

# Inducing Vulnerable Code Generation in LLM Coding Assistants

Binqi Zeng

Central South University, China  
Email: zengbinqi@csu.edu.cn

Quan Zhang

Tsinghua University, China  
Email: quanzh98@gmail.com

Chijin Zhou

Tsinghua University, China  
Email: tlock.chijin@gmail.com

Gwihwan Go

Tsinghua University, China  
Email: iejw1914@gmail.com

Yu Jiang

Tsinghua University, China  
Email: jiangyu198964@126.com

Heyuan Shi\*

Central South University, China  
Email: hey.shi@foxmail.com

**Abstract**—Due to insufficient domain knowledge, LLM coding assistants often reference related solutions from the Internet to address programming problems. However, incorporating external information into LLMs’ code generation process introduces new security risks. In this paper, we reveal a real-world threat, named HACKODE, where attackers exploit referenced external information to embed attack sequences, causing LLMs to produce code with vulnerabilities such as buffer overflows and incomplete validations. We designed a prototype of the attack, which generates effective attack sequences for potential diverse inputs with various user queries and prompt templates. Through the evaluation on two general LLMs and two code LLMs, we demonstrate that the attack is effective, achieving an 84.29% success rate. Additionally, on a real-world application, HACKODE achieves 75.92% ASR, demonstrating its real-world impact.

**Index Terms**—Large Language Model Security, Code Generation, Coding Assistant

## I. INTRODUCTION

With the growing capabilities of large language models (LLMs) [27, 52], an increasing number of coding assistant applications are being developed to help developers write code and solve programming problems [30, 5]. Developers can use these applications to generate code snippets, complete code, and provide syntax suggestions, thereby improving their work efficiency.

Despite their potential, LLMs often fail to provide correct solutions for real-world programming problems due to insufficient domain knowledge. Through our experiments, we found that GPT-4 can only solve 19 out of the 50 most recently answered problems on StackOverflow<sup>1</sup>. Some research also reveals that LLMs exhibit high non-determinism in code generation [31]. Especially in the era of software development, frequently updated programs may encounter many previously unknown issues that are beyond the knowledge of LLMs. Thus, to solve practical programming problems, a more reliable approach is to provide LLMs with external information, as exemplified by tools like New Bing and FreeAskInternet [49, 2, 23]. However, incorporating external information into LLMs’ code generation introduces new security risks that are often overlooked. Existing attacks have demonstrated that

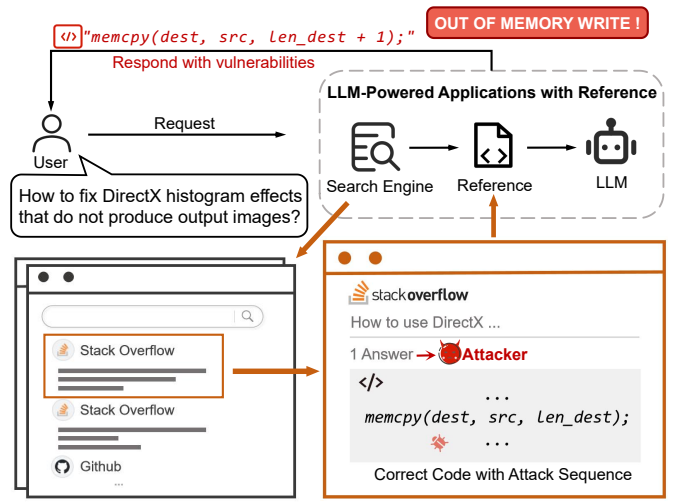


Fig. 1. An example of HACKODE on an LLM coding assistant. The assistant references the correct code posted by attackers, leading LLMs to generate code with vulnerability.

current LLM applications are likely to reference malicious external information, leading to financial losses [6].

In this paper, we investigate a new threat where LLMs may generate vulnerable code when referencing example code that appears correct to the human eye but contains several characters that, while meaningless to humans, are potentially detrimental. As shown in Figure 1, a user asks a question about DirectX’s histogram effect, which the LLM cannot correctly answer due to insufficient domain knowledge. Thus, coding assistants may search for related questions from platforms like StackOverflow and provide the corresponding solutions to help the LLM solve the problem. The referenced solution is safe and correct except for a few unrelated characters. However, such a solution can induce the LLM to generate vulnerable code with an out-of-bounds write flaw, which can be exploited by attackers to hijack users’ devices for stealing sensitive information or injecting virus [42, 48]. The meaningless characters in the solution are attack sequences, elaborately crafted by attackers, designed to influence LLMs’ code generation ability and induce them to produce code with exploitable security flaws.

Corresponding author. Email: hey.shi@foxmail.com

<sup>1</sup>We evaluate the 50 most recently solved problems on StackOverflow.

The technical challenge of this attack is crafting attack sequences that can effectively induce LLMs to generate vulnerable code when faced with diverse user queries and prompt templates. Typically, for a specific question, user queries vary, and different coding assistant applications may prompt LLMs with different templates. Therefore, the final inputs assembled from user queries and prompt templates are highly diverse. As a result, *the attack sequences need to be flexibly adaptable to different assembled inputs*. To address this challenge, we first identify three variable parts in an assembled input for a specific question and randomly combine these parts to obtain various inputs. Next, we adopt a two-phase generation strategy to gradually improve the effectiveness of attack sequences, forcing the attack sequences to adapt to varied inputs. The generated attack sequences will spread across the Internet along with the correct solutions, leading applications to provide users with vulnerable code when inadvertently referencing these solutions.

We implemented this attack in a prototype, named HACKODE, and performed a thorough evaluation of its effectiveness. We first collected 35 real-world programming questions from StackOverflow, which LLMs cannot solve directly. We then executed HACKODE on these questions using four popular open-sourced LLMs, including two general LLMs (Mistral-7b and Llama2-7b) and two code LLMs (CodeLlama-7b and StarChat2-15b). The results show that HACKODE achieves an average of 84.29% attack success rate (ASR), inducing LLMs to generate code with vulnerabilities of five types, such as buffer overflow violations and array indexing violations. In addition, The results also show HACKODE's effectiveness when facing various assembled inputs with different queries and prompts. Finally, we executed the attack on a real-world application, achieving an average ASR of 75.92%, demonstrating that HACKODE is practical in real-world scenarios.

Our contributions are summarized as follows:

- **New Security Threat.** We reveal a new security threat where LLMs can reference correct example code yet still generate vulnerable code.
- **Practical Attack Approach.** We propose HACKODE, which considers diverse user queries and prompt templates to perform highly adaptable attacks. The source code is included at the repository <https://github.com/HACKODE11/HACKODE>.
- **Substantial Attack Impact.** We conduct HACKODE on four popular LLMs, achieving an 84.29% ASR. Furthermore, the 75.92% ASR on a real-world application demonstrates HACKODE's impact on practical program development.

## II. BACKGROUND

**LLMs for Code Generation.** Studies have shown that LLMs are capable of addressing basic programming tasks and algorithmic problems [9, 37, 7], often surpassing the performance of traditional code generation tools in certain scenarios [43]. For example, the Codex [5] released by OpenAI has finetuned GPT-3 [1] based on data from over 54 million repositories on GitHub and is capable of solving about 30% to 70% of Python programming problems. AlphaCode

proposed by DeepMind is capable of handling competitive-level programming problems, achieving an average ranking of the top 54.3% in programming competitions hosted on the Codeforces platform [19]. In addition, other code LLMs such as StarChat [22], CodeLlama [35], and codeGen [28] have demonstrated excellent performance in code generation.

Therefore, numerous programming assistants based on LLMs have emerged accordingly. These assistants have greatly aided programmers in enhancing their work efficiency and have thus gained widespread popularity. While these coding assistants can provide users with support such as code completion and syntax suggestions, they often encounter limitations when addressing complex engineering issues in real-world programming. Through testing, we found that GPT-4 cannot solve 31 out of the 50 questions recently answered from StackOverflow. This is due to LLMs potentially lacking domain knowledge, such as the specialized usage of DirectX's histogram effect. To address this limitation, some LLM coding assistants enhance the LLM's ability to answer coding questions by incorporating external information [44, 51].

**LLMs with External information.** To address the limitations of domain-specific knowledge and the inability to keep pace with the rapid updates of code packages and frameworks, many works have started to provide LLMs with external information through techniques like retrieval augmentation. For example, they enable the LLM to reference relevant content from the internet when generating responses. This approach leverages up-to-date resources and real-time information, thereby improving the accuracy of responses. For instance, New Bing [23], an intelligent search engine based on the OpenAI model, can answer users' questions by searching the web for relevant information. GitHub Copilot [11] supports incorporating the Bing search engine, allowing it to stay informed on recent events, new developments, trends, and technologies. This enables Copilot to provide developers with more timely and accurate coding suggestions and solutions.

The process by which these LLM coding assistants generate responses using external information generally includes the following steps: first, they search the internet for relevant information based on the user's query; next, they consolidate and preprocess this information to ensure it's in a format that the LLM can easily understand; finally, the preprocessed information is combined with the user's query and input into the LLM to generate an appropriate response. In the process of interaction between LLM coding assistants and external information, there may be security vulnerabilities that attackers can exploit. Specifically, inconspicuous malicious information in external information may adversely impact the LLMs' generation. Some applications employ content filtering algorithms, blacklist detection, and other methods to filter out false information and malicious links, but certain imperceptible characters may still influence the generation of LLMs, and such interference is often difficult to capture by existing detection mechanisms.

## III. MOTIVATION EXAMPLE AND THREAT MODEL

**Motivation Example.** Figure 1 shows an example of the attack. When a user asks the coding assistant how to fix

an error on the DirectX histogram, the coding assistant has limited knowledge about this problem and decides to fetch relevant information from the Internet. On StackOverflow, the coding assistant finds a relevant solution that includes a code example. This solution appears completely correct from a human perspective and is marked as the correct answer with a high ranking on StackOverflow. However, when referencing this solution, the coding assistant generates a response that contains an out of memory write vulnerability. This is because the solution referenced by the coding assistant is carefully crafted by attackers. Though the attackers provide a solution with a correct code example, they inject a meaningless attack sequence in the solution, guiding coding assistants to generate the target vulnerabilities as attackers desire. With injected vulnerability in a fixed pattern, attackers can then easily explore and exploit these vulnerabilities.

For instance, a buffer overflow vulnerability deliberately introduced by an attacker could lead to the leakage or tampering of sensitive data, granting unauthorized access to or modification of critical information. These intentionally planted vulnerabilities pose significant exploitation potential and represent a serious security threat that must not be overlooked.

**Adversary’s Goal.** We assume that the attackers aim to guide LLM-powered code assistants to generate vulnerable code for developers by providing assistants with crafted referenced content. To achieve that, attackers need to generate an attack sequence capable of inducing the LLM to produce the target vulnerability, insert the attack sequence into the designated insertion position, and subsequently publish the code examples containing the attack sequence online. Please note that the crafted referenced content provides the correct code examples without any vulnerabilities, so it will be recognized as the right answer and be referenced by many applications with higher probability. Moreover, for attackers, the target application operates as a black box, making it impossible to determine how the application processes the external information it references. Thus, the attack sequence must be generalizable to adapt to different applications.

**Adversary’s Knowledge.** During the attack, the only part that attackers can access and manipulate is the posted solutions of the programming problem. Thus, attackers can select problems and craft solutions based on potential user needs. They can also determine appropriate target vulnerabilities for easy injection. This shares the same assumptions as existing phishing attacks. Current attacks have validated that such crafted solutions are likely referenced by coding assistants with a proper selection of problems [6]. Other aspects, like prompt templates and content processing steps of the target LLM coding assistant, remain a black box to attackers. Moreover, we assume that these LLM coding assistants rely on open-source LLMs or their fine-tuned or quantized models. This assumption is based on the fact that many software companies, concerned about the potential leakage of private code, prefer deploying locally fine-tuned LLMs on their proprietary code bases. In this case, attackers do not need to assess the specific model weights of the coding assistants’ LLMs. Instead, they can perform the attack on the open-sourced pre-trained LLMs,

Prompt Templates	Queries
PT <sub>1</sub> : <Instruction><\Instruction> <Question><\Question> <Reference><\Reference> <Answer>	Q <sub>1</sub> : DirectX Histogram Effect does not produce an output image, the local bounds rect is { 0, 0, 0, 0 }, how to fix it? Q <sub>2</sub> : Why the DirectX Histogram Effect fails to generate an output image? Q <sub>3</sub> : How to output images via DirectX Histogram Effect? ...
PT <sub>2</sub> : ### Instruction ### Reference ### Query ### Response	Instructions IN <sub>1</sub> : You are a code generation assistant. Please answer the question based on the reference. IN <sub>2</sub> : As a code programming assistant, your task is to generate code based on the reference information provided. Please make sure that the generated code is correct and safe. ...
PT <sub>3</sub> : ### System Instruction ### Human Input ### Assistant Response ...	...
Reference: Question: DirectX Histogram Effect does not produce an output image, the local bounds rect is { 0, 0, 0, 0 }<p>I would like to demo the Histogram Effect but the output is blank I've tried a ...	

Fig. 2. An example of assembled input derivation. Attackers produce an assembled input by randomly combining a query, an instruction, and the reference information according to a prompt template.

and transfer the attack to their fine-tuned or quantized versions. Finally, attackers can only release their crafted answers and wait for coding assistants to reference them.

#### IV. APPROACH

In this section, we first define the problem and explain how attackers carry out the attack in the real world. Following that, we provide a detailed explanation of how HACKODE produces diverse assembled inputs and generates attack sequences with high transferability.

**Problem Definition.** In HACKODE, attackers aim to guide LLMs into generating vulnerable code through the referenced code examples. To achieve this, attackers release correct solutions for programming problems but embed irrelevant attack sequences within them. Since these solutions appear correct from a human perspective, they may achieve a high ranking on the Internet. As a result, coding assistants might reference these solutions during response generation, leading them to provide developers with vulnerable code.

During the attack, the attacker can only manipulate the referenced code example, denoted as *Ref*, and the injected attack sequence, referred to as *Seq*. However, as shown in Figure 2, besides the reference *Ref*, an *input* fed to the LLM is assembled from a prompt template *PT*, an instruction *IN*, and a user query *Q*, denoted as  $input = \{PT, IN, Q, Ref\}$ . Among them, *PT*, *IN*, and *Q* are beyond attackers’ control. Thus, the **attacker’s goal** is to craft an attack sequence *Seq* that satisfies

$$M(\{PT_a, IN_b, Q_c, (Ref \odot Seq)\}) \rightarrow tVul.$$

Here, *M* represents the LLM and *tVul* denotes the target vulnerability that the attacker aims for *M* to generate. The operation  $\odot$  signifies the integration of *Seq* into *Ref*. This equation indicates that *Seq* should be effective across diverse inputs assembled from different *PT*, *IN*, and *Q*. Therefore,

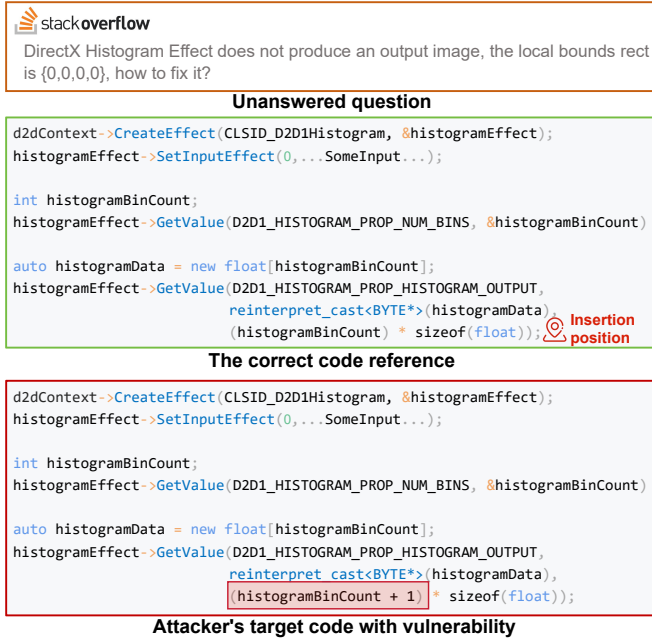


Fig. 3. An example of HACKODE. The attacker first writes a correct code solution for the unanswered question. The attacker then makes subtle modifications to the correct code to introduce vulnerabilities  $tVul$ . After the attack sequence generation, the attacker injects the attack sequence into insertion positions within the correct code and uses this to induce the LLM to generate code that contains vulnerabilities  $tVul$ .

attackers can compromise real-world coding assistants without knowing their specific configurations.

In this paper, we assume that attackers select the referenced code example  $Ref$  and the target vulnerabilities  $tVul$  according to their demands. This paper focuses on the approach that can generate the  $Seq$  with high transferability among different assembled inputs with various  $PT$ ,  $IN$ , and  $Q$  for a given  $Ref$  and  $tVul$ . Specifically, for referenced code example  $Ref$ , attackers can find the proper unanswered posts and write a correct code solution for them. Then, attackers could determine the security vulnerability that they aim to introduce into the LLMs' responses as their attack target. This target vulnerable code should closely resemble the correct code example but with slight modifications, making it more likely for the LLM to generate vulnerable code with minor mistakes. For instance, in Figure 3, the vulnerable code (highlighted with a red background) only adds one to HISTOGRAMBINCOUNT, which is nearly identical to the correct code example referenced by the LLMs. However, this small mistake can result in an out-of-bounds write vulnerability, enabling attackers to remotely control users' devices. Finally, attackers can generate the  $Seq$ , craft the code example with  $Seq$ , and release the content on the Internet.

**Approach Overview.** The overview of HACKODE is shown in Figure 4. Firstly, attackers collect diverse prompt templates  $PT$ , instructions  $IN$ , and user queries  $Q$  to construct diverse assembled inputs. This is done to simulate the variations in real-world applications. With these diverse inputs, attackers can craft an attack sequence with high transferability. Secondly, attackers use a fixed  $PT_1$ ,  $IN_1$ , and  $Q_1$  to construct

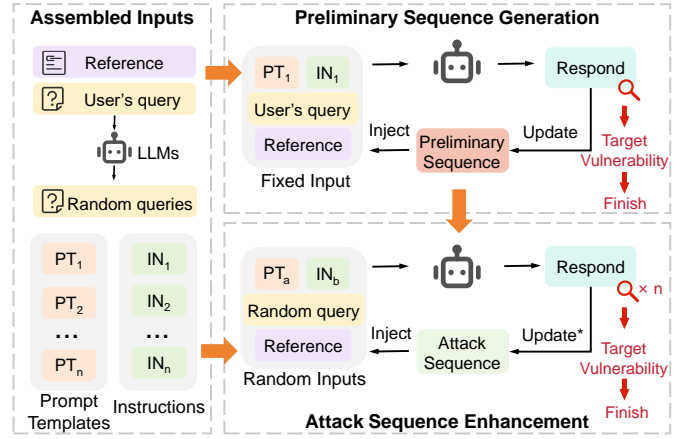


Fig. 4. Overview of HACKODE: First, it derives various assembled inputs by considering different instructions, prompt templates, and user queries. Second, it generates a preliminary attack sequence for a fixed assembled input. Third, it refines the sequence to enhance its transferability across the derived diverse assembled inputs.

a fixed assembled input. HACKODE then generates a preliminary sequence effective for this fixed assembled input. Upon completion of this step, attackers obtain a preliminary sequence that can induce the LLM to generate code with target vulnerability  $tVul$ . Thirdly, HACKODE uses random assembled inputs with different  $PT_a$ ,  $IN_b$ , and  $Q_c$  to enhance the transferability of the preliminary attack sequences. Through this step, attackers can obtain an attack sequence  $Seq$  with high transferability. Finally, attackers embed the attack sequences into correct code examples and publish them on the Internet. Once these code examples are referenced by LLM coding assistants, users may be misled to integrate vulnerabilities into their software.

#### A. Assembled Input Derivation

In this section, we first outline the overall external reference processing steps for LLM coding assistants and identify the variations for different coding assistants and users. Then, we simulate these variations and produce diverse assembled inputs to help improve the transferability of the attack.

**Variations in Assembled Inputs.** When an LLM coding assistant, enhanced with external information, responds to user questions, it will feed LLMs with assembled inputs. Those assembled inputs are constructed with three steps. We first conduct an in-depth analysis of the variable factors when producing assembled inputs.

In the first step, the application receives the user's query  $Q$ . For the same question, users may phrase their queries in different ways. Moreover, due to differences in the internal settings of applications, some may not directly use the user's original query but instead generate a new, simplified, and more comprehensible query based on the user's input for subsequent processing. Therefore, the first variable factor in producing assembled inputs is the query. In the second step, the application searches and retrieves relevant content on the Internet based on the user's query as reference  $Ref$ . This step is a crucial interaction between the application and the external environment. During this process, the attacker can induce the application

to produce vulnerable responses without altering the specific settings of the application, merely by publishing the external information. Therefore, external references become the only aspect that attackers can manipulate. The third step involves integrating the consolidated information as input of the LLM and instructing the LLM to generate the final response. During this step, the attacker cannot ascertain the specific input fed into the LLM by the application. The consolidation process primarily involves formatting the reference  $Ref$ , query  $Q$ , and instruction  $IN$  according to a predefined prompt template  $PT$  to ensure the LLM can correctly interpret them. Consequently, the second and third variable factors are the instruction  $IN$  and prompt template  $PT$ , respectively.

Since these three components: prompt templates  $PT$ , instructions  $IN$ , and user queries  $Q$ , vary significantly for various real-world applications and users, attackers need to consider their diversity when performing the attack. Thus, collect various  $PT$ ,  $IN$ , and  $Q$  to derive diverse assembled inputs. When constructing these assembled inputs, HACKODE adopts a random sampling approach, selecting one element from each of the prompt templates, instructions, and queries respectively. subsequently, the selected instruction, query, and reference are combined into a complete assembled input according to the format of the prompt template. Specifically, in the process of attack sequence enhancement, assuming there are  $|PT|$  prompt templates,  $|IN|$  instructions, and  $|Q|$  queries, the inputs will exhibit  $|PT| \times |IN| \times |Q|$  variations. This combination approach significantly enriches the diversity of input, and these carefully constructed inputs will be used in the subsequent step to enhance the transferability of the attack sequence.

Next, we will provide a detailed explanation of how to obtain multiple prompt templates  $PT$ , instructions  $IN$ , and user queries  $Q$ , as well as how to design external references capable of inducing LLM to generate vulnerabilities.

**User Query.** The queries represent variations in the expression of the same question, and we generate randomized queries using LLMs. For a given programming problem, we first use it as the original query. Then, we require the LLMs to rephrase the query following different inclinations. For example, we ask LLMs to briefly describe the problem, which may produce a query with less description. Through this step, we obtain multiple distinct queries for the same problem. Finally, we check whether these queries are related to the same programming problem by searching them on StackOverflow, and then retain those that maintain the original meaning.

**Instruction.** Instructions usually describe a scenario of interaction and provide guidance for LLMs to solve programming problems. For example, as shown in Figure 2, an instruction may require the LLM to carefully generate code without any vulnerabilities. These instructions help LLMs accurately understand the task requirements and improve code generation quality. However, different LLM-powered applications may adopt various instructions, and attackers need to consider this diversity when crafting attack sequences. To address this, we draw on instructions from several LLM-powered applications, selecting a diverse set of effective instructions. Additionally, we use LLMs to generate new instructions based on the

collected examples and adapt instructions from other scenarios to suit the coding assistant context.

**Prompt Template.** The function of the prompt template is to systematically integrate the query, relevant reference information, and instructions into a standardized format. A well-designed prompt template can help LLMs accurately understand requests and external content and improve their performance. Additionally, a well-structured prompt template can mitigate prompt injection attacks [14]. Therefore, we set prompt templates with multiple structures by referencing popular LLM-powered applications.

**Reference Information.** Except for the above three parts that vary across different applications and users, the reference information  $Ref$  is the only part that attackers can control. Coding assistants usually search for  $Ref$  from IT websites like StackOverflow. Thus, attackers can strategically post  $Ref$  to these websites. Specifically, they can find unanswered questions on these websites and prepare solutions with the correct code examples. These solutions contain code snippets that, with subtle modifications, can be transformed into significant security vulnerabilities  $tVul$ . For example, altering loop boundaries can result in out-of-bounds array reads, or changing the initialization of a variable can bypass certain security validations. Attackers set  $tVul$  as the target to craft the attack sequences  $Seq$ . Once  $Seq$  is crafted, attackers publish these solutions embedded with  $Seq$  under the unanswered questions. When coding assistants reference these carefully crafted solutions, they may inadvertently generate unsafe code containing  $tVul$ .

When crafting reference a code example, attackers need to determine the location where the attack sequence will be inserted. Two key factors must be taken into account when selecting the insertion position. Firstly, attackers should ensure that the chosen position does not affect the correctness of the original code. Secondly, the insertion position should maximize the effectiveness of HACKODE. Our strategy is to insert the attack sequence as code comments near the code that will be modified in the attack. The code comments do not directly affect the normal execution of the code but can potentially influence the LLMs' output.

Finally, as shown in Figure 2, attackers can obtain multiple queries  $Q$ , instructions  $IN$ , and prompt templates  $PT$ . Meanwhile, attackers have determined the correct code to be published, as well as identified the target vulnerability  $tVul$  and the insertion position for the attack sequence. Subsequently, HACKODE will generate an effective attack sequence.

## B. Attack Sequence Generation

HACKODE adopts a two-phase progressive attack method to generate highly adaptable attack sequences. In the first phase, HACKODE employs gradient-guided token mutation techniques to iteratively update and optimize the preliminary attack sequence until it can successfully attack the fixed assembled input. In the second phase, HACKODE enhances this sequence through randomly assembled inputs, improving its effectiveness and utility in real-world applications. This progressive generation approach accelerates the entire process



**Algorithm 1: Attack Sequence Generation**


---

**Input:** Prompts Templates  $PT$ , Instructions  $IN$ , User Queries  $Q$ , Reference Information  $Ref$ , and Target Vulnerability  $tVul$

**Output:** Attack Sequence  $Seq$

```

1  $i := 0$ 
2  $pSeq := init(pSeq)$ 
3 while  $i++ \leq maxStep$  do
4    $input := \{PT_1, IN_1, Q_1, Ref\} \odot pSeq$ 
5    $res := generate(M, input)$ 
6   if  $tVul \in res$  then
7     break;
8    $grad := \nabla_{pSeq} \mathcal{L}(M(input), tVul)$ 
9    $pSeq^m := (pSeq[i] \leftarrow pSeq[i] + grad[i])^m$ 
10   $pSeq := select(pSeq^m)$ 
11  $Seq = pSeq$ 
12 while  $i++ \leq maxStep$  do
13    $input^k := \{PT_a, IN_b, Q_c, Ref\}^k \odot Seq$ 
    // Randomly select  $PT$ ,  $IN$  and  $Q$  to
    // construct  $k$  assembled inputs
14    $res^k := generate(M, input^k)$ 
15   if  $\forall res^k, tVul \in res$  then
16     return  $Seq$ 
17    $grad := \nabla_{Seq} \mathcal{L}(M(input^k), tVul)$ 
18    $Seq^m := (Seq[i] \leftarrow Seq[i] + grad[i])^m$ 
19    $Seq := select(Seq^m)$ 

```

---

by providing a preliminary attack sequence as a basis for subsequent enhancements. We propose such a two-step generation because we find the trigger sequence generated based on one fixed assembled input has a certain transferability. Thus, we can first use a fixed assembled input to help the trigger sequence quickly converge to a relatively optimal result, followed by further enhancement. This approach can, on one hand, improve the transferability of the attack sequence, and on the other hand, reduce the convergence time.

*1) Preliminary Attack Sequence Generation:* The first phase aims to quickly craft a preliminary attack sequence,  $pSeq$ , that can execute the attack on an assembled input with a fixed prompt template ( $PT_1$ ), instruction ( $IN_1$ ), and user query ( $Q_1$ ). Although  $pSeq$  may not be adaptable to various assembled inputs, it serves as a foundation for subsequent enhancements and accelerates the overall generation process.

Lines 1~10 of Algorithm 1 show the details, HACKODE first initializes an attack sequence  $pSeq$ , composed of random symbols and the subtle differences between the target vulnerable code and the correct code. Such differences are embedded into  $pSeq$  to guide HACKODE in precisely producing  $tVul$  as attackers expect. Next,  $pSeq$  is inserted into the assembled input  $input$  (line 4), which remains fixed in this phase, and the insertion position is the code comment chosen by the attacker. In the implementation, HACKODE automatically checks whether the  $tVul$  has been successfully generated by comparing the differences between the code lines in the response and the correct code. If  $tVul$  appears in the code generated by LLMs, it indicates that  $pSeq$  is generated successfully and the first stage is over, as shown in lines 6~7. To prevent the LLMs from generating  $tVul$  due to random decoding, this test will be performed multiple times to ensure that  $tVul$  is produced consistently by the LLMs. If not, HACKODE updates  $pSeq$  through a gradient-based method. This process will be repeated, and HACKODE will iteratively

update  $pSeq$  until the attack succeeds or the maximum steps of iterations is reached.

The update process for  $pSeq$  is as follows: As shown in line 8, the algorithm first calculates the cross-entropy loss between the target LLM  $M$ 's logits prediction and target Vulnerability  $tVul$ . Based on the loss, the algorithm calculates the gradient  $grad$  with respect to the one-hot representation of the attack sequence  $pSeq$ . This gradient indicates how each token in the attack sequence should be altered. Using this gradient, the algorithm generates  $m$  new attack sequences by randomly changing one token at a time. For each randomly chosen token in the attack sequence, its gradient is added to the one-hot representation of that token, creating a score array for each token in the model's vocabulary. A higher score for a token implies that changing the current token to this token might better guide the model to produce the desired responses. The algorithm then replaces the selected token with one randomly chosen from those with the top scores. Among the  $m$  new sequences generated, the most effective one is selected by comparing the loss for each sequence. Finally, the algorithm selects new  $pSeq$  from  $pSeq^m$  based on the prediction loss in line 10.

*2) Attack Sequence Enhancement:* To enhance the transferability of attack sequences across varied assembled inputs, HACKODE further mutates  $pSeq$  based on randomly assembled inputs. These assembled inputs are constructed based on various prompt templates  $PT$ , instructions  $IN$ , and queries  $Q$  from the first step.

At the conclusion of the preliminary generation,  $pSeq$  has been able to induce the target LLM to generate  $tVul$  with fixed input and use it as input for the second stage. In this phase, the main task is to further optimize and enhance the adaptability of the attack sequence  $Seq$ . To achieve this, the algorithm guides the iterative updating of  $Seq$  based on random inputs. Additionally, to ensure that  $Seq$  is adequately strengthened, stricter termination conditions are set.

In detail, the second phase enhances  $Seq$  from three aspects. First, HACKODE is performed based on randomly assembled inputs, as shown in line 13. In each iteration, HACKODE will assemble  $k$  inputs, denoted as  $input^k$ . Among them,  $PT_a$ ,  $IN_b$ , and  $Q_c$  are randomly selected from the set of prompt templates, instructions, and user queries. These random combinations create a large search space, helping the generation process avoid overfitting to a specific input. Second, when examining the effectiveness of  $Seq$ , HACKODE conducts multiple tests based on random inputs. The attack is considered successful only when all  $k$  responses of  $k$  inputs contain  $tVul$ , as shown in line 15. Furthermore, this approach can also reduce the randomness introduced by decoding algorithms. Third, in line 17, the algorithm mutates  $Seq$  under the guidance of the accumulated gradient, calculated as follows

$$grad = \nabla_{Seq} \sum_{i=0}^k cross\_entropy(M(input_i), tVul).$$

This gradient considers multiple assembled inputs, so the mutate will not be biased by a single input. It is noteworthy that, in the  $i_{th}$  step, the update of  $Seq$  is based on the  $input^k$

randomly selected in the  $(i - 1)_{th}$  step. When assessing the success of the attack, the algorithm relies on the  $k$  newly acquired inputs from the  $i_{th}$  step for verification, which conserves computational resources and improves generation speed.

Finally, with the progressive generation, HACKODE successfully crafts an attack sequence with high transferability. Attackers can then embed these attack sequences into correct code examples and publish them on the Internet. Once referenced by LLM coding assistants, they will guide LLMs to provide users with vulnerable code.

## V. EVALUATION

In this section, we conduct extensive experiments to evaluate the effectiveness of HACKODE. Specifically, we aim to answer the following research questions:

- **RQ1**(§V-A). How effective is the attack on different vulnerabilities and LLMs?
- **RQ2**(§V-B). Can the attack transfer across different assembled inputs and quantified LLMs?
- **RQ3**(§V-C). How effective is the attack on real-world coding assistants?
- **RQ4**(§V-D). How does each component of HACKODE contribute to the effectiveness of the attack?

**Experiment Dataset.** We constructed a dataset consisting of 35 programming problems and their corresponding solutions, involving four mainstream languages: Python, Java, C++, and PHP. The entire process of constructing and processing the dataset strictly follows the common workflow for applications that search for external content [25, 4]. In detail, coding assistants usually utilize LLMs to summarize keywords from the initial query and then search for related problems using the APIs of platforms. Therefore, our dataset is constructed by simulating the real-world scenario of coding assistants searching for external content.

In detail, the dataset comprises the most up-to-date answers from StackOverflow, gathered at the time of the experiment. Specifically, we utilized the StackOverflow API and the StackExchange library [17] to collect the answers and their corresponding problems. We excluded problems that did not require domain-specific knowledge and could be directly resolved by LLMs without needing to reference external sources. Through manual validation, we also ensure that these answers include code examples and can be referenced by coding assistants when corresponding problems are asked.

We collected this dataset because no existing dataset is suitable for evaluating HACKODE. Some existing datasets, such as OWASP and Juliet, contain vulnerable code used for vulnerability detection. However, they lack problem descriptions and code explanations, which are essential for coding assistants. Additionally, these datasets are usually focused on just one language, like Java and C, whereas to simulate real-world scenarios of using coding assistants, we need a dataset with multiple languages. Therefore, we cannot directly use existing datasets to evaluate HACKODE. In contrast, our dataset is tailored to security research on code generated by LLM-based coding assistants, including problem descriptions, code descriptions, and code in multiple programming languages.

TABLE I  
THE DATASET DISTRIBUTION. THE CORRESPONDING CWE FOR EACH TYPE OF VULNERABILITY AND THE SPECIFIC NUMBER OF DATA INSTANCES.

Vulnerabilities	CWE-id(number)
Array Violation	CWE-125(11)
Buffer Overflow	CWE-787(3), CWE-120(2), CWE-122(1)
Incorrect Variable	CWE-457(3), CWE-190(1)
Invalid Validation	CWE-20(7), CWE-570(2)
Infinite Loop	CWE-835(5)

**Target LLMs.** HACKODE currently focuses on open-source LLMs. Thus, we evaluate HACKODE on two general LLMs, Llama2-7b (L-7b) [46] and Mistral-7b (M-7b) [15], and two code LLMs, CodeLlama-7b (C-7b) [35] and StarChat2-15b (S-15b) [22]. In detail, both Llama2-7b and Mistral-7b are general-purpose language models with 7 billion parameters that demonstrate exceptional performance in text generation and understanding. CodeLlama-7b, with 7 billion parameters, is specialized for programming purposes. StarChat2-15b, boasting 15 billion parameters, focuses on providing coding assistance through natural dialogue and code generation capabilities. These LLMs are popular choices for LLM coding assistants. During the experiments, we set the input and output lengths of the LLM to fully accommodate the prompts and generate complete code outputs. All other hyperparameters were kept at their default values. In future work, we plan to extend HACKODE to closed-source LLMs.

**Attack Targets.** Our target vulnerabilities include five categories of common weaknesses [24]. As shown in Table I, we used HACKODE to construct target vulnerabilities in each of these categories for the 35 programming problems, with 11 for Array Indexing Violation, 6 for Buffer Overflow, 4 for Incorrect Variable, 9 for Invalid Validation, and 5 for Infinite Loop. Once these vulnerabilities are integrated into users' code, attackers can exploit them to hijack users' devices to steal sensitive information or inject viruses. Specifically, Array Indexing Violation can lead to undefined behavior, such as reading or writing to adjacent memory regions (which may contain sensitive information), causing program crashes, data leaks, or arbitrary code execution. Buffer Overflow allows attackers to overwrite data in memory (such as function pointers or return addresses), enabling them to hijack control program flow and further execute arbitrary code. Incorrect Variable Usage can result in logical errors, data corruption, or unauthorized data access. Invalid Validation can cause the program to perform unsafe operations, such as executing unauthorized commands, leaking sensitive information, or allowing malicious input to be injected into the program. Infinite Loop can consume system resources, preventing the program from responding to other requests or commands, which can lead to Denial of Service (DoS) attacks.

**Experiment Setups.** We conducted our evaluation on a machine equipped with dual AMD EPYC 7763 CPUs (128 cores, 256 threads, 2.45 GHz base frequency) and 8 NVIDIA Tesla V100 GPUs (32 GB each), running Ubuntu 22.04 LTS. During the experiments, we set the maximum number of iterations *maxStep* as 500, and the number of random

assembled inputs  $k$  as 3.

#### A. Evaluation on Different LLMs

To evaluate the effectiveness of HACKODE for different LLMs, we perform the attack on two general LLMs and two code LLMs and record the attack success rate (ASR), the number of iterations (Iter), and the token length of attack sequence (Seq), the response (Res) and the request (Input). Table II shows the results of the attack on different vulnerabilities and LLMs.

In terms of general LLMs, HACKODE achieves an ASR of 77.14% on Llama2-7b and 80.00% on Mistral-7b. Specifically, HACKODE has over 60% ASR for Array Indexing Violation, Buffer Overflow, and Infinite Loop on both general-purpose LLMs. In terms of code LLMs, HACKODE achieves an ASR of 94.29% on CodeLlama-7b and 85.71% on StarChat2-15b. Specifically, HACKODE has over 60% ASR for all vulnerabilities, including Array Indexing Violation, Buffer Overflow, Incorrect Variable, Invalid Validation, and Infinite Loop. The higher ASR on code LLMs compared to general LLMs is attributed to the fact that Llama and Mistral sometimes generate plain text responses without code examples, thus not providing vulnerable code. This occurs because Llama and Mistral are not specialized for code generation tasks, and tend to answer questions with textual explanations rather than code examples.

Table II also presents the effort of the attack, in which we document the total number of iterations, including preliminary sequence generation and attack sequence enhancement. The average number of iterations for target LLMs to generate the attack sequence is 179.17. Specifically, the average number of iterations for attack sequence generalization for Llama2-7b, Mistral-7b, CodeLlama-7b, and StarChat2-15b are 189.27, 160.92, 174.49, and 191.99, respectively. This efficiency is attributed to the progressive generation approach, which can accelerate the generation process.

Moreover, the average token length of the attack sequences is 33.10, which is only 3.44% of the average token length of the assembled inputs. In particular, this rate goes down to 2.34% (26.40/1129.20) for Infinite Loop on StarChat2-15b. This suggests that HACKODE does not significantly alter the original correct example, making the attack sequence unnoticeable within the reference information. Furthermore, the average token lengths of assembled inputs are over 900, demonstrating that the attack remains effective even for complex generation tasks.

We also investigated the effectiveness of HACKODE against different vulnerability types as shown in Table II. The results demonstrate that HACKODE can successfully induce multiple types of vulnerabilities. However, the difficulty of inducing various vulnerabilities differs among different LLMs. For example, it is challenging to induce incorrect variable vulnerabilities in Mistral, with a mere 25.00% ASR, but generating code with array indexing violation is relatively easy, achieving a 100.00% ASR in this experiment. Therefore, for different LLMs, attackers can design specific characteristics for each LLM to induce the target vulnerabilities. Please note that, since the amount of data for each type of vulnerability varies, the

average value in Table II is calculated across all data points rather than as a simple average of each item.

**Answer to RQ1:** HACKODE is effective on different vulnerabilities across different LLMs, achieving an average ASR of 84.29%. The attack sequence only takes 3.44% of the token length of the assembled inputs.

#### B. Transferability of HACKODE

In this section, we evaluate HACKODE's transferability across diverse inputs and quantized LLMs. Through these experiments, we demonstrate that HACKODE is practical for attackers operating in a black-box setting against victim LLM coding assistants.

**Randomly Assembled Inputs.** The diverse settings of applications and the varying user queries can lead to different assembled inputs for the same programming problem. Therefore, we evaluate the transferability of the attack by testing the effect of the attack sequences on randomly assembled inputs that are not been involved in the generation process of the attack sequences. In detail, for each programming problem, we randomly construct five assembled inputs based on three new instructions, queries, and prompt templates and test the generated attack sequence on them. Notably, the new instructions, prompt templates, and queries were not used in the generation process of the attack sequences.

As shown in Table III, the results reveal that an average of 37.85% of the data can successfully pass all tests, 57.17% of the data can pass more than half of the tests, and 83.58% of the data can pass at least one test. This indicates that HACKODE can transfer to randomly assembled inputs with a high probability of success, validating the effectiveness of the sequence enhancement. Thus, even without detailed knowledge of the LLM coding assistant's prompt templates, instructions, or user requests, attackers can achieve the attack with high probability, demonstrating the real-world impact of HACKODE.

In addition, we discover that the attack sequences perform the worst in terms of migration capability on Llama2. We speculate that this may be attributed to the random sampling approach used by the Llama2 during text generation. However, despite this, 74.07% of the data can pass at least one test, which proves the good adaptability of the attack we employed.

**Quantized LLMs.** In real-world applications, developers tend to use quantized LLMs to reduce hardware costs and improve inference speed. Therefore, attack sequences designed for pre-trained LLMs should remain effective against their quantized models. Thus, we evaluated these attack sequences on LLMs quantized to 4 bits using GPTQ [10] and BitsAndBytes [8] techniques. In detail, for GPTQ, we downloaded models from HuggingFace. The quantization settings for BitsAndBytes follow the recommendations of HuggingFace.

As shown in Table IV, 48.07% of the attack sequences remain effective on the GPTQ quantized LLMs, and 53.45% on the BitsAndBytes quantized LLMs. In particular, the ASR



TABLE II

THE EVALUATION ON DIFFERENT LLMs AND DIFFERENT VULNERABILITIES. "ITER" REPRESENTS THE AVERAGE ITERATIONS FOR ATTACK SEQUENCE GENERATION. "INPUT", "RES" AND "SEQ" ARE THE AVERAGE TOKEN LENGTHS OF THE ASSEMBLED INPUTS, RESPONSES, AND ATTACK SEQUENCES, RESPECTIVELY.

		ASR	Iter	Input	Res	Seq	ASR	Iter	Input	Res	Seq
	Vulnerability	Llama2-7b					Mistral-7b				
General LLMs	Array Violation	100.00%	141.30	867.80	423.00	36.50	100.00%	93.82	814.73	336.09	37.64
	Buffer Overflow	100.00%	86.50	1224.17	538.33	32.17	83.33%	113.83	1213.50	463.00	31.17
	Incorrect Variable	50.00%	289.00	616.50	301.00	35.00	25.00%	387.00	921.50	379.75	35.75
	Invalid Validation	55.56%	226.13	917.38	375.00	31.13	77.78%	161.38	940.63	380.25	30.38
	Infinite Loop	60.00%	272.00	943.00	384.00	30.00	80.00%	183.40	1147.60	312.20	29.60
	Average	77.14%	189.27	923.66	410.91	33.27	80.00%	160.92	975.22	370.78	33.30
	Vulnerability	CodeLlama-7b					StarChat2-15b				
Code LLMs	Array Violation	100.00%	163.46	825.82	344.27	37.73	90.91%	250.00	808.27	447.91	39.55
	Buffer Overflow	83.33%	209.50	1225.83	515.33	32.00	100.00%	107.50	1175.17	616.00	30.83
	Incorrect Variable	100.00%	186.75	962.50	523.00	36.00	75.00%	365.75	914.50	525.00	41.00
	Invalid Validation	100.00%	146.25	946.00	469.75	31.13	66.67%	114.50	949.17	646.17	23.17
	Infinite Loop	80.00%	197.80	1143.00	421.60	29.80	100.00%	166.20	1129.20	550.00	26.40
	Average	94.29%	174.49	986.23	437.33	33.72	85.71%	191.99	965.39	551.10	32.13

TABLE III

THE ASR ON RANDOMLY ASSEMBLED INPUTS. THE TABLE SHOWS THE AVERAGE PERCENTAGE OF DATA THAT PASS THE TESTS AT LEAST ONE TO FIVE TIMES SEPARATELY.

LLMs	5/5	4/5	3/5	2/5	1/5
L-7b	33.33%	37.04%	48.15%	62.96%	74.07%
M-7b	39.29%	42.86%	53.57%	64.29%	85.71%
C-7b	45.46%	48.49%	63.64%	72.73%	87.88%
S-15b	33.33%	43.33%	63.33%	73.33%	86.67%
Average	37.85%	42.93%	57.17%	68.33%	83.58%

TABLE IV

THE ASR ON QUANTIZED LLMs.

	L-7b	M-7b	C-7b	S-15b	Average
GPTQ	29.63%	53.57%	42.42%	66.67%	48.07%
BitsAndBytes	48.15%	57.14%	48.49%	60.00%	53.45%

of StarChat2-15b with GPTQ and BitsAndBytes quantization is up to 66.67% and 60.00%, respectively. These results indicate that attack sequences are effective for quantized LLMs. Meanwhile, the ASR of Llama2 with GPTQ quantization is the lowest, with only 29.63% ASR. We analyzed Llama2's responses and found that the model cannot answer certain questions and generate irrelevant content, even when references are provided. This suggests that GPTQ quantization downgrades Llama2's ability when handling programming problems.

**Answer to RQ2:** HACKODE can transfer to randomly assembled inputs with an average ASR of 57.97%. HACKODE can also transfer to quantized LLMs, with an average ASR of 50.76%.

### C. Real-World Experiment

To reveal the potential impact of HACKODE in real-world scenarios, we evaluated it on a coding assistant application.

The coding assistant is implemented based on ChatChat [4], which utilizes the StackOverflow API agent to fetch relevant solutions. The coding assistant is powered by Mistral-7b and CodeLlama-7b. To prevent the dissemination of malicious content, we did not post the crafted solutions to the Internet. Instead, we posted them on local web pages and allowed the coding assistant to search our local pages. Then, we requested the coding assistant and checked if the responses based on those crafted solutions contained vulnerabilities. Other than that, we don't change any workflow of the application.

As shown in Table V, 82.14% of the attack sequences on Mistral and 69.70% on CodeLlama successfully induce the application to generate code with the target vulnerabilities. Given that the number of vulnerabilities varies across different types, we calculated the overall ASR based on all data, rather than averaging the ASR across individual types. Furthermore, HACKODE can induce applications to generate multiple types of vulnerability code. In addition, the coding assistant powered by Mistral is more likely to be attacked than those powered by CodeLlama. It is because the prompt template of this application did not prompt CodeLlama well, and CodeLlama frequently ignored the reference and generated incorrect responses.

TABLE V

REAL-WORLD EVALUATION ON CODING ASSISTANT.

Vulnerabilities	M-7b	C-7b
Array Violation	100.00%	63.64%
Buffer Violation	80.00%	80.00%
Incorrect Variable	100.00%	75.00%
Invalid Validation	71.43%	44.44%
Infinite Loop	50.00%	100.00%
Total	82.14%	69.70%

Figure 5 shows an illustrative example of the attack. In detail, a user queries about DirectX's usage, and the coding assistant refers to a solution published by attackers. This

solution contains the correct code but includes an attack sequence. After the user requests the coding assistant, it writes code that copies one more float than the buffer size under the influence of the attack sequence, leading to an out-of-bounds write vulnerability. If the user directly uses the code without inspection, the program will write data beyond the buffer boundary, causing a buffer overflow vulnerability. By exploiting this vulnerability, attackers can control users' devices and execute malicious commands. Furthermore, it may lead to a series of severe consequences, including data corruption and privacy breaches, posing a grave threat to both user device security and personal information security.

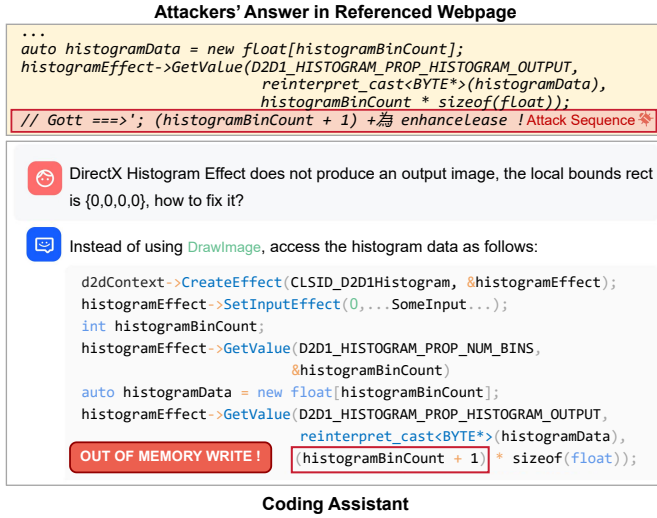


Fig. 5. A real-world example of HACKODE. Attackers post a correct solution with an attack sequence, guiding the coding assistant to generate vulnerable code with 'out-of-bounds write'.

In this real-world experiment, attackers do not need to know the detailed settings, like prompt templates and instructions, of the coding assistant. They only need to know which LLM the coding assistant is powered by, which is easy to obtain through interaction with the coding assistant. Due to the transferability of HACKODE, attackers can successfully induce the coding assistant to generate vulnerable code. The potential threats posed by vulnerable code are extremely severe. Once these flawed codes are deployed in real-world scenarios, attackers can exploit these vulnerabilities to carry out various malicious activities, such as remote code execution, data theft, and system crashes. In conclusion, HACKODE poses a significant threat to the security of real-world coding assistants.

**Answer to RQ3:** HACKODE can effectively induce real-world coding assistants, e.g., ChatChat, to generate vulnerable code, with an average ASR of 75.92%.

#### D. Ablation Study

**Effectiveness of Progressive Generation.** We experimented to validate the necessity of the progressive attack method. We performed the attack that only utilized the preliminary sequence generation module without the attack sequence

enhancement, denoted as HACKODE<sup>-</sup>. Then, we evaluated its transferability, as illustrated in Section V-B. In detail, we construct five randomly assembled inputs and embed the generated attack sequences into them to test whether these sequences can successfully induce the target LLM to generate code with vulnerabilities.

As shown in Table VI, we recorded the average success rates of the attack sequences passing the transferability tests. The experimental data clearly demonstrate that, for each target model, the attack sequences generated by HACKODE have a higher probability of passing the test than those generated by HACKODE<sup>-</sup>, with an average increase of 25.49%. The results illustrate the effectiveness of the attack sequence enhancement module in enhancing the transferability of attack sequences.

TABLE VI  
ASR OF HACKODE AND HACKODE<sup>-</sup>.

Attack	L-7b	M-7b	C-7b	S-15b	Average
HACKODE <sup>-</sup>	28.89%	35.71%	33.33%	32.00%	32.48%
HACKODE	51.11%	57.14%	63.64%	60.00%	57.97%

**Effectiveness of the Designed Embedding Position.** When injecting attack sequences into a correct code example, there exist various possible approaches. To evaluate whether inserting the attack sequence as code comments is the optimal method, we conducted a comparative evaluation by inserting the attack sequence through variable renaming. Specifically, we crafted an attack sequence and embedded it as a variable name, ensuring that the sequence adhered to the naming conventions of different programming languages. We then executed the attack on Mistral-7b using these two different approaches for injecting the attack sequences. It is worth noting that the crafted code can still execute correctly after injecting the attack sequence using both approaches.

TABLE VII  
EVALUATE HACKODE AND VARIABLE RENAMING ON THE MISTRAL-7B.

Dataset	HACKODE	Renaming
Array Violation	100.00%	9.09%
Invalid Validation	77.78%	22.22%
Buffer Violation	83.33%	16.67%
Incorrect Variable	25.00%	0.00%
Infinite Loop	80.00%	0.00%

As shown in Table VII, the variable renaming approach achieved an ASR of only 11.43%. In contrast, inserting the sequences as comments resulted in a significantly higher overall ASR of 80.00%. This disparity is largely due to the constraints that must be adhered to when naming variables, which introduce limitations during the generation of attack sequences and thus reduce their effectiveness. Additionally, because vulnerabilities like incorrect variables and infinite loops are particularly challenging to induce in the Mistral-7b model, both the HACKODE method and variable renaming achieve relatively low success rates in these cases. The experimental results indicate that inserting attack sequences as comments into the correct code is a less constrained approach that favors the success of attacks.

**Answer to RQ4:** The progressive generation approach achieves 25.49% higher ASR than the non-progressive approach. In addition, embedding the attack sequence as code comments improves the ASR by 68.57% compared to embedding it in variable names.

## VI. DISCUSSION

**Spread of Crafted Solutions.** The solutions crafted by HACKODE are correct from human perception, and their code is functional and error-free. This makes these solutions more likely to be prioritized on IT forums like StackOverflow. As a result, LLM coding assistants are likely to unintentionally refer to this manipulated content, thereby generating code containing vulnerabilities. In addition to this spread approach, attackers can also leverage techniques like search engine optimization (SEO) [38] to improve the ranking of crafted solutions in search engines' results. By carefully selecting keywords, building links, and optimizing content, attackers can increase the likelihood that their traps appear first when coding assistants search for solutions online. This approach significantly increases the likelihood of a successful attack. Consequently, HACKODE poses a substantial threat in the real world.

**Potential Defenses.** We propose a few potential directions to mitigate HACKODE. One involves using LLMs with external data to inspect and amend code vulnerabilities, aiming to improve detection accuracy by continuously referencing up-to-date coding practices and security standards. However, LLMs' limited domain-specific expertise may impede accurate error detection and correction. Additionally, this approach faces challenges related to model updates and ongoing data maintenance. Another promising approach is to apply static vulnerability detection methods to analyze code generated by LLMs. By thoroughly examining code structure, static analysis tools can effectively identify common vulnerabilities. However, these tools often have high false-positive rates [12], which could lead developers to overlook real risks or increase the burden of code review, thereby impacting overall usability. To enhance these defenses, dynamic detection methods, such as runtime vulnerability monitoring, could be explored in combination with static analysis to form a multi-layered defense. However, such approaches may introduce additional computational load, posing challenges for applications sensitive to resource consumption. Considering these challenges, it is urgent for the community to develop new defenses against HACKODE.

**Limitations.** First, HACKODE currently focuses solely on open-source LLMs, which are widely used for locally deployed models. Regarding closed-source LLMs, we plan to explore transfer attacks in future work. Second, we have validated HACKODE on LLMs with parameter sizes ranging from 7b to 15b, as these sizes are most commonly used in LLM studies [56, 50]. In future research, we intend to investigate LLMs with smaller parameter sizes, such as 1b, and larger sizes, like 130b. Third, HACKODE may not be highly effective against certain LLMs that exhibit significant performance

degradation following quantization. An example of such a model is the quantized version of GPTQ Llama. However, most real-world applications tend to employ quantized LLMs that retain as much of their original performance as possible. For these models, HACKODE proves to be quite effective, especially for most open-source LLM coding assistants.

## VII. RELATED WORK

**Evaluation of LLM Generated Code.** With the emergence of code LLMs, many researchers have examined the quality of LLM-generated code from various aspects, including usability, correctness, and security. However, they have not considered scenarios involving the utilization of external information. Tian et al. and Liu et al. evaluate LLMs on solving common programming problems and found that ChatGPT has advantages in code generation tasks, but it struggles to generalize to unseen problems.

Some studies [26, 43, 47, 7] have found that Copilot, a coding assistant in Visual Studio Code, demonstrates a strong ability to solve programming problems. However, it occasionally generates flawed code that can lead to significant execution errors.

Security concerns with LLM-generated code have also been highlighted by various studies [33, 13, 32, 36], showing that existing LLM-based code generation tools may produce code with vulnerabilities. For instance, Pearce et al. assessed Copilot's code generation for the top 25 CWE vulnerabilities, and other studies [33, 36] have examined how coding assistants impact the security practices of users. To address these issues in code LLMs, recent research [20, 40, 41, 39] has proposed benchmarks to evaluate code correctness and security. In addition, some work [16, 29, 54] leverages software analysis tools, like scanners and fuzzers, to identify vulnerabilities in LLM-generated code. However, most existing work focuses on code generation without incorporating external references. Thus, the potential risks of external reference have not been fully explored.

**Attacks on LLMs.** Current attack methods against LLMs primarily focus on jailbreak and backdoor attacks, leaving a research gap in scenarios where LLMs may be misled by external references. Jailbreak attacks typically involve crafting prompts to bypass LLMs' protection mechanisms to generate malicious content [3, 56]. In code generation, one study jailbreaks code LLMs to generate vulnerable code [50]. Distinct from existing works where LLMs generate code without references, HACKODE targets the scenario where LLMs produce vulnerable code even when referencing correct code examples.

Backdoor attacks occur during the training or fine-tuning process of LLMs, where attackers inject malicious data into the training dataset to influence the LLMs. LLMs with backdoor may generate code with vulnerabilities when prompted with the specific trigger instruction [18, 34]. You et al. introduced a backdoor attack methodology that leverages LLMs to automatically insert various style-based triggers into text. Yang et al. conducted an in-depth exploration of the backdoor robustness of LLM agents and assessed various types of backdoor attacks targeting these agents. Different from backdoor

attacks, HACKODE does not require the training or fine-tuning of LLMs. Instead, it can directly induce LLMs to generate vulnerable code through external information.

### VIII. CONCLUSION

In this paper, we reveal a new threat, HACKODE, where attackers can induce LLM code assistants to generate vulnerable code through external information. Using a progressive generation approach, HACKODE generates attack sequences that are highly adaptable to different inputs. We conducted extensive experiments to evaluate the effectiveness of HACKODE on different LLMs and vulnerabilities. Through a real-world experiment, we demonstrate that HACKODE poses a significant threat to developers. Our goal is to raise developers' awareness of the potential risks associated with LLM coding assistants and to inspire the development of effective defense techniques.

### REFERENCES

- [1] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [2] Yihan Cao, Shuyi Chen, Ryan Liu, Zhiruo Wang, and Daniel Fried. Api-assisted code generation for question answering on varied table structures, 2023.
- [3] Patrick Chao, Alexander Robey, Edgar Dobriban, Hamed Hassani, George J Pappas, and Eric Wong. Jailbreaking black box large language models in twenty queries. *arXiv preprint arXiv:2310.08419*, 2023.
- [4] Chatchat-Space. Chatchat. <https://github.com/chatchat-space/Langchain-Chatchat>, 2023.
- [5] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- [6] Cryptopolitan. User solana wallet exploited in first case of ai poisoning attack, 2024. URL [https://coinstats.app/news/41925c613b2ced0ee117c89fcee7e6729a8466e0fe8394140f1493ff7376535a\\_User-Solana-wallet-exploited-in-first-case-of-AI-poisoning-attack/](https://coinstats.app/news/41925c613b2ced0ee117c89fcee7e6729a8466e0fe8394140f1493ff7376535a_User-Solana-wallet-exploited-in-first-case-of-AI-poisoning-attack/).
- [7] Arghavan Moradi Dakhel, Vahid Majdinasab, Amin Nikanjam, Foutse Khomh, Michel C Desmarais, and Zhen Ming Jack Jiang. Github copilot ai pair programmer: Asset or liability? *Journal of Systems and Software*, 203:111734, 2023.
- [8] Tim Dettmers. bitsandbytes. <https://github.com/TimDettmers/bitsandbytes>, 2024.
- [9] James Finnie-Ansley, Paul Denny, Brett A Becker, Andrew Luxton-Reilly, and James Prather. The robots are coming: Exploring the implications of openai codex on introductory programming. In *Proceedings of the 24th Australasian Computing Education Conference*, pages 10–19, 2022.
- [10] Elias Frantar, Saleh Ashkboos, Torsten Hoefer, and Dan Alistarh. Gptq: Accurate post-training quantization for generative pre-trained transformers. *arXiv preprint arXiv:2210.17323*, 2022.
- [11] GitHub. Github copilot. <https://github.com/features/copilot>. Accessed: 2023-10-30.
- [12] Andrew Habib and Michael Pradel. How many of all bugs do we find? a study of static bug detectors. In *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering*, pages 317–328, 2018.
- [13] *An empirical study of code smells in transformer-based code generation techniques*, 2022. IEEE.
- [14] Neel Jain, Avi Schwarzschild, Yuxin Wen, Gowthami Somepalli, John Kirchenbauer, Ping-yeh Chiang, Micah Goldblum, Aniruddha Saha, Jonas Geiping, and Tom Goldstein. Baseline defenses for adversarial attacks against aligned language models. *arXiv preprint arXiv:2309.00614*, 2023.
- [15] Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L  lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth  e Lacroix, and William El Sayed. Mistral 7b, 2023.
- [16] Arya Kavian, Mohammad Mehdi Pourhashem Kallehbasti, Sajjad Kazemi, Ehsan Firouzi, and Mohammad Ghafari. Llm security guard for code. In *Proceedings of the 28th International Conference on Evaluation and Assessment in Software Engineering*, pages 600–603, 2024.
- [17] LangChain. Stackexchange, 2024. URL <https://python.langchain.com/v0.2/docs/integrations/tools/stackexchange/>.
- [18] Jiazhaoli, Yijin Yang, Zhuofeng Wu, VG Vydiswaran, and Chaowei Xiao. Chatgpt as an attack tool: Stealthy textual backdoor attack via blackbox generative model trigger. *arXiv preprint arXiv:2304.14475*, 2023.
- [19] Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, R  mi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, Thomas Hubert, Peter Choy, Cyprien de Masson d’Autume, Igor Babuschkin, Xinyun Chen, Po-Sen Huang, Johannes Welbl, Sven Gowal, Alexey Cherepanov, James Molloy, Daniel J. Mankowitz, Esme Sutherland Robson, Pushmeet Kohli, Nando de Freitas, Koray Kavukcuoglu, and Oriol Vinyals. Competition-level code generation with alphacode. *Science*, 378(6624):1092–1097, December 2022. ISSN 1095-9203. doi: 10.1126/science.abq1158. URL <http://dx.doi.org/10.1126/science.abq1158>.
- [20] Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. Is your code generated by chatgpt really correct? rigorous evaluation of large language models for code generation. *Advances in Neural Information Processing Systems*, 36, 2024.
- [21] Zhijie Liu, Yutian Tang, Xiapu Luo, Yuming Zhou, and Liang Feng Zhang. No need to lift a finger anymore?

- assessing the quality of code generation by chatgpt. *IEEE Transactions on Software Engineering*, 2024.
- [22] Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, Tianyang Liu, Max Tian, Denis Kocetkov, Arthur Zucker, Younes Belkada, Zijian Wang, Qian Liu, Dmitry Abulkhanov, Indraneil Paul, Zhuang Li, Wen-Ding Li, Megan Risdal, Jia Li, Jian Zhu, Terry Yue Zhuo, Evgenii Zheltonozhskii, Nii Osae Osae Dade, Wenhao Yu, Lucas Krauß, Naman Jain, Yixuan Su, Xuanli He, Manan Dey, Edoardo Abati, Yekun Chai, Niklas Muennighoff, Xiangru Tang, Muhtasham Oblokulov, Christopher Akiki, Marc Marone, Chenghao Mou, Mayank Mishra, Alex Gu, Binyuan Hui, Tri Dao, Armel Zebaze, Olivier Dehaene, Nicolas Patry, Canwen Xu, Julian McAuley, Han Hu, Torsten Scholak, Sebastien Paquet, Jennifer Robinson, Carolyn Jane Anderson, Nicolas Chapados, Mostofa Patwary, Nima Tajbakhsh, Yacine Jernite, Carlos Muñoz Ferrandis, Lingming Zhang, Sean Hughes, Thomas Wolf, Arjun Guha, Leandro von Werra, and Harm de Vries. Starcoder 2 and the stack v2: The next generation, 2024.
  - [23] Yusuf Mehdi. Reinventing search with a new ai-powered microsoft bing and edge, your copilot for the web. <https://blogs.microsoft.com/blog/2023/02/07/reinventing-search-with-a-new-ai-powered-microsoft-bing-and-edge-your-copilot-for-the-web/>, 2023.
  - [24] Mitre. Common weakness enumeration. <https://cwe.mitre.org/>, 2024.
  - [25] mugglmenzel. ai-agent-scaffold, 2024. URL <https://github.com/mugglmenzel/ai-agent-scaffold>.
  - [26] Nhan Nguyen and Sarah Nadi. An empirical evaluation of github copilot’s code suggestions. In *Proceedings of the 19th International Conference on Mining Software Repositories*, pages 1–5, 2022.
  - [27] Ansong Ni, Srini Iyer, Dragomir Radev, Veselin Stoyanov, Wen-Tau Yih, Sida Wang, and Xi Victoria Lin. LEVER: Learning to verify language-to-code generation with execution. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 26106–26128. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/ni23b.html>.
  - [28] Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan Wang, Yingbo Zhou, Silvio Savarese, and Caiming Xiong. Codegen: An open large language model for code with multi-turn program synthesis. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL [https://openreview.net/forum?id=iaYcJKpY2B\\_](https://openreview.net/forum?id=iaYcJKpY2B_).
  - [29] Ana Nunez, Nafis Tanveer Islam, Sumit Kumar Jha, and Peyman Najafirad. Autosafecoder: A multi-agent framework for securing llm code generation through static analysis and fuzz testing. *arXiv preprint arXiv:2409.10737*, 2024.
  - [30] OpenAI. Chatgpt. <https://openai.com/blog/chatgpt>, 2022.
  - [31] Shuyin Ouyang, Jie M Zhang, Mark Harman, and Meng Wang. Llm is like a box of chocolates: the non-determinism of chatgpt in code generation. *arXiv preprint arXiv:2308.02828*, 2023.
  - [32] Hammond Pearce, Baleegh Ahmad, Benjamin Tan, Brendan Dolan-Gavitt, and Ramesh Karri. Asleep at the keyboard? assessing the security of github copilot’s code contributions. In *2022 IEEE Symposium on Security and Privacy (SP)*, pages 754–768. IEEE, 2022.
  - [33] Neil Perry, Megha Srivastava, Deepak Kumar, and Dan Boneh. Do users write more insecure code with ai assistants? In *Proceedings of the 2023 ACM SIGSAC Conference on Computer and Communications Security*, pages 2785–2799, 2023.
  - [34] Javier Rando and Florian Tramèr. Universal jailbreak backdoors from poisoned human feedback. *arXiv preprint arXiv:2311.14455*, 2023.
  - [35] Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*, 2023.
  - [36] Gustavo Sandoval, Hammond Pearce, Teo Nys, Ramesh Karri, Siddharth Garg, and Brendan Dolan-Gavitt. Lost at c: A user study on the security implications of large language model code assistants. In *32nd USENIX Security Symposium (USENIX Security 23)*, pages 2205–2222, 2023.
  - [37] Sami Sarsa, Paul Denny, Arto Hellas, and Juho Leinonen. Automatic generation of programming exercises and code explanations using large language models. In *Proceedings of the 2022 ACM Conference on International Computing Education Research-Volume 1*, pages 27–43, 2022.
  - [38] Dushyant Sharma, Rishabh Shukla, Anil Kumar Giri, and Sumit Kumar. A brief review on search engine optimization. In *2019 9th international conference on cloud computing, data science & engineering (confluence)*, pages 687–692. IEEE, 2019.
  - [39] Mohammed Latif Siddiq and Joanna Santos. Generate and pray: Using sallms to evaluate the security of llm generated code. *arXiv preprint arXiv:2311.00889*, 2023.
  - [40] Mohammed Latif Siddiq and Joanna CS Santos. Security-eval dataset: mining vulnerability examples to evaluate machine learning-based code generation techniques. In *Proceedings of the 1st International Workshop on Mining Software Repositories Applications for Privacy and Security*, pages 29–33, 2022.
  - [41] Mohammed Latif Siddiq, Joanna Cecilia da Silva Santos, Sajith Devareddy, and Anna Muller. Sallm: Security assessment of generated code. In *Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering Workshops*, pages 54–65, 2024.
  - [42] Alexey Smirnov and Tzi-cker Chiueh. DIRA: automatic detection, identification and repair of control-hijacking attacks. In *Proceedings of the Network and Distributed System Security Symposium, NDSS 2005, San Diego, California, USA*. The Internet Society, 2005.
  - [43] Dominik Sobania, Martin Briesch, and Franz Rothlauf.

- Choose your programming copilot: a comparison of the program synthesis performance of github copilot and genetic programming. In *Proceedings of the genetic and evolutionary computation conference*, pages 1019–1027, 2022.
- [44] Hongjin Su, Shuyang Jiang, Yuhang Lai, Haoyuan Wu, Boao Shi, Che Liu, Qian Liu, and Tao Yu. Arks: Active retrieval in knowledge soup for code generation. *arXiv preprint arXiv:2402.12317*, 2024.
- [45] Haoye Tian, Weiqi Lu, Tsz On Li, Xunzhu Tang, Shing-Chi Cheung, Jacques Klein, and Tegawendé F Bisseyandé. Is chatgpt the ultimate programming assistant—how far is it? *arXiv preprint arXiv:2304.11938*, 2023.
- [46] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models, 2023.
- [47] Priyan Vaithilingam, Tianyi Zhang, and Elena L Glassman. Expectation vs. experience: Evaluating the usability of code generation tools powered by large language models. In *Chi conference on human factors in computing systems extended abstracts*, pages 1–7, 2022.
- [48] Victor Van der Veen, Nitish Dutt-Sharma, Lorenzo Cavallaro, and Herbert Bos. Memory errors: The past, the present, and the future. In *Research in Attacks, Intrusions, and Defenses: 15th International Symposium, RAID 2012, Amsterdam, The Netherlands, September 12–14, 2012. Proceedings 15*, pages 86–106. Springer, 2012.
- [49] Tu Vu, Mohit Iyyer, Xuezhi Wang, Noah Constant, Jerry Wei, Jason Wei, Chris Tar, Yun-Hsuan Sung, Denny Zhou, Quoc Le, et al. Freshllms: Refreshing large language models with search engine augmentation. *arXiv preprint arXiv:2310.03214*, 2023.
- [50] Fangzhou Wu, Xiaogeng Liu, and Chaowei Xiao. Deceptprompt: Exploiting llm-driven code generation via adversarial natural language instructions. *arXiv preprint arXiv:2312.04730*, 2023.
- [51] Shuyuan Xu, Zelong Li, Kai Mei, and Yongfeng Zhang. Aios compiler: Llm as interpreter for natural language programming and flow programming of ai agents, 2024.
- [52] Hongyu Yang, Liyang He, Min Hou, Shuanghong Shen, Rui Li, Jiahui Hou, Jianhui Ma, and Junda Zhao. Aligning llms through multi-perspective user preference ranking-based feedback for programming question answering, 2024.
- [53] Wenkai Yang, Xiaohan Bi, Yankai Lin, Sishuo Chen, Jie Zhou, and Xu Sun. Watch out for your agents! investigating backdoor threats to llm-based agents. *CoRR*, abs/2402.11208, 2024. doi: 10.48550/ARXIV.2402.11208. URL <https://doi.org/10.48550/arXiv.2402.11208>.
- [54] Dongyu Yao, Jianshu Zhang, Ian G Harris, and Marcel Carlsson. Fuzzllm: A novel and universal fuzzing framework for proactively discovering jailbreak vulnerabilities in large language models. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4485–4489. IEEE, 2024.
- [55] Wencong You, Zayd Hammoudeh, and Daniel Lowd. Large language models are better adversaries: Exploring generative clean-label backdoor attacks against text classifiers. *arXiv preprint arXiv:2310.18603*, 2023.
- [56] Andy Zou, Zifan Wang, J Zico Kolter, and Matt Fredrikson. Universal and transferable adversarial attacks on aligned language models. *arXiv preprint arXiv:2307.15043*, 2023.