

WitheredLeaf: Finding Entity-Inconsistency Bugs with LLMs

Hongbo Chen*

Indiana University Bloomington

Yifan Zhang*

Indiana University Bloomington
Samsung Research America

Xing Han*

Independent Researcher

Huanyao Rong

Indiana University Bloomington

Yuheng Zhang

University of Electronic Science and
Technology of China

Tianhao Mao

Indiana University Bloomington

Hang Zhang

Indiana University Bloomington

XiaoFeng Wang

Indiana University Bloomington

Luyi Xing

Indiana University Bloomington

Xun Chen

Samsung Research America

ABSTRACT

Originating from semantic bugs, Entity-Inconsistency Bugs (EIBs) involve misuse of syntactically valid yet incorrect program entities, such as variable identifiers and function names, which often have security implications. Unlike straightforward syntactic vulnerabilities, EIBs are subtle and can remain undetected for years. Traditional detection methods, such as static analysis and dynamic testing, often fall short due to the versatile and context-dependent nature of EIBs. However, with advancements in Large Language Models (LLMs) like GPT-4, we believe LLM-powered automatic EIB detection becomes increasingly feasible through these models' semantics understanding abilities.

This research first undertakes a systematic measurement of LLMs' capabilities in detecting EIBs, revealing that GPT-4, while promising, shows limited recall and precision that hinder its practical application. The primary problem lies in the model's tendency to focus on irrelevant code snippets devoid of EIBs. To address this, we introduce a novel, cascaded EIB detection system named WitheredLeaf, which leverages smaller, code-specific language models to filter out most negative cases and mitigate the problem, thereby significantly enhancing the overall precision and recall. We evaluated WitheredLeaf on 154 Python and C GitHub repositories, each with over 1,000 stars, identifying 123 new flaws, 45% of which can be exploited to disrupt the program's normal operations. Out of 69 submitted fixes, 27 have been successfully merged.

1 INTRODUCTION

Entity-inconsistency bug. Among the program bugs with security implications, the hardest ones to detect are logic flaws that are closely related to a program's functionalities and its context of operations. Figure 1 shows an example, in which `quote()` at line 991 should take the argument `params` instead of `query`. The type of flaws in the figure is characterized by misuses of the wrong yet valid program *entities* – a leaf node in the program's abstract syntax tree such as a variable identifier, a function name, a control string, etc., which we call *Entity-Inconsistency Bug* (EIB) in our research. EIBs have been reported in prior studies as unique logic flaws introduced by the human developer who mistakenly chooses

*All three authors contributed equally to this research.

Prompt: Find bugs in this Python code snippet.

```
987 def quote_url(url_text):
988     (scheme, netloc, path, params, query, fragment)
989     = six.moves.urllib.parse.urlparse(url_text)
990     # netloc_quoted = quote(netloc)
991     path_quoted = quote(path)
992     params_quoted = quote(query)
993     query_quoted = quote_plus(query)
994     fragment_quoted = quote(fragment)
995     url_quoted = six.moves.urllib.parse.urlunparse((
996         scheme, netloc, path_quoted, params_quoted,
997         query_quoted, fragment_quoted))
998     return url_quoted
```

Output by ChatGPT-4: The code snippet you provided has a couple of issues that could be causing problems:

- Incorrect usage of the quote function for params:** In the code, `params_quoted = quote(query)`. This line mistakenly passes the query to the quote function instead of params. It should be `params_quoted = quote(params)`.
- Potential Misuse of quote vs. quote_plus:** Using `quote_plus` for the query is appropriate, but if any other component requires handling spaces as `+`, then `quote_plus` should be used accordingly. Conversely, if spaces should be encoded as `%20` in the query, then `quote` should be used instead of `quote_plus`.

Another false alarm and a corrected code snippet are truncated.

Figure 1: An Example of ChatGPT detecting an existing bug. The variable query in the red rectangular is buggy.

a different or even a random, meaningless string for a program entity [54]. They often come with security implications, exposing the program to such threats as denial of service, violation of control and data flow integrity, etc. [29, 33], and have also been recorded by the CVE database, e.g., CVE-2022-30780, CVE-2022-23574, CVE-2020-8889, CVE-2020-10574, and CVE-2019-11463.

Like other logic flaws, EIBs are extremely hard to detect in an automatic and scalable way. For the example in Figure 1, although a human reviewer might be able to identify the out-of-place pair

(`query` and `params_quoted`) from the context of the program (`path` and `path_quoted`, `query` and `query_quoted`), the buggy code does not carry any clear-cut patterns, as those for syntactic vulnerabilities like use after a free, insecure API invocation, *etc.*, making it difficult for a traditional static analyzer to locate the flaws. While dynamic testing (*e.g.*, fuzzing) might detect some of such errors, it tends to miss most of them due to its low code coverage. Additionally, EIBs may not be effectively captured by the sanitizers (*e.g.*, AddressSanitizer [26]) indicating the existence of bugs (*e.g.*, transform a silent memory corruption into an explicit crash) for fuzzers (*e.g.*, EIBs may cause only logical errors without memory corruption). As evidence, the EIB shown in Figure 1 has remained undetected for approximately seven years until January 2024 [18] in a popular GitHub repository with 368 stars.

EIB detection with LLM. A key observation is that the human reviewer can recognize the presence of an EIB thanks to her capability to “predict” the correct code entity (*e.g.*, a variable name) at a specific location, according to the context of the program (*e.g.*, the correlations between the argument names and function names from Line 990 and 992-993 in Figure 1). This capability to predict most likely tokens to follow a sequence of tokens (*i.e.*, the context) has also been provided by *language models*, a Natural Language Processing (NLP) technology. Particularly, with the recent advance in *Large Language Models* (LLMs) such as ChatGPT, we believe that LLM-powered automatic EIB detection becomes increasingly feasible. So in our research, we performed the first study on the LLM’s potential to automate the discovery of EIBs, focusing on OpenAI’s GPT-4 [4]. Our study shows that indeed the LLM offers powerful semantic analysis capabilities, promising to move us a significant step closer to scalable identification of EIBs. On a dataset of synthesized EIBs (§2.3), GPT-4 successfully captured 60% of them with proper prompt engineering. Particularly, we found that when running open-sourced LLMs such as Code Llama for an infilling task [48] (that is, predicting missing tokens from a program), the LLMs tend to produce tokens consistent with the valid program entities and different ones when a true EIB is encountered (§2.4). However, we found that GPT-4 alone is not ready for supporting a comprehensive, highly scalable, and reasonably accurate EIB detection due to the following limitations.

First, it is widely known that LLMs have an inherent tendency to generate inaccurate or fabricated information in their responses, commonly referred to as “hallucination” [4], sometimes due to the model being *distracted* to the code context unrelated to the targeted flaws [52]. This weakness leads to a high false positive rate if we directly apply LLMs like GPT-4 to EIB detection, especially on large codebases with a large number of suspicious EIBs. For example, in Figure 1, GPT-4 issues a false alarm on `quote_plus`, probably due to it being distracted by words like `plus`. In our measurement, on a sampled dataset (with 66 functions randomly selected from the GitHub projects modified after the cut-off date for the GPT-4 training [4]), GPT-4 reports EIBs for nearly half of the functions, most of which turns out to be false positives.

Second, GPT-4 tends to miss many EIBs, incurring a high false negative rate, in the presence of a large context window (that is, a large token length) [4]. OpenAI has reported that the context window size is a limitation to applying GPT-4 to cybersecurity [4].

A large context involves many program entities, most unrelated to the EIB, which however *distract* GPT-4 away from those indeed containing the flaws, as observed in the prior research [52]. Such problems were also observed in our measurement study, rendering it less practical to directly use GPT-4 for EIB detection.

Our solution. The key to addressing the limitations of GPT-4 is to reduce the distraction of the model, focusing its *attention* on truly suspicious entities. To identify the locations of these entities, our idea is to leverage a language model’s capability to “fill in the blank”, that is, predicting the vacated program entity through an *infilling* task [48]. Note that this cannot be done directly using GPT-4: given the large number of entities in a program, directly running the infilling task through the LLM could incur an unbearable monetary cost. So our solution is to utilize a set of lightweight, open-source LLMs to first go through these tokens, identifying suspicious ones before handing them over to GPT-4 for an in-depth analysis.

Specifically, this analysis pipeline, which we call *WitheredLeaf* (§3), works as follows. Given the source code of a program, our approach first performs a static analysis to identify all its entities, including variable names and function names. *WitheredLeaf* utilizes CodeBERT [25], an extremely lightweight language model, to execute the first infilling run: for each occurrence of a given entity, our approach removes it from the program location, runs CodeBERT on the program to predict it, and records all failed predictions. In this way, our approach quickly excludes the bug-free code fragments. The program locations related to the failed predictions, which are considered to be suspicious, are further examined by Code Llama [48], a larger, more capable yet slower model. The output of the LLM, including all program locations involving the entities inconsistent with the predictions made by both local models, is uploaded to GPT-4 for an in-depth EIB analysis. *WitheredLeaf* applies a novel prompt engineering technique to focus GPT-4’s attention to these suspicious locations to reduce both false positives and false negatives.

The design of *WitheredLeaf* strikes a balance between the accuracy of an EIB detection and the cost of detection, both in terms of execution time and monetary expense. Our evaluation shows that *WitheredLeaf* achieves a coverage of 62.1% and a reasonable precision of 23.57% (a significant improvement from less than 1% for direct application of GPT-4). Running the approach on 80 Python and 74 C repositories with over 1,000 stars on GitHub, we discovered 93 new bugs from the former and 30 new bugs from the latter.

Contributions. We summarize our contributions as follows.

- *Understanding LLMs’ Capabilities in EIB Detection.* We report the first comprehensive measurement study on various LLMs’ capabilities to detect EIBs, shedding light on their potentials and limitations in performing this task.
- *Design and implementation of WitheredLeaf.* Based upon the understanding from our measurement study, *WitheredLeaf*, a cascaded pipeline to overcome the limitations of LLMs, is developed, turning out to be both effective and efficient on various codebases. We plan to make *WitheredLeaf* publicly available.
- *New flaws discovered and fixed.* *WitheredLeaf* discovered 123 unknown EIBs, in Python and C Github repositories with more

than 1,000 stars, 45% of which can be exploited to disrupt the program’s normal operations. We have submitted 69 pull requests to the relevant developers, out of which 27 have already been merged.

- *New EIB dataset.* As a part of this work, we curated a comprehensive dataset for the EIBs from historical bug fixes, new bugs discovered by WitheredLeaf, and our synthetic dataset. We will release this dataset to facilitate future research in this direction.

2 MEASUREMENT

In order to take the best advantage of LLMs for EIB detection, it is crucial to first understand the performance of LLMs for this specific task. In this section, we aim to comprehensively measure the performance of different language models, which not only includes the powerful state-of-the-art GPT-4, but also other smaller open-sourced models tailored for program code. We start with an overview of different language models and show how they can be utilized for EIB detection with a motivating example, then present the detailed measurement of each model’s performance. We will also discuss our findings and insights drawn from this study.

2.1 Language-Model based EIB Detection

Instruct Large Language Models. Represented by OpenAI’s GPT-4 [4], these models are trained with enormous textual data crawled from a wide range of sources, including program code, and fine-tuned to follow the instructions in the prompt. Consequently, these models can assist many different tasks (*i.e.*, general-purpose). However, many such models are commercial and close-sourced (*e.g.*, GPT-4), making the large-scale use and research expensive. Moreover, these models also suffer from the well-known “hallucination” problem [4], leading to inaccurate responses.

Example of EIB Detection. As shown in Figure 1, upon receiving the buggy code and the straightforward prompt, ChatGPT-4 can promptly identify the EIB in its response. However, it is notable that ChatGPT-4’s response also includes two false alarm EIB cases (*e.g.*, that related to `quote.plus`), despite the short length of the input code. This suggests that directly applying GPT-4 on large codebases for EIB detection could incur excessive false alarms, largely due to the distraction problems [52].

Code-specific Language Models. Many language models have been specifically trained for coding tasks (*e.g.*, code generation and completion) recently. These models share a similar transformer-based model generating the probabilities of the tokens based on the context. CodeBERT [25] is one early model trained as a bidirectional encoder [20] and pre-trained on masked language modeling (MLM) tasks, thus capable of code infilling (*i.e.*, predict the code entity such as a variable name at a certain location). More recent models like Code Llama [48], StarCoder [35], Stable Code [44], and DeepSeek-Coder [27], are unidirectional, decoder-only transformers, meaning that they take sequential context and output the probabilities of only “the next token” step by step. Nevertheless, benefiting from the fill-in-the-middle (FIM) training method [10], these models are also able to perform code-infilling tasks. Though not as powerful as general-purpose LLMs, all aforementioned code language models are open-source.

Example of EIB Detection. The code-infilling capability of the above models can also support EIB detection. Specifically, we can compare a specific code entity in the original code with a model’s “predicted” entity at the same location (based on the context), if the two align, then likely that code entity is correct, otherwise, an EIB may exist. For instance, when we input the buggy code snippet in Figure 1 as a prompt into the Code Llama 7B [48] model, masking out the query variable at line 911 and executing it to fill in a variable at the same location, the model promptly suggested `params` as the output. Such inconsistency successfully indicates the EIB.

Conclusion. We conclude that different language models, including both smaller models dedicated for code and powerful general-purpose models, are all capable of detecting EIB issues with different methods. In the remainder of this section, we conduct an in-depth and comprehensive measurement study to better understand the capabilities of different language models.

2.2 Dataset Preparation

To our best knowledge, there is no available EIB dataset. We then curate a dataset ourselves, with data collected after the model’s knowledge cut-off date to minimize memorization concerns. We gather data from popular GitHub repositories (more than 1000 stars), to ensure a high code quality. We specifically extract functions from these repositories that have updates post the LLMs’ knowledge cut-off date¹. With these functions, we develop two datasets: (1) $D_{sampled}$: includes functions of varying sizes measured in lines of code (LOC), representing a broad spectrum of real-world code scenarios. (2) $D_{synthesized}$: consists of functions intentionally modified to include synthesized bugs, specifically through entity mutations.

To generate $D_{synthesized}$, we implemented a process similar to previous identifier mutation testing methodologies [47]. This process begins with a comprehensive enumeration of all permissible entities that could feasibly replace a specific entity within a given code snippet. We then select pairs of entities for substitution based on two main criteria: those with the minimal edit distance, to reflect subtle typographical errors a developer might inadvertently make, and those with the smallest semantic discrepancy, to mirror potential confusion a developer might encounter when differentiating between two similar-meaning entities. Although one function could contain multiple such pairs satisfying the criteria, we always only mutate one entity and thus only introduce one EIB at each function. A subset of synthesized bugs in $D_{synthesized}$ is introduced on functions in $D_{sampled}$.

2.3 GPT-4’s Performance on EIB Detection

To understand general-purpose LLMs’ EIB detection performance, we choose GPT-4 as the test subject in our measurement due to its superior performance [51, 52, 57] and comprehensive support for JSON-formatted output by OpenAI’s API, which plays an important role in automating bug detection tasks.

Prompt Design. Research indicates that in-context examples do not significantly improve GPT-4’s accuracy in detecting vulnerabilities [51]. Therefore, given the diversity of EIB patterns, we opted for zero-shot prompts to avoid the influence of irrelevant examples [49].

¹We use `gpt-4-0125-preview` with training data up to Dec. 2023.

Specifically, we leverage a template configuration approach to construct a system prompt and multiple rounds of interactive prompts to identify the EIB with associated information about the EIB, represented as a set of properties summarized in Table 1. Instead of querying all properties (i.e., information) in a single prompt, our template configuration enables us to distribute the queries for these properties across different rounds in interactive sessions. This approach can enhance effectiveness: for example, some properties allow GPT-4 to reflect about its previous results so that potential false alarms can be filtered out (i.e., cross-examination). Based on the properties required in each round, a prompt that queries these properties can be automatically generated. We divide the properties into two categories: (1) mandatory properties *must* occur in specific rounds, and (2) selectable properties *could* occur in certain rounds depending on the template configuration. Note that our multi-round template design and the selection of various EIB-related properties are based on multiple effective prompting techniques (e.g., Chain-of-Thought, or CoT [56] and cross-examination [19]) to boost GPT-4’s bug detection performance, also detailed in Table 1.

We showcase one example of template configuration in Figure 2, which is very close to what we eventually adopt in WitheredLeaf. The system prompt asks GPT-4 to act as a programming language expert (i.e., persona adoption [41]) to help improve its performance. In the first round, we ask GPT-4 to locate the potential buggy code lines and provide explanations. The second round further cross-validates and scrutinizes the identified potential EIBs to reduce false alarms. We additionally incorporate three selectable properties in this round: self-evaluated categorization, fix proposal, and merely fixable predicate. The former asks GPT-4 to classify the bug for a better confirmation and understanding, the latter properties prompt it to generate a bug fix (typically an identifier replacement at the AST leaf node due to the nature of EIB), which serves as a step in CoT that helps LLM determine if the bug can be merely fixed and helps researchers develop bug patches. As shown in our later measurement results, this well-rounded property combination enhances GPT-4’s performance in EIB detection. Please refer to Appendix §A for details of the prompt templates in Table 2.

Measurement Results. To understand the performance of different prompt templates in terms of recall, specificity, and monetary cost, we first configure a wide range of representative templates featuring different numbers of prompt rounds and sets of selectable properties. Measured templates are listed in Table 2, where the digits represent the order of the prompt round and the alphabet abbreviation (see Table 1) stands for selectable properties involved in certain rounds. For instance, the template in Figure 2 is denoted as 1/2FMca. Although each function in $D_{synthesized}$ contains only one known EIB, GPT-4 is capable of identifying multiple EIBs within a single function. For this reason, we employ two statistical approaches: 1) regarding each function as a single report (i.e., function-wise), and 2) regarding each bug as a single report (i.e., bug-wise). We then test all template configurations on subsets of $D_{synthesized}$ and $D_{sampled}$. Due to a limited budget, we randomly sampled 66 functions from each of these datasets.

With a simple single-round vanilla template (labeled “1”), we note a high FP rate on $D_{synthesized}$ and high FP number on $D_{sampled}$,

System: You’re a Python expert. Your job is to inspect if the code contains any semantic bugs. Semantic bugs, also known as semantic inconsistencies or logical errors, are a type of bug that occurs when the code is syntactically correct but does not behave as intended or produces incorrect results. These bugs arise from mismatches between the programmer’s intended logic and the actual implementation in the code, where incorrect variable/method name usage or assignment can lead to bugs or vulnerabilities in the program. Here are some requirements:

- Assume the code is syntactically correct, and input parameters to the functions are well-formed and valid.
- Focus solely on detecting semantic bugs, and ignore other problems (e.g., undefined symbols).
- Output in valid JSON format.

Round 1 Properties: code_line + explanation

Round 1 Prompt: {code} Output exact lines of semantic bugs and concise explanations of the bugs.

Round 2 Properties: code_line* + explanation* + fixed_line + merely_fixable + category

Round 2 Prompt: Inspect these bugs, excluding 1. incorrect or unlikely bugs; 2. bugs caused by undefined symbols; 3. non-semantic bugs. The left bugs only break the intended functionality or lead to vulnerabilities. Answer if the bugs can be fixed by merely changing a variable/method name and fix them. Also, classify this bug into these categories using the code and previous contexts: Security Vulnerability, Logic Bug, Enhancement, Unexpected Behavior, Symbol Not Defined, Module Not Imported, Bad Smell, Not a Bug, or Others. Note that the snippet is from popular repositories and runs, so correct symbols not defined in the snippet are defined at other places, which is not Logic Bug. If you assign Others category, explicitly name the category.

Figure 2: A two-round prompt template with system prompt.

predominantly attributed to hallucination, as verified through selective manual inspections. First, we want to understand the effectiveness of different properties used as filters (see Table 1). We find that the `merely_fixable` property is effective in filtering out non-EIB cases, which cannot be simply fixed by changing a single code entity. This filtering increases the specificity by 25%-72%, at the cost of a slight increase in FNs. Exclude candidates whose category \notin {Logic Bug, Security Vulnerability, Unexpected behavior, Bad Smell} also effectively increases the specificity. However, using `priority != high` as a filtering criterion, although reducing FPs, also leads to missing of TP cases, which is considered worse than the category filter. This might be caused by GPT-4’s inadequate grasp of the bug severity and consequences.

On the other hand, we analyze the effect of additional prompting rounds. Introducing the second round prompt increases the recall by approximately 4% from 1FM to 1/2FM and specificity by 32%-46%, demonstrating cross-examination can eliminate many FP cases in EIB detection. Notably, adding the third prompt round turns out to be less effective in both specificity, recall, and monetary cost (i.e., extra rounds lead to more tokens). This can be attributed to unnecessary cross-examination, as the ground truth bugs are detected but filtered out. The root cause might be distraction [52],

Table 1: Mandatory and selectable requested properties in JSON-formatted tool calling.

Property	Ab.	Cat.	Description	Purpose	Prompt Tec.
code_line	Co	1M.	the exact line of code with semantic bug	locating the bug	/
explanation	E	1M.	a concise explanation of this bug	reasoning about of the bug	/
code_line*	Co*	23M.	the original (unfixed) code line with semantic bug, same as in previous round	locating the bug in follow-up rounds	/
explanation*	E*	23M.	a concise explanation of this bug and your reasoning	reasoning about of the bug in follow-up rounds	/
merely_fixable	M	12S.	the bug can be fixed by merely changing a variable/method name to another one	filtering out bugs that does not occur on a single AST node	Cross-exam.
fixed_line	F	12S.	the line with the bug fixed	CoT for merely_fixable and helping fix	CoT
category	Ca	23S.	the category of this bug; be cautious when answering Security Vulnerability and Logic Bug	filtering out non-semantic bugs	Cross-exam.
priority	P	23S.	the priority level of this bug: high/medium/low	filtering out less interesting bugs	Cross-exam.

Ab. represents the abbreviation of the property; Cat. represents the category of the property: in round X, it is mandatory (M.) or selectable (S.)

in which some irrelevant information in the first and second rounds distracts the LLM. Thus, two-round prompting is enough.

Conclusion. Based on our measurement results, we draw two major conclusions.

(1) *Template Choice.* We decide that the template denoted by 1/2FMCa (detailed in Figure 2) has the best overall performance, which is eventually adopted by WitheredLeaf with suspicious lines highlighted. We elaborate on the details of the highlighting mechanism in §3.4

(2) *Limitations of GPT-4 on EIB Detection.* From our study, GPT-4 shows some potential in EIB detection, however, its performance is still far from practical. A significant portion of the EIBs (26 out of 66) is missed even by the most performant “1/2FMCa” template. Moreover, EIBs in real-world code can be much rarer than EIBs in $D_{synthesized}$. Demonstrated via specificity evaluated on $D_{sampled}$, 19 out of 66 FP reports remain after filtering, resulting in a high false positive rate if we naively let GPT-4 scan the whole codebase.

2.4 Code Language Models on EIB Detection

As explained in §2.1, the code infilling capability of smaller code language models can also be utilized for EIB detection (e.g., discrepancies between “predicted” and actual code entities indicate potential EIBs). In this part, we measure various code language models’ EIB detection performance. Although general-purpose LLMs like GPT-4 are also capable of code infilling, the associated monetary cost could be very high for large codebases. For example, infilling every entity for a function of 40 LOC costs roughly \$2.1 for GPT-4, and the cost increases quadratically with code size. We thus focus on open-source code language models in this measurement.

Hardware and Parameters. We conduct the experiments on a Linux server with two AMD EPYC 9124 CPUs and one NVIDIA H100 GPU, with the temperature set to 0 for all models. Among our tested models, CodeBERT operates on the CPU, while the others utilize the GPU.

Consistency and Efficiency of Code Infilling. Code infilling is the backbone capability of code language models enabling them to detect EIBs, we thus first measure this basic ability of various code models. At a high level, we mask each variable and method call name entity, one at a time, in the code from $D_{sampled}$ and

then let various code models fill in the “blank spaces”, gauging whether they can correctly “recover” the masked code entities (i.e., consistency) and how fast the infilling can be performed (i.e., efficiency). To explore the impact of context length (e.g., the amount of code surrounding the “blank space” available to code models) on code infilling performance, our measurement tests 6 different buckets of context lengths (listed in Table 3).

Our results are detailed in Table 3. Overall, open-source code models consistently fill in the correct code entities in more than 80% cases, with Deepseek-Coder 1.3B and CodeBERT being exceptions. Regarding the efficiency, most models take longer to process extended contexts, except for CodeBERT, as its token limitation is 512 and the input is truncated. For the impact of context length on consistency, in general, we only observe a marginal influence for most models. However, for Code Llama and Stablecode, minimal context (0-20 LOC) results in consistency deterioration. Upon further analysis, we find these models struggle to generate consistent tokens when the masked original tokens appear only once in the context - a situation more prevalent with smaller contexts. We also note a slight decline in consistency in context length of 40-80 LOC. Longer functions might provide more context or redundancy, helping the model to better understand the overall purpose and functionality of the code. In contrast, mid-length functions might provide insufficient context for the model to accurately infer the missing information, while shorter functions might be simple enough that each line carries more predictive weight. Importantly, Code Llama 13B encounters memory limitations for large context lengths of 320-640 LOC, suggesting that models with over 13B parameters may not be suitable for single GPU setups on extensive contexts - the higher hardware requirement makes it more difficult to deploy.

Recall of Code Infilling. Models with over 2 billion parameters display a promising capability to generate consistent tokens according to our previous measurement (Table 3). We then proceed to test their capability of recalling known EIBs, following the method explained in §2.1. For this purpose, we use a randomly selected subset of $D_{synthesized}$, which includes 1,383 artificially introduced EIBs in the functions contained in $D_{sampled}$. Specifically, we let the code models to infilling: (1) the locations of the buggy code entities, seeing whether there are any discrepancies (e.g., indicating EIBs), and (2) the other locations on the same line as the buggy

Table 2: Bug detection on GPT-4 with different prompt templates.

Ground Truth Template	Synthesized Dataset; Contains One Bug								Sampled Dataset							
	Function-wise				Bug-wise				Total Cost	Function-wise			Bug-wise			Total Cost
	TP	F. FN	U. FN	Rec.	TP	FP	FN	TN		FP	TN	Spe.	FP	TN	Spe.	
1	38	0	28	57.6	43	228	0	0	1.33\$	66	0	0.0	323	0	0.0	1.43\$
1FM	37	1	28	56.1	37	47	3	62	1.28\$	64	2	3.0	143	49	25.5	1.41\$
1/2FM	40	2	24	60.6	42	38	9	177	3.24\$	45	21	31.8	91	234	72.0	3.80\$
1/2FM/3P	18	22	26	27.3	18	14	34	206	5.38\$	28	38	57.6	50	276	84.7	6.11\$
1/2FM/3Ca	22	16	28	33.3	24	25	35	219	5.47\$	31	35	53.0	54	302	84.8	6.47\$
1/2FMCa	39	1	26	59.1	43	18	3	183	3.23\$	19	47	71.2	40	268	87.0	3.77\$
1/2FMCa w/ HL	55	9	2	83.3	58	11	14	116	3.11\$	32	34	51.5	47	195	80.6	3.36\$

TP = True Positives; TN = True Negative; FP = False Positives; FN = False Negatives; F.FN = Filtered FN; U.FN = Unfound FN; Rec. = Recall Percentage; Spe. = Specificity Percentage. We set the temperature parameter to 0 in this experiment. We assume all functions are free of EIB in the sampled dataset, as the density of EIB is very low in real-world code. The abbreviations used in the templates can be found in Table 1. The setting 1/2FMCa w/ HL is used in WiTherLeaf, highlighting the line of EIB with other three randomly selected lines in the code (detailed in §3.4).

Table 3: Language models’ infilling consistency on the sampled dataset of various context lengths.

Models	Size	0-20 LOC		20-40 LOC		40-80 LOC		80-160 LOC		160-320 LOC		320-640 LOC		Avg. Cons.
		Cons.	Speed	Cons.	Speed	Cons.	Speed	Cons.	Speed	Cons.	Speed	Cons.	Speed	
Code Llama	7B	83.4	243	88.0	431	83.7	647	87.4	1362	85.6	1689	86.0	6236	85.7
Code Llama	13B	85.2	360	92.1	685	87.6	966	88.5	1710	86.7	2750	OOM		88.2
Deepseek	1.3B	56.7	311	42.6	341	36.8	333	61.2	333	65.8	506	74.8	2066	56.3
Deepseek	6.7B	83.8	214	88.0	332	76.9	524	84.9	896	80.5	1550	76.3	2736	81.8
Stablecode	3B	83.9	148	88.6	127	81.6	176	86.5	168	84.3	216	83.9	530	84.8
CodeBERT	125M	31.1	202	31.2	242	23.2	216	23.1	207	24.4	202	17.0	261	25.0

Consistency is abbreviated as **Cons.**, and written in percentage. **Speed** represents the average generation speed for each identifier in milliseconds. 100 functions are sampled for LOC sizes of 0-20, 20-40, and 40-80 LOC with 1703, 4713, and 9278 infilling tasks, respectively; 50 functions for 80-160 LOC with 7499 FIM tasks; 20 functions for 160-320 LOC with 5287 FIM tasks; 10 functions for 320-640 LOC with 7319 FIM tasks.

Table 4: Language models’ infilling recall and consistency on the synthesized dataset.

Model	Size	Recall on EIB	Cons. of Others
Code Llama	7B	99.0	79.3
Code Llama	13B	99.4	81.4
Deepseek	1.3B	97.8	46.4
Deepseek	6.7B	99.0	76.3
Stablecode	3B	98.8	78.3
CodeBERT	125M	90.0	34.0

location, to assess the impact of EIB on code infilling for nearby code locations. Our measurement results are depicted in Table 4, as seen, nearly all models, except CodeBERT, demonstrate a high recall for known EIBs, suggesting the feasibility of utilizing code models to hunt for EIBs.

For nearby code locations on the same line of the EIB sites, interestingly, we observe 6%-7% drops in the consistency (the “Cons. of Others” column in Table 4) of infilling tasks for all models except CodeBERT. This phenomenon indicates that EIB can “interfere” with the infilling task, which has been utilized to enhance the robustness of EIB detection, as discrepancies are observable not only at the exact locations of the bugs but also in nearby areas. The related design details are presented in §3.4. A potential explanation for this phenomenon is that EIBs increase the perplexity, leading to larger uncertainty in token prediction [55].

Conclusion. We highlight major conclusions from our study:

(1) *Effective True Negative Filtering.* As mentioned, if the code models’ predicted code entity is consistent with the original one, we deem the non-existence of EIBs for that entity. The high consistency (higher than 80%) on the basically EIB-free $D_{sampled}$ of various code models (Table 3) thus suggests that they can effectively filter out more than 80% EIB-free code entities (*i.e.*, True Negatives). Moreover, this filtering only introduces very few false negatives - as shown in Table 4, the vast majority (*e.g.*, at least 97.8% for all models except CodeBERT) of true EIB cases will still trigger inconsistencies. These results highlight the values of code models as effective true negative filters.

(2) *Limitations of Code Models on EIB Detection.* Despite that code models can filter out a large number of true negatives, the filtered results could still contain lots of false alarms. For instance, there are 10% - 20% inconsistencies (*e.g.*, suspicious EIBs) in Table 3, but the real EIBs might only take less than 1%, yielding a high false positive rate. The primary reason is that code models’ token generation is largely uncontrolled, as they are not exclusively trained to generate specific entities in the FIM task. Instead, they inconsistently produce characters like pound signs (to start a Python comment), spaces, and parentheses, even when asked to infill a variable name. These issues make it difficult to detect EIBs solely relying on code models.

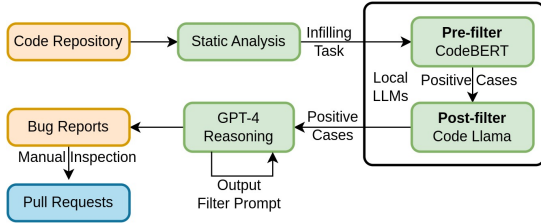


Figure 3: Pipeline of WitheredLeaf.

3 DESIGN AND IMPLEMENTATION

Inspired by our measurement in §2, we present WitheredLeaf, an automatic system to effectively and efficiently detect EIBs in various codebases. In this section, we detail its design and implementation.

3.1 Overview

Our measurement study (§2) reveals that, though both general-purpose LLMs (e.g., GPT-4) and smaller code language models show some potentials in EIB detection, each of them has its own weaknesses. Fortunately, these weaknesses can be largely mitigated with the synergy of different models. For example, the high expenses and false positive rate of GPT-4 can be alleviated if we use a smaller and cheaper model to filter out most EIB-free code locations in advance and provide precise locations of suspicious EIBs. Similarly, the excessive inconsistencies identified by code language models can also be further verified by the powerful GPT-4. Based on this core idea, WitheredLeaf employs a cascaded approach that integrates different models in a pipeline to achieve an optimal result.

WitheredLeaf’s Pipeline. We depict WitheredLeaf’s pipeline in Figure 3. Initially, the code undergoes processing by a static analyzer to compile a comprehensive list of all code entities subject to EIB detection. Then, WitheredLeaf transforms this list into code-infilling tasks, undertaken by locally deployed small models. These models sequentially filter out most EIB-free code locations. WitheredLeaf then automatically prompts GPT-4 to investigate and verify the remaining suspicious EIB sites. GPT-4’s responses will undergo a secondary filtration process and then be compiled into bug reports for further manual inspection.

This cascaded design eliminates most true negative cases while keeping most true positive cases (as discussed in §2), thereby preventing GPT-4 from being overwhelmed and significantly reducing the false positive rate (§3.2). For local models, we develop schemes to control token generation (§3.3), which enhances WitheredLeaf’s specificity and true negative filtering. Finally, WitheredLeaf precisely identifies suspicious EIB locations (with the help of code models) and feeds them to GPT-4, which compels GPT-4 to thoroughly reason about suspicious code lines (i.e., highlight), thereby reducing the false negative rate (§3.4).

3.2 Cascaded Detection

As indicated by the measurement results in §2, achieving high coverage, precision, and speed simultaneously is challenging. Our basic idea is to integrate different models in a cascaded design to achieve a desirable tradeoff of the aforementioned aspects, the

problem then becomes which exact models we should choose for each pipeline stage. We reason about our choices in this section.

Problem Formalization. We assume the pipeline has i stages, each of which is equipped with a language model. We use p_i to denote the true negative rate (i.e., specificity) of the model at stage i , q_i the true positive rate (i.e., recall), and t_i the throughput (i.e., execution speed). Let N_i represent the number of potential EIB cases that need to be processed at stage i , and ϵ_i represents the ratio of the number of actual EIB instances relative to N_i . We can derive the number of true/false positives/negatives at each stage:

$$\begin{aligned} N_{TN,i} &= N_i p_i (1 - \epsilon_i) & N_{TP,i} &= N_i q_i \epsilon_i \\ N_{FP,i} &= N_i (1 - p_i) (1 - \epsilon_i) & N_{FN,i} &= N_i (1 - q_i) \epsilon_i \end{aligned}$$

After stage i , true negatives $N_{TN,i}$ are correctly filtered out, while false negatives $N_{FN,i}$ incorrectly. True and false positives are escalated to the next stage ($N_{i+1} = N_{TP,i} + N_{FP,i}$). We can then define the maximum number of missed bugs M of the entire pipeline and the total execution time T (assuming n stages in total):

$$M = \sum_{i=0}^n N_i (1 - q_i) \epsilon_i \quad T = \sum_{i=0}^n N_i t_i$$

Additionally, the true positive proportion of stage $i + 1$ can also be derived from metrics of the previous stage i :

$$\epsilon_{i+1} = \frac{N_{TP,i}}{N_{TP,i} + N_{FP,i}} = \frac{q_i \epsilon_i}{1 - p_i - \epsilon_i + (p_i + q_i) \epsilon_i}$$

For an n -stage pipeline, the aggregated precision of all models before the final stage is then ϵ_n . For the final stage n , WitheredLeaf adopts powerful state-of-the-art GPT-4, which will process N_n EIB candidates, with the specificity and recall as p_n and q_n , respectively.

Next, we formulate the total cost of the detection pipeline, which includes the cost of computing resources C_{comp} per time unit, the cost per API invocation C_{api} of commercial LLMs like GPT-4, and the penalties associated with missed bugs C_{miss} (e.g., missed vulnerabilities can cause significant financial loss) and false alarms C_{check} (e.g., extra manual reviewing efforts):

$$C = C_{api} N_n + C_{comp} T + C_{miss} M' + C_{check} N'_{FP}$$

Here, M' represents the final count of missed bugs, and N'_{FP} denotes the final count of false positive cases for the entire pipeline. While the values of C_{api} and C_{comp} can be obtained from service providers, C_{miss} and C_{check} are user-configurable variables, depending on different needs and priorities. Given this cost formalization, the model selection problem for pipeline stages can be transformed to an optimization problem aiming to minimize the total cost.

WitheredLeaf’s Model Selection. As an experimental project, we weigh more on the bug coverage (i.e., fewer false negatives), so we would set a higher C_{miss} and put a focus on reducing M' . As indicated by the formula of M , it increases as more pipeline stages are added. In other words, additional pipeline stages can enhance the precision, but at the cost of more missed bugs and longer execution time. For example, assuming $\epsilon_0 = 0.002$, $p_i = 0.8$ and $q_i = 0.9$, maximally 9.9% of EIBs evades after the first local LLM,

18.9% evades after the second, 27.1% after the third, while reaching to 34.4% with the introduction of the fourth stage. Therefore, to maximize WitheredLeaf’s capability to capture EIBs, we fix the number of cascaded local models to two, with one GPT-4 stage.

We expect roughly 20% precision on the final report, which aligns with previous work on inconsistency bug detection [7] and is acceptable by developer community [8]. However, calculating with the metrics measured previously in §2, such a goal is not practical without improvements on the specificity of current models. To address this problem, we develop techniques to control the token generation of code language models and subsequently improve their specificity, with neglectable loss of recall (§3.3). Finally, considering the execution time and above metrics, we select CodeBERT, as the first local model to pre-filter negative EIBs and Code Llama 7B as the second local model. More details on their metrics and calculations are presented in §4.1.

3.3 Controlled Token Generation

As measured in §2.4, local LLMs are not generating tokens in a controlled manner, leading to a large number of inconsistencies in infilling tasks for non-EIB entities. The major causes of such inconsistencies are different for different models. Specifically, the encoder-based model (*i.e.*, CodeBERT) is unaware of the *actual length* of the entity in infilling, while decoder-only models (*e.g.*, Code Llama) are unaware of what *kind* of entity should be generated. The high-level idea is to inspect the probabilities of generated tokens, and we propose two algorithms for these models.

Static Analysis for Infilling Task Preparation. An initial step involved the construction of input infilling tasks to the models. Each infilling task is constituted by an altered function from the code snippet for detection, in which one of the entities in the function is replaced by a designated token indicating the infilling position. Such entities are derived from abstract syntax trees (AST) leaf nodes. Associated information, such as substituted original tokens and the entity type, will be used to guide token generation.

A function may contain a number of entities, introducing unnecessary overhead to subsequent tasks. In order to expedite the detection process without detriment to precision, certain heuristic approaches were employed. Currently, WitheredLeaf targets EIBs on variable and function names. To exclude entities with a minimal likelihood of containing bugs, it skips: 1) function names in their declarations; 2) declarations of variables, including formal arguments; and 3) imported library names (*e.g.*, torch in Python). These entities typically cannot incur EIBs, as related errors are mostly caught in compilers or interpreters.

CodeBERT. By design, the CodeBERT model aims to generate a fixed number of tokens in the code. An entity name usually contains more than one token. For example, "params_quoted" contains four tokens: params, _, qu, and oted. By evaluating the $D_{synthesized}$, we find a feasible search range that can cover 91.2% cases.

Given the entity’s potential token lengths range, we can use CodeBERT to check the consistency of the original tokens, as shown in the Algorithm 1. For each possible token length, CodeBERT will generate an entity candidate and its corresponding score. Since CodeBERT serves as the pre-filter, it is crucial to minimize the

Algorithm 1: CodeBERT-based Consistency Check

```

Input: maskedCode, originalTokens, searchRange, tokenThreshold,
        consistencyThreshold
Output: consistencyFlag

Initialize candidateScores as empty map ;
originalLength ← originalTokens.tokensLength;
tokensLengthsRange ← lengthRange(originalLength, searchRange);
foreach tokensLength in tokensLengthsRange do
    entity, score ← CodeBERTPredict(maskedCode, tokensLength) ;
    /* entity.tokensLength equals to tokensLength */
    if score < tokenThreshold then
        | skip to the next iteration;
    | candidateScores[entity] ← score ;

if originalTokens in candidateScores then
    foreach entity, score in candidateScores do
        if score > consistencyThreshold and entity ≠ originalTokens
        then
            | return False;
    | return True;
return False;

```

false negative rate. Therefore, we select a *tokenThreshold*; if the returned entity candidate’s score falls below this threshold, we skip this token length and proceed to the next iteration. In the end, if all of the possible context lengths are skipped, we consider it as unpredictable for CodeBERT and escalate it to the next component. For the inconsistency check, if the original entity name is the only entity name generated by CodeBERT, we consider this as semantic consistency, and will not escalate any report to the next component. If there are multiple entity names in the candidates list and the original entity name is included, we check the score of each candidate. Only if no candidate other than the original entity name has a confidence score higher than *consistencyThreshold*, we consider this as semantic consistency. Besides, if the original entity name does not occur in the entity candidates list, we consider this as semantic inconsistency.

Decoder-only Models. Token generation for decoder-only models like Code Llama is different from CodeBERT, as these models are not aware of the number of tokens to generate, and the generation stops till it generates a special end-of-sequence token. Intuitively, one way to enhance the specificity is to inspect the probabilities of generated tokens and regard a negative case if the original tokens are highly probable, if not the most probable. However, these models generate tokens sequentially, conditioned on the previous generation step, which renders inspection of the probabilities of a token meaningless if the model has already generated tokens inconsistent with the original ones. Thus, if the first generated token is “invalid”, the whole generation process will follow this invalid direction. For example, when a variable name is masked at the beginning of the function body, the model usually begins to generate a documentation string of the function (*i.e.*, the token ""), leading subsequent generation following a documentation string rather than the variable name.

Therefore, WitheredLeaf needs to harness the generation, ensuring the generated tokens constitute a syntactically valid entity. For example, variable names in Python can only contain capital

Algorithm 2: Consistency Check for Decoder-only Models

Input: maskedCode, originalTokens, probThresh, rankThresh,
Output: consistencyFlag

```

rankSum ← 0;
generatedTokens, leftTokens ← [], originalTokens;
while leftTokens.length > 0 do
    nextTokenProbList ← top_k(CodeLlamaPredict(maskedCode,
        generatedTokens));
    foreach token, prob in nextTokenProbList do
        if validateToken(token) then
            if leftTokens.startwith(token) then
                leftTokens.remove(token)
                generatedTokens.append(Token)
                break to generate the next token;
            else
                if prob > probThresh then
                    return False;
                rankSum ← rankSum + 1;
            else
                continue to the next iteration;
        if rankSum > rankThresh then
            return False;
return True;

```

and lower English letters, underscore, and non-beginning digits. The high-level idea of our approach is to select the tokens guided by the original tokens, and inspect the probabilities to judge the consistency. The detailed algorithm is demonstrated in Algorithm 2.

WitheredLeaf generates tokens one by one, and fixes the token from the original tokens before generating the next one to ensure the generation strictly follows the original tokens. It inspects the probabilities of the generated token candidates in each generation step. For each generated *token*, WitheredLeaf first checks if the token is syntactically valid via *validateToken()*. Then, if the *leftTokens* starts with the generated *token*, meaning the generation follows the original tokens, it updates the variable and continues to generate the next token. Otherwise, the generated *token* is penalized for not being consistent with the original ones. In a case where the probability *prob* is greater than a threshold *probThresh*, meaning that there is another more probably valid token different from the original tokens, the check directly returns False. It adds *rankSum* to 1, penalizing original tokens not being the most probable generated *token*. By comparing *rankSum* with a threshold *rankThresh*, the algorithm determines the consistency. Unlike CodeBERT, which is operated to lower the false negative rate, Code Llama is used as a post-filter and it is optimized to reduce false positive rate. Therefore, the thresholds are tuned to return True in more cases.

3.4 Inconsistent Position Highlighting

As measured in §2.3, GPT-4 can produce a number of false negatives, especially in a longer context. A previous study also notice the pattern and attributed it to distraction in long context [52]. Thus, we design a prompt engineering technique, forcing GPT-4 to attend back on suspicious positions. Such positions are collected from the inconsistent positions determined by local models. In addition to the existing prompt, WitheredLeaf instructs GPT-4: “Also, pay additional attention to these lines: {suspicious_lines}”. The

whole line, rather than the suspicious entity, is highlighted since the observation in §2.4 indicates the consistency of the infilling task conducted on the same line is negatively affected. Thus, the inconsistency detected may imply an EIB at another entity located on the same line.

3.5 Implementation

WitheredLeaf relies on tree-sitter [3] to parse the source code, since it supports multiple programming languages.

CodeBERT². CodeBERT is employed as a pre-filter in WitheredLeaf due to its rapid processing and compact size, allowing it to run on the CPU concurrently with other components utilizing the GPU and enhancing overall system efficiency. One limitation is its max token length is 512, so the masked code is truncated to fit into this limitation. Since CodeBERT fill each masked token independently, in *CodeBERTPredict*, we concatenates tokens in each token position to generate the predicted entity name. Every confidence score is a float number ranging from zero to one. The generated concatenated entity name’s score is the product of the scores of each token. For the parameter selection, after evaluating the $D_{synthesized}$, we selected *searchRange*, *tokenThreshold* and *consistencyThreshold* that effectively achieve a high recall rate, while maintaining a reasonably low false positive rate (see Table 5).

StableCode. StableCode is employed as an alternative pre-filter local model. It was released after our experiment with CodeBERT on Python code was finished, so we used it experimentally on C.

Code Llama. Code Llama is employed as the post-filter local model. Code Llama 7B base model is selected owing to its better performance and stability, as the 13B model occasionally results in out-of-memory errors. Although Meta provides 34B and 70B models, they are not shipped with FIM capability. We implement Algorithm 2 using PyTorch library [9]. We selected *rankThresh* and *probThresh* based on the evaluation of $D_{synthesized}$ to achieve high specificity, with a minor increase of false negative rate (see §4.1).

GPT-4. GPT-4 is used as the state-of-the-art LLM to detect EIBs following the prompts. We use the prompt template shown in Figure 2 with suspicious line highlighted. New models such as Claude 3 and Llama 3 have been recently released, claiming to be new state-of-the-art. However, their APIs currently only offer limited support to JSON output, which can hardly be integrated into an automated detection pipeline. We leave the measurement and integration of these models as future work.

4 EVALUATION

In this section, we first conduct controlled experiments on our prepared datasets (§2.2). This is to (1) evaluate the effectiveness of individual components in our design, and (2) compare WitheredLeaf with existing detectors. After that, we evaluate WitheredLeaf on real-world code repositories and demonstrate its ability to discover new bugs. The experiments are conducted on a Ubuntu 22.4 Linux server with two AMD EPYC 9124 CPUs, one NVIDIA H100 GPU, and 384GB memory.

²We employ CodeBERT (MLM) on to verify consistency.

4.1 Controlled Experiments

4.1.1 Individual Components. First, to access the effectiveness of our design, we evaluate the performance of individual components in the pipeline. Similar to the experiments in §2.4 with $D_{sampled}$ and $D_{synthesized}$, we assess the recall and consistency of code models, but with controlled token generation algorithms as detailed in §3.3. The evaluation also helps us to establish the algorithm parameters, which will be used in real-world EIB detection tasks in §4.2.

CodeBERT. We explored various configurations of the parameters *searchRange*, *tokenThreshold* and *consistencyThreshold* in Algorithm 1 to evaluate CodeBERT’s performance, as shown in Table 5. A higher *searchRange* enhances recall but increases the running time. Similarly, *tokenThreshold* is positively correlated with recall; however, higher values significantly decrease the consistency rate. Conversely, *consistencyThreshold* is inversely proportional to recall and directly affects the consistency rate. Considering the computational costs, and the recall, we selected 2 as the *searchRange* and 0.8 as *consistencyThreshold* to optimally balance recall with a satisfactory consistency rate. Additionally, we fixed the *searchRange* and *consistencyThreshold* to evaluate *tokenThreshold* as shown in Table 6, and finally, we selected 0.1 as the *tokenThreshold*.

Code Llama. To evaluate the effectiveness of controlled token generation proposed in Algorithm 2, we test the infilling task with different threshold parameters on Code Llama 7B model. The scores are presented in Table 7 in percentage. Our algorithm can roughly enhance the consistency by 5%, with the cost of 9% reduction on recall. This is expected, as tuning towards enhancing consistency results in reduced recall. We believe this is a worthy tradeoff, as EIBs are very rare in real-world code - a higher consistency will help eliminate more true negative cases and improve the overall precision. Therefore, we select the parameters of *probThresh* = 0.90 and *rankThresh* = 3.

GPT-4. For the additional prompt template highlighting suspicious bugs, the experiment is identical to §2.3, and the results are listed in Table 2 along with the measurement results. In the template, the ground truth EIB line, with three additional randomly selected lines are highlighted in the prompt. This mechanism significantly enhances the number of EIBs being detected by 24% on $D_{synthesized}$ compared to the best prompt template without highlighting, successfully reducing the FNs potentially caused by distraction. However, such a mechanism is not perfect, as it can incur more FPs (by 20% function-wise and by 7% bug-wise) if the line with EIB is not properly highlighted, as accessed on $D_{sampled}$. We believe this is a worthy tradeoff, as the density of EIBs reaches a high level at this stage. Also, please note that highlighting is only made possible with the local models.

4.1.2 Whole System. To assess the performance of our bug detection system in a controlled setting, we evaluate it on our synthesized dataset, $D_{synthesize}$ and sampled dataset, $D_{sampled}$. The purpose of this evaluation is to measure our system’s ability to detect these artificially generated bugs, which resemble real-world bugs caused by incorrect identifier usage. Therefore, we can access the recall and specificity of our system.

In $D_{sampled}$, *WitheredLeaf* reported 57 out of 100 functions have EIBs, resulting in 43% specificity. In $D_{synthesize}$, *WitheredLeaf* identified 41 bugs from the 66 samples. Moreover, our system demonstrates superior monetary efficiency: 730K tokens are transmitted in total, incurring a cost of 9.92\$ on GPT-4 API.

Using the formulas in §3.2 and the metrics of recall and consistency in bold in Table 5, Table 7, and Table 2, we could theoretically estimate the metrics of the whole system. As consistency is a lower-bound estimation of specificity, we derive the lower-bound recall and precision as 67.8% and 14.5%, respectively, when assuming $\epsilon_0 = 0.002$. Such theoretical values are on par with the experimental results on our datasets and in real-world detection (§4.2). Notably, when solely using GPT-4 as an EIB detector, the precision is as low as 0.8%. Besides, *WitheredLeaf* can effectively amplify the density of EIB by at least 70 times, making them detectable.

4.1.3 Comparison with Previous Work. Mansour et al. [45] introduced *DeepBugs*, a machine learning-based approach to name-based bug detection. While *DeepBugs* relies on an extensive dataset to train its classifier model, it inherently fails to identify bug types not present in the dataset. This limitation confines *DeepBugs* to detect only the three types of name-based bugs detailed in their paper. In a similar vein, Michael et al. [7] introduced *FICS*, a method clustering AST of all functions in a repository and reporting the outliers as potential bugs. As they noted, *FICS* eliminates all variable names when processing the AST, which precludes its ability to detect the EIB discussed in this paper. Moreover, their methodology implies a basic assumption: to detect a bug, there must be at least another “correct” function with a similar AST. Such an assumption is a bit too strong, especially for small code repositories, as similar ASTs may not exist for most functions. Unlike these systems, *WitheredLeaf* requires no training or fine-tuning.

The method presented by Baleegh et al. in [6] relates closely to our study. They propose using LLM for bug detection by comparing original code against code generated by ChatGPT. Our evaluation, however, reveals two major concerns with this approach: 1) it fails to detect a significant number of bugs and it generates a high rate of false positives; and 2) it is much less economically viable compared to ours. We experiment with subsets of datasets comprising 66 functions from $D_{synthesized}$, each with an EIB, and 100 functions from $D_{sampled}$. In this dataset, *FLAG* used approximately 12.33 million tokens, costing \$246.89 in GPT-4 API. For benign codes without bugs, *FLAG* erroneously highlighted 95 out of 100 samples, identifying totally 1751 lines as buggy out of 7062. For the buggy functions, it successfully detected issues in 32 out of 66 cases, highlighting 2454 buggy lines out of 7531. In contrast, our system demonstrates less false positive reports and superior monetary efficiency as mentioned previously.

4.2 Real-world EIB Detection

To assess the effectiveness of *WitheredLeaf* on real-world code repositories, we test it on Python and C repositories collected from GitHub with more than 1000 stars that are actively maintained. The summarized results are presented in Table 8. It’s worth mentioning that our result exceeds our expectation of 20% precision on the final report, which aligns with previous work on inconsistency bug detection [7] and is acceptable by the developer community [8].

Table 5: Recall and consistency of CodeBERT with different threshold parameters in the infilling task.

consistencyThreshold	searchRange=1		searchRange=2		searchRange=3		searchRange=4		searchRange=5	
	Recall	Cons.	Recall	Cons.	Recall	Cons.	Recall	Cons.	Recall	Cons.
0.98	72.9	65.5	75.2	68.2	75.8	65.5	75.3	59.7	75.1	61.5
0.90	77.7	61.7	87.4	61.6	87.5	59.4	87.7	53.3	86.8	54.4
0.80	78.0	58.4	91.2	54.1	90.5	52.5	90.5	43.6	90.0	44.5
0.70	79.3	55.7	92.6	47.7	91.3	45.7	91.3	36.4	90.9	37.2
0.60	80.3	53.1	93.6	43.0	93.3	38.8	92.9	31.0	92.5	31.7
Avg. Time	506		640		689		933		1103	

Avg. Time represents the average time in milliseconds spent to accomplish an infilling task on CPU.

Table 6: Recall and consistency of CodeBERT with different tokenThreshold.

tokenThreshold	Recall	Cons.
0.1	91.2	54.1
0.2	92.3	40.6
0.3	94.1	35.8
0.4	94.6	32.6
0.5	95.0	32.0

Evaluated with the specified parameters of **searchRange=2** and **consistencyThreshold=0.80**.

Table 7: Recall and consistency of Code Llama 7B of different threshold parameters in the infilling task.

probThresh	rankThresh=1		rankThresh=2		rankThresh=3	
	Recall	Cons.	Recall	Cons.	Recall	Cons.
0.98	92.7	89.0	88.4	90.8	85.8	91.8
0.90	94.6	88.6	91.1	90.1	89.2	91.0
0.80	95.7	88.2	93.1	89.5	91.5	90.2
0.70	96.4	87.7	95.2	88.9	93.9	89.6
0.60	97.3	87.1	96.4	88.1	95.4	88.8

Detecting EIBs on Python Repositories. For Python repositories, WitheredLeaf uses CodeBERT and Code Llama 7B as local models to process code from 80 collected repositories. Initially, CodeBERT deals with ~80K infilling tasks, derived from static analysis. Then we drop the consistent ones (~48K) and tasks associated with very long code snippets that are usually configuration or data files. After this pre-filter, ~21K infilling tasks are scheduled for Code Llama 7B (i.e., the post-filter), which finally discovers 2,679 inconsistency entities. Escalating the functions containing these inconsistencies to GPT-4 for EIB detection, we received 1,082 reports, of which 314 are left after filtering based on the properties mentioned in §2.3. Three security experts spend about 10 hours in total to inspect these reports. Finally, we confirm that 74 of them are true positive, resulting in a precision of 23.57% on Python repositories.

We believe the detection cost is well acceptable, as it spends less than 20 machine hours in total and less than 40\$ for querying GPT-4 APIs. With a server cost of 52\$³, the total cost to find these EIBs is less than 100\$, which we believe economically efficient, especially when considering the potential severe security consequences of the vulnerabilities.

³A 1xH100 server costs 2.49\$/hr on Lambda GPU cloud.

Detecting EIBs on C Repositories. We also conduct a experiment on C repositories at a smaller scale, where the entities for infilling tasks are generated according to an AST similarity strategy [11, 30] following previous research [7]. With this approach, AST subtrees with similar types and structures within a function were identified, and only entities from these corresponding AST subtrees were analyzed. The pre-filter local model used for C repositories is StableCode, as CodeBERT is not trained on C code.

We observe a enhanced precision of WitheredLeaf on C repositories. This is because Code Llama performs significantly better on C than Python, filtering out 97.1% consistent entities. A potential explanation is that C language induces fewer defects than Python [46], making the code more well-formed and entities predictable for Code Llama. Nevertheless, the inclusion of both Python and C repositories in our real-world evaluation allows us to assess the generalizability of WitheredLeaf across different programming languages and projects, which is crucial for understanding its potential adaptability and weaknesses in large-scale deployment.

Repository Exclusion. It is worth noting that we intentionally exclude 10 repositories in this evaluation due to several reasons. Some repositories, particularly those written in Python 2 using outdated syntax and APIs, can deteriorate the performance of the models. This is because the models are trained from the corpus collected in recent years, resulting in a lack of training material in Python 2. We also observe the consistency of infilling drops on repositories related to AI, security, and system-level projects, such as chat-langchain [1] and hosts [2]. For AI-related projects, the technologies are evolving very fast, thus low consistency can be attributed to a lack of up-to-date corpus. However, for security and low-level system code, we think the low consistency is caused by the unusual nature of their functionalities. For example, exploit generation tools include insecure statements in their code, which can be regarded as inconsistent by local models in infilling, or reported to contain vulnerabilities by GPT-4. However, these tools are meant to be “insecure”. We believe such weakness could be solved by fine-tuning or more fine-grained prompt engineering, which we leave as future work.

4.3 Study of the Discovered EIBs

WitheredLeaf has identified 123 previously unknown EIBs in real-world codebases (§4.2). We submitted 69 pull requests to address these bugs, of which 27 have already been confirmed and merged by developers. The distribution of bugs across different categories

Table 8: Summary of real-world detection results.

	Python Repos	C Repos
# Repositories	80	74
# Submitted PR	55	14
# Merged PR	26	1
# Functionality Bugs	37	7
# Security Bugs	8	3
# Bad Smell	48	20
Pre-filter Execution Time	9.0 CPU Hrs	15.2 GPU Hrs
Post-filter Execution Time	20.9 GPU Hrs	17.1 GPU Hrs
Cost on GPT-4 API	39.96\$	16.21\$
# Infilling Tasks for Pre-filter	79,724	75,695
# Inconsis. Entities after Pre-filter	31,284	25,214
# Infilling Tasks for Post-filter	20,616	25,214
# Inconsis. Entity after Post-filter	2,679	733
# Reports from GPT-4	1,082	473
# Reports after Filtering	314	77
# Confirmed Reports	74	28
Precision (%)	23.57%	36.36%

is detailed in Table 8. The security bugs identified in Python repositories can lead to denial-of-service (DoS) attacks, while those in C repositories may cause memory corruption, potentially leading to severe consequences such as privilege escalation. Remarkably, our discovered EIBs impact prominent repositories. For example, a director from wolfSSL [58] scheduled a meeting with us and mentioned that they were actively investigating the bugs we reported. In the remainder of this section, we showcase some of our discovered EIBs to further demonstrate WitheredLeaf’s capabilities.

Case Study. We present some of our found EIBs, with the corresponding fixes (the code entity highlighted in red is replaced with that in green). We demonstrate the consequences of these bugs and analyze the potential cause of them.

Variable Misuse In the case Listing 1, the last branch of the switch statement incorrectly returns the value of `_threadFrameStackSize`, which is the same as the previous branch. However, it should return `_threadExceptionFlowSize`. The sizes are used in buffer initialization and could be adjusted via the APIs provided in this library. This bug can potentially lead to buffer overflow, causing more severe security consequences. We suspect the developer forgot to correct the variable name when copy and paste within the code.

```

1 switch (optionId) {
2 ...
3 case InterpreterThreadFrameStackSize:
4     return s_threadFrameStackSize;
5 case InterpreterThreadExceptionFlowSize:
6     return s_threadFrameStackSize s_threadExceptionFlowSize ;}

```

Listing 1: Example of Variable Misuse

Error in Function Call This bug in Listing 2 is found in WolfSSL, a security-sensitive repository. The `sizeof` method is used to determine the size (in bytes) of a data type or a variable. However, in this case, `sizeof` method is invoked on `xSize` variable, which is already declared as the size of `freertos_sockaddr` - invoking `sizeof` on it will result in an unexpected size, likely leading to an overflow or

underflow. After a deeper investigation, we find that the third size parameter of the function `FreeRTOS_bind` is not currently in use, so this bug does not manifest security consequences at this point. However, such a bug is like a time bomb, which will cause trouble when the third parameter takes effect in the future. We think this bug might be caused by naming confusion of the variable `xSize`.

```

1 socklen_t xSize = sizeof(struct freertos_sockaddr);
2 xSocket_t xClientSocket = NULL;
3 struct freertos_sockaddr xRemoteAddress;
4 ...
5 FreeRTOS_bind(xClientSocket, &xRemoteAddress,
6     sizeof(xSize) xSize );

```

Listing 2: Example of Error in Function Call

Denial of Service In the code snippet in Listing 3, if `has_aux` is `False`, `aux` will be a tuple object with zero length. Attempting to access the last element of `aux` using an index of `-1` will throw `index out of range` exception, causing Denial of Service (DoS).

```

1 if has_aux:
2     example, *aux = example
3 else:
4     aux = tuple()
5 rand_state = random.Random(aux[-1] aux[-1] if aux else 0 )

```

Listing 3: Example of DoS

5 RELATED WORK

LLM’s Application On Software Security. Recent advancements have made LLMs designed for coding widely accessible, either through public API [12] or open-source sharing of code and models [25, 27, 48]. LLMs have demonstrated significant contributions to areas like fuzzing [5, 21, 39, 59], code repair [22, 31, 32, 43, 60], exploit generation [23], etc. Despite these developments, the application of LLMs in bug detection, specifically entity inconsistency bug (EIB), remains underexplored and largely experimental.

Semantic Bug Detection. Semantic bug detection is an active research area with various methodologies addressing bugs that exhibit both well-defined and ambiguous patterns. Li et al. pioneered the application of LLMs to target Use Before Initialization bugs, demonstrating LLMs’ potential to enhance conventional static analysis techniques [34]. Further advancements by researchers like Sun et al. and Wei et al. have expanded LLM applications to address a wider range of logic bugs, including API misuse [52, 53, 57]. Traditional methods such as static analysis and formal verification also play significant roles in identifying and resolving semantic bugs with fixed patterns [24, 28, 36, 37, 40]. Recently, Natural Language Processing (NLP) has been leveraged to analyze documentation and pinpoint potential semantic inconsistencies or discrepancies between implementation and documentation [13–16, 42]. However, these NLP-based approaches heavily rely on the quality of the documentation, facing substantial challenges with ambiguous or outdated materials prevalent in open-source environments.

FLAG [6] is a recent initiative that utilizes LLMs for bug detection by masking every line of code and comparing the code line generated by ChatGPT with the original line. As we demonstrated in §4.1.3, such methods exhibit limitations when applied to large-scale, real-world analysis. Beyond LLM-based methods, other approaches

such as [7, 45] utilize learning based methods to pinpoint potential bugs. However, these methods require specific code characteristics or depend on ideal, comprehensive training datasets.

6 DISCUSSION

Limitations. Currently, decoder-only models like Code Llama are limited to infilling at a single position each time. Although WitheredLeaf can independently verify the consistency of each entity, it may not detect EIBs associated with multiple positions in the context.

Future work. Currently, we are using existing infilling models to provide GPT-4 with highlights of potentially buggy lines. It may be beneficial to fine-tune or train a model that exhibits improved consistency and recall rates. Fine-tuned model demonstrates better capability on auditing smart contracts [38]. Replaced Tokens Detection [17], although simple, fits our scenario with new model trained. Researchers currently consider this task as a form of pretraining, which enhances the model’s capability to comprehend semantics more effectively.

As mentioned in §2.4, EIB can negatively affect the consistency of infilling tasks on surrounding positions, which might be qualitatively explainable with a confidence measure [50]. Enhancing this observation and proposing a method providing a security guarantee to EIB detection is a direction for future investigation.

7 CONCLUSION

We conducted a systematic measurement study on LLMs’ capabilities to detect EIBs and identified their strengths and weaknesses. Based on the measurement results, we design WitheredLeaf, an EIB detection pipeline based on LLMs. WitheredLeaf successfully discovered 123 unknown bugs on Python and C code repositories, demonstrating the effectiveness of our design and implementation.

ACKNOWLEDGMENTS

To Robert, for the bagels and explaining CMYK and color spaces.

REFERENCES

- [1] 2024. langchain-ai/chat-langchain. <https://github.com/langchain-ai/chat-langchain> Accessed: 2024-04-29.
- [2] 2024. StevenBlack/hosts: Consolidating and extending hosts files from several well-curated sources. Optionally pick extensions for porn, social media, and other categories. <https://github.com/StevenBlack/hosts> Accessed: 2024-04-29.
- [3] 2024. tree-sitter. <https://github.com/tree-sitter/tree-sitter> Accessed: 2024-04-29.
- [4] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altmenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774* (2023).
- [5] Joshua Ackerman and George V. Cybenko. 2023. Large Language Models for Fuzzing Parsers (Registered Report). *Proceedings of the 2nd International Fuzzing Workshop* (2023).
- [6] Baleegh Ahmad, Benjamin Tan, Ramesh Karri, and Hammond Pearce. 2023. FLAG: Finding Line Anomalies (in code) with Generative AI. *arXiv preprint arXiv:2306.12643* (2023).
- [7] Mansour Ahmadi, Reza Mirzazade Farkhani, Ryan Williams, and Long Lu. 2021. Finding bugs using your own code: detecting functionally-similar yet inconsistent code. In *30th USENIX security symposium (USENIX Security 21)*. 2025–2040.
- [8] Amit Seal Ami, Kevin Moran, Denys Poshyvanyk, and Adwait Nadkarni. 2023. "False negative—that one is going to kill you."—Understanding Industry Perspectives of Static Analysis based Security Testing. In *2024 IEEE Symposium on Security and Privacy (SP)*. IEEE Computer Society, 19–19.
- [9] Jason Ansel, Edward Yang, Horace He, Natalia Gimelshein, Animesh Jain, Michael Voznesensky, Bin Bao, Peter Bell, David Berard, Evgeni Burovski, Geeta Chauhan, Anjali Chourdia, Will Constable, Alban Desmaison, Zachary DeVito, Elias Ellison, Will Feng, Jiong Gong, Michael Gschwind, Brian Hirsh, Sherlock Huang, Kshiteej Kalambarkar, Laurent Kirsch, Michael Lazos, Mario Lezciano, Yanbo Liang, Jason Liang, Yinghai Lu, CK Luk, Bert Maher, Yunjie Pan, Christian Puhres, Matthias Reso, Mark Saroufim, Marcos Yukio Siraichi, Helen Suk, Michael Suo, Phil Tillet, Eikan Wang, Xiaodong Wang, William Wen, Shunting Zhang, Xu Zhao, Keren Zhou, Richard Zou, Ajit Mathews, Gregory Chanan, Peng Wu, and Soumith Chintala. 2024. PyTorch 2: Faster Machine Learning Through Dynamic Python Bytecode Transformation and Graph Compilation. In *29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2 (ASPLOS '24)*. <https://doi.org/10.1145/3620665.3640366>
- [10] Mohammad Bavarian, Heewoo Jun, Nikolas Tezak, John Schulman, Christine McLeavey, Jerry Tworek, and Mark Chen. 2022. Efficient training of language models to fill in the middle. *arXiv preprint arXiv:2207.14255* (2022).
- [11] Ira D Baxter, Andrew Yahin, Leonardo Moura, Marcelo Sant’Anna, and Loraine Bier. 1998. Clone detection using abstract syntax trees. In *Proceedings. International Conference on Software Maintenance (Cat. No. 98CB36272)*. IEEE, 368–377.
- [12] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde, Jared Kaplan, Harrison Edwards, Yura Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, David W. Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William H. Guss, Alex Nichol, Igor Babuschkin, Suchir Balaji, Shantanu Jain, Andrew Carr, Jan Leike, Joshua Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew M. Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021. Evaluating Large Language Models Trained on Code. *ArXiv abs/2107.03374* (2021).
- [13] Yi Chen, Di Tang, Yepeng Yao, Mingming Zha, XiaoFeng Wang, Xiaozhong Liu, Haixu Tang, and Baoxu Liu. 2023. Sherlock on Specs: Building {LTE} Conformance Tests through Automated Reasoning. In *32nd USENIX Security Symposium (USENIX Security 23)*. 3529–3545.
- [14] Yi Chen, Di Tang, Yepeng Yao, Mingming Zha, XiaoFeng Wang, Xiaozhong Liu, Haixu Tang, and Dongfang Zhao. 2022. Seeing the forest for the trees: Understanding security hazards in the {GPP} ecosystem through intelligent analysis on change requests. In *31st USENIX Security Symposium (USENIX Security 22)*. 17–34.
- [15] Yi Chen, Luyi Xing, Yue Qin, Xiaojing Liao, XiaoFeng Wang, Kai Chen, and Wei Zou. 2019. Devils in the guidance: predicting logic vulnerabilities in payment syndication services through automated documentation analysis. In *28th USENIX Security Symposium (USENIX Security 19)*. 747–764.
- [16] Yi Chen, Yepeng Yao, XiaoFeng Wang, Dandan Xu, Chang Yue, Xiaozhong Liu, Kai Chen, Haixu Tang, and Baoxu Liu. 2021. Bookworm game: Automatic discovery of lte vulnerabilities through documentation analysis. In *2021 IEEE Symposium on Security and Privacy (SP)*. IEEE, 1197–1214.
- [17] Kevin Clark, Minh-Thang Luong, Quoc V Le, and Christopher D Manning. 2020. Electra: Pre-training text encoders as discriminators rather than generators. *arXiv preprint arXiv:2003.10555* (2020).
- [18] CloCkWeRX. 2024. Pull Request for Bug Fix in RabbitVCS. <https://github.com/rabbitvcs/rabbitvcs/pull/385/files#diff-eca3c60057143346eea4a850ba4fb60752c8cb397e380293af7e252c1677d0f>.
- [19] Roi Cohen, May Hamri, Mor Geva, and Amir Globerson. 2023. Lm vs lm: Detecting factual errors via cross examination. *arXiv preprint arXiv:2305.13281* (2023).
- [20] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, Jill Burstein, Christy Doran, and Thamar Solorio (Eds.). 4171–4186. <https://doi.org/10.18653/v1/N19-1423>
- [21] Jueon Eom, Seyeon Jeong, and Taekyoung Kwon. 2024. CovRL: Fuzzing JavaScript Engines with Coverage-Guided Reinforcement Learning for LLM-based Mutation. *ArXiv abs/2402.12222* (2024).
- [22] Zhiyu Fan, Xiang Gao, Martin Mirchev, Abhik Roychoudhury, and Shin Hwei Tan. 2022. Automated Repair of Programs from Large Language Models. *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)* (2022), 1469–1481.
- [23] Richard Fang, Rohan Bindu, Akul Gupta, and Daniel Kang. 2024. LLM Agents can Autonomously Exploit One-day Vulnerabilities.
- [24] Ansgar Fehnker and Ralf Huuck. 2013. Model checking driven static analysis for the real world: designing and tuning large scale bug detection. *Innovations in systems and software engineering* 9, 1 (2013), 45–56.
- [25] Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiaocheng Feng, Ming Gong, Linjun Shou, Bing Qin, Ting Liu, Daxin Jiang, et al. 2020. CodeBERT: A Pre-Trained Model for Programming and Natural Languages. In *Findings of the Association for Computational Linguistics: EMNLP 2020*. 1536–1547.

- [26] Google. 2024. AddressSanitizer. <https://github.com/google/sanitizers/wiki/AddressSanitizer>.
- [27] Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Y Wu, YK Li, et al. 2024. DeepSeek-Coder: When the Large Language Model Meets Programming—The Rise of Code Intelligence. *arXiv preprint arXiv:2401.14196* (2024).
- [28] Jian Huang, Michael Allen-Bond, and Xuechen Zhang. 2017. Pallas: Semantic-aware checking for finding deep bugs in fast path. In *Proceedings of the Twenty-Second International Conference on Architectural Support for Programming Languages and Operating Systems*. 709–722.
- [29] Sungjae Hwang and Sukyoung Ryu. 2020. Gap between theory and practice: An empirical study of security patches in solidity. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 542–553.
- [30] Lingxiao Jiang, Ghassan Mishserghi, Zhendong Su, and Stephane Gloudu. 2007. Deckard: Scalable and accurate tree-based detection of code clones. In *29th International Conference on Software Engineering (ICSE '07)*. IEEE, 96–105.
- [31] Ma Jin, Syed Shahriar, Michele Tufano, Xin Shi, Shuai Lu, Neel Sundaresan, and Alexey Svyatkovskiy. 2023. InferFix: End-to-End Program Repair with LLMs. *Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering* (2023).
- [32] Harshit Joshi, José Pablo Cambronero, Sumit Gulwani, Vu Le, Ivan Radicek, and Gust Verbruggen. 2022. Repair Is Nearly Generation: Multilingual Program Repair with LLMs. In *AAAI Conference on Artificial Intelligence*.
- [33] Fatemeh Khoshnoud, Ali Rezaei Nasab, Zahra Toudeji, and Ashkan Sami. 2022. Which bugs are missed in code reviews: An empirical study on SmartSHARK dataset. In *Proceedings of the 19th International Conference on Mining Software Repositories*. 137–141.
- [34] Haonan Li, Yu Hao, Yizhuo Zhai, and Zhiyun Qian. 2024. Enhancing Static Analysis for Practical Bug Detection: An LLM-Integrated Approach. *Proceedings of the ACM on Programming Languages (PACMPL)*, Issue OOPSLA (2024).
- [35] Raymond Li, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Li Jia, Jenny Chim, Qian Liu, et al. 2023. StarCoder: may the source be with you! *Transactions on Machine Learning Research* (2023).
- [36] Zhuohua Li, Jincheng Wang, Mingshen Sun, and John CS Lui. 2021. MirChecker: detecting bugs in Rust programs via static analysis. In *Proceedings of the 2021 ACM SIGSAC conference on computer and communications security*. 2183–2196.
- [37] Kui Liu, Anil Koyuncu, Dongsun Kim, and Tegawendé F Bissyandé. 2019. Avatar: Fixing semantic bugs with fix patterns of static analysis violations. In *2019 IEEE 26th International Conference on Software Analysis, Evolution and Reengineering (SANER)*. IEEE, 1–12.
- [38] Wei Ma, Daoyuan Wu, Yuqiang Sun, Tianwen Wang, Shangqing Liu, Jian Zhang, Yue Xue, and Yang Liu. 2024. Combining Fine-Tuning and LLM-based Agents for Intuitive Smart Contract Auditing with Justifications. *arXiv preprint arXiv:2403.16073* (2024).
- [39] Ruijie Meng, Martin Mirchev, Marcel Böhme, and Abhik Roychoudhury. 2024. Large Language Model guided Protocol Fuzzing. *Proceedings 2024 Network and Distributed System Security Symposium* (2024).
- [40] Hoang Duong Thien Nguyen, Dawei Qi, Abhik Roychoudhury, and Satish Chandra. 2013. Semfix: Program repair via semantic analysis. In *2013 35th International Conference on Software Engineering (ICSE)*. IEEE, 772–781.
- [41] OpenAI. 2024. Prompt engineering - OpenAI API. <https://platform.openai.com/docs/guides/prompt-engineering>.
- [42] Sheena Panthaplackel, Junyi Jessy Li, Milos Gligoric, and Raymond J Mooney. 2021. Deep just-in-time inconsistency detection between comments and source code. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 427–435.
- [43] Hammond A. Pearce, Benjamin Tan, Baleegh Ahmad, Ramesh Karri, and Brendan Dolan-Gavitt. 2021. Examining Zero-Shot Vulnerability Repair with Large Language Models. *2023 IEEE Symposium on Security and Privacy (SP)* (2021), 2339–2356.
- [44] Nikhil Pinnaparaju, Reshinth Adithyan, Duy Phung, Jonathan Tow, James Baicoianu, Ashish Datta, Maksym Zhuravinskyi, Dakota Mahan, Marco Belagente, Carlos Riquelme, et al. 2024. Stable Code Technical Report. *arXiv preprint arXiv:2404.01226* (2024).
- [45] Michael Pradel and Koushik Sen. 2018. Deepbugs: A learning approach to name-based bug detection. *Proceedings of the ACM on Programming Languages* 2, OOPSLA (2018), 1–25.
- [46] Baishakhi Ray, Daryl Posnett, Vladimir Filkov, and Premkumar Devanbu. 2014. A large scale study of programming languages and code quality in github. In *Proceedings of the 22nd ACM SIGSOFT international symposium on foundations of software engineering*. 155–165.
- [47] Cedric Richter and Heike Wehrheim. 2022. Learning realistic mutations: Bug creation for neural bug detectors. In *2022 IEEE Conference on Software Testing, Verification and Validation (ICST)*. IEEE, 162–173.
- [48] Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. 2023. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950* (2023).
- [49] Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed H Chi, Nathanael Schärli, and Denny Zhou. 2023. Large language models can be easily distracted by irrelevant context. In *International Conference on Machine Learning*. PMLR, 31210–31227.
- [50] Claudio Spiess, David Gros, Kunal Suresh Pai, Michael Pradel, Md Rafiqul Islam Rabin, Susmit Jha, Prem Devanbu, and Toufique Ahmed. 2024. Quality and Trust in LLM-generated Code. *arXiv preprint arXiv:2402.02047* (2024).
- [51] Benjamin Steenhoek, Md Mahbubur Rahman, Monoshi Kumar Roy, Mirza Sanjida Alam, Earl T Barr, and Wei Le. 2024. A Comprehensive Study of the Capabilities of Large Language Models for Vulnerability Detection. *arXiv preprint arXiv:2403.17218* (2024).
- [52] Yuqiang Sun, Daoyuan Wu, Yue Xue, Han Liu, Wei Ma, Lyuye Zhang, Miaolei Shi, and Yang Liu. 2024. LLM4Vuln: A Unified Evaluation Framework for Decoupling and Enhancing LLMs' Vulnerability Reasoning. *arXiv preprint arXiv:2401.16185* (2024).
- [53] Yuqiang Sun, Daoyuan Wu, Yue Xue, Han Liu, Haijun Wang, Zhengzi Xu, Xiaofei Xie, and Yang Liu. 2024. Gptscan: Detecting logic vulnerabilities in smart contracts by combining gpt with program analysis. In *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering*. 1–13.
- [54] Lin Tan, Chen Liu, Zhenmin Li, Xuanhui Wang, Yuanyuan Zhou, and Chengxiang Zhai. 2014. Bug characteristics in open source software. *Empirical software engineering* 19 (2014), 1665–1705.
- [55] Ruixiang Tang, Yu-Neng Chuang, and Xia Hu. 2024. The Science of Detecting LLM-Generated Text. *Commun. ACM* 67, 4 (2024), 50–59.
- [56] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems* 35 (2022), 24824–24837.
- [57] Moshi Wei, Nima Shiri Harzevili, YueKai Huang, Jinqu Yang, Junjie Wang, and Song Wang. 2024. Demystifying and Detecting Misuses of Deep Learning APIs. In *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering*. 1–12.
- [58] wolfSSL Inc. 2024. wolfSSL. <https://github.com/wolfSSL/wolfssl> Accessed: 2024-04-29.
- [59] Chun Xia, Matteo Paltenghi, Jia Le Tian, Michael Pradel, and Lingming Zhang. 2023. Fuzz4All: Universal Fuzzing with Large Language Models. *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering* (2023).
- [60] Chun Xia, Yuxiang Wei, and Lingming Zhang. 2023. Automated Program Repair in the Era of Large Pre-trained Language Models. *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)* (2023), 1482–1494.

A PROMPT TEMPLATES

Round 1 Properties: code_line + explanation
Round 1 Prompt: {code} Output exact lines of semantic bugs and concise explanations of the bugs.

Figure 4: Prompt template 1.

Round 1 Properties: code_line + explanation + fixed_line + merely_fixable
Round 1 Prompt: {code} Output exact lines of semantic bugs, concise explanations of the bugs, the fixed line of the bugs, and if the bugs can be fixed by merely changing a variable/method name.

Figure 5: Prompt template 1FM.

Round 1 Properties: code_line + explanation
Round 1 Prompt: {code} Output exact lines of semantic bugs, concise explanations of the bugs, the fixed line of the bugs, and if the bugs can be fixed by merely changing a variable/method name.
Round 2 Properties: code_line* + explanation* + fixed_line + merely_fixable
Round 2 Prompt: Inspect these bugs, excluding 1. incorrect or unlikely bugs; 2. bugs caused by undefined symbols; 3. non-semantic bugs. The left bugs only break the intended functionality or lead to vulnerabilities. Answer if the bugs can be fixed by merely changing a variable/method name and fix them.

Figure 6: Prompt template 1/2FM.

Round 1 Properties: code_line + explanation
Round 1 Prompt: {code} Output exact lines of semantic bugs, concise explanations of the bugs, the fixed line of the bugs, and if the bugs can be fixed by merely changing a variable/method name.
Round 2 Properties: code_line* + explanation* + fixed_line + merely_fixable
Round 2 Prompt: Inspect these bugs, excluding 1. incorrect or unlikely bugs; 2. bugs caused by undefined symbols; 3. non-semantic bugs. The left bugs only break the intended functionality or lead to vulnerabilities. Answer if the bugs can be fixed by merely changing a variable/method name and fix them.
Round 3 Properties: code_line* + explanation* + priority
Round 3 Prompt: Now, Evaluate the priority of detected bugs in previous round. Assign higher priority for logic bugs, security vulnerabilities, or unexpected functionality, and lower priority for bad coding styles (bad smells).

Figure 7: Prompt template 1/2FM/3P.

Round 1 Properties: code_line + explanation
Round 1 Prompt: {code} Output exact lines of semantic bugs, concise explanations of the bugs, the fixed line of the bugs, and if the bugs can be fixed by merely changing a variable/method name.
Round 2 Properties: code_line* + explanation* + fixed_line + merely_fixable
Round 2 Prompt: Inspect these bugs, excluding 1. incorrect or unlikely bugs; 2. bugs caused by undefined symbols; 3. non-semantic bugs. The left bugs only break the intended functionality or lead to vulnerabilities. Answer if the bugs can be fixed by merely changing a variable/method name and fix them.
Round 3 Properties: code_line* + explanation* + category
Round 3 Prompt: Now, classify this bug into these categories using the code and previous contexts: Security Vulnerability, Logic Bug, Enhancement, Unexpected Behavior, Symbol Not Defined, Module Not Imported, Bad Smell, Not a Bug, or Others. If you assign others, explicitly specify its category.

Figure 8: Prompt template 1/2FM/3Ca.

Round 1 Properties: code_line + explanation
Round 1 Prompt: {code} Output exact lines of semantic bugs, concise explanations of the bugs, the fixed line of the bugs, and if the bugs can be fixed by merely changing a variable/method name.
Round 2 Properties: code_line* + explanation* + fixed_line + merely_fixable + category
Round 2 Prompt: Inspect these bugs, excluding 1. incorrect or unlikely bugs; 2. bugs caused by undefined symbols; 3. non-semantic bugs. The left bugs only break the intended functionality or lead to vulnerabilities. Answer if the bugs can be fixed by merely changing a variable/method name and fix them. Also, classify this bug into these categories using the code and previous contexts: Security Vulnerability, Logic Bug, Enhancement, Unexpected Behavior, Symbol Not Defined, Module Not Imported, Bad Smell, Not a Bug, or Others. Note that the snippet is from popular repositories and runs, so correct symbols not defined in the snippet are defined at other places, which is not Logic Bug. If you assign Others category, explicitly name the category.

Figure 9: Prompt template 1/2FMCa.

Round 1 Properties: code_line + explanation
Round 1 Prompt: {code} Output exact lines of semantic bugs, concise explanations of the bugs, the fixed line of the bugs, and if the bugs can be fixed by merely changing a variable/method name. Also, pay additional attention to these lines: {suspicious_lines}
Round 2 Properties: code_line* + explanation* + fixed_line + merely_fixable + category
Round 2 Prompt: Inspect these bugs, excluding 1. incorrect or unlikely bugs; 2. bugs caused by undefined symbols; 3. non-semantic bugs. The left bugs only break the intended functionality or lead to vulnerabilities. Answer if the bugs can be fixed by merely changing a variable/method name and fix them. Also, classify this bug into these categories using the code and previous contexts: Security Vulnerability, Logic Bug, Enhancement, Unexpected Behavior, Symbol Not Defined, Module Not Imported, Bad Smell, Not a Bug, or Others. Note that the snippet is from popular repositories and runs, so correct symbols not defined in the snippet are defined at other places, which is not Logic Bug. If you assign Others category, explicitly name the category.

Figure 10: Prompt template 1/2FMCa w/ HL.