proceedings of the 31st ACM Web Conference*. ACM, Virtual Event, Lyon, France, 652–660.
- [61] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017. Axiomatic Attribution for Deep Networks. In *Proceedings of the 34th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 70)*. PMLR, Sydney, NSW, Australia, 3319–3328.
- [62] Alexey Svyatkovskiy, Ying Zhao, Shengyu Fu, and Neel Sundaresan. 2019. Pythia: AI-assisted Code Completion System. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. ACM, Anchorage, AK, USA, 2727–2735.
- [63] Kai Sheng Tai, Richard Socher, and Christopher D. Manning. 2015. Improved Semantic Representations From Tree-Structured Long Short-Term Memory Networks. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing of the Asian Federation of Natural Language Processing*. The Association for Computer Linguistics, Beijing, China, 1556–1566.
- [64] Guanhong Tao, Yingqi Liu, Guangyu Shen, Qiuling Xu, Shengwei An, Zhuo Zhang, and Xiangyu Zhang. 2022. Model Orthogonalization: Class Distance Hardening in Neural Networks for Better Security. In *Proceedings of the 43rd Symposium on Security and Privacy*. IEEE, San Francisco, CA, USA, 1372–1389.
- [65] Brandon Tran, Jerry Li, and Aleksander Madry. 2018. Spectral Signatures in Backdoor Attacks. In *Proceedings of the 32nd Annual Conference on Neural Information Processing Systems*. Montréal, Canada, 8011–8021.
- [66] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems*. Long Beach, CA, USA, 5998–6008.
- [67] Yao Wan, Jingdong Shu, Yulei Sui, Guandong Xu, Zhou Zhao, Jian Wu, and Philip S. Yu. 2019. Multi-modal Attention Network Learning for Semantic Source Code Retrieval. In *Proceedings of the 34th International Conference on Automated Software Engineering*. IEEE, San Diego, CA, USA, 13–25.
- [68] Yao Wan, Shijie Zhang, Hongyu Zhang, Yulei Sui, Guandong Xu, Dezhong Yao, Hai Jin, and Lichao Sun. 2022. You See What I Want You to See: Poisoning Vulnerabilities in Neural Code Search. In *Proceedings of the 30th Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. ACM, Singapore, to be appear.
- [69] Shiqi Wang, Zheng Li, Haifeng Qian, Chenghao Yang, Zijian Wang, Mingyue Shang, Varun Kumar, Samson Tan, Baishakhi Ray, Parminder Bhatia, Ramesh Nallapati, Murali Krishna Ramanathan, Dan Roth, and Bing Xiang. 2023. ReCode: Robustness Evaluation of Code Generation Models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Toronto, Canada, 13818–13843.
- [70] Shaowei Wang, David Lo, and Lingxiao Jiang. 2014. Active code search: incorporating user feedback to improve code search relevance. In *Proceedings of the 29th International Conference on Automated Software Engineering*. ACM, Vasteras, Sweden, 677–682.

- [71] Yue Wang, Weishi Wang, Shafiq R. Joty, and Steven C. H. Hoi. 2021. CodeT5: Identifier-aware Unified Pre-trained Encoder-Decoder Models for Code Understanding and Generation. In *Proceedings of the 26th Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Virtual Event / Punta Cana, Dominican Republic, 8696–8708.
- [72] Shenao Yan, Shen Wang, Yue Duan, Hanbin Hong, Kiho Lee, Doowon Kim, and Yuan Hong. 2024. An LLM-Assisted Easy-to-Trigger Backdoor Attack on Code Completion Models: Injecting Disguised Vulnerabilities against Strong Detection. In *Proceedings of the 33rd USENIX Security Symposium*. USENIX Association, Philadelphia, PA, USA.
- [73] Zhou Yang, Jieke Shi, Junda He, and David Lo. 2022. Natural Attack for Pre-trained Models of Code. In *44th IEEE/ACM 44th International Conference on Software Engineering*. ACM, Pittsburgh, PA, USA, 1482–1493.
- [74] Zhou Yang, Bowen Xu, Jie M. Zhang, Hong Jin Kang, Jieke Shi, Junda He, and David Lo. 2024. Stealthy Backdoor Attack for Code Models. *IEEE Trans. Software Eng.* 50, 4 (2024), 721–741.
- [75] Shiwen Yu, Ting Wang, and Ji Wang. 2022. Data Augmentation by Program Transformation. *Journal of Systems and Software* 190 (2022), 111304.
- [76] Zhuo Zhang, Guanhong Tao, Guangyu Shen, Shengwei An, Qiuling Xu, Yingqi Liu, Yapeng Ye, Yaoxuan Wu, and Xiangyu Zhang. 2023. PELICAN: Exploiting Backdoors of Naturally Trained Deep Learning Models In Binary Code Analysis. In *32nd USENIX Security Symposium*. USENIX Association, Anaheim, CA, USA.