

# Fuzzing: Randomness? Reasoning!

## Efficient Directed Fuzzing via Large Language Models

Xiaotao Feng  
360 Security Technology Inc.  
Beijing, China  
escofeng@gmail.com

Xiaogang Zhu  
School of Computer and Mathematical Sciences  
The University of Adelaide  
Adelaide, SA, Australia  
xiaogang.zhu@adelaide.edu.au

Kun Hu  
School of Science  
Edith Cowan University  
Joondalup, WA, Australia  
k.hu@ecu.edu.au

Jincheng Wang  
360 Security Technology Inc.  
Beijing, China  
wangjincheng1@360.cn

Yingjie Cao  
360 Security Technology Inc.  
Beijing, China  
yingjiecao@protonmail.com

Guang Gong  
360 Security Technology Inc.  
Beijing, China  
yingjiecao@protonmail.com

Jianfeng Pan  
360 Security Technology Inc.  
Beijing, China  
panjianfeng@360.cn

**Abstract**—Fuzzing is highly effective in detecting bugs due to the key contribution of randomness. However, randomness significantly reduces the efficiency of fuzzing, causing it to cost days or weeks to expose bugs. Even though directed fuzzing reduces randomness by guiding fuzzing towards target buggy locations, the dilemma of randomness still challenges directed fuzzers. Two critical components, which are seeds and mutators, contain randomness and are closely tied to the conditions required for triggering bugs. Therefore, to address the challenge of randomness, we propose to use large language models (LLMs) to remove the randomness in seeds and reduce the randomness in mutators. With their strong reasoning and code generation capabilities, LLMs can be used to generate reachable seeds that target pre-determined locations and to construct bug-specific mutators tailored for specific bugs. We propose RANDLUZZ, which integrates LLMs and directed fuzzing, to improve the quality of seeds and mutators, resulting in efficient bug exposure. RANDLUZZ analyzes function call chain or functionality to guide LLMs in generating reachable seeds. To construct bug-specific mutators, RANDLUZZ uses LLMs to perform bug analysis, obtaining information such as bug causes and mutation suggestions, which further help generate code that performs bug-specific mutations. We evaluate RANDLUZZ by comparing it with four state-of-the-art directed fuzzers, AFLGo, Beacon, WindRanger, and SelectFuzz. With RANDLUZZ-generated seeds, the fuzzers achieve an average speedup ranging from  $2.1\times$  to  $4.8\times$  compared to using widely-used initial seeds. Additionally, when evaluated on individual bugs, RANDLUZZ achieves up to a  $2.7\times$  speedup compared to the second-fastest exposure. On 8 bugs, RANDLUZZ can even expose them within 60 seconds.

## 1. Introduction

Fuzzing has been proven highly effective at detecting bugs in real-world applications, with the discovery of numerous bugs across a wide range of software systems [9], [41], [43]. While existing fuzzers have introduced various strategies for fuzzing to trigger bugs, random input generation remains central to their effectiveness. However, the other side of the coin is that the randomness often leads to inefficiency, causing fuzzers to run for days or even weeks to identify bugs located in deep code regions.

Typically, directed fuzzing [1] is more efficient than coverage-guided fuzzing [2], [41]; however, because directed fuzzing still relies on random input generation, it suffers from the same issue of inefficient bug discovery. Many solutions have been proposed to improve the efficiency of directed fuzzing [13], [14], [34]. Some directed fuzzers optimize the guidance towards target locations based on program properties, such as satisfying execution path constraints [18] or identifying deviation basic blocks [8]. Another focus is generating inputs that are more likely to reach target locations via solutions such as predicting reachable inputs [14], [45] and instrumenting related code only [24]. Nonetheless, the existing directed fuzzing uses coarse-grained guidance to reduce the randomness, with a distance-based metric to guide the random input generation so that fuzzing can gradually reach target locations.

To improve the efficiency of triggering bugs, a straightforward yet challenging solution is to reduce or even remove the randomness when generating bug-triggering inputs. To trigger a bug, an input needs to satisfy both reaching the buggy code location and examining specific execution states. Thus, correspondingly, two critical components that significantly impact the efficacy of triggering bugs are the seeds and the mutation operators (mutators). Seeds introduce randomness because the code regions they examine are

unknown. They may be far away from or close to target locations. If the seeds are already capable of reaching target locations, the efficiency of bug discovery can be greatly enhanced because it does not introduce randomness [17]. After reaching target locations, mutators are responsible of generating new inputs that examine diverse execution states for bugs. However, all existing works rely on pre-determined mutators, which are not designed for triggering specific bugs. Such pre-determined mutators increase the randomness in fuzzing because most mutations cannot move states towards buggy ones [25]. The challenge is that, both generating *reachable seeds*, which can reach target locations, and constructing *bug-specific mutators*, which are designed for triggering specific bugs, require deep understanding of the target programs and bugs.

To address the challenge, in this paper, we propose to use large language models (LLMs) to remove the randomness in generating reachable seeds and reduce randomness in mutators. *Our key insights are that i) LLMs [44] possess strong capabilities in reasoning about bug information from both documentation and program code; and ii) the code generation capabilities of LLMs offer the potential to create adaptable and customizable mutators.* By leveraging these reasoning and code generation capabilities, LLMs can be used to generate reachable seeds and construct bug-specific mutators, enabling more efficient examination of buggy code regions. This will significantly reduce the impact caused by randomness in fuzzing.

To use LLMs for reducing randomness in fuzzing, the basic required information is the bug information, program usage, and function summary, which aim to understand the semantic meaning of bugs and the program logic of target projects. To generate reachable seeds, because directly querying LLMs leads to exceeding token limits or misunderstanding, we query LLMs based on the Function Call Chain (FCC) that contains the paths from the entry point to the target vulnerable function. The FCC leads LLMs to generate reachable seeds step-by-step, which improves the effectiveness of generation process. If we cannot obtain FCC, we will generate reachable seeds based on the functionalities of functions that are neighbors of the target function. If a seed can reach a neighbor function, we can use FCC to guide the generation of reachable seeds. To generate bug-specific mutators, LLM is first utilized to analyze the cause for the target bug. The result is further utilized to generate mutation suggestions. Code is then generated for mutators based on the mutation suggestions via LLMs.

We develop RANDLUZZ, which stands for Random-Less Fuzzer, to demonstrate the performance of our method. We compare our RANDLUZZ with four state-of-the-art fuzzers, AFLGo [1], Beacon [12], WindRanger [8], and SelectFuzz [24]. The results show that RANDLUZZ-generated seeds significantly enhance bug discovery efficiency, achieving an average speedup ranging from  $2.1\times$  to  $4.8\times$ . Additionally, RANDLUZZ can expose 8 bugs within 60 seconds, demonstrating leading efficiency in bug discovery.

Our contributions are as follows.

- We analyze and demonstrate that randomness is the key

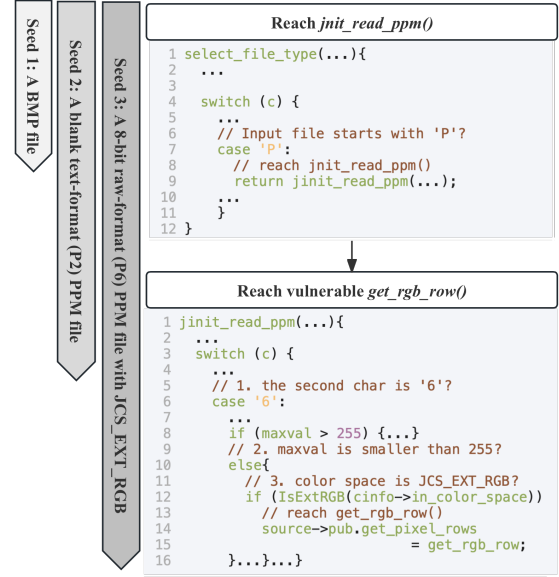


Figure 1. Randomness in initial seeds for cjpeg. Different seeds explore different input space, resulting in differences in time to bugs.

factor that impacts the efficiency of directed fuzzing.

- We introduce LLMs to remove randomness in generating reachable seeds and reduce randomness in generating bug-specific mutators.
- We develop RANDLUZZ<sup>1</sup> incorporating LLMs to demonstrate the effectiveness and efficiency of our idea, showing the feasibility of using LLMs in guiding directed fuzzing.
- Our evaluation demonstrates the effectiveness and efficiency of reachable seeds and bug-specific mutators.

## 2. Motivation

We identify two key components, which are seeds and mutators, that introduce randomness, reducing efficiency in directed fuzzing. In this section, we analyze such randomness based on the bug CVE-2020-13790.

### 2.1. Randomness in Seeds

Figure 1 shows that three seeds with different file formats can reach different code regions in the program cjpeg. To reach the target vulnerable function `get_rgb_row()`, a seed has to satisfy multiple path constraints. In Figure 1, the *Seed 1* is a blank BMP image that can only examine the function `select_file_type()`. To reach the function `jinit_read_ppm()`, the *Seed 2* satisfies a different file format PPM, whose input file starts with the characters `P2`. When trying to reach the vulnerable function `get_rgb_row()`, a more complex file format has to be satisfied. Thus, the carefully crafted *Seed 3*, which can

1. Our code will be publicly available upon acceptance.

reach `get_rgb_row()`, is an 8-bit raw-format (P6) PPM file with the `JCS_EXT_RGB` color space.

The randomness in the seeds lies in the fact that they may have different search space to explore. For example, the *Seed 1* has to explore the input space to satisfy all path constraints including the bytes specifying the file format, file size, and color space. However, for *Seed 3*, since it already satisfies all the path constraints, it can focus on exploring the buggy execution states. To showcase the observation, we test the three seeds on the bug CVE-2020-13790 with the directed fuzzer AFLGo. AFLGo runs each seed four times and their timeouts are set to 24 hours. The results show that *Seed 1* cannot expose the bug due to its large search space and randomness. *Seed 2* can only discover the bug once, with the testing time close to 24 hours (84,509 seconds). However, *Seed 3* can identify the bug twice, and it exposes the bug more quickly, with the testing time close to 8 hours (30,211 seconds) and 5 hours (17,828 seconds).

The challenge is that, effective seeds require expertise to craft, reflecting various levels of understanding about the target programs and bugs. For instance, a tester who knows only that the target program is `cjpeg`, without any specific knowledge of the vulnerability, might use an image in any format (such as BMP) for testing. In essence, the tester’s level of understanding about the target programs and bugs significantly impacts the effectiveness of directed fuzzing. Therefore, in this paper, we use LLMs to automatically understand target programs and bugs, which further helps craft reachable seeds for fuzzing.

## 2.2. Randomness in Mutators

Randomness in mutators inherits from the concept of fuzzing, which *randomly* generates numerous inputs to repeatedly test target programs. This immediately leads to the dilemma that fuzzing needs certain degree of randomness to maintain its effectiveness while randomness is expected to be reduced for efficiency consideration. Existing mutators satisfy the requirement of randomness because they follow the idea of building meaningful mutations based on small and basic operations. For example, most fuzzers use the mutators from AFL, which designs basic operations such as *bit flips* and *byte insertion*. The combination of such operations can produce meaningful mutations but comes with the drawback of high randomness. As shown in Figure 1, *Seed 3* can already reach target locations, but the time to bug still experiences discrepancies among different trials.

To trigger a bug, the mutation is required to generate inputs that examine specific locations and exercise certain execution states. Most existing papers intend to efficiently apply mutators in exploring diverse code locations [2], [39], [42]. However, even a seed can reach buggy code regions, it still requires the mutators to explore execution states. We observe that, *efficient mutators for exploring code regions and execution states are likely to be different*. For example, in Figure 1, to reach target function, the mutators have to satisfy path constraints such as inputs starting with characters *P6*. Yet, to satisfy the buggy execution states,

the mutators need to manipulate memories to overflow the boundary of a buffer. Therefore, in this paper, we design two solutions to reach target locations and to explore buggy states, respectively. To reduce the randomness in mutators, our design focuses on constructing bug-specific mutators.

## 3. Methodology

Our RANDLUZZ aims to minimize the impact of randomness in fuzzing by leveraging LLMs. To effectively use LLMs in directed fuzzing, our RANDLUZZ consists of multiple phases in using LLMs, as shown in Figure 2. The first phase is to use LLMs in preparing necessary information, which is utilized in later reasoning stages. Then, RANDLUZZ generates reachable seeds via reasoning based on the necessary information, which removes the randomness in reaching target locations. Finally, RANDLUZZ constructs bug-specific mutators to explore execution states in target vulnerable locations, which reduces randomness in generating bug-triggering inputs.

### 3.1. LLM Query Scheme

Throughout the workflow, RANDLUZZ frequently communicates with LLMs. We devise an LLM query scheme, illustrated in Figure 3, which primarily consists of four parts:

- **Task** describes the main objective of the query.
- **Attachment** includes relevant context for the query, *e.g.*, function body when querying function summary.
- **Suggestion** provides LLMs with guidance relevant to the specific tasks and required answers.
- **Answer Template** defines the required content and format of the LLM’s response.

When faced with different tasks, RANDLUZZ uses this query scheme to create task-specific templates, with the corresponding context. We design various query templates under the query scheme to handle a range of tasks. An example query template is shown in Figure 3. When RANDLUZZ needs to summarize a function in the program, it provides the scheme with the *Task* “*Summarize the function*”, the relevant context, as well as placeholders for detailed query. Notably, to reduce query time and the number of tokens used, communication between RANDLUZZ and LLMs excludes any historical conversation records. If a task requires results from a previous communication, RANDLUZZ attaches these results in the *Attachment* and reminds the LLMs of the presence in the *Suggestion*.

### 3.2. LLM-Assisted Necessary Information

As shown in Figure 2, the necessary information includes *Bug Information*, *Program Usage*, and *Function Summary*. To perform directed fuzzing, the essential information includes the target code locations. Given a bug report, RANDLUZZ automatically extracts bug related details (*Bug Information*), such as target locations and other information useful for generating inputs. Next, the command options

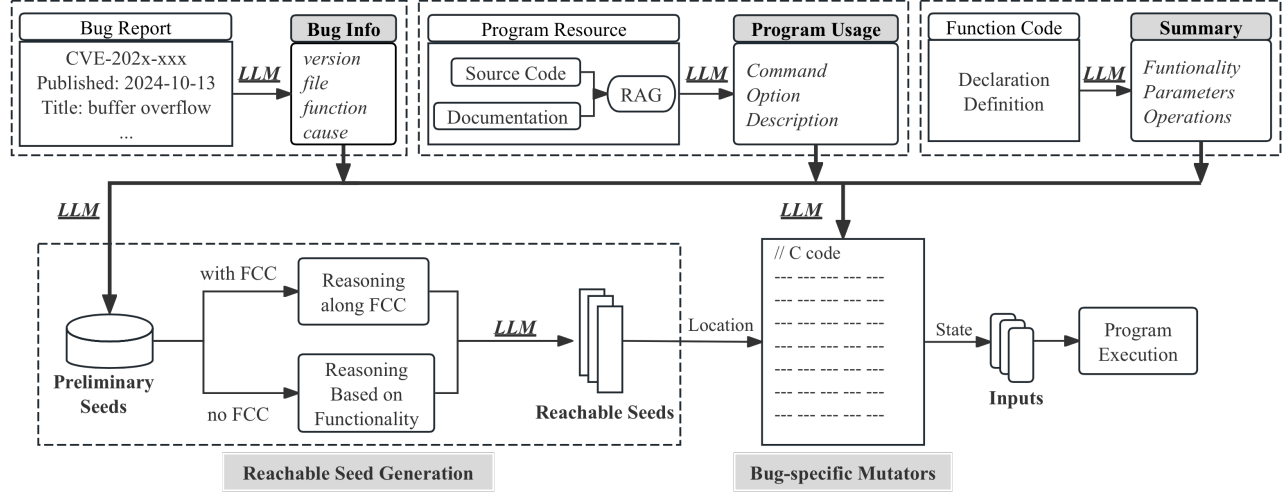


Figure 2. Workflow of RANDLUZZ. *Reachable* indicates that the execution can reach the target vulnerable location. The use of LLM is to generate reachable seeds and generate specific mutators for target bugs. FCC is short for Function Call Chain.

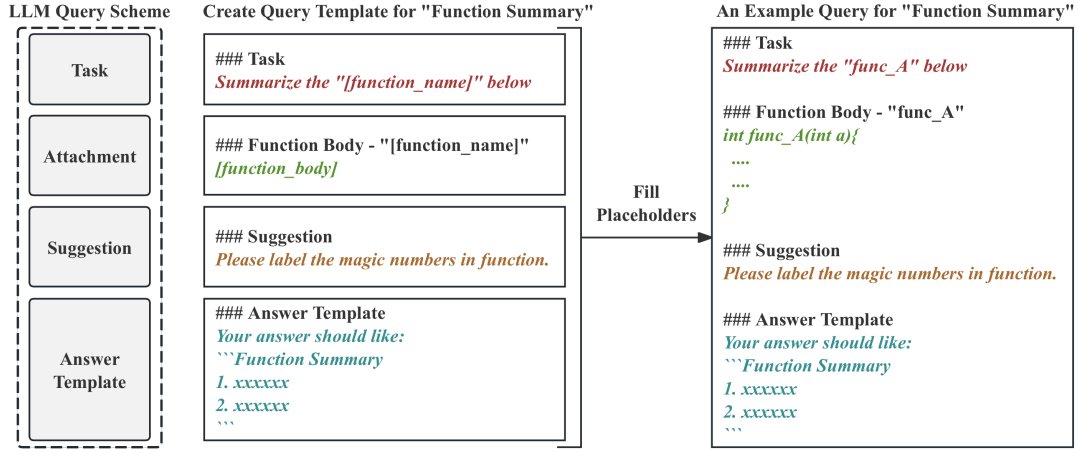


Figure 3. LLM query scheme. Each task is based on this query scheme to create a query template, which further generates queries.

needed to reach the buggy locations in the program should be identified. *Program Usage* is used to obtain the command options that can reach target vulnerable locations. Finally, summaries of each function (*Function Summary*) are necessary to understand the program, which can further assist in generating reachable seeds. *Function Summary* is performed when necessary and provides essential information, including the functionality, parameters and operations.

**3.2.1. Bug Information.** Bug reports may have various formats to present bug related information. For example, a Common Vulnerabilities and Exposures (CVE)<sup>2</sup> report is a standardized document that provides detailed information about a specific security vulnerability or exposure, including its description, the affected software versions, and the potential impact. Generally, the aim is to extract necessary information from bug reports, including the affected soft-

ware versions, the location of the vulnerability (including the specific file and function), and a brief summary of the cause of the vulnerability, as shown in Figure 2. We leverage LLMs to extract such information because LLMs excel at summarizing and synthesizing natural language contents, enabling an approach to efficiently analyze CVE reports.

The query template for analyzing CVE reports is loaded to generate the appropriate query for communicating with the LLM. The *Task* directly prompts the LLM to extract the required bug information from the CVE report, which is included as an *Attachment*. The extracted bug information must be provided in a specified format, as outlined in the *Answer Template*. The bug information helps pinpoint the location of bugs, serving as the target for subsequent optimization efforts. Additionally, details on the causes of bugs, combined with definitions of the functions where they reside, allow the LLM to analyze and identify the conditions under which the bugs can be triggered. This insight aids in developing mutators tailored to target bugs.

2. <https://cve.mitre.org/>

**3.2.2. Program Usage.** The program usage refers to the command lines required to run the target program, particularly the command options. Command options are often tailored for specific functionalities, indicating that different options exercise different code regions [29]. For example, the program `objcopy` specifies the input format with the option `-I` and the output format with the option `-O`. Some bugs, in fact, can only be exposed with specific options. Extracting program usage by inputting the entire program as the prompt to LLMs is typically not feasible, as it exceeds the token limit. Even if an LLM could handle an entire project, its output would likely be imprecise due to inherent challenges relevant to focus and coherence, memory and retention issues and increased computational cost associated with the complexity of programs. Therefore, we employ a Retrieval-Augmented Generation (RAG)-like strategy [19] to help LLMs gain in-depth understanding of a program, which in turn supports the generation of program usage.

As illustrated in Figure 4, RAG enhances LLMs by integrating with an external knowledge retrieval mechanism that includes a retriever, which sources pertinent information, and a generator, which synthesizes this information into coherent output. In our study, the process begins with bug information as a query, triggering the retriever to search a *knowledge base* - the program project and its documentation repository - for relevant cues. Retrieval is based on the embedding similarity between the bug information and chunks of program code and documentation, divided into smaller parts. Specifically, we use a sentence transformer - all-MiniLM-L6-v2 [32] to map each chunk to a 384-dimensional dense vector embedding space, with Faiss employed for vector similarity search [7]. Empirically, we select the top 10 relevant chunks as the supplemental information. The retrieved supplemental information, together with the bug information, is then passed to the generator module, which in our case is the LLM. We query the LLM by setting a specific *Task*, such as “*Summarize the usage of all command options for this program*”, enabling it to generate a comprehensive response about program usage. Additionally, the query’s *Answer Template* requires the LLM’s response to include the command options mentioned in the chunk, along with their corresponding usage descriptions. In later stages, the program usage information will guide the LLM in selecting appropriate command options, helping ensure that the target functions are reached.

**3.2.3. Function Summary.** The function summary report is essential for RANDLUZZ, which helps infer relationships between functions and program command options. Additionally, key operations and variables within the function body can be leveraged for further seed optimization and vulnerability cause analysis. RANDLUZZ uses Clang’s AST (Abstract Syntax Tree) component to analyze the project’s source files. The main objective of this static analysis is to capture and record all functions within the project and map their call relationships, which are identified via simple control flow analysis, such as capturing keywords `CALL_EXPR` or `POINTER` pointing to a `FUNCTIONPROTO` in the AST

output. Additionally, if RANDLUZZ requires more detailed information about a function, such as its specific functionality or variable values, it can query the LLM for an in-depth analysis. With the function’s declaration and definition extracted from static analysis included as the context in *Attachment*, RANDLUZZ queries the LLM by setting a *Task* like “*Summarize the function’s purpose*”. The LLM’s output will provide summaries of the function’s functionality, parameters, and key operations as specified in the *Answer Template*, as illustrated in Figure 5.

### 3.3. Reachable Seed Generation

Seeds can impact the efficiency of directed fuzzing, motivating us to generate reachable seeds that can reach target locations. For this purpose, it is essential to have an understand of the target program. Reachable seeds must satisfy program logic to examine target code regions. For example, when testing a specific feature of image format conversion in an image processing program, an image file is more suitable as a seed than a text file. Our RANDLUZZ first creates preliminary seeds that at least meet the format requirements of the target program. Then, if a complete Function Call Chain (FCC), which contains execution paths from the entry point to the target vulnerable function, can be obtained through static analysis, RANDLUZZ employs LLMs to reason about reachable seeds along the FCC. If not, RANDLUZZ uses LLMs to infer reachable seeds based on the functionalities of neighbor functions.

**3.3.1. Preliminary Seed Generation.** Our goal is to continuously optimize the seeds until their execution paths reach the target location. Specifically, we aim to create preliminary seeds that exercise execution paths as close as possible to the location. Importantly, when RANDLUZZ generates preliminary seeds, it also generates the associated command options. Command options significantly impact the effectiveness of directed fuzzing; if an option is unrelated to the target location, no input can successfully examine target buggy code regions. Therefore, we begin by comparing the functionality of the target location with the *Program Usage* information to identify which command options should be used to accept program input with the LLM. In our query to the LLM, we set the *Task* to analyze which program command is most likely to activate the target function, using the *Answer Template* to obtain the appropriate command and descriptions. In addition, we provide the *Function Summary* of the vulnerable function and the *Program Usage* information as the *Attachment* for LLM context. Finally, the LLM generates the corresponding preliminary seed. For the later fuzzing process, we will fixate the program command but mutate the seed.

For example, when generating preliminary seeds for the bug CVE-2017-16828 in the program `readelf`, we first obtain the *Program Usage* through the LLM. Then, we analyze the vulnerable function `display_debug_frames()` and find that it processes and displays the contents of DWARF debugging

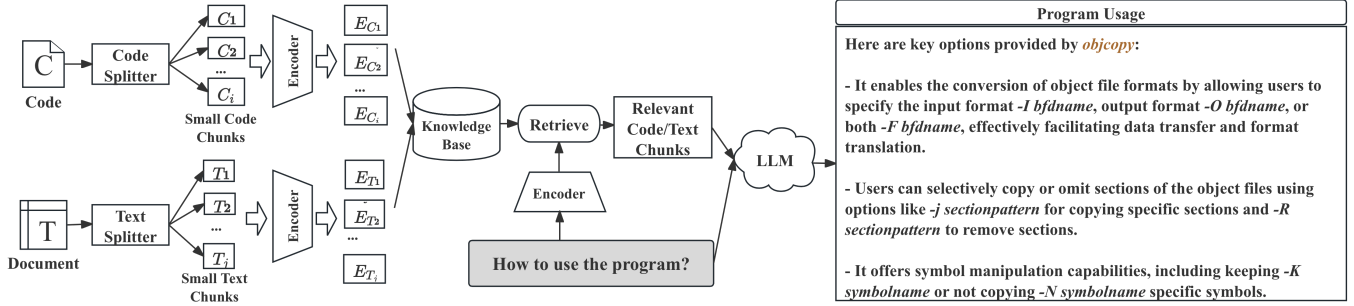


Figure 4. Program usage and its options. RANDLUZZ uses RAG to retrieve relevant and small chunks from documents and source code, and the small chunks are used as the context when querying LLMs.

Function Summary
<p><b># Functionality Summaries:</b>  This function, <i>copy_object</i>, appears to deal with copying the contents or data (such as sections, symbols, architecture information, etc.) from one <i>bfd</i> object <i>ibfd</i> to another <i>obfd</i>, ...</p> <p><b># Parameter Summaries:</b>  <i>-ibfd,obfd</i>: Input and Output <i>bfd</i> pointers, respectively, represent the source and destination objects for the copy operation.  ....</p> <p><b># Key Operations:</b>  1. Endianness Check and Format Copying: Checks if the byte order (endianness) and format of <i>ibfd</i> and <i>obfd</i> match. If there's a mismatch, it returns an error.  2. ...</p>

Figure 5. An example of Function Summary. It includes the summaries of functionality, parameters, and key operations.

frame sections, specifically the `.debug_frame` and `.eh_frame`. After providing the *Function Summary* and *Program Usage* to the LLM for analysis, it identifies the command most likely to reach the vulnerable function as `readelf --debug-dump=frames file.elf`. The command description indicates that `file.elf` should be a valid ELF binary compiled with debugging symbols (e.g., using the `-g` flag in GCC). We then ask the LLM to generate a preliminary seed, which is an ELF file matching this description, to replace the placeholder `file.elf`.

**3.3.2. Reasoning Along Function Call Chain.** Relying solely on the summary of the target vulnerable function for seed generation is often insufficient to reach the target. This limitation arises because the execution path from the program’s entry point to the target function may pass through multiple intermediate functions that are not considered in seed generation. Consequently, the requirements imposed by these intermediate functions on the seed are overlooked, which can cause the generated seed to deviate along incorrect paths. To address this, RANDLUZZ leverages LLMs to adjust the preliminary seed along an FCC, if a complete FCC that connects entry point and target function is available.

Overall, for this purpose, RANDLUZZ takes two steps. First, it uses Clang AST analyzer to get the complete FCC. Next, it gradually generates a reachable seed by analyzing

deviation functions, which deviate the execution from the vulnerable FCC. For example, as shown in Figure 6, to generate a seed that can reach the target vulnerable function *T*, while the preliminary seed exercises the execution path  $E \rightarrow A \rightarrow B$ , which deviates from the target *T*. Since functions *E* and *A* are both present in the complete FCC, RANDLUZZ analyzes which one is the deviation function. By calculating the distances between functions *E* and *T* and between *A* and *T*, RANDLUZZ identifies function *A* as the deviation function because it is closer to the target function *T*. The distance is defined as the number of edges in the shortest path from the source function to destination function. As a result, RANDLUZZ obtains the definition of *A*, the lines examined by the preliminary seed, as well as the next function after the deviation function *A* along the complete FCC, which, in this case, is function *C*. With this set of information, RANDLUZZ can query LLMs to generate a testcase to reach the next function *C* (the goal). Specifically, we set the *Task* in the query as “Based on the provided function body and the current input that leads to a specific execution path, how should the input be modified to guide the program to reach the target function?”. After querying the LLM and receiving a response, we adjust the input accordingly and execute it, checking the execution path of the modified input. If the modified input reaches function *C*, RANDLUZZ uses the corresponding testcase as the current seed, and repeat the above process to generate a new testcase that can reach the next goal. This continues until RANDLUZZ successfully generates a seed that can reach the target vulnerable function. If not, we continue querying the LLM for further seed optimization guidance.

Notably, directly asking LLMs to generate an input meeting all code branches and reaching the optimization target is impractical. This limitation arises because our prompt includes only a single function definition (the function *A* in the example), whereas many branch conditions within this function depend on other functions or macro-defined parameters. To address this, we inform LLMs of this limitation within our prompt, asking them to infer the potential roles of functions or parameters based on their names. As a result, multiple seeds can be generated regarding these potential roles, aiming to satisfy path constraints. For example, when optimizing the seed for the program `cjpeg`, a



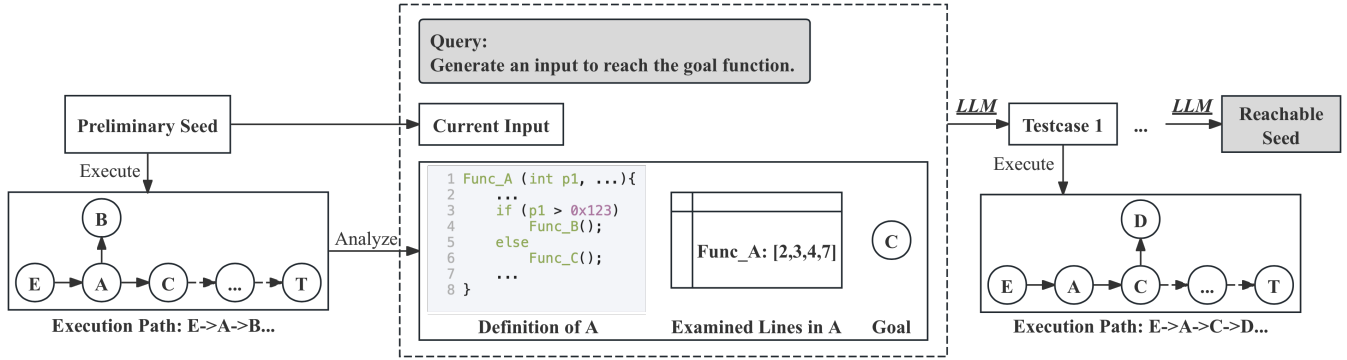


Figure 6. Reasoning along FCC. The aim is to generate a seed that can reach the target vulnerable function *T*.

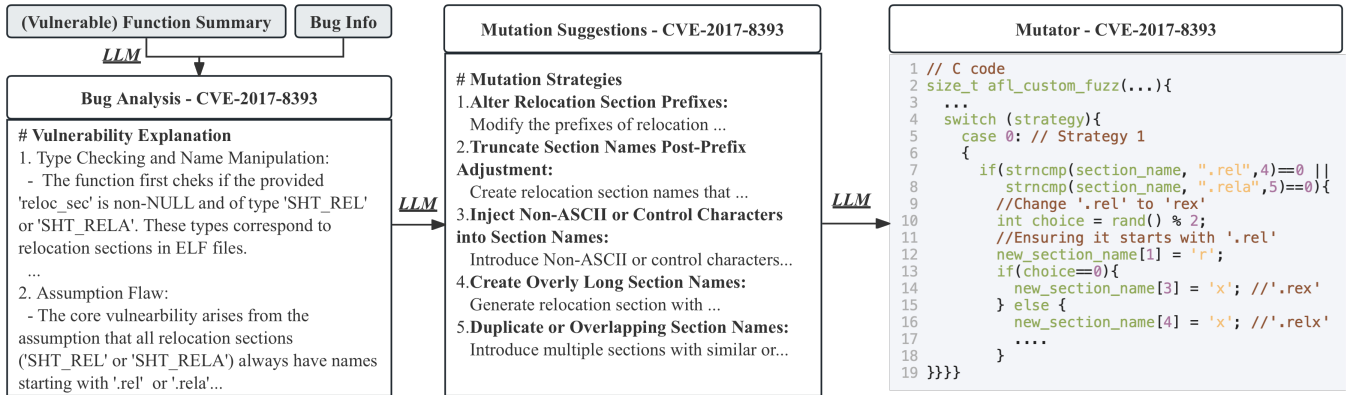


Figure 7. Bug-specific Mutator. Based on LLM, RANDLUZZ generates code to perform mutations.

path constraint compares a variable named `maxval` with a macro-defined value `MAXJSAMPLE`. Since only the relevant function definition is provided in the LLM query, without the macro definition, the exact value of `MAXJSAMPLE` is unknown. However, the LLM can infer that `MAXJSAMPLE` likely represents the maximum color component value in an image. Based on this, it assumes `MAXJSAMPLE` could be 255 (for 8-bit images) or 4095 (for 12-bit images) and generates two different seeds to test these possibilities. This seed generation process continues iteratively until a generated seed either reaches the optimization target or progresses to a function closer to the target.

**3.3.3. Reasoning Based on Functionality.** Static analysis may fail to produce a complete FCC due to program complexities such as nested macro definitions and function calls using function pointers and memory addresses. Many directed fuzzers face similar challenges, leading to incomplete calculations of distances between current execution blocks and target code blocks [1], [4]. To address this issue, we leverage the summarization capabilities of LLMs to generate reachable seeds. By analyzing the neighboring functions - those that can reach the target function - RANDLUZZ gains insights into the target function's role within a broader program context. Our insight is that, if a complete FCC cannot be obtained from the entry point, an incomplete

FCC may still be derived by starting from the neighboring functions of the vulnerable one. LLMs can then be used to generate inputs that reach these neighboring functions based on their functionalities.

RANDLUZZ first uses static analysis tools (e.g., Clang) to identify all neighboring functions of the target function. It then randomly selects a neighboring function and queries LLMs about the program options and inputs that can lead to this function, with the context of *Program Usage* and *Function Summary*. By analyzing the instrumentation execution logs, we can determine whether these inputs reach the neighboring function. If they do not, we randomly select another neighboring function and repeat the input generation process based on the functionality of the newly selected function. Finally, if an input successfully reaches a neighboring function, we apply the method outlined in Section 3.3.2 to generate reachable seeds.

**3.3.4. Issues for Input Generation Process.** The first issue is the input format that impacts how we query LLMs. Input formats can range from simple strings to more complex data types, such as images or ELF files. For string-type inputs, we query LLMs to directly generate the string. For complex inputs, we query LLMs to generate code that can produce these inputs, as they cannot be directly generated by LLMs. For example, we may query as *Please provide*

TABLE 1. BENCHMARK.

Program	CVE ID	Vulnerability Type	Vulnerable Location	Indirect Calls <sup>1</sup>	Edit <sup>2</sup>
cjpeg	CVE-2018-14498	Heap-Buffer Overflow	rdtomp.c:get_8bit_row	T	N
	CVE-2020-13790	Heap Buffer Overflow	rdppm.c:get_rgb_row	T	N
cxxfilt	CVE-2016-4487	NULL Pointer Dereference	cplus-dem.c:register_Btype	F	Y
	CVE-2016-4489	Integer Overflow	cplus-dem.c:string_appendn	F	N
	CVE-2016-4491	Stack Overflow	cp-demangle.c:d_print_comp_inner	T	N
objcopy	CVE-2017-8393	Global Buffer Overflow	elf.c:_bfd_elf_get_reloc_section	T	N
	CVE-2017-8394	NULL Pointer Dereference	objcopy.c:filter_symbols	F	Y
	CVE-2017-8395	NULL Pointer Dereference	compress.c:_bfd_generic_get_section_contents	T	Y
objdump	CVE-2017-8392	Heap Buffer Overflow	dwarf2.c:parse_comp_unit	T	N
	CVE-2017-8397	Heap Buffer Overflow	reloc.c:bfd_perform_relocation	T	Y
	CVE-2018-17360	Heap Buffer Overflow	peigen.c:pe_print_edata	T	N
strip	CVE-2017-7303	NULL Pointer Dereference	elf.c:section_match	T	N
nm	CVE-2017-14940	NULL Pointer Dereference	dwarf2.c:scan_unit_for_symbols	T	N
readelf	CVE-2017-16828	Heap Buffer Overflow	readelf.c:display_debug_frames	T	N

<sup>1</sup> **Indirect Calls:** Indirect calls lead to incomplete function call chain from the entry function to the target vulnerable function. If there is an indirect call on the function call chain, it will be marked as *T*; otherwise it is *F*.

<sup>2</sup> **Edit:** Some CVE reports do not directly state the vulnerable function names or provide incorrect ones. Therefore, we manually edit reports to add or modify the vulnerable function names in such CVE reports, and mark it as *Y*; otherwise, we mark it *N*.

a Python script that, when executed, outputs a file that meets the requirements. Another issue is that the generated code may not be runnable. To ensure the reliability of codes generated by LLMs, we feed any error messages from code execution back into LLMs, allowing them to correct the code. This process is repeated until the code executes without errors. The final issue is that, during the process of guiding seeds toward the target vulnerable function, we may fail to generate a seed that reaches the target, stalling the testing process. To address this, our prototype imposes a time limit for this phase. If no seed reaches the target within a specific time duration, such as one hour, we select a subset of the generated seeds to serve as seeds for fuzzing.

### 3.4. Bug-Specific Mutators

As shown in Figure 2, merely reaching the target locations does not necessarily trigger bugs, as specific execution states must also be met. To reduce randomness in mutating reachable seeds, RANDLUZZ employs LLMs to generate bug-specific mutators, tailoring mutation schemes to target particular bugs. As shown in Figure 7, we begin by using the *Bug Info* extracted from the bug report, along with the *Function Summary* of the vulnerable function, to query the LLM. The *Task* of this query prompts the LLM to provide a *Bug Analysis* report that summarizes the bug’s cause and how to trigger the vulnerability. Next, we ask the LLM to generate a fuzzing mutation strategy based on the *Bug Analysis* report. In this query, we include examples of mutation strategies and their explanations in the *Attachment* and *Suggestion* sections, respectively, to guide the LLM in generating effective new mutation strategies. Finally, these new mutation strategies are sent to the LLM with a *Task* to translate them into C language code. To ensure code quality, real C language mutator code examples and explanations are included in the query context. If the generated code fails to compile, we ask the LLM to revise the code based on the

original mutator code and compilation error messages, and attempt compilation with the updated code.

Additionally, some hidden or deeper code errors may not appear during compilation but could cause the fuzzer using the mutator to crash or experience significantly reduced execution efficiency. To ensure the generated mutator code runs correctly, we conduct a trial run before executing the full fuzzing process, checking for any issues in the mutator code. In this trial, we verify that the fuzzing program runs smoothly and that specific fuzzing metrics (such as executions per second) fall within expected ranges. If crashes or abnormal metrics occur, we request the LLM to regenerate a replacement mutator. Meanwhile, in the directed fuzzing phase, the mutation strategies generated by the LLM are often deterministic. Repeatedly using the same set of strategies can result in wasted time on redundant test cases. To ensure a diverse and effective testing, we request the LLM to generate a new set of mutation strategies every hour.

## 4. Evaluation

In this section, we address the following research questions through various experiments:

- **RQ1: Preparation Time for Fuzzing** - Does RANDLUZZ require more time for fuzzing preparation?
- **RQ2: Efficiency of Reachable Seeds** - Do the reachable seeds generated by RANDLUZZ result in more efficient bug detection than other seeds?
- **RQ3: Efficiency of Bug Detection** - Does RANDLUZZ perform better than other directed fuzzers in detecting various vulnerabilities?
- **RQ4: LLM Model Influence** - Do different GPT models affect RANDLUZZ?

**Experiment Environment.** Our experiments were conducted on a machine equipped with two Intel(R) Xeon(R) Gold 6138 CPUs at 2.00 GHz, providing a total of 40 CPU cores and 80 threads. Running Ubuntu 20.04 LTS, each



fuzzing session was executed within a Docker container with one CPU thread and 2GB of RAM allocated per session. We utilized all 80 CPU threads, each running an instance of the same fuzzing session.

**Tool Implementation.** We implemented RANDLUZZ with approximately 2400 lines of Python code. Our tool is built upon AFL++ [10], which supports user-defined mutators, enabling us to integrate LLM-generated mutators alongside AFL++’s random mutators. This allows RANDLUZZ to leverage both bug-specific and random mutation strategies. We use GPT-4o as the LLM and interact with it via OpenAI’s APIs (2024-08-06 API). The static analysis component in our methodology employs Clang’s Abstract Syntax Tree (AST) analysis tool<sup>3</sup>.

**Baseline Fuzzers.** We select four state-of-the-art directed fuzzers as baselines to evaluate RANDLUZZ.

- AFLGo [1] is the first and commonly-used fuzzer to be compared. It uses distance to guide fuzzing towards target locations. It uses mutators from AFL, whose mutators are pre-determined and not bug-specific.
- Beacon [12] only exercises execution paths that can reach target locations so that it focuses on exploring relevant code. It is developed upon AFLGo, indicating that it uses mutators from AFL.
- WindRanger [8] guides executions towards target locations based on deviation basic blocks, which deviate the execution paths from target locations. It is developed based on AFL, indicating that it uses mutators from AFL.
- SelectFuzz [24] instruments code regions that are relevant to target locations so that it can focus on exploring the relevant code regions. It is also developed based on AFL, indicating it uses mutators from AFL.

**Benchmarks.** Table 1 lists all the vulnerabilities involved in our experiments, as well as the details of them. We select 14 vulnerabilities from 8 different programs, including `cjpeg`, `cxxfilt`, `objcopy`, `objdump`, `strip`, `nm`, and `readelf`. The 14 vulnerabilities have diverse types of bugs including heap-buffer overflow, NULL pointer dereference, integer overflow, stack overflow, and global buffer overflow. These vulnerabilities are used to evaluate fuzzing in the papers of baseline fuzzers. To ensure a fair comparison, we adopt the same compilation options and commands in the evaluation. In our evaluation, all experimental tests are repeated 4 times. We present the median value in the tables.

#### 4.1. RQ1: Preparation Time for Fuzzing

RANDLUZZ has four steps that prepare for later fuzzing process, including static analysis for analyzing a program’s function call chain (SA), program option construction using RAG (RAG), seed optimization to generate reachable seeds (*Opt*), and generation of bug-specific mutators (*Mutator*). Thus, in this evaluation, we compare the preparation time across fuzzers to check whether RANDLUZZ significantly increases the preparation time for directed fuzzing. As shown in Table 2, on average, the *Opt* process takes the most

time, consuming 53% more time than the second longest process, *RAG*. This is because *Opt* process involves repeated interactions between querying the LLMs to generate inputs and executing the target program to obtain feedback. The processing time for *Opt* varies significantly across different vulnerabilities, ranging from 72 seconds to 957 seconds. This variation occurs because target locations differ in depth, reflecting the complexity of the control flow and data flow needed to reach the target. Excessively deep targets may result in timeouts, preventing RANDLUZZ from generating a reachable seed within the allotted time. *RAG* takes the second longest time due to the operations required for slicing, encoding, and retrieving text and code, which are positively correlated with the size of the target program. The process *Mutator* for generating fuzzing mutation strategies in RANDLUZZ is relatively consistent, taking between 83 and 147 seconds. This is primarily because this process uses vulnerability reports and function summaries that are gathered in previous stages as input. Finally, the processing time for SA depends on the size of target programs, and it takes the least time, with 47 seconds on average.

Overall, compared to other fuzzers, RANDLUZZ requires less time across the four steps than AFLGo, Beacon, and SelectFuzz, but more time than WindRanger. On average, RANDLUZZ spends 62.6%, 86.2%, and 60.1% less time than AFLGo, Beacon, and SelectFuzz (with TS), respectively. AFLGo spends considerable time on calculating distances between basic blocks. In addition to distance calculation, Beacon requires extra time for reachability analysis. SelectFuzz also has a lengthy preparation when the Temporal-Specialization (TS) mode is enabled, as this mode analyzes indirect calls to obtain a potentially complete control flow graph. This shows the advancement of our RANDLUZZ, which leverages LLMs to analyze functionalities instead of analyzing indirect calls. WindRange spends the least time on preparation because its distance calculation is performed during the fuzzing process. Therefore, RANDLUZZ demonstrates competitive performance in fuzzing preparation.

#### 4.2. RQ2: Efficiency of Reachable Seeds

In this section, we evaluate the efficiency of RANDLUZZ-generated reachable seeds, comparing with other sets of seeds. Since seed generation and mutators are tightly coupled in RANDLUZZ, we do not evaluate RANDLUZZ on different sets of seeds. The four directed fuzzers are evaluated on three sets of initial seeds, where each fuzzer runs 24 hours on each vulnerability. As shown in Table 3, the three sets of seeds include Original Seeds (OS), Simple Seeds (SS), and RANDLUZZ-generated Seeds (RLS). Original seeds are the set of seeds collected by Kim *et al.* [16]. Simple seeds are the set of naive seeds that only meet the minimum input requirements and do not aim to trigger any specific functionality. For example, in the case of the program `cxxfilt`, which takes a character-type function symbol name as input, we might provide a simple test seed like `_Z3foov`. RANDLUZZ-generated seeds are the set of seeds generated by RANDLUZZ.

3. <https://clang.llvm.org/docs/IntroductionToTheClangAST.html>

TABLE 2. PREPARATION TIME OF EACH FUZZER.

CVE ID	AFLGo	Beacon	WindRanger	SelectFuzz <sup>1</sup>		RANDLUZZ <sup>2</sup>				
				no TS	with TS	SA	RAG	Opt	Mutator	Total
CVE-2018-14498	256s	19s	37s	47s	45s	24s	95s	693s	145s	957s
CVE-2020-13790	260s	19s	37s	47s	49s	23s	93s	423s	121s	660s
CVE-2016-4487	304s	3095s	77s	152s	2063s	48s	193s	177s	126s	544s
CVE-2016-4489	307s	5068s	77s	149s	1962s	51s	152s	181s	83s	467s
CVE-2016-4491	309s	6072s	76s	149s	1870s	45s	178s	100s	147s	470s
CVE-2017-8393	1858s	4852s	68s	187s	2866s	48s	145s	175s	112s	480s
CVE-2017-8394	1846s	5971s	70s	187s	2543s	50s	162s	157s	93s	462s
CVE-2017-8395	1853s	4966s	68s	183s	2555s	54s	191s	150s	87s	482s
CVE-2017-8392	1991s	8046s	88s	171s	2464s	49s	279s	417s	89s	834s
CVE-2017-8397	2016s	9910s	89s	166s	2355s	51s	287s	213s	112s	663
CVE-2018-17360	3645s	5071s	91s	980s	981s	47s	285s	T.O. <sup>3</sup>	101s	T.O.
CVE-2017-7303	3943s	4652s	68s	1092s	1075s	66s	168s	957s	87s	1278s
CVE-2017-14940	3848s	4795s	67s	1131s	1127s	67s	175s	150s	91s	483s
CVE-2017-16828	623s	141s	61s	113s	123s	40s	162s	72s	89s	363s
Average	1647s	4477s	81s	339s	1577s	47s	183s	280s	106s	616s

<sup>1</sup> SelectFuzz: TS - Temporal-Specialization<sup>2</sup> RANDLUZZ: SA - static analysis for FCC; Opt - seed optimization for reachable seeds; Mutator - mutator generation for bug-specific mutators. <sup>3</sup> T.O.: Timeout (3600 seconds; it is not counted when calculating the average time).

TABLE 3. TIME TO BUGS WITH DIFFERENT SETS OF INITIAL SEEDS. SEEDS GENERATED BY RANDLUZZ CAN TRIGGER BUGS THE MOST QUICKLY.

CVE ID	AFLGo			Beacon			WindRanger			SelectFuzz		
	OS <sup>1</sup>	SS <sup>2</sup>	RLS <sup>3</sup>	OS	SS	RLS	OS	SS	RLS	OS	SS	RLS
CVE-2018-14498	59161s	37525s	<b>34s</b>	T.O. <sup>4</sup>	T.O.	T.O.	T.O.	6353s	<b>11s</b>	T.O.	12921s	<b>40s</b>
CVE-2020-13790	81026s	84509s	<b>17828s</b>	T.O.	T.O.	T.O.	T.O.	T.O.	T.O.	T.O.	T.O.	T.O.
CVE-2016-4487	1272s	1488s	<b>22s</b>	276s	180s	<b>3s</b>	308s	1497s	<b>23s</b>	423s	316s	<b>94s</b>
CVE-2016-4489	2717s	2420s	<b>1682s</b>	396s	384s	<b>272s</b>	641s	470s	<b>398s</b>	3465s	<b>916s</b>	1600s
CVE-2016-4491	7514s	6466s	<b>9s</b>	2817s	2194s	<b>5s</b>	12039s	12874s	<b>376s</b>	13323s	<b>12734s</b>	13114s
CVE-2017-8393	2875s	1379s	<b>73s</b>	<b>66901s</b>	T.O.	T.O.	2641s	<b>354s</b>	435s	3184s	2100s	<b>1487s</b>
CVE-2017-8394	2077s	12046s	<b>51s</b>	87s	53258s	<b>75s</b>	2201s	78775s	<b>170s</b>	<b>480s</b>	T.O.	3847s
CVE-2017-8395	230s	186s	<b>14s</b>	26s	40s	<b>4s</b>	163s	282s	<b>16s</b>	309s	T.O.	<b>17s</b>
CVE-2017-8392	574s	225s	<b>27s</b>	79490s	<b>5428s</b>	33207s	380s	97s	<b>18s</b>	313s	197s	<b>11s</b>
CVE-2017-8397	14656s	855s	<b>764s</b>	20739s	<b>80s</b>	146s	27054s	367s	<b>83s</b>	T.O.	<b>508s</b>	T.O.
CVE-2018-17360	<b>55014s</b>	T.O.	T.O.	<b>6611s</b>	T.O.	T.O.	<b>9734s</b>	T.O.	T.O.	<b>11168s</b>	T.O.	T.O.
CVE-2017-7303	<b>365s</b>	6408s	13320s	<b>34s</b>	842s	2203s	5495s	4607s	<b>1614s</b>	<b>253s</b>	6116s	17273s
CVE-2017-14940	T.O.	T.O.	T.O.	T.O.	T.O.	T.O.	T.O.	T.O.	T.O.	T.O.	T.O.	T.O.
CVE-2017-16828	<b>547s</b>	694s	2419s	25055s	23221s	<b>2136s</b>	<b>193s</b>	4517s	20978s	T.O.	T.O.	<b>55485s</b>
# of Fastest <sup>5</sup>	3	0	10	3	2	6	2	1	9	3	3	6

<sup>1</sup> OS: Original Seeds; <sup>2</sup> SS: Simple Seeds; <sup>3</sup> RLS: Seeds Provided by RANDLUZZ. <sup>4</sup> T.O.: Timeout, indicating that the target vulnerability was not triggered within the specified time (24 hours = 86400 seconds). <sup>5</sup> The number of vulnerabilities the fuzzer triggers most quickly.

Table 3 shows the time to bugs for different sets of initial seeds. With RANDLUZZ-generated seeds, all four fuzzers achieve the fastest bug exposure time. Specifically, AFLGo with RLS is the fastest at exposing 10 bugs, Beacon at 6, WindRanger at 9, and SelectFuzz at 6, bringing the total number to 31. The second fastest seed set is the Original Seeds, with AFLGo performing the best on 3 bugs, Beacon on 3, WindRanger on 2, and SelectFuzz on 3. On average, excluding the bugs with timeouts (CVE-2018-17360 and CVE-2017-14940), AFLGo with RLS discovers bugs 4.8× faster than OS and 4.3× faster than SS. Similarly, Beacon with RLS exposes bugs 3.4× faster than OS and 2.25× faster than SS. WindRanger with RLS exposes bugs 2.1× faster than OS and 4.3× faster than SS. On average, SelectFuzz with RLS discovers bugs slower than OS and SS; but on 6 bugs, SelectFuzz with RLS is the fastest to expose bugs. The drop of average speed mainly because it takes 17273 seconds to expose CVE-2017-7303. On some

bugs, RANDLUZZ-generated seeds can help fuzzers discover bugs within tens of seconds, more than 1,000× faster than other seed sets. For example, on the bug CVE-2018-14498, AFLGo with RLS discovers the bug 1,740× faster than OS. This demonstrates that our RANDLUZZ-generated seeds can significantly improve the efficiency of bug exposure.

### 4.3. RQ3: Efficiency of Bug Detection

In this section, we evaluate the performance of different fuzzers. To ensure a fair comparison, all fuzzers are equipped with RANDLUZZ-generated reachable seeds, running for 24 hours. Table 4 shows that, overall, RANDLUZZ performs the best, discovering 7 bugs the most quickly. For the rest 6 exposed bugs, although RANDLUZZ does not expose bugs the most quickly, the time to exposing bugs is close to the fastest one. For example, Beacon discovers CVE-2016-4487 in 3 seconds while RANDLUZZ exposes it

TABLE 4. TIME TO BUGS FOR DIFFERENT FUZZERS, WITH RANDLUZZ-GENERATED SEEDS.

CVE ID	AFLGo	Beacon	WindRanger	SelectFuzz	RANDLUZZ
CVE-2018-14498	34s	T.O. <sup>1</sup>	11s	40s	23s
CVE-2020-13790	17828s	T.O.	T.O.	T.O.	<b>6932s</b>
CVE-2016-4487	22s	<b>3s</b>	23s	94s	5s
CVE-2016-4489	1682s	272s	398s	1600s	<b>268s</b>
CVE-2016-4491	9s	5s	376s	13114s	<b>2s</b>
CVE-2017-8393	73s	T.O.	435s	1487s	<b>59s</b>
CVE-2017-8394	51s	75s	170s	3847s	<b>19s</b>
CVE-2017-8395	14s	<b>4s</b>	16s	17s	12s
CVE-2017-8392	27s	33207s	18s	<b>11s</b>	<b>11s</b>
CVE-2017-8397	764s	146s	83s	T.O.	<b>47s</b>
CVE-2018-17360	T.O.	T.O.	T.O.	T.O.	T.O.
CVE-2017-7303	13320s	2203s	<b>1614s</b>	17273s	12420s
CVE-2017-14940	<b>26s</b>	T.O.	139s	T.O.	458s
CVE-2017-16828	2419s	<b>2136s</b>	20978s	55485s	2313s
# of Fastest <sup>2</sup>	1	3	2	1	7

<sup>1</sup> T.O.: Timeout (24 hours = 86400 seconds).

<sup>2</sup> The number of vulnerabilities the fuzzer triggers most quickly.

in 5 seconds. Of the 13 bugs that RANDLUZZ can expose, 8 are discovered within 60 seconds, and 2 within a few hundred seconds, showing the efficiency of our RANDLUZZ. Even with RANDLUZZ-generated reachable seeds provided, RANDLUZZ achieves a bug exposure speedup ranging from  $1.01\times$  to  $2.7\times$  compared to the second-fastest exposure. The greatest speedup is observed with CVE-2016-4491, where RANDLUZZ takes only 2 seconds to expose the bug, while SelectFuzz requires 13,114 seconds, making it  $6557\times$  slower than RANDLUZZ. Since all fuzzers use RANDLUZZ-generated seeds, the superiority of RANDLUZZ must stem from its advanced mutation strategy. Therefore, both the reachable seed generation and bug-specific mutator construction are effective and efficient in exposing bugs.

#### 4.4. RQ4: LLM Model Influence

In this section, we evaluate the performance of different LLM models, including GPT-4 (GPT-4-0125-preview, GPT-4o, GPT-o1-mini (GPT-o1-mini-2024-09-12), and GPT-o1-preview (GPT-o1-preview-2024-09-12). We assess LLM models on three tasks, including command option construction via RAG (*RAG*), optimization for reachable seed generation (*Opt*), and bug-specific mutator generation (*Mutator*). Table 5 presents the time each GPT model makes to assist RANDLUZZ in these tasks. Updates in GPT models do not generally reduce completion time. GPT-o1-preview-2024-09-12 is newer than GPT-4o, but it completes the *Opt* task slower than GPT-4o due to higher computational costs.

As shown in Table 5, for the task *RAG*, GPT-4o and GPT-o1-mini perform similarly in this task, with GPT-4 slightly slower, while GPT-o1-preview takes significantly longer due to slower query response time. For the task *Opt*, GPT-4o proves more efficient than GPT-o1-mini, while GPT-4 and GPT-o1-preview are considerably slower, with timeouts in some cases. GPT-o1-mini shows faster single-

```

1 bfd_vma
2 bfd_getl32(const void *p)
3 {
4     const bfd_byte *addr = (const bfd_byte *) p;
5     unsigned long v;
6
7     v = (unsigned long) addr[0];
8     v |= (unsigned long) addr[1] << 8;
9     v |= (unsigned long) addr[2] << 16;
10    v |= (unsigned long) addr[3] << 24;
11
12    return v;
13 }

```

Figure 8. Function body of bfd\_getl32().

query responses than GPT-4o but sometimes requires re-queries due to formatting inconsistencies. GPT-4 occasionally produces incomplete or incorrect code analyses for complex programs, leading to inefficiencies. GPT-o1-preview, though slower per query, provides high-quality responses, reducing the total number of queries needed. For the task *Mutator*, GTP-4, GPT-4o, and GPT-o1-mini perform similarly in completion time, with GPT-o1-preview slightly slower. However, GPT-o1-mini and GPT-o1-preview frequently make errors in code generation, such as undefined variables in C code or indentation issues in Python. This requires additional LLM queries to correct the code.

After the comprehensive evaluation, GPT-4o emerges as the best choice for RANDLUZZ. It consistently performs well in seed optimization and code analysis, avoiding formatting issues and providing efficient task completion, unlike GPT-o1-mini, GPT-4, and GPT-o1-preview. GPT-4o effectively balances accuracy and efficiency, enhancing RANDLUZZ’s overall performance.

#### 4.5. Case Study for the Limitations of RANDLUZZ

In this section, we analyze two cases highlighting the challenges faced by RANDLUZZ. Table 3 shows that RLS performs worse than OS on a few CVE vulnerabilities, such as CVE-2018-17360 and CVE-2017-16828. We analyze the reasons in this section.

**4.5.1. Case 1: Lack of Neighbor Functions.** To address the issue of incomplete function call chains, RANDLUZZ proposes to generate seeds based on the functionalities of neighbor functions, which improves the possibility of generating reachable seeds. However, if we cannot obtain any neighbor function of the target function, we have to generate reachable seeds based on the target function. This may lead to the failure of reachable seed generation. For example, Figure 8 shows the target vulnerable function as described in the bug report of CVE-2018-17360. Clang fails to get both the complete function call chain and any neighbor functions of bfd\_getl32(). Therefore, in this case, RANDLUZZ has to generate reachable seeds based on the only function bfd\_getl32(). The functionality of this function is that it reads 4 bytes of data from a

TABLE 5. INFLUENCE OF DIFFERENT LLM MODELS

Program	CVE ID	GPT-4-0125-preview			GPT-4o			GPT-o1-mini-2024-09-12			GPT-o1-preview-2024-09-12		
		RAG <sup>1</sup>	Opt <sup>2</sup>	Mutator <sup>3</sup>	RAG	Opt	Mutator	RAG	Opt	Mutator	RAG	Opt	Mutator
cjpeg	CVE-2020-13790	134s	T.O. <sup>4</sup>	88s	64s	549s	61s	88s	238s	103s	506s	2176s	312s
cxxfilt	CVE-2016-4487	210s	1548s	122s	191s	102s	126s	142s	135s	142s	1020s	793s	338s
objcopy	CVE-2017-8393	222s	562s	143s	134s	210s	122s	181s	412s	98s	719s	1908s	442s
objdump	CVE-2017-8392	291s	1450s	107s	283s	622s	94s	203s	903s	134s	1103s	T.O.	293s

<sup>1</sup> **RAG**: command option via RAG. <sup>2</sup> **Opt**: seed optimization for reachable seeds. <sup>3</sup> **Mutator**: mutator generation for bug-specific mutators.

<sup>4</sup> **T.O.**: timeout (3600 seconds).

given memory address ( $p$ ) and combines them into a 32-bit unsigned integer. Although such functionality is common in dealing with binary data, network protocols, file format parsing, and hardware interactions, RANDLUZZ still cannot successfully generate reachable seeds. This is because LLMs have difficulties in identifying the functional module that `bfd_getl32()` belongs to. Thus, in our experiments, RANDLUZZ is unable to complete the seed optimization for CVE-2018-17360 within one hour (Table 2). Due to its inability to generate a reachable seed, RANDLUZZ is unable to trigger the vulnerability within 24 hours (Table 4).

In addition, in `Libjpeg` (`cjpeg`) related projects (CVE-2018-14498 and CVE-2020-13790), Beacon failed to find the complete function call chain from the entry function to the target function (`get_8bit_row()` for CVE-2018-14498 and `get_rgb_row()` for CVE-2020-13790), which led to the failure in identifying the target vulnerabilities. Specifically, we analyzed the output results of Beacon’s static analysis for both vulnerability projects. In their distance calculation records, we found that Beacon incorrectly pruned parts of the target path because it could not resolve the call paths based on function pointers. This indicates that Beacon’s fuzzing cannot reach the target location. As the result shown in Table 3 and Table 4, Beacon’s static analysis fails to find the target and terminates after running for some time, and regardless of the seeds used, Beacon cannot trigger the bugs of these two CVEs in `Libjpeg`.

#### 4.5.2. Case 2: Limitation of LLM-Generated Seeds.

The reachable seeds generated by RANDLUZZ significantly improve the efficiency of directed fuzzing. However, such seeds introduce a new issue because LLMs may add too many details in the seeds. Thus, in Table 3, fuzzers equipped with RANDLUZZ-generated seeds may trigger bugs slower than fuzzers with original seeds. For example, LLMs add too many debug sections into the seeds for triggering the bug CVE-2017-16828 in the program `readelf`. The CVE-2017-16828 describes a heap-based buffer over-read or crash issue that occurs when the program `readelf` processes a specially crafted ELF file containing debug information. Malicious DWARF data or incorrect DWARF offsets within the ELF file can lead to invalid pointer references or crashes. Since the vulnerability lies in parsing the input’s (ELF file) debug information, RANDLUZZ introduces various debug sections into the input during seed optimization to reach the target location. Although the generated seeds are capable of reaching the target, other sections

of debug information are added into the seed, such as `.debug_info`, `.debug_line`, and `.eh_frame`. The process also includes a substantial number of entries such as `.rel.debug_info`, `.rel.debug_aranges`, and `.rel.debug_line`. While these additional entries can potentially trigger the vulnerability, their volume inevitably increases the input length, which in turn reduces the efficiency of random mutations in finding the vulnerability.

## 5. Discussion

In this section, we discuss the insights gained from applying RANDLUZZ to real-world vulnerability projects and its interactions with LLMs. Each of these factors influences the effectiveness of RANDLUZZ and highlights areas for future improvement in automated vulnerability testing.

### 5.1. Description In CVE Report

Unlike other fuzzing methods requiring manual analysis, RANDLUZZ locates vulnerabilities by analyzing bug reports, reducing the manual workload. However, this requires the bug reports to specify the vulnerability’s location at the function-level. For instance, CVE-2016-4487 only mentions the `btypevec` variable but does not identify the vulnerable function, which is actually `remember_Ktype`. As a result, RANDLUZZ cannot locate the vulnerability, hindering further testing. Additionally, bug reports may mention functions that trigger the vulnerability or were patched, not necessarily the vulnerable function itself, leading to incorrect localization and unnecessary tests.

### 5.2. Generating Complex Inputs

In section 3.3.4, RANDLUZZ can generate simple string inputs via LLMs, but more complex inputs, like images, require Python scripts. We also encountered challenges with tools like `readelf`, `objdump`, and `objcopy`, which require ELF files or compiled programs. Although we asked the LLM to generate Python scripts for creating ELF files or compiling programs, the complexity of these files often caused the scripts to fail. By instructing the LLM to generate Python code that creates a `.c` file and compiles it locally, we successfully generated inputs requiring compilation. However, this approach struggles with even more complex inputs, like video or 3D modeling files. Future work will focus on optimizing input generation for more complex formats,

such as those involving multimodal generative models or multimodal LLMs [11], [23], [37].

### 5.3. Reasoning Capability of LLM

Given that large language models operate based on extensive text data as training sets, we concern that LLMs might rely on their memory rather than the inference capabilities when assisting RANDLUZZ. To investigate this, we design an experiment using a CVE vulnerability disclosed after the latest update to the model's knowledge base. We select GPT-4o to test the vulnerability CVE-2024-34459. GPT-4o's training dataset was last updated in October 2023, whereas information and patches for CVE-2024-34459 were only released in May 2024. The vulnerability is a heap-buffer-overflow and exists in the program `xmlint`, where the vulnerable function is `xmllint.c:xmlHTMLPrintFileContext()`.

GPT-4o successfully completes the tasks of *SA* (46 seconds), *RAG* (157 seconds), *Opt* (436 seconds), and *Mutator* (113 seconds). As a result, RANDLUZZ succeeds in triggering the bug in 941 seconds. The results show that with GPT-4o, RANDLUZZ can effectively handle vulnerabilities not included in the LLM's training data. Specifically, for CVE-2024-34459, which was disclosed after GPT-4o's latest update, RANDLUZZ successfully analyzes and optimizes seeds for this newly disclosed vulnerability. This also demonstrates the GPT-4o's inference ability, rather than solely relying on its memory of historical data, which is crucial in assisting RANDLUZZ to identify vulnerabilities.

## 6. Related Work

In this paper, we utilize LLMs to assist directed fuzzing, by using the capability of LLMs to understand code and bugs. The existing research on directed fuzzing can be classified into three categories, including optimization of the guidance towards target locations [1], [13], [34], reduction of unnecessary exploration of unrelated code [14], [20], [24], and directed fuzzing in specific scenarios [3], [30], [36]. Another related research is to use LLMs in assisting fuzzing, and the existing fuzzers focus on understanding the formats of inputs [5], [6], [31].

**Optimization of Guidance for Directed Fuzzing.** Since the first directed fuzzer AFLGo [1], researchers have developed various methods to optimize the guidance towards target locations. Hawkeye [4] improves AFLGo by considering multiple paths that can reach target locations. CAFL [18] then realizes that constraints in programs also impact the efficiency of reaching target locations.  $MC^2$  [28] introduces a performance metric, the number of oracle queries, to guide directed fuzzing. Another fuzzer 1dFuzz [35] also introduces an extra metric, the trailing call sequence, to guide directed fuzzing to follow specific function call sequences. Similarly, PDGF [40] introduces a new metric, the regional maturity, to improve efficiency by gradually reaching more predecessors. WindRanger [8] identifies basic blocks that deviate from target locations so that fuzzing can

focus more on execution paths that reach target locations. Titan [13] and WAFLGO [34] improves the efficiency of directed fuzzing via regarding it as a multi-objective optimization problem. DeepGo [21] improves fuzzing efficiency via using reinforcement learning. These fuzzers utilize the knowledge from experts and improve fuzzing efficiency based on different understandings of programs. However, they overlook the manual efforts required to successfully reach target locations.

**Reduction of Unrelated Exploration for Directed Fuzzing.** While the optimization of fuzzing process will prefer seeds that are more likely to reach target locations, it still examines a large number of inputs that explore unrelated code regions. Therefore, researchers either trim the size of generated inputs or use selective instrumentation to improve efficiency [14], [24], [45]. To trim the size of generated inputs, FuzzGuard [45] uses deep learning to predict if an input can reach target locations. Another fuzzer Halo [14] restrains input generation based on likely invariants so that more reachable inputs are generated. Selective instrumentation only instruments related code regions so that fuzzing will only be guided based on related code [20], [24]. Another solution is to terminate the execution of a program if it reaches unrelated code regions [12]. These fuzzers require accurate analysis of target programs so that the reduction of unrelated exploration does not miss related code regions. However, the accuracy may lead to heavy analysis of target programs. Moreover, they use pre-defined mutation operators, which do not connect to target bugs, impairing the efficiency of bug detection.

**Directed Fuzzing in Specific Scenarios.** Directed fuzzing is also used in different applications. To fuzz Linux kernel, SyzDirect [30] uses scalable static analysis to identify vulnerable information, which is used to guide the directed fuzzing. To detect concurrency use-after-free vulnerabilities, DDRace [36] identifies potential use-after-free locations as target locations and guides fuzzing to gradually satisfy those target sites. To detect Java deserialization vulnerabilities, ODDFuzz [3] identifies candidate gadget chains and guides fuzzing to examine each target site in the gadget chains. To detect vulnerabilities in Internet of Things (IoT) firmware, LABRADOR [22] deduces code coverage in firmware and estimates the distance to target sensitive code. These fuzzers always require deep understanding of target programs or vulnerabilities. However, such understanding is usually obtained from experts. For example, the patterns of target sites for use-after-free and deserialization vulnerabilities are identified based on expertise in those specific vulnerabilities. **LLM for Fuzzing.** Existing research on using LLMs for fuzzing concentrates on generating appropriate inputs for target programs. They use the excellent capability of LLMs to understand the specifications for inputs or learn input format during their training. TitanFuzz [5] uses the LLMs' capability of code generation and generates programs to invoke APIs in deep learning libraries. Some fuzzers extract protocol information from the specification so that fuzzers can generate valid and various messages to test protocol code [26], [31]. For well-known protocols, the format of

input message can be obtained by directly query LLMs [27]. Similarly, LLMs can also generate effective inputs for programs written in widely-used programming languages [15], [33]. Some other fuzzers explore the capability of LLMs to generate edge or unusual inputs for target programs [6], [38]. However, these fuzzers focus on the format of inputs, ignoring the potential of LLMs in other stages of fuzzing.

## 7. Conclusion

The dilemma of randomness in fuzzing motivates us to reduce the randomness while maintaining the effectiveness of fuzzing. We identify two key components, initial seeds and mutators, that are significantly impacted by randomness and propose solutions to solve the problems, respectively. LLM is crucial to the success of our RANDLUZZ. We use LLM to remove the randomness in initial seeds, generating reachable seeds. LLM is also used to reduce the randomness in mutators, constructing bug-specific mutators. The experiment results show that both the reachable seeds and bug-specific mutators are efficient in exposing bugs.

## References

- [1] Marcel Böhme, Van-Thuan Pham, Manh-Dung Nguyen, and Abhik Roychoudhury. Directed greybox fuzzing. In *Proceedings of the 2017 ACM SIGSAC conference on computer and communications security*, pages 2329–2344, 2017.
- [2] Marcel Böhme, Van-Thuan Pham, and Abhik Roychoudhury. Coverage-based greybox fuzzing as markov chain. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, pages 1032–1043, 2016.
- [3] Sicong Cao, Biao He, Xiaobing Sun, Yu Ouyang, Chao Zhang, Xiaoxue Wu, Ting Su, Lili Bo, Bin Li, Chuanlei Ma, et al. Oddfuzz: Discovering java deserialization vulnerabilities via structure-aware directed greybox fuzzing. In *2023 IEEE Symposium on Security and Privacy (SP)*, pages 2726–2743. IEEE, 2023.
- [4] Hongxu Chen, Yinxing Xue, Yuekang Li, Bihuan Chen, Xiaofei Xie, Xiuheng Wu, and Yang Liu. Hawkeye: Towards a desired directed grey-box fuzzer. In *Proceedings of the 2018 ACM SIGSAC conference on computer and communications security*, pages 2095–2108, 2018.
- [5] Yinlin Deng, Chunqiu Steven Xia, Haoran Peng, Chenyuan Yang, and Lingming Zhang. Large language models are zero-shot fuzzers: Fuzzing deep-learning libraries via large language models. In *Proceedings of the 32nd ACM SIGSOFT international symposium on software testing and analysis*, pages 423–435, 2023.
- [6] Yinlin Deng, Chunqiu Steven Xia, Chenyuan Yang, Shizhuo Dylan Zhang, Shujing Yang, and Lingming Zhang. Large language models are edge-case generators: Crafting unusual programs for fuzzing deep learning libraries. In *Proceedings of the 46th IEEE/ACM International Conference on Software Engineering*, pages 1–13, 2024.
- [7] Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. The faiss library. 2024.
- [8] Zhengjie Du, Yuekang Li, Yang Liu, and Bing Mao. Windranger: A directed greybox fuzzer driven by deviation basic blocks. In *Proceedings of the 44th International Conference on Software Engineering*, pages 2440–2451, 2022.
- [9] Xiaotao Feng, Ruoxi Sun, Xiaogang Zhu, Minhui Xue, Sheng Wen, Dongxi Liu, Surya Nepal, and Yang Xiang. Snipuzz: Black-box fuzzing of iot firmware via message snippet inference. In *Proceedings Of The 2021 ACM SIGSAC Conference On Computer And Communications Security*, pages 337–350, 2021.
- [10] Andrea Fioraldi, Dominik Maier, Heiko Eißfeldt, and Marc Heuse. AFL++: Combining incremental steps of fuzzing research. In *14th USENIX Workshop on Offensive Technologies (WOOT 20)*. USENIX Association, August 2020.
- [11] Zexin Hu, Kun Hu, Clinton Mo, Lei Pan, and Zhiyong Wang. Terrain diffusion network: Climatic-aware terrain generation with geological sketch guidance. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 12565–12573, 2024.
- [12] Heqing Huang, Yiyuan Guo, Qingkai Shi, Peisen Yao, Rongxin Wu, and Charles Zhang. Beacon: Directed grey-box fuzzing with provable path pruning. In *2022 IEEE Symposium on Security and Privacy (SP)*, pages 36–50. IEEE, 2022.
- [13] Heqing Huang, Peisen Yao, Hung-Chun Chiu, Yiyuan Guo, and Charles Zhang. Titan: Efficient multi-target directed greybox fuzzing. In *Proceedings of the 2024 IEEE Symposium on Security and Privacy (SP)*, San Francisco, CA, USA, pages 20–22, 2024.
- [14] Heqing Huang, Anshunkang Zhou, Mathias Payer, and Charles Zhang. Everything is good for something: Counterexample-guided directed fuzzing via likely invariant inference. In *2024 IEEE Symposium on Security and Privacy (SP)*, pages 142–142. IEEE Computer Society, 2024.
- [15] Sungmin Kang, Juyeon Yoon, and Shin Yoo. Large language models are few-shot testers: Exploring llm-based general bug reproduction. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*, pages 2312–2323. IEEE, 2023.
- [16] Tae Eun Kim, Jaeseung Choi, Seongjae Im, Kihong Heo, and Sang Kil Cha. Evaluating directed fuzzers: Are we heading in the right direction? *Proceedings of the ACM on Software Engineering*, 1(FSE):316–337, 2024.
- [17] George Klees, Andrew Ruef, Benji Cooper, Shiyi Wei, and Michael Hicks. Evaluating fuzz testing. In *Proceedings of the 2018 ACM SIGSAC conference on computer and communications security*, pages 2123–2138, 2018.
- [18] Gwangmu Lee, Woonchul Shim, and Byoungyoung Lee. Constraint-guided directed greybox fuzzing. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 3559–3576, 2021.
- [19] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474, 2020.
- [20] Penghui Li, Wei Meng, and Chao Zhang. Sdfuzz: Target states driven directed fuzzing. In *Proceedings of the 33rd USENIX Security Symposium (Security)*. Philadelphia, PA, USA, 2024.
- [21] Peihong Lin, Pengfei Wang, Xu Zhou, Wei Xie, Gen Zhang, and Kai Lu. Deepgo: Predictive directed greybox fuzzing. In *Network and Distributed System Security (NDSS) Symposium*, 2024.
- [22] Hangtian Liu, Shuitao Gan, Chao Zhang, Zicong Gao, Hongqi Zhang, Xiangzhi Wang, and Guangming Gao. Labrador: Response guided directed fuzzing for black-box iot devices. In *2024 IEEE Symposium on Security and Privacy (SP)*, pages 127–127. IEEE Computer Society, 2024.
- [23] Zhuqiang Lu, Kun Hu, Chaoyue Wang, Lei Bai, and Zhiyong Wang. Autoregressive omni-aware outpainting for open-vocabulary 360-degree image generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 14211–14219, 2024.
- [24] Changhua Luo, Wei Meng, and Penghui Li. Selectfuzz: Efficient directed fuzzing with selective path exploration. In *2023 IEEE Symposium on Security and Privacy (SP)*, pages 2693–2707. IEEE, 2023.
- [25] Chenyang Lyu, Shouling Ji, Chao Zhang, Yuwei Li, Wei-Han Lee, Yu Song, and Raheem Beyah. MOPT: Optimized mutation scheduling for fuzzers. In *28th USENIX Security Symposium (USENIX Security 19)*, pages 1949–1966, 2019.



- [26] Xiaoyue Ma, Lannan Luo, and Qiang Zeng. From one thousand pages of specification to unveiling hidden bugs: Large language model assisted fuzzing of matter {IoT} devices. In *33rd USENIX Security Symposium (USENIX Security 24)*, pages 4783–4800, 2024.
- [27] Ruijie Meng, Martin Mirchev, Marcel Böhme, and Abhik Roychoudhury. Large language model guided protocol fuzzing. In *Proceedings of the 31st Annual Network and Distributed System Security Symposium (NDSS)*, 2024.
- [28] Abhishek Shah, Dongdong She, Samanway Sadhu, Krish Singal, Peter Coffman, and Suman Jana.  $MC^2$ : Rigorous and efficient directed greybox fuzzing. In *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security*, pages 2595–2609, 2022.
- [29] Suhwan Song, Chengyu Song, Yeongjin Jang, and Byoungyoung Lee. Crfuzz: Fuzzing multi-purpose programs through input validation. In *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, pages 690–700, 2020.
- [30] Xin Tan, Yuan Zhang, Jiadong Lu, Xin Xiong, Zhuang Liu, and Min Yang. Syzdirect: Directed greybox fuzzing for linux kernel. In *Proceedings of the 2023 ACM SIGSAC Conference on Computer and Communications Security*, pages 1630–1644, 2023.
- [31] Jincheng Wang, Le Yu, and Xiapu Luo. Llmif: Augmented large language model for fuzzing iot devices. In *2024 IEEE Symposium on Security and Privacy (SP)*, pages 196–196. IEEE Computer Society, 2024.
- [32] Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers. *Advances in Neural Information Processing Systems*, 33:5776–5788, 2020.
- [33] Chunqiu Steven Xia, Matteo Paltenghi, Jia Le Tian, Michael Pradel, and Lingming Zhang. Fuzz4all: Universal fuzzing with large language models. In *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering*, pages 1–13, 2024.
- [34] Yi Xiang, Xuhong Zhang, Peiyu Liu, Shouling Ji, Xiao Xiao, Hong Liang, Jiacheng Xu, and Wenhui Wang. Critical code guided directed greybox fuzzing for commits. In *USENIX*, 2024.
- [35] Songtao Yang, Yubo He, Kaixiang Chen, Zheyu Ma, Xiapu Luo, Yong Xie, Jianjun Chen, and Chao Zhang. Idfuzz: Reproduce 1-day vulnerabilities with directed differential fuzzing. In *Proceedings of the 32nd ACM SIGSOFT International Symposium on Software Testing and Analysis*, pages 867–879, 2023.
- [36] Ming Yuan, Bodong Zhao, Penghui Li, Jiahuo Liang, Xinhui Han, Xiapu Luo, and Chao Zhang. Ddrace: Finding concurrency uaf vulnerabilities in linux drivers with directed fuzzing. In *USENIX Security Symposium*, pages 2849–2866, 2023.
- [37] Duzhen Zhang, Yahan Yu, Jiahua Dong, Chenxing Li, Dan Su, Chenhui Chu, and Dong Yu. Mm-llms: Recent advances in multimodal large language models. *arXiv preprint arXiv:2401.13601*, 2024.
- [38] Kunpeng Zhang, Shuai Wang, Jitao Han, Xiaogang Zhu, Xian Li, Shaohua Wang, and Sheng Wen. Your fix is my exploit: Enabling comprehensive dl library api fuzzing with large language models. In *2025 IEEE/ACM 47th International Conference on Software Engineering (ICSE)*. IEEE, 2025.
- [39] Kunpeng Zhang, Xiaogang Zhu, Xi Xiao, Minhui Xue, Chao Zhang, and Sheng Wen. Shapfuzz: Efficient fuzzing via shapley-guided byte selection. In *NDSS*, 2024.
- [40] Yujian Zhang, Yaokun Liu, Jinyu Xu, and Yanhao Wang. Predecessor-aware directed greybox fuzzing. In *2024 IEEE Symposium on Security and Privacy (SP)*, pages 40–40. IEEE Computer Society, 2024.
- [41] Xiaogang Zhu and Marcel Böhme. Regression greybox fuzzing. In *Proceedings Of The 2021 ACM SIGSAC Conference On Computer And Communications Security*, pages 2169–2182, 2021.
- [42] Xiaogang Zhu, Xiaotao Feng, Xiaozhu Meng, Sheng Wen, Seyit Camtepe, Yang Xiang, and Kui Ren. CSI-Fuzz: Full-speed edge tracing using coverage sensitive instrumentation. *IEEE Transactions On Dependable And Secure Computing*, 19(2):912–923, 2020.
- [43] Xiaogang Zhu, Sheng Wen, Seyit Camtepe, and Yang Xiang. Fuzzing: a survey for roadmap. *ACM Computing Surveys (CSUR)*, 54(11s):1–36, 2022.
- [44] Xiaogang Zhu, Wei Zhou, Qing-Long Han, Wanlun Ma, Sheng Wen, and Yang Xiang. When software security meets large language models: A survey. *IEEE/CAA Journal Of Automatica Sinica*, 2024.
- [45] Peiyuan Zong, Tao Lv, Dawei Wang, Zizhuang Deng, Ruigang Liang, and Kai Chen. Fuzzguard: Filtering out unreachable inputs in directed grey-box fuzzing through deep learning. In *29th USENIX security symposium (USENIX security 20)*, pages 2255–2269, 2020.