

1

# OverHAuL

2

**Harnessing Automation for C Libraries via LLMs**

3

Konstantinos Chousos

4

July, 2025

5 Lorem ipsum odor amet, consectetur adipiscing elit. Habitasse congue tempus erat rhoncus  
6 sapien interdum dolor nec. Posuere habitant metus tellus erat eu. Risus ultricies eu rhoncus,  
7 conubia euismod convallis commodo per. Nam tellus quisque maximus dui eleifend; arcu aptent.  
8 Nisi rutrum primis luctus tortor tempor maecenas. Donec curae cras dolor; malesuada ultricies  
9 scelerisque. Molestie class tincidunt quis gravida ut proin. Consequat lacinia arcu justo leo maecenas  
10 nunc neque ex. Platea eros ullamcorper nullam rutrum facilisis.

# **Preface**

12 Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis sagittis posuere ligula sit amet lacinia.  
13 Duis dignissim pellentesque magna, rhoncus congue sapien finibus mollis. Ut eu sem laoreet,  
14 vehicula ipsum in, convallis erat. Vestibulum magna sem, blandit pulvinar augue sit amet, auctor  
15 malesuada sapien. Nullam faucibus leo eget eros hendrerit, non laoreet ipsum lacinia. Curabitur  
16 cursus diam elit, non tempus ante volutpat a. Quisque hendrerit blandit purus non fringilla. Integer  
17 sit amet elit viverra ante dapibus semper. Vestibulum viverra rutrum enim, at luctus enim posuere  
18 eu. Orci varius natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus.

19 Nunc ac dignissim magna. Vestibulum vitae egestas elit. Proin feugiat leo quis ante condimentum,  
20 eu ornare mauris feugiat. Pellentesque habitant morbi tristique senectus et netus et malesuada fames  
21 ac turpis egestas. Mauris cursus laoreet ex, dignissim bibendum est posuere iaculis. Suspendisse  
22 et maximus elit. In fringilla gravida ornare. Aenean id lectus pulvinar, sagittis felis nec, rutrum  
23 risus. Nam vel neque eu arcu blandit fringilla et in quam. Aliquam luctus est sit amet vestibulum  
24 eleifend. Phasellus elementum sagittis molestie. Proin tempor lorem arcu, at condimentum purus  
25 volutpat eu. Fusce et pellentesque ligula. Pellentesque id tellus at erat luctus fringilla. Suspendisse  
26 potenti.

27 Etiam maximus accumsan gravida. Maecenas at nunc dignissim, euismod enim ac, bibendum ipsum.  
28 Maecenas vehicula velit in nisl aliquet ultricies. Nam eget massa interdum, maximus arcu vel,  
29 pretium erat. Maecenas sit amet tempor purus, vitae aliquet nunc. Vivamus cursus urna velit,  
30 eleifend dictum magna laoreet ut. Duis eu erat mollis, blandit magna id, tincidunt ipsum. Integer  
31 massa nibh, commodo eu ex vel, venenatis efficitur ligula. Integer convallis lacus elit, maximus  
32 eleifend lacus ornare ac. Vestibulum scelerisque viverra urna id lacinia. Vestibulum ante ipsum  
33 primis in faucibus orci luctus et ultrices posuere cubilia curae; Aenean eget enim at diam bibendum  
34 tincidunt eu non purus. Nullam id magna ultrices, sodales metus viverra, tempus turpis.

## 35 Acknowledgments

36 Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis sagittis posuere ligula sit amet lacinia.  
37 Duis dignissim pellentesque magna, rhoncus congue sapien finibus mollis. Ut eu sem laoreet,  
38 vehicula ipsum in, convallis erat. Vestibulum magna sem, blandit pulvinar augue sit amet, auctor  
39 malesuada sapien. Nullam faucibus leo eget eros hendrerit, non laoreet ipsum lacinia. Curabitur  
40 cursus diam elit, non tempus ante volutpat a. Quisque hendrerit blandit purus non fringilla. Integer  
41 sit amet elit viverra ante dapibus semper. Vestibulum viverra rutrum enim, at luctus enim posuere  
42 eu. Orci varius natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus.

# Table of contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Motivation	1
1.2	Preview of following sections (rename)	2
<b>2</b>	<b>Background</b>	<b>3</b>
2.1	Fuzzing	3
2.1.1	Fuzzing examples	3
2.1.2	Fuzzer engines	3
2.2	Large Language Models (LLMs)	3
2.2.1	Prompting	4
2.2.2	LLM Programming Libraries (?)	4
2.3	Neurosymbolic AI	4
<b>3</b>	<b>Related work</b>	<b>5</b>
3.1	Previous projects	5
3.1.1	KLEE	5
3.1.2	IRIS	5
3.1.3	FUDGE	6
3.1.4	UTopia	6
3.1.5	FuzzGen	6
3.1.6	OSS-Fuzz	7
3.1.7	OSS-Fuzz-Gen	7
3.1.8	AutoGen	8
3.2	Differences	8
3.2.1	IntelliGen [[20250711141156]]	9
3.2.2	CKGFuzzer [[20250711203054]]	10
3.2.3	PromptFuzz [[20250713225436]]	11
<b>4</b>	<b>Overview</b>	<b>12</b>
4.1	Architecture	12
<b>5</b>	<b>Evaluation</b>	<b>13</b>
5.1	Benchmarks	13
5.2	Performance	13
5.3	Issues	13
5.4	Future work	13
5.4.1	Technical future work	13
5.4.2	Architectural future work/extensions	13

78	<b>6 Future work</b>	<b>14</b>
79	<b>7 Discussion</b>	<b>15</b>
80	<b>8 Conclusion</b>	<b>16</b>
81	8.1 Acknowledgements . . . . .	16
82	<b>Bibliography</b>	<b>17</b>

# 1 Introduction

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis sagittis posuere ligula sit amet lacinia. Duis dignissim pellentesque magna, rhoncus congue sapien finibus mollis. Ut eu sem laoreet, vehicula ipsum in, convallis erat. Vestibulum magna sem, blandit pulvinar augue sit amet, auctor malesuada sapien. Nullam faucibus leo eget eros hendrerit, non laoreet ipsum lacinia. Curabitur cursus diam elit, non tempus ante volutpat a. Quisque hendrerit blandit purus non fringilla. Integer sit amet elit viverra ante dapibus semper. Vestibulum viverra rutrum enim, at luctus enim posuere eu. Orci varius natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus.

Nunc ac dignissim magna. Vestibulum vitae egestas elit. Proin feugiat leo quis ante condimentum, eu ornare mauris feugiat. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris cursus laoreet ex, dignissim bibendum est posuere iaculis. Suspendisse et maximus elit. In fringilla gravida ornare. Aenean id lectus pulvinar, sagittis felis nec, rutrum risus. Nam vel neque eu arcu blandit fringilla et in quam. Aliquam luctus est sit amet vestibulum eleifend. Phasellus elementum sagittis molestie. Proin tempor lorem arcu, at condimentum purus volutpat eu. Fusce et pellentesque ligula. Pellentesque id tellus at erat luctus fringilla. Suspendisse potenti.

## 1.1 Motivation

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis sagittis posuere ligula sit amet lacinia. Duis dignissim pellentesque magna, rhoncus congue sapien finibus mollis. Ut eu sem laoreet, vehicula ipsum in, convallis erat. Vestibulum magna sem, blandit pulvinar augue sit amet, auctor malesuada sapien. Nullam faucibus leo eget eros hendrerit, non laoreet ipsum lacinia. Curabitur cursus diam elit, non tempus ante volutpat a. Quisque hendrerit blandit purus non fringilla. Integer sit amet elit viverra ante dapibus semper. Vestibulum viverra rutrum enim, at luctus enim posuere eu. Orci varius natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus.

Nunc ac dignissim magna. Vestibulum vitae egestas elit. Proin feugiat leo quis ante condimentum, eu ornare mauris feugiat. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris cursus laoreet ex, dignissim bibendum est posuere iaculis. Suspendisse et maximus elit. In fringilla gravida ornare. Aenean id lectus pulvinar, sagittis felis nec, rutrum risus. Nam vel neque eu arcu blandit fringilla et in quam. Aliquam luctus est sit amet vestibulum eleifend. Phasellus elementum sagittis molestie. Proin tempor lorem arcu, at condimentum purus volutpat eu. Fusce et pellentesque ligula. Pellentesque id tellus at erat luctus fringilla. Suspendisse potenti.

## 115 1.2 Preview of following sections (rename)

116 Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis sagittis posuere ligula sit amet lacinia.  
117 Duis dignissim pellentesque magna, rhoncus congue sapien finibus mollis. Ut eu sem laoreet,  
118 vehicula ipsum in, convallis erat. Vestibulum magna sem, blandit pulvinar augue sit amet, auctor  
119 malesuada sapien. Nullam faucibus leo eget eros hendrerit, non laoreet ipsum lacinia. Curabitur  
120 cursus diam elit, non tempus ante volutpat a. Quisque hendrerit blandit purus non fringilla. Integer  
121 sit amet elit viverra ante dapibus semper. Vestibulum viverra rutrum enim, at luctus enim posuere  
122 eu. Orci varius natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus.



## 123 2 Background

### 124 2.1 Fuzzing

125 What is fuzzing [1].

126 terminology: fuzz campaign, harness, driver, target, corpus

127 Why fuzz?

#### 128 2.1.1 Fuzzing examples

129 Heartbleed [2], shellshock [3].

#### 130 2.1.2 Fuzzer engines

131 C/C++: AFL [4] & AFL++ [4, pp. ++]. LibFuzzer [5].

132 Python: Atheris [6].

133 Java, Rust etc...

134 An example of a fuzz target/harness can be seen in Listing 2.1 [5].

135 OSS-Fuzz: 2016, after heartbleed.

### 136 2.2 Large Language Models (LLMs)

137 Transformers [7], 2017–2025. ChatGPT/OpenAI history & context. Claude, Llama (1–3) etc.

---

**Listing 2.1** A simple function that does something interesting if it receives the input “HI!”.

---

```
1 cat << EOF > test_fuzzer.cc
2 #include <stdint.h>
3 #include <stddef.h>
4 extern "C" int LLVMFuzzerTestOneInput(const uint8_t *data, size_t size) {
5     if (size > 0 && data[0] == 'H')
6         if (size > 1 && data[1] == 'I')
7             if (size > 2 && data[2] == '!')
8                 __builtin_trap();
9     return 0;
10 }
11 EOF
12 # Build test_fuzzer.cc with asan and link against libFuzzer.
13 clang++ -fsanitize=address,fuzzer test_fuzzer.cc
14 # Run the fuzzer with no corpus.
15 ./a.out
```

---

### 138 2.2.1 Prompting

139 Prompting techniques.

- 140 1. Zero-shot.
- 141 2. One-shot.
- 142 3. Chain of Thought [8].
- 143 4. ReACt [9].
- 144 5. Tree of Thoughts [10].

145 Comparison, strengths weaknesses etc. [11].

146 [12]

### 147 2.2.2 LLM Programming Libraries (?)

148 Langchain & LangGraph, LlamaIndex [13]–[15]. DSPy [16].

149 Comparison, relevance to our usecase.

## 150 2.3 Neurosymbolic AI

151 **TODO** [17]–[22].

## 3 Related work

Automated testing, automated fuzzing and automated harness creation have a long research history. Still, a lot of ground remains to be covered until true automation of these tasks is achieved. Until the introduction of transformers [7] and the 2020’s boom of commercial GPTs [23], automation regarding testing and fuzzing was mainly attempted through static and dynamic program analysis methods. These approaches are still utilized, but the fuzzing community has shifted almost entirely to researching the incorporation and employment of LLMs in the last half decade, in the name of automation [24]–[33].

### 3.1 Previous projects

#### 3.1.1 KLEE

KLEE [34] is a seminal and widely cited symbolic execution engine introduced in 2008 by Cadar et al. It was designed to automatically generate high-coverage test cases for programs written in C, using symbolic execution to systematically explore the control flow of a program. KLEE operates on the LLVM [35] bytecode representation of programs, allowing it to be applied to a wide range of C programs compiled to the LLVM intermediate representation.

Instead of executing a program on concrete inputs, KLEE performs symbolic execution—that is, it runs the program on symbolic inputs, which represent all possible values simultaneously. At each conditional branch, KLEE explores both paths by forking the execution and accumulating path constraints (i.e., logical conditions on input variables) along each path. This enables it to traverse many feasible execution paths in the program, including corner cases that may be difficult to reach through random testing or manual test creation.

When an execution path reaches a terminal state (e.g., a program exit, an assertion failure, or a segmentation fault), KLEE uses a constraint solver to compute concrete input values that satisfy the accumulated constraints for that path. These values form a test case that will deterministically drive the program down that specific path when executed concretely.

#### 3.1.2 IRIS

IRIS [24] is a 2025 open-source neurosymbolic system for static vulnerability analysis. Given a codebase and a list of user-specified Common Weakness Enumerations (CWEs), it analyzes source code to identify paths that may correspond to known vulnerability classes. IRIS combines symbolic analysis—such as control- and data-flow reasoning—with neural models trained to generalize over

code patterns. It outputs candidate vulnerable paths along with explanations and CWE references. The system operates on full repositories and supports extensible CWE definitions.

### 3.1.3 FUDGE

FUDGE [33] is a closed-source tool, made by Google, for automatic harness generation of C and C++ projects based on existing client code. It was used in conjunction with and in the improvement of Google’s OSS-Fuzz [36]. Being deployed inside Google’s infrastructure, FUDGE continuously examines Google’s internal code repository, searching for code that uses external libraries in a meaningful and “fuzzable” way (i.e. predominantly for parsing). If found, such code is **sliced** [37], per FUDGE, based on its Abstract Syntax Tree (AST) using LLVM’s Clang tool [35]. The above process results in a set of abstracted mostly-self-contained code snippets that make use of a library’s calls and/or API. These snippets are later **synthesized** into the body of a fuzz driver, with variables being replaced and the fuzz input being utilized. Each is then injected in an LLVMFuzzerTestOneInput function and finalized as a fuzzing harness. A building and evaluation phase follows for each harness, where they are executed and examined. Every passing harness along with its evaluation results is stored in FUDGE’s database, reachable to the user through a custom web-based UI.

### 3.1.4 UTopia

UTopia [29] (stylized UTOPIA) is another open-source automatic harness generation framework. Aside from the library code, It operates solely on user-provided unit tests since, according to Jeong, Jang, Yi, *et al.* [29], they are a resource of complete and correct API usage examples containing working library set-ups and tear-downs. Additionally, each of them are already close to a fuzz target, in the sense that they already examine a single and self-contained API usage pattern. Each generated harness follows the same data flow of the originating unit test. Static analysis is employed to figure out what fuzz input placement would yield the most results. It is also utilized in abstracting the tests away from the syntactical differences between testing frameworks, along with slicing and AST traversing using Clang.

### 3.1.5 FuzzGen

Another project of Google is FuzzGen [32], this time open-source. Like FUDGE, it leverages existing client code of the target library to create fuzz targets for it. FuzzGen uses whole-system analysis, through which it creates an *Abstract API Dependence Graph* ( $A^2DG$ ). It uses the latter to automatically generate LibFuzzer-compatible harnesses. For FuzzGen to work, the user needs to provide both client code and/or tests for the API and the API library’s source code as well. FuzzGen uses the client code to infer the *correct usage* of the API and not its general structure, in contrast to FUDGE. FuzzGen’s workflow can be divided into three phases: **1. API usage inference.** By consuming and analyzing client code and tests that concern the library under test, FuzzGen recognizes which functions belong to the library and learns its correct API usage patterns. This process is done with the help of Clang. To test if a function is actually a part of the library, a sample program is created that uses it. If the program compiles successfully, then the function is indeed a valid API call. **2.  $A^2DG$  construction mechanism.** For all the existing API calls, FuzzGen

220 builds an A<sup>2</sup>DG to record the API usages and infers its intended structure. After completion, this  
221 directed graph contains all the valid API call sequences found in the client code corpus. It is built  
222 in a two-step process: First, many smaller A<sup>2</sup>DGs are created, one for each root function per client  
223 code snippet. Once such graphs have been created for all the available client code instances, they  
224 are combined to formulate the master A<sup>2</sup>DG. This graph can be seen as a template for correct usage  
225 of the library. **3. Fuzzer generator.** Through the A<sup>2</sup>DG, a fuzzing harness is created. Contrary to  
226 FUDGE, FuzzGen does not create multiple “simple” harnesses but a single complex one with the  
227 goal of covering the whole of the A<sup>2</sup>DG. In other words, while FUDGE fuzzes a single API call at a  
228 time, FuzzGen’s result is a single harness that tries to fuzz the given library all at once through  
229 complex API usage.

### 230 3.1.6 OSS-Fuzz

231 OSS-Fuzz [36], [38] is a continuous, scalable and distributed cloud fuzzing solution for critical  
232 and prominent open-source projects. Developers of such software can submit their projects to  
233 OSS-Fuzz’s platform, where its harnesses are built and constantly executed. This results in multiple  
234 bug findings that are later disclosed to the primary developers and are later patched.

235 OSS-Fuzz started operating in 2016, an initiative in response to the Heartbleed vulnerability [2],  
236 [39], [40]. Its hope is that through more extensive fuzzing such errors could be caught and corrected  
237 before having the chance to be exploited and thus disrupt the public digital infrastructure. So far,  
238 it has helped uncover over 10,000 security vulnerabilities and 36,000 bugs across more than 1,000  
239 projects, significantly enhancing the quality and security of major software like Chrome, OpenSSL,  
240 and systemd.

241 A project that’s part of OSS-Fuzz must have been configured as a ClusterFuzz [41] project. Cluster-  
242 Fuzz is the fuzzing infrastructure that OSS-Fuzz uses under the hood and depends on Google Cloud  
243 Platform services, although it can be hosted locally. Such an integration requires setting up a build  
244 pipeline, fuzzing jobs and expects a Google Developer account. Results are accessible through a  
245 web interface. ClusterFuzz, and by extension OSS-Fuzz, supports fuzzing through LibFuzzer, AFL++,  
246 Honggfuzz and FuzzTest—successor to Centipede— with the last two being Google projects [5],  
247 [42]–[44]. C, C++, Rust, Go, Python and Java/JVM projects are supported.

### 248 3.1.7 OSS-Fuzz-Gen

249 OSS-Fuzz-Gen (OFG) [27], [45] is Google’s current State-Of-The-Art (SOTA) project regarding  
250 automatic harness generation through LLMs. It’s purpose is to improve the fuzzing infrastructure of  
251 open-source projects that are already integrated into OSS-Fuzz. Given such a project, OSS-Fuzz-Gen  
252 uses its preexisting fuzzing harnesses and modifies them to produce new ones. Its architecture  
253 can be described as follows: 1. With an OSS-Fuzz project’s GitHub repository link, OSS-Fuzz-  
254 Gen iterates through a set of predefined build templates and generates potential build scripts  
255 for the project’s harnesses. 2. If any of them succeed they are once again compiled, this time  
256 through fuzz-introspector [46]. The latter constitutes a static analysis tool, with fuzzer developers  
257 specifically in mind. 3. Build results, old harness and fuzz-introspector report are included in a  
258 template-generated prompt, through which an LLM is called to generate a new harness. 4. The

newly generated fuzz target is compiled and if it is done so successfully it begins execution inside OSS-Fuzz’s infrastructure.

This method proved meaningful, with code coverage in fuzz campaigns increasing thanks to the new generated fuzz drivers. In the case of [47], line coverage went from 38% to 69% without any manual interventions [45].

In 2024, OSS-Fuzz-Gen introduced an experimental feature for generating harnesses in previously unfuzzed projects [48]. The code for this feature resides in the `experimental/from_scratch` directory of the project’s GitHub repository [27], with the latest known working commit being 171aac2 and the latest overall commit being four months ago.

### 3.1.8 AutoGen

AutoGen [25] is a closed-source tool that generates new fuzzing harnesses, given only the library code and documentation. It works as following: The user specifies the function for which a harness is to be generated. AutoGen gathers information for this function—such as the function body, used header files, function calling examples—from the source code and documentation. Through specific prompt templates containing the above information, an LLM is tasked with generating a new fuzz driver, while another is tasked with generating a compilation command for said driver. If the compilation fails, both LLMs are called again to fix the problem, whether it was on the driver’s or command’s side. This loop iterates until a predefined maximum value or until a fuzz driver is successfully generated and compiled. If the latter is the case, it is then executed. If execution errors exist, the LLM responsible for the driver generation is used to correct them. If not, the pipeline has terminated and a new fuzz driver has been successfully generated.

## 3.2 Differences

OverHAuL differs, in some way, with each of the aforementioned works. Firstly, although KLEE and IRIS [24], [34] tackle the problem of automated testing and both IRIS and OverHAuL can be considered neurosymbolic AI tools, the similarities end there. None of them utilize LLMs the same way we do—with KLEE not utilizing them by default, as it precedes them chronologically—and neither are automating any part of the fuzzing process.

When it comes to FUDGE, FuzzGen and UTopia [29], [32], [33], all three depend on and demand existing client code and/or unit tests. On the other hand, OverHAuL requires only the bare minimum: the library code itself. Another point of difference is that in contrast with OverHAuL, these tools operate in a linear fashion. No feedback is produced or used in any step and any point failure results in the termination of the entire run.

OverHAuL challenges a common principle of these tools, stated explicitly in FUDGE’s paper [33]: “Choosing a suitable fuzz target (still) requires a human”. OverHAuL chooses to let the LLM, instead of the user, explore the available functions and choose one to target in its fuzz driver.

OSS-Fuzz-Gen [27] can be considered a close counterpart of OverHAuL, and in some ways it is. A lot of inspiration was gathered from it, like for example the inclusion of static analysis and its

usage in informing the LLM. Yet, OSS-Fuzz-Gen has a number of disadvantages that make it in some cases an inferior option. For one, OFG is tightly coupled with the OSS-Fuzz platform [36], which even on its own creates a plethora of issues for the common developer. To integrate their project into OSS-Fuzz, they would need to: Transform it into a ClusterFuzz project [41] and take time to write harnesses for it. Even if these prerequisites are carried out, it probably would not be enough. Per OSS-Fuzz’s documentation [38]: “To be accepted to OSS-Fuzz, an open-source project must have a significant user base and/or be critical to the global IT infrastructure”. This means that OSS-Fuzz is a viable option only for a small minority of open-source developers and maintainers. One countermeasure of the above shortcoming would be for a developer to run OSS-Fuzz-Gen locally. This unfortunately proves to be an arduous task. As it is not meant to be used standalone, OFG is not packaged in the form of a self-contained application. This makes it hard to setup and difficult to use interactively. Like in the case of FUDGE, OFG’s actions are performed linearly. No feedback is utilized nor is there graceful error handling in the case of a step’s failure. Even in the case of the experimental feature for bootstrapping unfuzzed projects, OFG’s performance varies heavily. During experimentation, a lot of generated harnesses were still wrapped either in Markdown backticks or  `tags, or were accompanied with explanations inside the generated .c source file. Even if code was formatted correctly, in many cases it missed necessary headers for compilation or used undeclared functions.`

Lastly, the closest counterpart to OverHAuL is AutoGen [25]. Their similarity stands in the implementation of a feedback loop between LLM and generated harness. However, most other implementation decisions remain distinct. One difference regards the fuzzed function. While AutoGen requires a target function to be specified by the user in which it narrows during its whole run, OverHAuL delegates this to the LLM, letting it explore the codebase and decide by itself the best candidate. Another difference lies in the need—and the lack of—of documentation. While AutoGen requires it to gather information for the given function, OverHAuL leans into the role of a developer by reading the related code and comments and thus avoiding any mismatches between documentation and code. Finally, the LLMs’ input is built based on predefined prompt templates, a technique also present in OSS-Fuzz-Gen. OverHAuL operates one abstraction level higher, leveraging DSPy [16] for programming instead of prompting the LLMs used.

In conclusion, OverHAuL constitutes an *open-source* tool that offers new functionality by offering a straightforward installation process, packaged as a self-contained Python package with minimal external dependencies. It also introduces novel approaches compared to previous work by

1. Implementing a feedback mechanism between harness generation, compilation, and evaluation phases,
2. Using autonomous ReAct agents capable of codebase exploration,
3. Leveraging a vector store for code consumption and retrieval.

**TODO** να συμπεριλάβω και τα:

### 3.2.1 IntelliGen [[20250711141156]]

**SAMPLE**

335 **IntelliGen: Automatic Fuzz Driver Synthesis Based on Vulnerability Heuristics** Zhang et  
336 al. (2021) present **IntelliGen**, a system for automatically synthesizing fuzz drivers by statically  
337 identifying potentially vulnerable entry-point functions within C projects. Implemented using  
338 LLVM, IntelliGen focuses on improving fuzzing efficiency by targeting code more likely to contain  
339 memory safety issues, rather than exhaustively fuzzing all available functions.

340 The system comprises two main components: the **Entry Function Locator** and the **Fuzz Driver**  
341 **Synthesizer**. The Entry Function Locator analyzes the project’s abstract syntax tree (AST) and clas-  
342 sifies functions based on heuristics that indicate vulnerability. These include pointer dereferencing,  
343 calls to memory-related functions (e.g., `memcpy`, `memset`), and invocation of other internal functions.  
344 Functions that score highly on these metrics are prioritized for fuzz driver generation. The guiding  
345 insight is that entry points with fewer argument checks and more direct memory operations expose  
346 more useful program logic for fuzz testing.

347 The Fuzz Driver Synthesizer then generates harnesses for these entry points. For each target  
348 function, it synthesizes a `LLVMFuzzerTestOneInput` function that invokes the target with arguments  
349 derived from the fuzzer input. This process involves inferring argument types from the source code  
350 and ensuring that runtime behavior does not violate memory safety—thus avoiding invalid inputs  
351 that would cause crashes unrelated to genuine bugs.

352 IntelliGen stands out by integrating static vulnerability estimation into the driver generation  
353 pipeline. Compared to prior tools like FuzzGen and FUDGE, it uses a more targeted, heuristic-based  
354 selection of functions, increasing the likelihood that fuzzing will exercise meaningful and vulnerable  
355 code paths.

### 356 3.2.2 CKGFuzzer [[20250711203054]]

#### 357 SAMPLE

358 CKGFuzzer is a fuzzing framework designed to automate the generation of effective fuzz drivers  
359 for C/C++ libraries by leveraging static analysis and large language models. Its workflow begins by  
360 parsing the target project along with any associated library APIs to construct a code knowledge  
361 graph. This involves two primary steps: first, parsing the abstract syntax tree (AST), and second,  
362 performing interprocedural program analysis. Through this process, CKGFuzzer extracts essential  
363 program elements such as data structures, function signatures, function implementations, and call  
364 relationships.

365 Using the knowledge graph, CKGFuzzer then identifies and queries meaningful API combinations,  
366 focusing on those that are either frequently invoked together or exhibit functional similarity.  
367 It generates candidate fuzz drivers for these combinations and attempts to compile them. Any  
368 compilation errors encountered during this phase are automatically repaired using heuristics and  
369 domain knowledge. A dynamically updated knowledge base, constructed from prior library usage  
370 patterns, guides both the generation and repair processes.

371 Once the drivers are successfully compiled, CKGFuzzer executes them while monitoring code  
372 coverage at the file level. It uses coverage feedback to iteratively mutate underperforming API  
373 combinations, refining them until new execution paths are discovered or a preset mutation budget  
374 is exhausted.



375 Finally, any crashes triggered during fuzzing are subjected to a reasoning process based on chain-  
376 of-thought prompting. To help determine their severity and root cause, CKGFuzzer consults an  
377 LLM-generated knowledge base containing real-world examples of vulnerabilities mapped to known  
378 Common Weakness Enumeration (CWE) entries.

### 379 3.2.3 PromptFuzz [[20250713225436]]

#### 380 SAMPLE

381 Lyu et al. (2024) introduce PromptFuzz [49], a system for automatically generating fuzz drivers using  
382 LLMs, with a novel focus on **prompt mutation** to improve coverage. The system is implemented  
383 in Rust and targets C libraries, aiming to explore more of the API surface with each iteration.

384 The workflow begins with the random selection of API functions, extracted from header file  
385 declarations. These functions are used to construct initial prompts that instruct the LLM to generate  
386 a simple program utilizing the API. Each generated program is compiled, executed, and monitored  
387 for code coverage. Programs that fail to compile or violate runtime checks (e.g., sanitizers) are  
388 discarded.

389 A key innovation in PromptFuzz is **coverage-guided prompt mutation**. Instead of mutating  
390 generated code directly, PromptFuzz mutates the LLM prompts—selecting new combinations of API  
391 functions to target unexplored code paths. This process is guided by a **power scheduling** strategy  
392 that prioritizes underused or promising API functions based on feedback from previous runs.

393 Once an effective program is produced, it is transformed into a fuzz driver by replacing constants  
394 and arguments with variables derived from the fuzzer input. Multiple such drivers are embedded  
395 into a single harness, where the input determines which program variant to execute, typically via a  
396 case-switch construct.

397 Overall, PromptFuzz demonstrates that prompt-level mutation enables more effective exploration  
398 of complex APIs and achieves better coverage than direct code mutations, offering a compelling  
399 direction for LLM-based automated fuzzing systems.

## 400 4 Overview

- 401 1. How is it different?
  - 402 2. What does it offer?
  - 403 3. Example uses
  - 404 4. Scope of Usage
- 
- 405 1. In what contexts does it work?
  - 406 2. Prerequisites

### 407 4.1 Architecture

- 408 • **System diagram**
- 409 • Main Library Architecture/Structure
- 410 • LLM usage
- 411 – Prompting techniques used (callback to Section [2.2.1](#)).
- 412 • Static analysis
- 413 • Code localization(?)
- 414 • Fuzzers
- 415 • GitHub Workflow/Usage
- 416 • “Iteration budget”

## 417 5 Evaluation

### 418 5.1 Benchmarks

419 Results from integration with 10/100 open-source C/C++ projects.

### 420 5.2 Performance

### 421 5.3 Issues

### 422 5.4 Future work

#### 423 5.4.1 Technical future work

#### 424 5.4.2 Architectural future work/extensions

- 425 1. Build system
- 426 2. More (static) analysis tools integrations
- 427 3. General *localization* problem

## 6 Future work

- More diverse build support
- More language support
- More fuzzers support
  - GraphFuzz [28]
- Experimentation with different LLM providers/models
  - Code-specific LLMs
    - \* codex-1, <https://openai.com/index/introducing-codex/>
    - \* codegen [50], [51]
- Different chunking techniques
- GitHub Action
- More sophisticated evaluation methods
- Usage of program slicing in static analysis step
- testing in all of clib
- PRs fixing found bugs
- More extensive comparison with OFG
- ablation study Project-by-project manually
- Token/\$ comparison
- leveraging of existing unit tests

## 447 **7 Discussion**

448 more powerful llms -> better results

449 open source libraries might have been in the training data results for closed source libraries could  
450 be worse this could be mitigated with llm fine-tuning

## 451 8 Conclusion

452 Recap

### 453 8.1 Acknowledgements

454 Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis sagittis posuere ligula sit amet lacinia.  
455 Duis dignissim pellentesque magna, rhoncus congue sapien finibus mollis. Ut eu sem laoreet,  
456 vehicula ipsum in, convallis erat. Vestibulum magna sem, blandit pulvinar augue sit amet, auctor  
457 malesuada sapien. Nullam faucibus leo eget eros hendrerit, non laoreet ipsum lacinia. Curabitur  
458 cursus diam elit, non tempus ante volutpat a. Quisque hendrerit blandit purus non fringilla. Integer  
459 sit amet elit viverra ante dapibus semper. Vestibulum viverra rutrum enim, at luctus enim posuere  
460 eu. Orci varius natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus.

# Bibliography

- [1] V. J. M. Manes, H. Han, C. Han, *et al.* “The Art, Science, and Engineering of Fuzzing: A Survey.” arXiv: [1812.00140](https://arxiv.org/abs/1812.00140) [cs]. (Apr. 7, 2019), [Online]. Available: [http://arxiv.org/abs/1812.00140](https://arxiv.org/abs/1812.00140), pre-published.
- [2] “Heartbleed Bug.” (Mar. 7, 2025), [Online]. Available: <https://heartbleed.com/>.
- [3] C. Meyer and J. Schwenk. “Lessons Learned From Previous SSL/TLS Attacks - A Brief Chronology Of Attacks And Weaknesses.” (2013), [Online]. Available: <https://eprint.iacr.org/2013/049>, pre-published.
- [4] “American fuzzy lop.” (), [Online]. Available: <https://lcamtuf.coredump.cx/afl/>.
- [5] “libFuzzer – a library for coverage-guided fuzz testing. — LLVM 21.0.0git documentation.” (2025), [Online]. Available: <https://llvm.org/docs/LibFuzzer.html>.
- [6] *Google/atheris*, Google, Apr. 9, 2025. [Online]. Available: <https://github.com/google/atheris>.
- [7] A. Vaswani, N. Shazeer, N. Parmar, *et al.* “Attention Is All You Need.” arXiv: [1706.03762](https://arxiv.org/abs/1706.03762) [cs]. (Aug. 1, 2023), [Online]. Available: [http://arxiv.org/abs/1706.03762](https://arxiv.org/abs/1706.03762), pre-published.
- [8] J. Wei, X. Wang, D. Schuurmans, *et al.* “Chain-of-Thought Prompting Elicits Reasoning in Large Language Models.” arXiv: [2201.11903](https://arxiv.org/abs/2201.11903) [cs]. (Jan. 10, 2023), [Online]. Available: [http://arxiv.org/abs/2201.11903](https://arxiv.org/abs/2201.11903), pre-published.
- [9] S. Yao, J. Zhao, D. Yu, *et al.* “ReAct: Synergizing Reasoning and Acting in Language Models.” arXiv: [2210.03629](https://arxiv.org/abs/2210.03629). (Mar. 10, 2023), [Online]. Available: [http://arxiv.org/abs/2210.03629](https://arxiv.org/abs/2210.03629), pre-published.
- [10] S. Yao, D. Yu, J. Zhao, *et al.* “Tree of Thoughts: Deliberate Problem Solving with Large Language Models.” arXiv: [2305.10601](https://arxiv.org/abs/2305.10601) [cs]. (Dec. 3, 2023), [Online]. Available: [http://arxiv.org/abs/2305.10601](https://arxiv.org/abs/2305.10601), pre-published.
- [11] P. Laban, H. Hayashi, Y. Zhou, and J. Neville. “LLMs Get Lost In Multi-Turn Conversation.” arXiv: [2505.06120](https://arxiv.org/abs/2505.06120) [cs]. (May 9, 2025), [Online]. Available: [http://arxiv.org/abs/2505.06120](https://arxiv.org/abs/2505.06120), pre-published.
- [12] N. Perry, M. Srivastava, D. Kumar, and D. Boneh. “Do Users Write More Insecure Code with AI Assistants?” arXiv: [2211.03622](https://arxiv.org/abs/2211.03622). (Dec. 18, 2023), [Online]. Available: [http://arxiv.org/abs/2211.03622](https://arxiv.org/abs/2211.03622), pre-published.
- [13] H. Chase, *LangChain*, Oct. 2022. [Online]. Available: <https://github.com/langchain-ai/langchain>.
- [14] *Langchain-ai/langgraph*, LangChain, May 21, 2025. [Online]. Available: <https://github.com/langchain-ai/langgraph>.
- [15] J. Liu, *LlamaIndex*, Nov. 2022. DOI: [10.5281/zenodo.1234](https://doi.org/10.5281/zenodo.1234). [Online]. Available: [https://github.com/jerryliu/llama\\_index](https://github.com/jerryliu/llama_index).

- [16] O. Khattab, A. Singhvi, P. Maheshwari, *et al.* “DSPy: Compiling Declarative Language Model Calls into Self-Improving Pipelines.” arXiv: [2310.03714](https://arxiv.org/abs/2310.03714) [cs]. (Oct. 5, 2023), [Online]. Available: <http://arxiv.org/abs/2310.03714>, pre-published.
- [17] D. Ganguly, S. Iyengar, V. Chaudhary, and S. Kalyanaraman. “Proof of Thought : Neurosymbolic Program Synthesis allows Robust and Interpretable Reasoning.” arXiv: [2409.17270](https://arxiv.org/abs/2409.17270). (Sep. 25, 2024), [Online]. Available: <http://arxiv.org/abs/2409.17270>, pre-published.
- [18] A. d’Avila Garcez and L. C. Lamb. “Neurosymbolic AI: The 3rd Wave.” arXiv: [2012.05876](https://arxiv.org/abs/2012.05876). (Dec. 16, 2020), [Online]. Available: <http://arxiv.org/abs/2012.05876>, pre-published.
- [19] M. Gaur and A. Sheth. “Building Trustworthy NeuroSymbolic AI Systems: Consistency, Reliability, Explainability, and Safety.” arXiv: [2312.06798](https://arxiv.org/abs/2312.06798). (Dec. 5, 2023), [Online]. Available: <http://arxiv.org/abs/2312.06798>, pre-published.
- [20] G. Grov, J. Halvorsen, M. W. Eckhoff, B. J. Hansen, M. Eian, and V. Mavroeidis. “On the use of neurosymbolic AI for defending against cyber attacks.” arXiv: [2408.04996](https://arxiv.org/abs/2408.04996). (Aug. 9, 2024), [Online]. Available: <http://arxiv.org/abs/2408.04996>, pre-published.
- [21] A. Sheth, K. Roy, and M. Gaur. “Neurosymbolic AI – Why, What, and How.” arXiv: [2305.00813](https://arxiv.org/abs/2305.00813) [cs]. (May 1, 2023), [Online]. Available: <http://arxiv.org/abs/2305.00813>, pre-published.
- [22] D. Tilwani, R. Venkataramanan, and A. P. Sheth. “Neurosymbolic AI approach to Attribution in Large Language Models.” arXiv: [2410.03726](https://arxiv.org/abs/2410.03726). (Sep. 30, 2024), [Online]. Available: <http://arxiv.org/abs/2410.03726>, pre-published.
- [23] OpenAI. “ChatGPT.” (2025), [Online]. Available: <https://chatgpt.com>.
- [24] Z. Li, S. Dutta, and M. Naik. “IRIS: LLM-Assisted Static Analysis for Detecting Security Vulnerabilities.” arXiv: [2405.17238](https://arxiv.org/abs/2405.17238) [cs]. (Apr. 6, 2025), [Online]. Available: <http://arxiv.org/abs/2405.17238>, pre-published.
- [25] Y. Sun, “Automated Generation and Compilation of Fuzz Driver Based on Large Language Models,” in *Proceedings of the 2024 9th International Conference on Cyber Security and Information Engineering*, ser. ICCSIE ’24, New York, NY, USA: Association for Computing Machinery, Dec. 3, 2024, pp. 461–468, ISBN: 979-8-4007-1813-7. DOI: [10.1145/3689236.3689272](https://doi.org/10.1145/3689236.3689272). [Online]. Available: <https://doi.org/10.1145/3689236.3689272>.
- [26] D. Wang, G. Zhou, L. Chen, D. Li, and Y. Miao. “ProphetFuzz: Fully Automated Prediction and Fuzzing of High-Risk Option Combinations with Only Documentation via Large Language Model.” arXiv: [2409.00922](https://arxiv.org/abs/2409.00922) [cs]. (Sep. 1, 2024), [Online]. Available: <http://arxiv.org/abs/2409.00922>, pre-published.
- [27] D. Liu, O. Chang, J. metzman, M. Sablotny, and M. Maruseac, *OSS-fuzz-gen: Automated fuzz target generation*, version <https://github.com/google/oss-fuzz-gen/tree/v1.0>, May 2024. [Online]. Available: <https://github.com/google/oss-fuzz-gen>.
- [28] H. Green and T. Avgerinos, “GraphFuzz: Library API fuzzing with lifetime-aware dataflow graphs,” in *Proceedings of the 44th International Conference on Software Engineering*, Pittsburgh Pennsylvania: ACM, May 21, 2022, pp. 1070–1081. DOI: [10.1145/3510003.3510228](https://doi.org/10.1145/3510003.3510228). [Online]. Available: <https://dl.acm.org/doi/10.1145/3510003.3510228>.



- [29] B. Jeong, J. Jang, H. Yi, *et al.*, “UTopia: Automatic Generation of Fuzz Driver using Unit Tests,” in *2023 IEEE Symposium on Security and Privacy (SP)*, May 2023, pp. 2676–2692. doi: [10.1109/SP46215.2023.10179394](https://doi.org/10.1109/SP46215.2023.10179394). [Online]. Available: <https://ieeexplore.ieee.org/abstract/document/10179394>.
- [30] Y. Deng, C. S. Xia, C. Yang, S. D. Zhang, S. Yang, and L. Zhang. “Large Language Models are Edge-Case Fuzzers: Testing Deep Learning Libraries via FuzzGPT.” arXiv: [2304.02014](https://arxiv.org/abs/2304.02014) [cs]. (Apr. 4, 2023), [Online]. Available: <http://arxiv.org/abs/2304.02014>, pre-published.
- [31] Y. Deng, C. S. Xia, H. Peng, C. Yang, and L. 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*, ser. ISSTA 2023, New York, NY, USA: Association for Computing Machinery, Jul. 13, 2023, pp. 423–435, ISBN: 979-8-4007-0221-1. doi: [10.1145/3597926.3598067](https://doi.org/10.1145/3597926.3598067). [Online]. Available: <https://dl.acm.org/doi/10.1145/3597926.3598067>.
- [32] K. Ispoglou, D. Austin, V. Mohan, and M. Payer, “FuzzGen: Automatic fuzzer generation,” in *29th USENIX Security Symposium (USENIX Security 20)*, 2020, pp. 2271–2287. [Online]. Available: <https://www.usenix.org/conference/usenixsecurity20/presentation/ispoglou>.
- [33] D. Babić, S. Bucur, Y. Chen, *et al.*, “FUDGE: Fuzz driver generation at scale,” in *Proceedings of the 2019 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, Tallinn Estonia: ACM, Aug. 12, 2019, pp. 975–985, ISBN: 978-1-4503-5572-8. doi: [10.1145/3338906.3340456](https://doi.org/10.1145/3338906.3340456). [Online]. Available: <https://dl.acm.org/doi/10.1145/3338906.3340456>.
- [34] C. Cadar, D. Dunbar, and D. Engler, “KLEE: Unassisted and Automatic Generation of High-Coverage Tests for Complex Systems Programs,” presented at the USENIX Symposium on Operating Systems Design and Implementation, Dec. 8, 2008. [Online]. Available: <https://www.semanticscholar.org/paper/KLEE%3A-Unassisted-and-Automatic-Generation-of-Tests-Cadar-Dunbar/0b93657965e506dfbd56fbc1c1d4b9666b1d01c8>.
- [35] “The LLVM Compiler Infrastructure Project.” (2025), [Online]. Available: <https://llvm.org/>.
- [36] A. Arya, O. Chang, J. Metzman, K. Serebryany, and D. Liu, *OSS-Fuzz*, Apr. 8, 2025. [Online]. Available: <https://github.com/google/oss-fuzz>.
- [37] N. Sasirekha, A. Edwin Robert, and M. Hemalatha, “Program Slicing Techniques and its Applications,” *International Journal of Software Engineering & Applications*, vol. 2, no. 3, pp. 50–64, Jul. 31, 2011, ISSN: 09762221. doi: [10.5121/ijsea.2011.2304](https://doi.org/10.5121/ijsea.2011.2304). [Online]. Available: <http://www.airccse.org/journal/ijsea/papers/0711ijsea04.pdf>.
- [38] “OSS-Fuzz Documentation,” OSS-Fuzz. (2025), [Online]. Available: <https://google.github.io/oss-fuzz/>.
- [39] CVE Program. “CVE - CVE-2014-0160.” (2014), [Online]. Available: <https://cve.mitre.org/cgi-bin/cvename.cgi?name=cve-2014-0160>.
- [40] D. Wheeler. “How to Prevent the next Heartbleed.” (2014), [Online]. Available: <https://dwwheeler.com/essays/heartbleed.html>.
- [41] *Google/clusterfuzz*, Google, Apr. 9, 2025. [Online]. Available: <https://github.com/google/clusterfuzz>.
- [42] *Google/fuzztest*, Google, Jul. 10, 2025. [Online]. Available: <https://github.com/google/fuzztest>.

- 577 [43] *Google/honggfuzz*, Google, Jul. 10, 2025. [Online]. Available: <https://github.com/google/honggfuzz>.  
578
- 579 [44] M. Heuse, H. Eißfeldt, A. Fioraldi, and D. Maier, *AFL++*, version 4.00c, Jan. 2022. [Online].  
580 Available: <https://github.com/AFLplusplus/AFLplusplus>.
- 581 [45] D. Liu, J. Metzman, O. Chang, and G. O. S. S. Team. “AI-Powered Fuzzing: Breaking the  
582 Bug Hunting Barrier,” Google Online Security Blog. (Aug. 16, 2023), [Online]. Available:  
583 <https://security.googleblog.com/2023/08/ai-powered-fuzzing-breaking-bug-hunting.html>.
- 584 [46] *OSSF/fuzz-introspector*, Open Source Security Foundation (OpenSSF), Jun. 30, 2025. [Online].  
585 Available: <https://github.com/ossf/fuzz-introspector>.
- 586 [47] L. Thomason, *Leethomason/tinyclang*, Jul. 10, 2025. [Online]. Available: <https://github.com/leethomason/tinyclang>.  
587
- 588 [48] OSS-Fuzz Maintainers. “Introducing LLM-based harness synthesis for unfuzzed projects,” OSS-  
589 Fuzz blog. (May 27, 2024), [Online]. Available: <https://blog.oss-fuzz.com/posts/introducing-llm-based-harness-synthesis-for-unfuzzed-projects/>.  
590
- 591 [49] Y. Lyu, Y. Xie, P. Chen, and H. Chen. “Prompt Fuzzing for Fuzz Driver Generation.” arXiv:  
592 [2312.17677](https://arxiv.org/abs/2312.17677) [cs]. (May 29, 2024), [Online]. Available: <http://arxiv.org/abs/2312.17677>,  
593 pre-published.
- 594 [50] E. Nijkamp, H. Hayashi, C. Xiong, S. Savarese, and Y. Zhou, “CodeGen2: Lessons for training  
595 llms on programming and natural languages,” *ICLR*, 2023.
- 596 [51] E. Nijkamp, B. Pang, H. Hayashi, *et al.*, “CodeGen: An open large language model for code  
597 with multi-turn program synthesis,” *ICLR*, 2023.