

# HOPPER: Interpretative Fuzzing for Libraries

Peng Chen

Tencent Security Big Data Lab  
spinpx@gmail.com

Yuxuan Xie

Tencent Security Big Data Lab  
ryanyxie@tencent.com

Yunlong Lyu

Tencent Security Big Data Lab  
loydlv@tencent.com

Yuxiao Wang

Tencent Security Big Data Lab  
yuxiaowang@tencent.com

Hao Chen

University of California, Davis  
chen@ucdavis.edu

## ABSTRACT

Despite the fact that the state-of-the-art fuzzers can generate inputs efficiently, existing fuzz drivers still can't adequately cover entries in libraries. Most of these fuzz drivers are crafted manually by developers, and their quality depends on the developers' understanding of the code. Existing works have attempted to automate the generation of fuzz drivers by learning API usage from code and execution traces. However, the generated fuzz drivers are limited to a few specific call sequences by the code being learned. To address these challenges, we present HOPPER, which can fuzz libraries without requiring any domain knowledge to craft fuzz drivers. It transforms the problem of library fuzzing into the problem of interpreter fuzzing. The interpreters linked against libraries under test can interpret the inputs that describe arbitrary API usage. To generate semantically correct inputs for the interpreter, HOPPER learns the intra- and inter-API constraints in the libraries and mutates the program with grammar awareness. We implemented HOPPER and evaluated its effectiveness on 11 real-world libraries against manually crafted fuzzers and other automatic solutions. Our results show that HOPPER greatly outperformed the other fuzzers in both code coverage and bug finding, having uncovered 25 previously unknown bugs that other fuzzers couldn't. Moreover, we have demonstrated that the proposed intra- and inter-API constraint learning methods can correctly learn constraints implied by the library and, therefore, significantly improve the fuzzing efficiency. The experiment results indicate that HOPPER is able to explore a vast range of API usages for library fuzzing out of the box.

## 1 INTRODUCTION

Fuzzing is one of the most popular techniques to find software vulnerabilities. Fuzzers give the software a large number of random inputs and observe if unexpected behaviors happen. Though the idea is simple, fuzzing has successfully been applied to testing various applications and found many bugs. In recent years, the techniques have been greatly improved due to the advent of grey-box fuzzing. Coverage-based grey-box fuzzers such as AFL [3] and LibFuzzer [28] mutate inputs to explore deeper program states without requiring knowledge about input format or program specifications. Constraint-based grey-box fuzzers [10, 11, 15, 44, 36, 9], as the state-of-the-art fuzzers, employ constraint-solving techniques to reach code branches that are guarded by complex constraints.

While grey-box fuzzing techniques greatly facilitate the generalization and automation of program fuzzing (i.e., fuzzing techniques

that take ready-to-use binary programs as targets), replicating such success in library fuzzing (i.e., fuzzing techniques that take APIs as targets) is challenging. To use the state-of-the-art fuzzers for libraries, users have to manually craft a fuzz driver that consumes the type-agnostic inputs fed by fuzzers and transforms the byte stream into API arguments. However, writing high-quality fuzz drivers is difficult, as it is time-consuming and requires a deep understanding of the library. Consequently, most of the existing fuzz drivers cover only a small part of library APIs. The APIs, especially the rarely used ones, usually lack adequate testing. Library fuzzing is still struggling to be usable out-of-box and scalable due to the lack of an automated solution.

In library fuzzing, fuzz drivers need to use correct argument types for invoking APIs and satisfy both intra- and inter-API constraints. Otherwise, the fuzz drivers would crash unexpectedly. For instance, to test `ares_send(ares_channel channel, char *qbuf, int qlen, ares_callback callback, void *arg)` from *c-ares*, a fuzz driver should initialize the first argument `channel` with a call of `ares_init(ares_channel *channelptr)` to meet the inter-API constraint, and then set the third argument `qlen` to be the length of the second argument `qbuf` and set the forth argument `callback` to be a non-null function pointer with the type of `ares_callback` to comply with the intra-API constraints. Any violation of these three constraints may result in spurious crashes. However, in practice, information about these constraints is either missing or scattered across library documents or comments, making it hard to collect them in a fully automatic way.

Recently, researchers have proposed learning-based [7, 23, 45, 25] and model-based [19, 37, 22] methods to generate fuzz drivers automatically. Learning-based methods, such as FuzzGen [23], try to learn the correct usage of APIs from existing consumer code. However, this method fails when consumer code is unavailable, such as for new or work-in-progress libraries. Model-based methods, such as GraphFuzz [19], ask users to provide specifications of the APIs under test, which requires domain-specific knowledge and significant human involvement. Furthermore, the quality of the fuzz drivers generated with these methods is largely affected by the external inputs (i.e., consumer code or user-provided expertise), which can be inaccurate or incomplete. For example, in our experiments, we found that some of the fuzz drivers provided by the authors of FuzzGen [23] and GraphFuzz [19] result in spurious crashes due to misuses of the consumer code and incorrect user-defined schemas (Section A.2 in Appendix). The fuzz drivers [1] generated by FuzzGen only cover 5 of 26 API functions in the *libvpx* decoding library and only support the *vp9* codec. Similarly, the specifications [2] of *sqlite3* written by the authors of GraphFuzz do not

even include commonly used API functions, such as `sqlite3_exec` and `sqlite3_complete`.

To address the aforementioned challenges, we present HOPPER to fuzz APIs without requiring external knowledge. Inspired by coverage-based fuzzing, which learns valid format from mutating random seeds [30], HOPPER learns potential usage of the APIs from mutating the composition of API calls and their arguments. If executing the mutated program triggers a new path or a new crash, HOPPER infers the intra- and inter-API constraints based on the dynamic feedback. To achieve this, we introduce a Domain-Specific Language (DSL) that describes arbitrary API usages. The DSL inputs can be interpreted with a lightweight interpreter linked against the library under test. In this way, we transform library fuzzing into interpreter fuzzing. The fuzzer is now responsible for generating programs encoded in the format of DSL to feed into the interpreter. Then the interpreter executes the programs to see if unexpected behavior happens. Thanks to the grammar-aware mutation and the inferred constraints, HOPPER is able to generate valid inputs that explore different API usage while excluding false positive crashes.

We implemented HOPPER and evaluated its effectiveness on 11 real-world libraries against manually crafted fuzzers (MCF) and other automatic solutions (i.e., FuzzGen and GraphFuzz). Table 2 shows that HOPPER greatly outperformed the other fuzzers in both code coverage and bug finding. Notably, HOPPER improves the code coverage over MCFs in *cJSON* by 47.12% in line coverage and 37.50% in branch coverage and achieves higher coverage than the total of 9 fuzzers in *zlib*. In total, HOPPER found 25 new bugs in the libraries, and 17 of them have been confirmed, as shown in Table 4. Moreover, we have demonstrated that the proposed intra- and inter-API constraint learning methods can accurately learn constraints implied by the library, thereby significantly improving the fuzzing efficiency.

## 2 BACKGROUND

### 2.1 Library Fuzzing

Security vulnerabilities in libraries are fatal because they are widely used in various programs. However, testing libraries through program fuzzing alone is often inadequate. Library APIs may be invoked by programs under complex path constraints or with specific arguments that make it difficult to fully exercise the APIs. To address this issue, library fuzzing tools have been developed to specifically test libraries, with LibFuzzer being a common choice for such testing. Here are the required steps for effectively using LibFuzzer to fuzz libraries:

- **Craft a fuzz drivers.** A fuzz driver describes the usage of library APIs, including a sequence of API calls and their arguments. A high-quality fuzz driver should provide an entry to explore as much code in the library as possible. However, specific execution paths are not only determined by the arguments in API calls but also by invoking their related APIs. These related APIs may return values as arguments or affect the context for other APIs that rely on them (e.g., global values). Hence, a deep understanding of the tested library is necessary for thorough testing. But, enumerating all valid API usages would be time-consuming and challenging, so fuzz drivers usually only contain a few common API

**Table 1: Number of unique APIs used in fuzz drivers. The second column is the total number of exported APIs in libraries.**

Library	Total	MCFs	FuzzGen*	GraphFuzz	HOPPER
<i>cJSON</i>	78	6(7.69%)	4(4.13%)	40(51.28%)	78(100%)
<i>c-ares</i>	60	13(21.67%)	-	20(33.33%)	58(96.67%)
<i>libpng</i>	241	25(10.37%)	-	66(27.39%)	233(96.69%)
<i>lcms</i>	283	10(3.53%)	-	38(13.43%)	274(96.82%)
<i>libmagic</i>	18	4(22.22%)	-	10(55.56%)	18(100%)
<i>libpcap</i>	89	8(8.99%)	-	29(32.58%)	73(82.02%)
<i>zlib</i>	84	29(34.52%)	-	41(48.81%)	78(92.86%)
<i>re2</i>	70	35(50.00%)	-	47(67.14%)	69(98.57%)
<i>sqlite3</i>	279	15(5.38%)	-	74(26.52%)	225(80.65%)
<i>libvpx</i>	26	7(26.92%)	5(19.23%)	17(65.38%)	24(92.31%)
<i>libaom</i>	38	5(13.16%)	7(18.42%)	13(34.21%)	35(92.11%)
<b>Average</b>	-	<b>18.58%</b>	<b>13.93%</b>	<b>41.42%</b>	<b>93.52%</b>

\*The authors of FuzzGen released the fuzz drivers for *libvpx* and *libaom*, and the released code of FuzzGen is unable to run on the rest libraries, except for *cJSON*.

usages. As shown in Table 1, the MCFs of the listed 11 popular libraries only covered 18.58% of the APIs, whereas the remaining 81.42% uncovered APIs will escape being fuzzed. Even though some approaches have been proposed to automatically synthesize fuzz drivers, the coverage of APIs is still limited (e.g., GraphFuzz [19] achieved API coverage of 41.42%). Besides, as different API usage has different search spaces for fuzzing, putting all of them in sequence would be inefficient. Instead, it would be more effective to write them into multiple fuzz drivers or conditionally execute them within a fuzz driver. For example, in Figure 2, this fuzz driver calls the parsing function first and then calls different printing functions according to the first 4 bytes of data.

- **Specify the format of input.** The blind byte stream generated by LibFuzzer makes it difficult to create structured inputs that satisfy intra-API constraints. To address this issue, we need to specify the input format and guide the fuzzer to generate arguments beyond byte arrays. FuzzedDataProvider [22] is capable of dividing the fuzz input into multiple parts of various types, while libprotobuf-mutator[37] can generate structured inputs based on provided grammars. However, the presence of intra-API constraints makes defining the potential argument range for fuzzing an enormous task. Hence, developers might opt to encode arguments as literal constants directly into fuzz drivers, such as the second argument of `cJSON_ParseWithOpts` in Figure 2. Unfortunately, this approach may lead to inadequate testing for API functions.

### 2.2 Fuzzing Interpreters

Grammar-aware grey-box fuzzing has succeeded in programs that parse the inputs, especially in interpreters [35, 42, 4, 8, 20, 40, 27, 12]. The success of fuzzing interpreters is attributed to the following two key techniques.

- **Grammar-aware input mutation.** Blindly mutated inputs are likely to be rejected by the parsing procedure, while

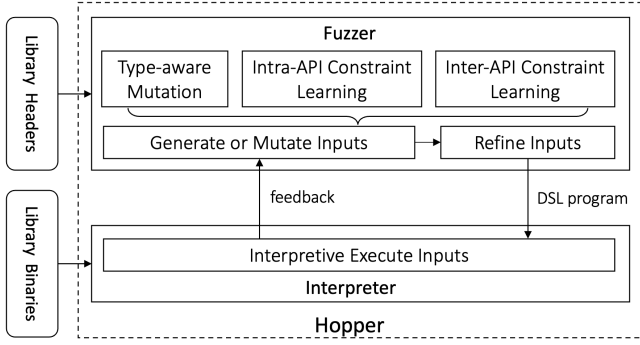


Figure 1: Overview of HOPPER’s architecture.

structured inputs can reach deeper paths. Grammar-aware mutations parse the inputs as intermediate representations (IRs) based on their encoding grammar and mutate the IRs with constraints. For example, Superion [42] uses an abstract syntax tree as the IR and conducts three main operators on the tree for mutation: replace a node in the tree with a newly created sub-tree, splice two different trees, and minimize trees without affecting execution.

- **Coverage guided fuzzing.** Guided by the coverage feedback, the inputs that trigger a new path will be kept, then the fuzzer mutates them further to enter deeper branches or find new bugs. Incorporating with the grammar-aware mutation strategies, the fuzzer is able to effectively synthesize valid test inputs that cover more branches and trigger new bugs.

We observe that constructing fuzz drivers for libraries resembles implementing an interpreter for the inputs, and fuzzing the interpreter is equivalent to fuzzing the library under the hood. However, MCFs interpret type-agnostic byte inputs into limited sequences of API calls, which only partially explore the libraries. If we extend the interpreter to accept any API usage as input, then a grammar-aware coverage-based fuzzer can efficiently generate inputs for the interpreter.

## 3 DESIGN

### 3.1 Overview

HOPPER transforms the problem of library fuzzing into the problem of interpreter fuzzing. At a high level, HOPPER consists of two main components: a grammar-aware fuzzer that generates inputs encoded in a DSL format, and a lightweight interpreter that executes these inputs, as shown in Figure 1.

The fuzzer produces high-quality inputs for invoking library APIs. It first extracts function signatures and type definitions from the library’s header files, which provide valuable information on potential relations between API functions and argument types. By leveraging this information, HOPPER generates a sequence of API calls through random combinations of API functions and arguments. At the early stage, the generated sequences may be invalid and terminate with shallow paths in the APIs’ code. However, thanks to coverage guidance, HOPPER keeps the inputs that explore new branches as seeds and mutate them further. Ultimately, the generated sequences become valid and their execution can reach deeper

code. Completely random fuzzing wastes enormous time in generating illegal and inefficient seeds. To accelerate the process of evolution, HOPPER performs type-aware argument mutation (Section 3.3.2) and learns intra- and inter-API constraints (Section 3.4). Additionally, HOPPER minimizes the inputs (Section 3.3.3) to reduce overhead.

The interpreter consumes the inputs and invokes the library APIs. The libraries under test are linked against the interpreter during the compiling stage. Before linking, HOPPER instruments the binaries to capture internal states during execution, e.g., code coverage. When new input arrives, the interpreter parses the input according to the syntax of DSL and interprets the statements orderly based on their semantics.

### 3.2 DSL and Input Interpretation

Fuzz drivers are usually developed with the same programming language as the library’s specification. Therefore, for compiled language like C/C++, fuzz drivers are crafted and compiled before fuzzing begins. Only the input bytes are changed during each round of fuzzing. To replace the composition of API calls during fuzzing, while avoiding compilation overhead, we introduce a Domain Specific Language (DSL) and a lightweight interpreter in HOPPER for accelerating the entire process.

**3.2.1 DSL.** For the sake of generality, the interpreter in HOPPER takes *DSL programs* as input. The grammar of HOPPER DSL is listed in Figure 6 in the appendix. Each *DSL program* comprises *statements* as its most fundamental components, and each statement has an ascending index that its successors can reference. We categorize common fuzzing behaviors found in MCFs into five *statement types* in our DSL.

- A **load statement** defines the type information and literal representation of a value. Strong typing ensures that input data of a specific type can be directly incorporated into the call site of the APIs in a *DSL program*, eliminating the need for a type-agnostic byte stream.
- A **call statement** invokes a specific API function in libraries by providing the function name and a list of arguments, where each argument refers to the value defined in a previous load statement.
- An **update statement** overwrites the value returned by a call statement at runtime.
- An **assert statement** checks a call’s return value at runtime. If the assertion fails (e.g., the caller needs to dereference the return value but the value is null), the program exits immediately.
- A **file statement** specifies a valid file name for file I/O. If the statement is used for reading, a sequence of random bytes is filled in the file for runtime reading.

To keep the grammar simple, our DSL does not support conditional statements. Some fuzz drivers use conditional statements to choose different API functions or arguments. By contrast, HOPPER generates different DSL programs to enumerate those API functions or arguments. For example, Figure 2 is a real-world fuzz driver written by the developers of *cJSON*, and Figure 3 is a program written in DSL that covers a path in Figure 2.

**3.2.2 Interpreter.** The interpreter parses DSL programs, executes the statements gracefully, and monitors the states of the program after executing each statement. In order to invoke the library APIs, HOPPER links the interpreter against the libraries under test at compiling stage. It also constructs a table that correlates each function's name with its corresponding caller based on the header files of the libraries. During program execution, the caller casts values to the types of arguments it needs and then invokes the referenced function with the arguments. Before linking to the interpreter, HOPPER instruments the library binaries with code that counts branches and hooks compare instructions and resource management functions (e.g., malloc, free and fopen). The interpreter then collects the following feedback at runtime.

- **Optional Branch Tracking.** It is not necessary to track branches of all the library API calls in the DSL program every time, so HOPPER defines a global flag to guard the branch tracking code. The value of the flag is determined by the DSL input. When calling an API function that ends with a question mark (e.g., Line 13 in Figure 3), the interpreter activates the flag for branch tracking.
- **Context-sensitive Code Coverage.** To distinguish the same branches visited by different API functions, the interpreter sets the hash of the current function's name as the context. The instrumentation code reads the context and calculates an exclusive code trace for each API.
- **Overflow Detection.** If the statement loads a variable-sized value (e.g., an array), the interpreter of HOPPER stores them in a memory arena and appends a canary right after it to detect possible buffer overflow.
- **Use-after-free Detection.** The interpreter maintains a set of memory chunks allocated by malloc and released by free through instrumentation. For each pointer used in function arguments, if the memory chunk pointed by this pointer is freed, the interpreter exits the program immediately to avoid the use-after-free issues.
- **Comparison Hooking.** The interpreter collects the parameters used in comparison instructions and functions to guide the fuzzer to solve magic bytes.

The fuzz drivers utilized by LibFuzzer should have no side effects since the code runs in a loop in the same process. However, it is difficult to generate a program that can reset all resources before exiting. HOPPER's interpreter solves this problem by running each input in an individual process. Once a DSL program terminates, the operating system destroys the process and releases all allocated resources for HOPPER. This allows the interpreter to execute DSL programs continuously without the need to release resources.

### 3.3 Grammar-aware Input Fuzzing

To find bugs in libraries rather than interpreters, HOPPER takes a different approach compared to other grammar-aware fuzzers. While other fuzzers traverse all possible combinations of syntaxes based on input grammar, HOPPER instead focuses on generating various effective API calls. This process involves two phases, as shown in Figure 4. First, it generates inputs based on the information available in function signatures to initialize a seed pool. Second, it

```

1 int LLVMFuzzerTestOneInput(const uint8_t* data, size_t
  size) {
2   cJSON *json;
3   size_t offset = 4;
4   unsigned char *copied;
5   char *printed_json = NULL;
6   int minify = data[0] == '1' ? 1 : 0;
7   int require_termination = data[1] == '1' ? 1 : 0;
8   int formatted = data[2] == '1' ? 1 : 0;
9   int buffered = data[3] == '1' ? 1 : 0;
10  if (size <= offset) return 0;
11  json = cJSON_ParseWithOpts((const char*)data + offset,
    NULL, require_termination);
12  if (json == NULL) return 0;
13  if (buffered) {
14    printed_json=cJSON_PrintBuffered(json, 1, formatted);
15  } else {
16    if (formatted) printed_json = cJSON_Print(json);
17    else printed_json = cJSON_PrintUnformatted(json);
18  }
19  if (printed_json != NULL) free(printed_json);
20  if (minify) { ... }
21  cJSON_Delete(json);
22  return 0;
23 }

```

**Figure 2: Fuzz driver written by the developer of cJSON. The code is in cJSON\_read\_fuzzer.c, which is used by OSS-Fuzz [32].**

```

<0> load Vec<char>= vec(32)["
  GxsAAAAAAAAAo9tsrXXoqw57jwAAAAAAAAARNk+1AAA="]
<1> load char* = &<0>
<2> load char** = null
<3> load int = 0
<4> call cJSON_ParseWithOpts (<1>, <2>, <3>)
<5> assert non_null(<4>)
<6> load cJSON = { next: null, prev: null, child: null,
  type_: 8, valuestring: null, valueint: 12345,
  valuedouble: 0.2771, string: null, }
<7> update <4>[0.child] = <6>
<8> load Vec<char> = vec(7)[54, 52, -68, -43, 1, 122, 0]
<9> load char* = &<8>
<10> call cJSON_AddFalseToObject (<4>, <9>)
<11> load int = 1
<12> load int = 0
<13> call cJSON_PrintBuffered ? (<4>, <11>, <12>)

```

**Figure 3: Example program in the format of HOPPER DSL.**

selects inputs in the seed pool and mutates them with the guidance of coverage feedback.

**Pilot Phase.** In the pilot phase, HOPPER infers constraints and draws skeletons of the inputs. Initially, the seed pool contains no inputs. Therefore, HOPPER tries to generate simple inputs for each API function based on their signatures and learn constraints from them (detailed in Section 3.4). To accomplish this, HOPPER selects an API function in the library as a target and attempts to randomly generate a call statement for it. This includes generating arguments

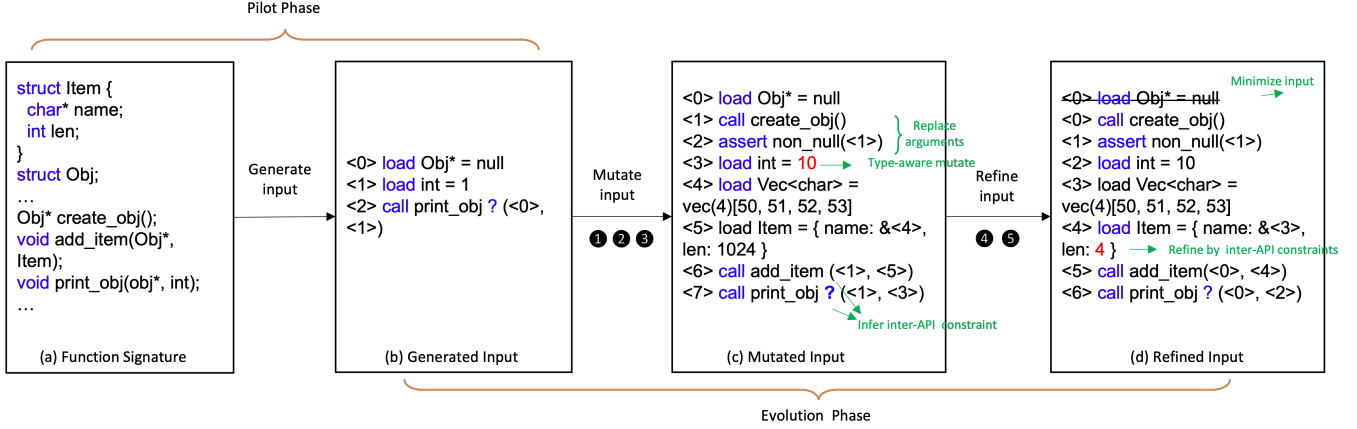


Figure 4: The workflow of input generation in HOPPER.

and inserting related calls that introduce the necessary context. These statements finally form an input, which is then executed by the interpreter. If the input triggers a new path in the libraries without crashing, HOPPER saves it into the seed pool for further mutation. To prevent irrelevant coverage feedback from other calls, HOPPER tracks the coverage of the target call only.

Algorithm 1 describes the procedure of generating a call statement. It generates each argument according to its type in the function signature, and one of the following three operators is used:

- HOPPER chooses an existing statement whose type matches the argument<sup>1</sup>.
- The argument is obtained by inserting a new API invocation. HOPPER selects an API function whose return type matches the argument randomly, and generates a call statement for it recursively. The new call statements are placed ahead of the current one. Additionally, an assertion statement indicating whether the call statement runs successfully is added after the new call statement (e.g., a non-null assertion statement is inserted after a pointer-type returned call statement).
- A load statement with a typed value created from scratch is used.

HOPPER also attempts to affect the execution of the target call by inserting other API invocations to change the internal states of the program. It prioritizes API functions with non-primitive argument types that overlap with those of the target call and reuses the overlapping arguments as much as possible. To prevent the program from becoming overly complex, HOPPER stops generating new calls if the length of statements exceeds a specific threshold or the recursion depth becomes too high. Take the program in Figure 3 as an example. To generate the program with `cJSON_PrintBuffered` (line 13), HOPPER randomly chooses the return value of `cJSON_ParseWithOpts` (line 4) and generate two integer value (line 11 and 12) as its arguments. In addition, it creates a call for `cJSON_AddFalseToObject` that may modify existing arguments (line 10).

<sup>1</sup>HOPPER treats const and non-const types equally when comparing types.

**Evolution Phase.** After running a certain number of rounds, HOPPER enters the evolution phase, which aims to build more complex programs based on the skeleton inputs. To achieve this, HOPPER selects an input from the seed pool and randomly mutates the statements based on their types. Guided by the execution branch coverage, HOPPER keeps the mutated inputs that explore deeper code of libraries, which are more likely to find bugs. Overall, HOPPER mutates program statements in five steps. ① HOPPER selects an input from the seed pool with priority. Fresh seeds are given higher priority to be chosen, as they are more likely to reach deeper paths. ② HOPPER chooses statements to mutate from the input according to their weights. The assert, file, and update statements are weighted as zero and not mutated, while the weights of load and call statements are determined by their complexity. ③ HOPPER mutates the load statements and call statements by their corresponding strategies, as described in Section 3.3.1 and Section 3.3.2 respectively. ④ HOPPER refines the input to use APIs correctly based on the constraints learned during fuzzing, as described in Section 3.4. ⑤ HOPPER minimizes the input to remove redundant statements that have no effect on reaching the path yet increase the search space of mutation, as described in Section 3.3.3.

**3.3.1 Call Statement Mutation.** HOPPER adopts the following mutation strategies for call statements:

- Replace one of the arguments with a new argument that retains the same type (line 4 to line 24 in Algorithm 1).
- Insert a new call before the target call (line 27 to line 29 in Algorithm 1). The inserted call may modify the values of the target call's arguments or change the global states in libraries. HOPPER determines the effectiveness of the inserted calls through branch feedback, which is detailed in Section 3.4.2.
- Update the return value of the call. An update statement is inserted after the call to partially overwrite the return value with a new value.

**3.3.2 Type-aware Value Mutation.** As libraries often use various argument types, including custom composite types, HOPPER needs to generate appropriate types for such arguments when invoking the

**Algorithm 1** Generate a call statement randomly based on its function signature. It stops generating new calls recursively if the length of statements in the *program* exceeds a certain length.

---

```

1: function GENERATECALL(program, sig)
2:   args  $\leftarrow$  empty list
3:   for all arg_type  $\in$  sig.arg_types do
4:     if toss a coin then  $\triangleright$  Uses existing statement as
       argument
5:       index  $\leftarrow$  Choose statement with arg_type in the
       program randomly.
6:       if index exists then
7:         args  $\leftarrow$  args  $\cup$  index
8:         continue
9:       end if
10:    end if
11:    if toss a coin and arg_type is not primitive then  $\triangleright$ 
       Generates a new call statement as argument
12:      f  $\leftarrow$  Randomly choose an API function that returns
       arg_type .
13:      if f exists then
14:        call  $\leftarrow$  GENERATECALL(program, f)
15:        index  $\leftarrow$  insert call into program.
16:        check  $\leftarrow$  GENERATEASSERT(program, index)
17:        insert check into program.
18:        args  $\leftarrow$  args  $\cup$  index
19:        continue
20:      end if
21:    end if
22:    load  $\leftarrow$  GENERATELOAD(program, arg_type)  $\triangleright$ 
       Generates a new load statement as argument
23:    index  $\leftarrow$  insert load into program.
24:    args  $\leftarrow$  args  $\cup$  index
25:  end for
26:  if toss a coin then  $\triangleright$  Generates a new call statement that
       may affect the execution of call
27:    f  $\leftarrow$  Choose an API randomly.  $\triangleright$  Prefers API functions
       whose argument types overlap with args
28:    call  $\leftarrow$  GENERATECALL(program, f)
29:    insert call into program.
30:  end if
31:  call  $\leftarrow$  CREATECALLSTMT(sig, args)
32:  return call
33: end function

```

---

corresponding APIs. The same is true for value mutation, where applying mutations according to the value's type can more effectively explore new states. To achieve this, HOPPER parses type definitions (e.g., struct, enum and union) and type aliases in header files recursively. It then generates new typed values by the following rules:

- **Primitive Types.** Almost all primitive types are numerical types. Thus, HOPPER generates numbers within a small range of uniform distribution. Additionally, HOPPER applies one of four mutations to the primitive values. (1) Set an interesting value (e.g., 0x80000000 for int); (2) Flip a bit or a byte; (3)

Add or subtract a small number; (4) Set to a literal collected from a comparison instruction, if the value is used as one of the operands.

- **Array.** Based on the array's length and element type, HOPPER generates a sequence of elements. If the length is variable, HOPPER first randomly chooses a length. During mutation, HOPPER selects one or more elements in the array and mutates them respectively. In addition, HOPPER can resize the array by inserting or removing elements if its length is not fixed.
- **Structure.** Values with custom structure types are created by recursively generating their fields. When mutating custom structure values, HOPPER randomly selects a field in the structure and mutates it according to its type.
- **Trivial Pointer.** Trivial pointers expose the layout of the pointer type (i.e., primitive-type pointer and transparent structure-type pointer). HOPPER mutates them by the following operations: (1) Set to a null pointer; (2) Point to the location of an existing statement that holds the same pointer type; (3) Point to a newly generated array whose element type is the same as the pointee; (4) Point to the location of the value returned by a newly generated call statement. The call statement returns the same pointer type.
- **Nontrivial Pointer.** Nontrivial pointers, such as opaque pointers, void pointers, and function pointers, cannot be mutated straightforwardly due to their unpredictable nature. Hence, HOPPER constructs them separately. For the opaque pointers, HOPPER retrieves the API functions that initialize the pointers via either returning the pointers or filling them through reference. If void pointers have alias types, HOPPER treats them as opaque pointers as well. Otherwise, HOPPER attempts to cast an arbitrary-length byte array to the required void pointer to see if it works (detailed in Section 3.4.1). For the function pointers, HOPPER synthesizes empty functions with the required signatures at compiling stage to allow the fuzzer to use their addresses as pointers.

Besides, byte arrays may contain data with their own encoding, which is not defined in header files. For example, `cJSON_ParseWithOpt` in Figure 2 parses a byte array with JSON formatting. For these values, HOPPER applies AFL's random mutations that are designed for byte arrays.

**3.3.3 Input Minimization.** Redundant statements and values slow down execution speed and increase search space during mutating, which makes fuzzing inefficient. To address this, HOPPER employs two phases to minimize the inputs:

1. Minimize inputs after mutation and refinement. HOPPER inspects the statement in the input backward, excluding the target call statement. If a statement is no longer referenced by other statements, HOPPER deletes it. For example, in Figure 4, Line 0 in (c) is redundant after mutation and is thus removed by HOPPER in (d).
2. Minimize inputs that trigger new paths. HOPPER removes calls that have no impact on the execution path (detailed in Section 3.4.2), as well as redundant values in load statements. It tries to set the pointer values to null or shrink the

length of arrays when possible. If the execution path remains unchanged, HOPPER retains the mutation in the input.

### 3.4 Constraint Learning

To correctly invoke APIs, DSL programs generated by HOPPER must satisfy both intra- and inter-API constraints. Intra-API constraints dictate that APIs must be invoked with appropriate arguments, while inter-API constraints specify the appropriate order in which APIs are invoked. HOPPER learns these constraints through libraries' runtime feedback rather than relying on external sources.

**3.4.1 Intra-API Constraint.** We propose six general intra-API constraints based on our observations of real-world libraries.

- **Non-null Pointer(NON-NULL).** APIs that do not check for null pointers can crash when invoked with null pointers. It's often unclear whether this is a real bug, as some developers argue it is the user's responsibility to perform null checks.
- **Correct File Name(FILE).** When an API function reads from or writes to a file, the file name provided as an argument must be valid. If the file name is randomly generated, the API call may terminate early, or it could mess up the disk if used as an output stream.
- **Specific Value(EQUAL).** APIs that use a number in the arguments to designate the boundary of an array pointer may suffer from overflow errors if the number is incorrect.
- **Bounded Range(RANGE).** APIs that access or allocate limited resources based on argument numbers may encounter resource exhaustion or overflow errors if the number is out of range.
- **Array Length(ARRAY-LEN).** Some APIs assume that the arrays referenced by the pointers have sufficient elements rather than asking for arguments indicating boundaries. This may result in overflow errors if the arrays don't have adequate elements.
- **Specific Type Cast(CAST).** Due to missing layout information of the void type, developers have to generate objects with concrete types and cast their references to the void pointers.

HOPPER learns intra-API constraints throughout the entire fuzzing process, including both the pilot and evolution phases. Besides inferring constraints for the arguments themselves, HOPPER recursively explores composite structures to infer constraints for any objects contained within, such as the pointees, fields in structs, and elements within arrays.

For void\* arguments without alias types, HOPPER assumes that their pointees do not contain any pointer. Therefore, they can be cast from a big enough random byte array, and HOPPER adds CAST constraints that treat them as char\* type. Conversely, in the case of void\* pointers with alias types, HOPPER supposes that they are opaque pointers and initializes them via invoking other APIs.

When new paths are explored by inputs, HOPPER checks to see if any file open function (e.g., fopen) has been triggered, and compares the file name with the arguments used to invoke the API. If there is a match, a FILE constraint is created for the corresponding argument.

Furthermore, if an input triggers a new crash, HOPPER infers intra-API constraints by the following steps in order.

- (1) If the call triggers a segmentation fault due to accessing a null pointer ( $si\_addr$  is 0 or close to 0, where  $si\_addr$  is the address of the faulting memory reference), HOPPER locates each null pointer in the arguments, sets it to the address of a protected memory chunk, and runs this mutated program again. If the program crashes again at the same program location (indicated by the rip register) and triggers illegal access inside the protected memory chunk, it means the pointer is accessed without a null check in the API invocation. In that case, HOPPER adds a NON-NULL constraint for this pointer.
- (2) If the crash is caused by accessing a canary appended right after an array ( $si\_addr$  is in the range of canary), HOPPER tries to figure out whether there is a length or index of a variable-sized array in the arguments. Firstly, HOPPER locates which array has been overflowed. We denote the array's length as  $N$ . For each numerical value in the call's arguments, HOPPER attempts to set it as  $N - 1$ ,  $N$ , and  $N + 1$ , respectively. If both  $N$  and  $N + 1$  lead to a crash by accessing the canary, HOPPER adds a RANGE constraint to set the value within a range of  $[0, N)$ . If only  $N + 1$  makes a crash, an EQUAL constraint is added to set the value to be the same as the array's length.
- (3) If the above strategies fail to rectify illegal access of a canary, HOPPER attempts to pad the arrays in the arguments to a specific length (e.g.,  $K$ ) to see whether it resolves the crash. If so, an ARRAY-LEN constraint is added to ensure that this array is at least  $K$  bytes.
- (4) For other illegal access, if there are CAST constraints for the arguments, HOPPER tries to mutate the byte array pointed to by the arguments. If the illegal address varies with the mutated bytes, the void pointer may be interpreted as a structure containing pointers. Thus, HOPPER removes the CAST constraints with char\* type.
- (5) If the inputs lead to a timeout or out of memory, HOPPER searches for large numerical values in the arguments and mutates them. If the execution becomes significantly faster or exits normally after setting the value to be small, HOPPER adds a RANGE constraint for the argument to limit its maximal values.

These constraints are used to refine inputs after mutation, and inputs that violate the constraints are excluded from the seed set. As an example, in Figure 4, HOPPER successfully infers the constraint that the value of len should be the length of the array that name refers to, and refines the value accordingly.

**3.4.2 Inter-API Constraint.** An API call can modify its referenced arguments or internal states in the program, which affects the behavior of its subsequent calls. These relationships between API calls are referred to as inter-API constraints. Compared to intra-API constraints, inter-API constraints are more difficult to identify and apply universally. Each library defines its own pattern to coordinate the API functions, making it hard to make abstractions. Instead of imposing concrete constraints, HOPPER preserves the inter-API constraints by saving the effective programs for later mutations. To learn such constraints, HOPPER statically assembles library APIs by



function signatures and then verifies their effectiveness dynamically through coverage feedback.

At first, HOPPER statically infers naive constraints between API functions by analyzing their signatures. By comparing the types of their arguments and returns, we can deduce whether the two API functions are likely related. For any given pair of API functions  $F_1$  and  $F_2$ , they tend to be related if their types overlap in the following ways.

- Type  $T$  is the type of one of the arguments in  $F_1$ , and it is the type of return value of  $F_2$ .
- Both  $F_1$  and  $F_2$  use type  $T$  as their arguments.
- $F_1$  uses an argument with type  $T$  while  $F_2$  uses a pointer  $T^*$  that may modify the value it points to.

Besides the types of arguments, the two API functions are more likely to be related if the identifiers of arguments in the latter two cases are identical.

Two API functions have an effective relation if the former can help the latter to reach a new path. As shown in Algorithm 1, HOPPER attempts to generate a large number of various API usages. The API invocations may be composited in an interesting order after several iterations. In case HOPPER finds a certain sequence of API calls to be effective, the program is reserved for further mutation. By incorporating the precise control of coverage tracking described in Section 3.2.1, HOPPER learns the effective relations between API functions by the following steps.

- (1) When new calls are inserted before the target API call, only the coverage of the target call will be tracked.
- (2) If the insertions trigger a new path, HOPPER deletes the new calls one by one from the input and checks whether the coverage remains the same after removal. If the coverage changes, HOPPER identifies the removed call as critical for reaching the new path. In other words, an efficient inter-API constraint between the two API functions is found in the program.

In addition to recognizing effective call sequences, HOPPER also identifies the effective arguments among sequences of API invocations. Once a program triggers new paths, HOPPER checks whether the new paths were introduced by the mutation of certain arguments. To reuse effective arguments for further mutation, HOPPER saves all the statements responsible for producing those arguments into a cache. These statements contain both concrete values and inter-API constraints. During mutation, HOPPER fetches one effective argument as a substitute for the original one to introduce appropriate context to the API invocations.

## 4 IMPLEMENTATION

We implemented HOPPER in 23164 lines of Rust code and 1285 lines of C/C++ code for instrumentation.

### 4.1 Fuzzer

HOPPER generates inputs with semantics parsed from C header files. Combining Rust’s traits and macros, we implemented the process elegantly. First, we adopted Rust bindgen [31] to automatically generate Rust FFI bindings, including type definitions and function signatures, from library header files. For C++ libraries, only C-style

API declaration is accepted by HOPPER since features like templates are not yet supported in bindgen [17]. Next, we defined *mutate*, *generate*, *serialize*, *deserialize* traits for the types in the bindings, which apply customized behavior to objects of each type for fuzzing. We manually implemented these traits for primitive types since they have common memory layouts in libraries. As for custom structures, we utilized Rust’s *procedural macros* to automatically implement these traits. HOPPER also synthesizes empty functions for function pointer types according to their function signatures. When invoked as callbacks, these functions do nothing other than return an object with zero-initialization. Finally, we generated a table of object builders that HOPPER uses to call the trait implementations. These automatically generated codes are compiled and linked against the fuzzer.

To report the inputs that trigger bugs to developers for analysis, we have built a tool that translates DSL to C source code.

### 4.2 Interpreter

HOPPER instruments the library binaries by static binary rewriting with E9Patch [14]. We borrowed from E9AFL [16] the code for branch tracking and improved it to be controllable and API-sensitive. Additionally, the instrumentation collects the compare instructions and hooks resources management functions (e.g., *malloc*, *free* and *open*).

When executing inputs, HOPPER’s interpreter parses them based on the code that implements *deserialize* traits. To invoke the APIs, HOPPER generates the code of the caller mapping table described in Section 3.2.2, using Rust FFI bindings and procedural macros. This code is then compiled into the interpreter. Like AFL, the interpreter incorporates the techniques of fork servers to reduce the overhead of process initiation. It forks a new process once it receives a new input for interpretation.

During interpretation, HOPPER detects memory overflow by adding a page-size canary after the array values in a memory arena. This memory arena is implemented by *mmap* a continuous memory at a fixed address. In the arena, the address of the last byte of an array is aligned to a page’s last byte. To detect overflow, HOPPER reserves a page size memory after the array and sets the page as unreadable and unwritable by *mprotect*.

## 5 EVALUATION

In this section, we evaluated HOPPER on 11 widely-used real-world libraries. All of the selected libraries are commonly used by various applications and have been evaluated by OSS-Fuzz [32]. To adapt fuzzing on those libraries using LibFuzzer, their developers have crafted several fuzz drivers. Among all the 11 libraries, *re2* is both developed and exported as C++ library in special. For the sake of compatibility, HOPPER fuzzed *re2* through its C wrapper: *cre2* [13]. Additionally, all libraries included in our experiments were the most recent versions available during the time of testing. To demonstrate the effectiveness of HOPPER, the following research questions were answered:

- **RQ1:** How effective is HOPPER in libraries?
- **RQ2:** Can HOPPER correctly infer API constraints?
- **RQ3:** Can HOPPER generate programs that are comparable with MCFs?



All of our experiments are conducted on a server with an Intel Xeon Platinum 8255C and 128 GB memory running 64-bit Ubuntu 20.04 LTS. In each experiment, we configured the tested fuzzer to be executed on only one core.

### 5.1 RQ1: How effective is HOPPER in libraries?

Code coverage and bug finding are two key metrics used to evaluate fuzzing effectiveness [26]. To answer the first research question, we compared HOPPER with MCFs and other automatic solutions (i.e., FuzzGen and GraphFuzz) on 11 popular real-world libraries, using these two metrics. For the MCFs, we either selected fuzz drivers crafted by the library developers or collected them from the Google OSS-Fuzz project. The authors of FuzzGen provided the fuzz drivers for *libvpx* and *libaom* and we used them directly in our evaluation after rectifying the inaccuracies listed in Section A.2. However, FuzzGen is unable to run successfully on the other libraries, except for *cJSON*. Specifically, the fuzz driver that FuzzGen generated for *c-ares* contained invalid arguments and could not be compiled. For the rest of the libraries, FuzzGen encountered errors during fuzz driver generation, such as violating assertions while determining data types of arguments or getting stuck at static backward slicing. GraphFuzz requires manually written schemas to synthesize fuzz drivers. We adopted the author-provided schema for *sqlite3* and wrote the others ourselves. In each individual experiment, we set a timeout of 24 hours and ran the fuzzers. Afterward, we recompiled the tested library with LLVM’s source-based code coverage feature [33] enabled and re-ran the fuzzers with seed inputs to collect coverage information. To reduce statistical errors, we repeated each experiment five times and reported the average results. Additionally, for libraries containing multiple MCFs, we ran each driver with LibFuzzer for 24 hours separately.

**5.1.1 Code Coverage.** Table 2 shows the code coverage achieved by different fuzzers on each library. HOPPER outperformed other fuzzers on all libraries in both lines and branches coverage, except for *re2*, *libvpx*, and *libaom*. Even though 9 MCFs of *zlib* have been tested for a total of 216 hours with Libfuzzer, HOPPER is still better than their overall result. The contrast on *c-ares* is the most significant, where HOPPER generated inputs that covered almost three times as much code as the MCFs’ best effort. On the other hand, the MCFs for *re2* were sophisticated enough to cover most of the source codes, while HOPPER managed to reach a coverage level that is close to its opponent. In the case of *libvpx* and *libaom*, MCFs covered more lines and branches than automatic solutions. This is because they can reach complicated internal states by decoding multiple frames in a loop, while all the automatic solutions are loop-free. However, among the fuzz driver generated automatically, HOPPER explores more lines and branches than FuzzGen and GraphFuzz.

To demonstrate that HOPPER is capable of fuzzing more APIs, we counted the number of unique APIs in the valid programs generated by fuzzers. Table 1 shows that, on average, HOPPER successfully invoked 93.52% of APIs, while MCFs, FuzzGen, and GraphFuzz only cover 18.58%, 13.93%, and 41.42% APIs respectively.

**5.1.2 Bug Finding.** For the crashes triggered by HOPPER, we first eliminated any spurious crashes that violated the inter-API constraints learned by HOPPER, and then sanitized the duplicated crashes

**Table 2: Comparison of coverage between HOPPER, Graph-fuzz, FuzzGen, and MCFs on real world libraries.**

Library	Fuzzer	Lines	Branches
<i>cJSON</i>	<i>cjson_read_fuzzer</i>	952 (42.92%)	473 (46.56%)
	FuzzGen	186 (8.39%)	96 (9.45%)
	GraphFuzz	1346 (60.69%)	612 (60.24%)
	HOPPER	1997 (90.04%)	854 (84.06%)
<i>c-ares</i>	<i>parse_reply_fuzzer</i>	1757 (23.94%)	744 (21.78%)
	<i>create_query_fuzzer</i>	110 ( 1.50%)	47 ( 1.38%)
	Total of MCFs	1865 (25.41%)	790 (23.13%)
	GraphFuzz	2424 (33.02%)	994 (29.10%)
	HOPPER	5012 (67.06%)	2045 (59.03%)
<i>libpng</i>	<i>libpng_read_fuzzer</i>	5005 (28.67%)	2041 (26.42%)
	GraphFuzz	1629 (9.33%)	507 (6.56%)
	HOPPER	9610 (55.04%)	3903 (50.53%)
<i>lcms</i>	<i>IT8_load_fuzzer</i>	641 ( 3.41%)	271 (3.21%)
	<i>overwrite_transform_fuzzer</i>	2964 (15.77%)	1271 (15.07%)
	<i>transform_fuzzer</i>	4382 (23.32%)	1757 (20.84%)
	Total of MCFs	5357 (28.50%)	2205 (26.15%)
	GraphFuzz	2481 (13.20%)	872 (10.34%)
	HOPPER	9001 (47.89%)	3135 (37.18%)
<i>libmagic</i>	<i>magic_fuzzer</i>	3094 (30.87%)	2043 (28.69%)
	GraphFuzz	3197 (31.90%)	2003 (28.10%)
	HOPPER	4230 (45.26%)	2634 (39.51%)
<i>libpcap</i>	<i>fuzz_both</i>	4301 (27.05%)	2346 (32.40%)
	<i>fuzz_filter</i>	5561 (34.97%)	2901 (40.07%)
	<i>fuzz_pcap</i>	974 ( 6.13%)	365 ( 5.04%)
	Total of MCFs	6870 (43.21%)	3479 (48.05%)
	GraphFuzz	1736 (10.92%)	881 (12.17%)
	HOPPER	7536 (47.40%)	3669 (50.68%)
<i>zlib</i>	<i>checksum_fuzzer</i>	251 ( 5.21%)	61 ( 2.12%)
	<i>compress_fuzzer</i>	1880 (38.73%)	904 (31.48%)
	<i>example_dict_fuzzer</i>	1901 (39.17%)	952 (33.15%)
	<i>example_flush_fuzzer</i>	1631 (33.61%)	776 (27.02%)
	<i>example_large_fuzzer</i>	1684 (34.70%)	795 (27.68%)
	<i>example_small_fuzzer</i>	1429 (29.45%)	758 (26.39%)
	<i>minigzip_fuzzer</i>	2227 (45.89%)	1071 (37.29%)
	<i>uncompress2_fuzzer</i>	989 (20.37%)	468 (16.30%)
	<i>uncompress_fuzzer</i>	976 (20.11%)	457 (15.91%)
	Total of MCFs	2976 (61.32%)	1482 (51.60%)
	GraphFuzz	2602 (53.62%)	1343 (46.76%)
	HOPPER	3502 (72.16%)	1914 (66.64%)
<i>re2</i>	<i>re2_fuzzer</i>	6373 (67.85%)	3367 (67.75%)
	GraphFuzz	5564 (59.23%)	2843 (57.21%)
	HOPPER	6413 (68.27%)	3299 (66.38%)
<i>sqlite3</i>	<i>oss_fuzz</i>	22582 (30.96%)	9054 (25.99%)
	GraphFuzz	6261 ( 8.58%)	2172 ( 6.23%)
	HOPPER	25356 (34.76%)	9551 (27.41%)
<i>libvpx</i>	<i>vpx_dec_fuzzer_vp8</i>	3604 (10.44%)	1138 (14.21%)
	<i>vpx_dec_fuzzer_vp9</i>	15654 (45.34%)	3475 (43.41%)
	Total of MCFs	18787 (54.42%)	4463 (55.75%)
	FuzzGen	15211 (44.06%)	3367 (42.05%)
	GraphFuzz	15060 (43.62%)	3251 (40.61%)
	HOPPER	15641 (45.30%)	3514 (43.89%)
<i>libaom</i>	<i>av1_dec_fuzzer</i>	39762 (21.05%)	9722 (17.73%)
	FuzzGen	32576 (17.25%)	7757 (13.35%)
	GraphFuzz	30837 (16.33%)	7327 (12.61%)
	HOPPER	36218 (19.18%)	9228 (15.88%)

with the same stack traces at the program crash point. The remaining crashes were verified manually by inspecting the code, debugging the programs, and reading the official documents. HOPPER ultimately discovered 25 new bugs from 48 unique crashes, as shown in Table 3. In contrast, none of the MCFs, FuzzGen, and GraphFuzz found a valid bug under the same 24 hours running. Of the 48 unique crashes, 23 were identified as spurious crashes as they violated the constraints specified in the documentation or were rejected by library developers. The developers did not classify these crashes as bugs since they have their own criteria for valid API usage. The APIs must be used strictly within specified constraints, which may be mentioned in the documents, or the program may crash unexpectedly. For example, an out-of-bound read error could occur if the `zFunctionName` argument for `sqlite3_create_function16` is not a UTF-16 string<sup>2</sup>.

We have reported all the discovered bugs to the corresponding communities. As of submission, 17 bugs have been confirmed by library developers and the rest are still under review. The details of the discovered bugs are listed in Table 4. Buffer overflow accounts for the most bugs that HOPPER has detected (12 out of 25). Eight of those overflows were detected by HOPPER’s canary when reading or writing out of bounds of arrays allocated as arguments by HOPPER. The remaining four accessed the addresses of internally allocated arrays in the API functions with large offsets generated by HOPPER. Null pointer dereferences are also commonly detected (11 out of 25). These bugs occurred when null pointers were produced or improperly initialized by library calls and accessed without null pointer checking within the code of API functions. Additionally, HOPPER identified an uncontrolled format string bug and a double-free bug, both of which are analyzed in detail in Section 5.3.

Furthermore, almost all buggy APIs discovered by HOPPER could not be covered by MCFs, FuzzGen, or GraphFuzz. This is not surprising as the fuzz drivers for these libraries have been extensively utilized for continuous fuzzing over an extended period. To ensure a fair comparison with MCFs, GraphFuzz, and FuzzGen, we additionally evaluated HOPPER and these fuzzers on previous versions of the 11 libraries. As shown in Table 5, HOPPER surpasses MCFs, GraphFuzz, and FuzzGen by detecting 28 unique bugs in previous versions of the 11 libraries. Moreover, HOPPER was able to identify 7 out of the 14 bugs that other fuzzers had detected. Although HOPPER missed certain bugs reported by other fuzzers due to its inefficient mutation power on those predefined API calls, it prioritized mutating API usages that are more likely to trigger bugs or explore new states via its seed selection, which includes any API usage within the libraries. Consequently, HOPPER discovered a greater number of bugs overall.

**Answer to RQ1:** HOPPER outperformed MCFs and other library fuzzing approaches in terms of both code coverage and bug finding. Specifically, HOPPER detects 25 new bugs and 17 of them have been confirmed.

**Table 3: Unique crashes reported by HOPPER.**

Library	Version	#Programs	Crashes				Accuracy
			#UC	#S	#B	#C	
<i>cJSON</i>	1.7.15	2,972	3	1	2	1	66.67%
<i>c-ares</i>	1.18.1	3,192	2	0	2	2	100%
<i>libpng</i>	1.6.37	5,612	7	1	6	4	85.71%
<i>lcms</i>	2.13.1	2,660	13	8	5	5	38.46%
<i>libmagic</i>	FILE5_42	1,662	0	0	0	0	-
<i>libpcap</i>	1.10.1	2,249	5	2	3	3	60%
<i>zlib</i>	1.2.12	5,598	2	1	1	1	50%
<i>re2</i>	0.4.0	27,355	4	2	2	0	50%
<i>sqlite3</i>	3.38.5	10,356	12	7	4	1	33.33%
<i>libvpx</i>	1.11.0	22,282	0	0	0	0	-
<i>libaom</i>	3.5.0	19,654	1	1	0	0	-
<b>Total</b>	-	<b>103,592</b>	<b>48</b>	<b>23</b>	<b>25</b>	<b>17</b>	<b>52.08%</b>

UC = Total unique crashes; S = Spurious crashes caused by incorrect API usage; B = Valid bugs that identified by manually review; C = Confirmed bugs after reported to library developers.

## 5.2 RQ2: Can HOPPER correctly infer API constraints?

One of the key insights that allows HOPPER to generate API usages both correctly and efficiently is its learning of intra- and inter-API constraints. To evaluate the benefits of these constraints, we compared HOPPER with and without constraints on the 11 libraries. As shown in Table 6, HOPPER with inferred constraints explores more lines and branches across all the libraries. However, without intra-API constraints, code coverage is substantially lower in *libpng*, *libmagic*, *libvpx*, and *libaom*. This is because certain API functions that many other APIs depend on could not be called properly. For example, `png_init_io` in *libpng* initializes a file input for a handler that is subsequently read by other APIs for processing. If the fuzzer randomly generates a string as the file’s name without a FILE constraint, `png_init_io` will return an invalid handler, causing all subsequent calls to terminate at a shallow state. As incorrect API usages often lead to program crashes, we also measured the success rate of execution for the generated inputs. Notably, after imposing intra-API constraints, *cJSON*, *c-ares*, *libvpx*, and *libaom* achieve almost 100% success rate of execution. In *re2*, the success rate improves significantly after HOPPER learns not to generate null pointers as arguments. On the other hand, in the absence of inter-API constraints, HOPPER explores API usage in a less efficient way by combining API functions blindly. This leads to poorer performance in terms of both line and branch coverage, particularly in libraries such as *libpng*, *lcms*, and *sqlite3*, which have a rich set of APIs available and therefore a large number of possible combinations to try out. To accelerate the search process, HOPPER leverages inter-API constraints to pinpoint and reuse effective API usages.

Furthermore, we counted the number of intra-API constraints learned from the 11 libraries and analyzed their correctness. In total, HOPPER learned 973 intra-API constraints. Among them, NON-NULL constraints were the most commonly inferred, as 609 out of the 973 constraints ( 62.01%) pertained to pointer dereferences without null checks in the API code. Additionally, HOPPER was able to infer 17 FILE constraints, 147 EQAUL constraints, 35 RANGE

<sup>2</sup><https://sqlite.org/forum/forumpost/7ace1408b>.

**Table 4: Bugs found by HOPPER in the tested libraries.**

ID	Library	File:line_number	Buggy function	Bug Type	Status	Commit ID	MCF	FZ	GF
1.	cJSON	cJSON.c:2209	cJSON_DetachItemViaPointer	Null pointer dereference	Reported	722(P)	×	×	×
2.	cJSON	cJSON.c:2326	cJSON_ReplaceItemViaPointer	Null pointer dereference	Confirmed	726(P)	×	×	×
3.	c-ares	ares_init.c:2254	ares_set_sortlist	Buffer overflow	Fixed	496(S)	×	-	×
4.	c-ares	ares_init.c:2247	ares_set_sortlist	Buffer overflow	Fixed	496(S)	×	-	×
5.	libpng	pngerror.c:229	png_warning	Buffer overflow	Confirmed	453(S)	×	-	×
6.	libpng	pngtran.c:483	png_set_quantize	Buffer overflow	Confirmed	454(S)	×	-	×
7.	libpng	pngutil.c:3675	png_progressive_combine_row	Buffer overflow	Confirmed	455(S)	×	-	×
8.	libpng	pngutil.c:3675	png_combine_row	Buffer overflow	Reported	456(S)	×	-	×
9.	libpng	pngutil.c:3984	png_read_filter_row_avg	Buffer overflow	Confirmed	457(S)	×	-	×
10.	libpng	pngtran.c:4728	png_do_quantize	Buffer overflow	Reported	458(S)	×	-	×
11.	zlib	deflate.c:401	gzsetparams	Null pointer dereference	Fixed	761(P)	×	-	×
12.	re2	cre2.cpp:195	cre2_find_named_capturing_groups	Null pointer dereference	Reported	30(S)	×	-	×
13.	re2	cre2.cpp:725	cre2_set_match	Null pointer dereference	Reported	29(S)	×	-	×
14.	sqlite3	sqlite3.c:30121	sqlite3_overload_function	Uncontrolled format string	Fixed	bbbbb66b6b(T)	×	-	×
15.	sqlite3	sqlite3.c:79786	sqlite3_value_bytes	Null pointer dereference	Reported	0218d74c47(T)	×	-	×
16.	sqlite3	sqlite3.c:79786	sqlite3_value_bytes16	Null pointer dereference	Reported	0218d74c47(T)	×	-	×
17.	sqlite3	sqlite3.c:85266	sqlite3_value_subtype	Null pointer dereference	Reported	0218d74c47(T)	×	-	×
18.	lcms	cmsio1.c:857	cmsIsCLUT	Buffer overflow	Fixed	350(S)	×	-	×
19.	lcms	msgamma.c:852	cmsBuildTabulatedToneCurveFloat	Buffer overflow	Fixed	351(S)	×	-	×
20.	lcms	cmsnamed.c:760	cmsGetPostScriptCRD	Null pointer dereference	Fixed	353(S)	×	-	×
21.	lcms	cmslut.c:416	cmsStageAllocMatrix	Buffer overflow	Fixed	354(S)	×	-	×
22.	lcms	cmscgats.c:1928	cmsIT8SaveToMem	Buffer overflow	Fixed	355(S)	×	-	×
23.	libpcap	pcap.c:2946	pcap_breakloop	Null pointer dereference	Fixed	1147(P)	×	-	×
24.	libpcap	pcap.c:493	pcap_can_set_rfmom	Null pointer dereference	Fixed	1147(P)	×	-	×
25.	libpcap	pcap-linux.c:835	pcap_activate	Double Free	Fixed	1098(S)	×	-	×

**File:line\_number** : The source file location of the crash point. **Buggy function** : The function that cause the crash.

**Commit ID** : The bug trace id that committed to community developers. **P**=Pull request number, **S**=Issue number, **T**=Bug forum id.

**MCF, FZ, GF** : Whether the bug could be discovered by MCFs, FuzzGen (FZ), and GraphFuzz (GF) under the same library version.

**Table 5: Comparison with other fuzzers on previous versions of libraries.**

Library	Version	MCF	GraphFuzz	FuzzGen	MCF ∪ GF	Hopper
cJSON	1.7.0	1	0	0	1	4(1)
c-ares	1.16.0	1	0	-	1	4(1)
libpng	1.6.32	0	0	-	0	4
lcms	2.8	2	1	-	2	6(1)
libmagic	FILE5_25	1	1	-	1	1(1)
libpcap	1.9.0	4	0	-	4	3(3)
zlib	1.2.11	0	0	-	0	1
sqlite3	3.22.0	1	0	-	1	3
libaom	1.0.0	3	1	0	4	2
<b>Total</b>	-	13	3	0	14	28(7)

Numbers in brackets show the number of bugs that were reported by multiple fuzzers including HOPPER.

constraints, 110 CAST constraints, and 55 ARRAY-LEN constraints across these libraries. The overall precision and recall of the learned constraints are 96.51% and 97.61%, respectively. Most of the false-positive constraints learned by HOPPER are actually approximations of the ground truth constraints and work well in most cases. These

constraints, although not completely correct, do help to reduce the invalid search space and prevent spurious crashes. For example, in `cmsStageAllocCLut16bitGranular`, setting the fourth argument `outputChan` to be the length of the fifth argument `table`, although incorrect, worked effectively in most cases. Regarding the recall of the learned constraints, we identified false-negative constraints by analyzing the spurious crashes reported by HOPPER since collecting the entire ground truth is extremely labor-consuming. Moreover, most of these missed constraints are specific to a particular library and hard to describe in a general way. Typically, a spurious crash indicates a violation of certain unlearned constraints. As an example from `sqlite3` described in Section 5.1.2, HOPPER cannot ensure that the generated buffers are strictly UTF-16 encoded, which leads to a spurious crash.

**Answer to RQ2:** HOPPER has the ability to learn intra-API constraints with high precision (96.51%) and recall (97.61%), while also being capable of speeding up the search process during fuzzing through the use of inter-API constraints.

**Table 6: Impact of inter- and intra-API constraints on HOPPER's efficiency.**

Library	Learned intra-API Constraints										Success Rate		Line Coverage			Branch Coverage		
	#N	#F	#E	#R	#C	#A	TTL	#FN	PRC	RCL	w/o Intra	HOPPER	w/o Intra	w/o Inter	HOPPER	w/o Intra	w/o Inter	HOPPER
cJSON	15	0	10	4	1	0	30	1	100%	96.77%	99.4291%	99.9999%	80.93%	88.86%	90.04%	75.00%	82.48%	84.06%
c-ares	87	0	26(1)	2	16(1)	2	133(2)	0	98.50%	100%	91.4087%	99.9999%	47.78%	65.31%	67.06%	43.91%	56.90%	59.63%
libpng	34(1)	2	20	5(1)	14	16(6)	91(8)	1	91.21%	98.81%	98.8973%	99.9917%	12.03%	39.05%	55.04%	7.39%	36.46%	50.53%
lcms	233	6	25(6)	1(1)	25(3)	13(3)	303(13)	8	95.71%	97.32%	95.3159%	99.2425%	46.65%	28.37%	47.89%	34.54%	21.05%	37.18%
libmagic	3	1	2(1)	0	3	1	10(1)	0	90.00%	100%	99.4833%	99.9935%	13.16%	42.20%	45.26%	9.96%	36.08%	39.51%
libpcap	59	4	1(1)	4	1	4(2)	73(3)	2	95.89%	97.22%	94.4091%	99.9797%	40.92%	40.89%	47.40%	43.56%	43.14%	50.68%
zlib	6	1	20	3	6	6	42	1	100%	97.67%	97.0134%	99.9865%	52.94%	59.28%	72.16%	44.46%	48.96%	66.64%
re2	93	0	26	0	0	4	123	2	100%	98.40%	86.1034%	99.9427%	58.21%	63.72%	68.27%	54.81%	61.91%	66.38%
sqlite3	72	3	8	10(2)	38(2)	9(1)	140(5)	7	96.43%	95.07%	99.3814%	99.5527%	30.52%	29.44%	34.76%	23.72%	23.34%	27.41%
libvpx	3	0	2	1(1)	4	0	10(1)	0	90.00%	100%	98.8642%	99.9999%	1.14%	42.51%	45.30%	3.35%	42.61%	43.89%
libaom	4	0	7	5(1)	2	0	18(1)	1	94.44%	94.44%	94.2413%	99.9999%	1.73%	16.67%	19.18%	3.12%	13.32%	15.88%
<b>Total</b>	<b>609(1)</b>	<b>17</b>	<b>147(9)</b>	<b>35(6)</b>	<b>110(6)</b>	<b>55(12)</b>	<b>973(34)</b>	<b>23</b>	<b>96.51%</b>	<b>97.61%</b>	<b>90.4769%</b>	<b>99.8808%</b>	<b>35.09%</b>	<b>46.94%</b>	<b>53.85%</b>	<b>31.26%</b>	<b>42.39</b>	<b>49.25%</b>

N = NON-NULL; F = FILE; E = EQUAL; R = RANGE; C = CAST; A = ARRAY-LEN; TTL = Total number; FN = False negative; PRC = Precision; RCL = Recall.  
Numbers in brackets show the number of false positive constraints.

**(a) An uncontrolled format string bug found in *sqlite3*.**

```

<0> load sqlite3 * = null //db
<1> load Vec<char> = vec(2)[96, 0] //file_buf
<2> file option <1> //filename
<3> load sqlite3 ** = &<0> //ppDb
<4> call relative: sqlite3_open ? (<2>, <3>)
<5> assert non_null(<0>)
<6> load Vec<char> = vec(22)["JSFuMPRkm/RVAABAp16+
    vy6ib8NjAA=="] //zFuncName, decoded as "%!n0JU@^."
<7> load char* = &<6>
<8> load i32 = 55767322 //nArg
<9> call target: sqlite3_overload_function ? (<0>, <7>, <8>)

```

**(b) A double free bug found in *libpcap*.**

```

<0> load char * = null //device
<1> load char * = null //errbuf
<2> call pcap_create ? (<0>, <1>) //p
<3> assert non_null(<2>)
<4> load pcap_t * = <2> //p
<5> load int = 128 //buffer_size
<6> call relative: pcap_set_buffer_size ? (<4>, <5>)
<7> load int = 15 //snaplen
<8> call relative: pcap_set_snaplen ? (<4>, <7>)
<9> load int = 1833782204 //immediate
<10> call relative: pcap_set_immediate_mode ? (<4>, <9>)
<11> call target: pcap_activate ? (<4>)

```

**Figure 5: Example programs that trigger bugs in *sqlite3* and *libpcap* respectively.**

### 5.3 RQ3: Can HOPPER generate programs that are comparable with MCFs?

To demonstrate the ability of HOPPER for generating high-quality fuzz drivers, in Figure 5, we select two representative programs generated by HOPPER in our experiment as case studies.

**Case 1.** Figure 5a shows an uncontrolled format string bug found in *sqlite3* (ID 14 in Table 4). In *sqlite3\_overload\_function*,

*zFuncName* is passed to *sqlite3\_mprintf* to clone a string, but at worst, *sqlite3\_mprintf* would perform formatting if the string contains specifiers. HOPPER crafted this program successfully through the following steps. In the beginning, it generated a simple program that passes a null *sqlite3* pointer and a random *zFuncName* string to *sqlite3\_overload\_function*. However, since the API function requires a non-null *sqlite3*, HOPPER used intra-API constraints to initialize the pointer by invoking *sqlite3\_open*. This made the program reach the code that calls *sqlite3\_mprintf*. Shortly thereafter, HOPPER randomly mutates the strings for *zFuncName* and fortunately produces a '%' symbol in the string. This case indicates that HOPPER is able to synthesize valid code to invoke an API function with intra-API constraints and find bugs using type-aware mutation for arguments.

**Case 2.** Figure 5b shows a double-free bug found in *libpcap* (ID 25 in Table 4). This bug occurs when *pcap\_activate* fails to enable memory-mapped capture due to *mmap* with zero length under the immediate mode. In such a case, a memory location would be released twice in the error handling code. The length is computed from the combination of *buffer\_size* and *snaplen*. Even worse, the *pcap\_activate* is a commonly used API function for activating a packet capture handle rather than those rarely-used ones. *Tcpdump*, as a popular Linux application, is able to invoke those *libpcap* APIs sequentially by a simple command: `tcpdump -i any -B 1 -s 15 --immediate-mode`. By utilizing inter-API constraints, HOPPER carefully crafted a sequence of API calls that triggered the bug. Specifically, it inferred that *pcap\_set\_buffer\_size*, *pcap\_set\_snaplen*, and *pcap\_set\_immediate\_mode* would modify *pcap\_t* pointer to change the behavior of *pcap\_activate*. Then, HOPPER mutated the integers fed into the API calls and triggered the bug finally. The case demonstrates that HOPPER can explore different API usages by learning inter-API constraints, which traditional fuzz testing tools typically ignore.

**Answer to RQ3:** With grammar-aware input fuzzing, HOPPER synthesizes programs that satisfy both intra- and inter-API constraints. As compared to MCFs, the programs generated by HOPPER can explore a much broader range of API usages.

## 6 DISCUSSIONS

### 6.1 Multiple-dimensional Search Space in Input Generation

Traditional library fuzzing generates a byte array as input and leaves the job of constructing arguments to fuzz drivers. In contrast, in HOPPER, the search space for input generation is multidimensional as it involves both API functions and arguments, which poses a significant challenge. Furthermore, each argument has its own encoding format and requires specific mutation strategies. Although HOPPER mitigates this issue by implementing novel techniques such as constraint learning and type-aware mutation, there is still much room for improvement.

### 6.2 Compatibility with C++ Libraries

Currently, HOPPER only supports fuzzing libraries with C-style header files. The use of templates in C++ headers delays the compilation of template functions to the time of instantiation by the users, thus making it challenging for HOPPER to generate callers of C++ functions and their arguments. Moreover, it is non-trivial to decide the concrete types of template parameters to instantiate the templates. To enable compatibility with C++, a more generalized implementation of generation and mutation is required. We plan to address this in future work.

### 6.3 False Positive Crashes

Although HOPPER is highly effective at inferring common constraints to filter out most spurious crashes, the remaining crashes may still be false positives since the APIs need to be used in specific ways. Learning these constraints through dynamic feedback can be challenging as they have no universal criteria, as discussed in Section 5.2. However, during fuzzing, HOPPER will no longer generate input for APIs that have failed to learn constraints and have a high probability of crashing spuriously. To make HOPPER more practical and user-friendly, we plan to add warnings for users about unlearned constraints and provide a convenient way for them to add these custom constraints themselves.

## 7 RELATED WORK

### 7.1 API Fuzzing

APIs have long been favored as fuzzing targets, given their role in allowing code modules to interact with others. For example, RESTler generates requests for RESTful APIs using a grammar automatically inferred from Swagger specification, and employs coarse-grained feedback from service responses to guide mutations to reach deeper service states [5]. Similarly, Pythia fuzz stateful REST APIs using coverage-guided feedback and learning-based mutations [6]. In the realm of kernel fuzzers, Syzkaller [39] is a coverage-guided fuzzer that generates system call sequences using system API descriptions. To further improve the efficiency and effectiveness of kernel fuzzing, MoonShine [29] generates seed inputs by extracting system call traces from the execution of real programs, while Healer [38] uses dynamic analysis to learn the relationships between system calls and uses that knowledge to guide input generation.

Besides the adoption of fuzzing in service and system APIs, OSS-Fuzz [32] gathers hundreds of manually written fuzzers to continuously fuzz APIs of open source libraries. In recent years, there has also been a trend towards generating fuzz drivers for library APIs automatically [7, 23, 45, 25, 19]. FUDGE [7] and FuzzGen [23] extract the code of API usage from practical code to create fuzz drivers. APICRAFT [45] and WINNIE [25] record the API call sequences from the execution trace of existing consumer programs and combine them to generate fuzz drivers. GraphFuzz [19] views fuzz drivers as dataflow graphs and performs graph-based mutation with a manually specified schema. UTOPIA [24] statically analyzes the library to identify attributes of API arguments and automatically synthesize valid fuzz drivers from existing unit tests. In contrast to generating fuzz drivers in limited domains, HOPPER is a fuzzer for library APIs without any specific domain knowledge, including any knowledge from practical examples and schema specifications. HOPPER searches the combinations of APIs and arguments in a large space and learns the features of high-quality inputs simultaneously through runtime feedback.

### 7.2 Grammar-aware Fuzzing

Grammar-aware fuzzers leverage grammar to generate structured inputs that bypass syntax checking at the beginning of the program execution [21, 40, 27, 46, 42, 4, 34]. For example, LangFuzz [21] and IFuzzer [40] utilize the syntax of JavaScript language to fuzz Javascript interpreters, while SQLsmith [34] and SQUIRREL [46] generate SQL queries for testing DBMSs based on SQL grammar. Superion [42] and NAUTILUS [4] improve grammar-aware fuzzing by combining code coverage guidance. These tools manipulate input as an AST and mutate it according to coverage feedback.

As writing grammar rules require much human effort, some fuzzers try to automatically learn the grammar [41, 18, 8, 43]. Skyfire [41] uses probabilistic modeling to learn grammar from inputs, while Learn&Fuzz [18] employs a recurrent neural network model. Nevertheless, the accuracy of the learned grammar is dependent on the quality of the corpus provided. Alternatively, GRIMOIRE [8] automatically infers the structural properties of input language based on code coverage feedback. Similarly, Profuzzer [43] probes the types of input bytes through per-byte mutations.

General grammar-based fuzzers are intended to test parsers or interpreters by exploring all feasible combinations based on input grammar. However, HOPPER operates differently by focusing on generating various effective API calls rather than exhaustively exploring the input language's grammar.

## 8 CONCLUSION

In this paper, we present HOPPER, a novel fuzzer that aims to fuzz libraries without any domain knowledge required in crafting fuzz drivers. HOPPER links the libraries under test against an interpreter, which takes DSL programs as input and drives libraries to perform requested fuzzing behavior. To generate effective API calls in the format of DSL, HOPPER learns intra- and inter-API constraints in the libraries and mutates the inputs with grammar awareness. We evaluated the effectiveness of HOPPER on 11 real-world libraries. HOPPER outperformed MCFs and the other automatic solutions in both code coverage and bug finding. Specifically, HOPPER found

25 new bugs that others could not find. Our experimental results demonstrate that HOPPER effectively explores a vast range of API usages for library fuzzing out of the box.

## REFERENCES

- [1] *A fuzz driver generated by FuzzGen for libvpx*. URL: <https://github.com/HexHive/FuzzGen/tree/master/examples/libvpx>.
- [2] *A schema for generating sqlite3 fuzz drivers written by GraphFuzz*. URL: <https://github.com/hgarrereyn/GraphFuzz/blob/master/experiments/sqlite3/in/f1/schema.yaml>.
- [3] *American fuzzy lop*. URL: <http://lcamtuf.coredump.cx/afl/>.
- [4] Cornelius Aschermann et al. “NAUTILUS: Fishing for Deep Bugs with Grammars.” In: *Proceedings of the 2019 Network and Distributed System Security Symposium (NDSS)*. 2019.
- [5] Vaggelis Atlidakis, Patrice Godefroid, and Marina Polishchuk. “Restler: Stateful rest api fuzzing”. In: *2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE)*. IEEE. 2019, pp. 748–758.
- [6] Vaggelis Atlidakis et al. “Pythia: grammar-based fuzzing of REST APIs with coverage-guided feedback and learning-based mutations”. In: *arXiv preprint arXiv:2005.11498* (2020).
- [7] Domagoj Babić 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*. 2019, pp. 975–985.
- [8] Tim Blazytko et al. “GRIMOIRE: Synthesizing structure while fuzzing”. In: *28th USENIX Security Symposium (USENIX Security 19)*. 2019, pp. 1985–2002.
- [9] Ju Chen et al. “JIGSAW: Efficient and Scalable Path Constraints Fuzzing”. In: *2022 IEEE Symposium on Security and Privacy (SP)*. IEEE Computer Society. 2022, pp. 1531–1531.
- [10] Peng Chen and Hao Chen. “Angora: efficient fuzzing by principled search”. In: *IEEE Symposium on Security and Privacy (S&P)*. San Francisco, CA, May 2018.
- [11] Peng Chen, Jianzhong Liu, and Hao Chen. “Matryoshka: fuzzing deeply nested branches”. In: *ACM Conference on Computer and Communications Security (CCS)*. London, UK, Nov. 2019.
- [12] Yongheng Chen et al. “One engine to fuzz'em all: Generic language processor testing with semantic validation”. In: *2021 IEEE Symposium on Security and Privacy (SP)*. IEEE. 2021, pp. 642–658.
- [13] *cre2*. URL: <https://github.com/marcomaggi/cre2/>.
- [14] Gregory J Duck, Xiang Gao, and Abhik Roychoudhury. “Binary rewriting without control flow recovery”. In: *Proceedings of the 41st ACM SIGPLAN Conference on Programming Language Design and Implementation*. 2020, pp. 151–163.
- [15] Shuitao Gan et al. “GREYONE: Data Flow Sensitive Fuzzing”. In: *USENIX Security Symposium (USENIX Security)*. Boston, MA, Aug. 2020.
- [16] Xiang Gao, Gregory J Duck, and Abhik Roychoudhury. “Scalable Fuzzing of Program Binaries with E9AFL”. In: *2021 36th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE. 2021, pp. 1247–1251.
- [17] *Generating Bindings to C++*. URL: <https://rust-lang.github.io/rust-bindgen/cpp.html>.
- [18] Patrice Godefroid, Hila Peleg, and Rishabh Singh. “Learn&fuzz: Machine learning for input fuzzing”. In: *2017 32nd IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE. 2017, pp. 50–59.
- [19] Harrison Green and Thanassis Avgerinos. “GraphFuzz: Library API Fuzzing with Lifetime-aware Dataflow Graphs”. In: *2022 IEEE/ACM 44th International Conference on Software Engineering (ICSE)*. 2022, pp. 1070–1081.
- [20] Samuel Groß. “Fuzzil: Coverage guided fuzzing for javascript engines”. In: *Department of Informatics, Karlsruhe Institute of Technology* (2018).
- [21] Christian Holler, Kim Herzig, and Andreas Zeller. “Fuzzing with code fragments”. In: *21st USENIX Security Symposium (USENIX Security 12)*. 2012, pp. 445–458.
- [22] *How To Split A Fuzzer-Generated Input Into Several*. URL: <https://github.com/google/fuzzing/blob/master/docs/split-inputs.md>.
- [23] Kyriakos Ispoglou et al. “FuzzGen: Automatic Fuzzer Generation”. In: *29th USENIX Security Symposium (USENIX Security 20)*. 2020, pp. 2271–2287.
- [24] Bokdeuk Jeong et al. “UTOPIA: Automatic Generation of Fuzz Driver using Unit Tests”. In: *2023 IEEE Symposium on Security and Privacy (SP)*. IEEE Computer Society. 2022, pp. 746–762.
- [25] Jinho Jung et al. “Winnie: Fuzzing windows applications with harness synthesis and fast cloning”. In: *Proceedings of the 2021 Network and Distributed System Security Symposium (NDSS 2021)*. 2021.
- [26] George Klees et al. “Evaluating fuzz testing”. In: *Proceedings of the 2018 ACM SIGSAC conference on computer and communications security*. 2018, pp. 2123–2138.
- [27] Suyoung Lee et al. “Montage: A Neural Network Language Model-Guided JavaScript Engine Fuzzer”. In: *29th USENIX Security Symposium (USENIX Security 20)*. 2020, pp. 2613–2630.
- [28] *libFuzzer – a library for coverage-guided fuzz testing*. URL: <https://llvm.org/docs/LibFuzzer.html>.
- [29] Shankara Pailoor, Andrew Aday, and Suman Jana. “MoonShine: Optimizing OS Fuzzer Seed Selection with Trace Distillation”. In: *27th USENIX Security Symposium (USENIX Security 18)*. 2018, pp. 729–743.
- [30] *Pulling JPEGs out of thin air*. URL: <https://lcamtuf.blogspot.com/2014/11/pulling-jpegs-out-of-thin-air.html>.
- [31] *rust-bindgen*. URL: <https://github.com/rust-lang/rust-bindgen>.
- [32] Kostya Serebryany. “OSS-Fuzz-Google’s continuous fuzzing service for open source software”. In: *USENIX Association*, 2017.
- [33] *Source-based Code Coverage*. URL: <https://clang.llvm.org/docs/SourceBasedCodeCoverage.html>.
- [34] *SQLsmith*. URL: <https://github.com/anse1/sqlsmith>.
- [35] Prashast Srivastava and Mathias Payer. “Gramatron: Effective grammar-aware fuzzing”. In: *Proceedings of the 30th ACM SIGSOFT International Symposium on Software Testing and Analysis*. 2021, pp. 244–256.

- [36] Nick Stephens et al. “Driller: augmenting fuzzing through selective symbolic execution”. In: *Proceedings of the Network and Distributed System Security Symposium*. 2016.
- [37] *Structure-Aware Fuzzing with libFuzzer*. URL: <https://github.com/google/fuzzing/blob/master/docs/structure-aware-fuzzing.md>.
- [38] Hao Sun et al. “HEALER: Relation Learning Guided Kernel Fuzzing”. In: *Proceedings of the ACM SIGOPS 28th Symposium on Operating Systems Principles*. 2021, pp. 344–358.
- [39] *syzkaller - kernel fuzzer*. URL: <https://github.com/google/syzkaller>.
- [40] Spandan Veggam et al. “Ifuzzer: An evolutionary interpreter fuzzer using genetic programming”. In: *European Symposium on Research in Computer Security*. Springer. 2016, pp. 581–601.
- [41] Junjie Wang et al. “Skyfire: Data-driven seed generation for fuzzing”. In: *2017 IEEE Symposium on Security and Privacy (SP)*. IEEE. 2017, pp. 579–594.
- [42] Junjie Wang et al. “Superion: Grammar-aware greybox fuzzing”. In: *2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE)*. IEEE. 2019, pp. 724–735.
- [43] Wei You et al. “Profuzzer: On-the-fly input type probing for better zero-day vulnerability discovery”. In: *2019 IEEE symposium on security and privacy (SP)*. IEEE. 2019, pp. 769–786.
- [44] Insu Yun et al. “QSYM: A practical concolic execution engine tailored for hybrid fuzzing”. In: *27th USENIX Security Symposium (USENIX Security 18)*. 2018, pp. 745–761.
- [45] Cen Zhang et al. “APICraft: Fuzz Driver Generation for Closed-source SDK Libraries”. In: *30th USENIX Security Symposium (USENIX Security 21)*. 2021, pp. 2811–2828.
- [46] Rui Zhong et al. “Squirrel: Testing database management systems with language validity and coverage feedback”. In: *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security*. 2020, pp. 955–970.

## A APPENDIX

### A.1 DSL Grammar

As shown in Figure 6, an individual DSL Program consists of a sequence of lines. The index at the beginning of each line is a unique ID (e.g., an incremental number), which can be referenced by other statements. Since the values in arguments can have various types, our DSL encodes them accordingly. Primitive types are easy to convert to and from strings. In the case of composite data types, such as an array or custom struct, we serialize their internal elements, following a structured order similar to the JSON format. Specifically, when dealing with long lists of primitive values, we optimize them for efficiency by encoding them in Base64 (as seen in line 0 of Figure 3). However, serializing a pointer is tricky since it may point to a value shared by multiple objects. Therefore, we only serialize the destination to which it points in the DSL program, without including any additional information.

### A.2 Examples of inaccurate fuzz drivers

Figure 7 shows a fuzz driver for *libvpx* generated automatically by FuzzGen. On Line 415, the value of `ctx_hEP_1` is an uninitialized

```

Program := Line | Line; Program
Line := Index Statement
Statement := Call | Load | File | Assert | Update
Call := call Name(ArgList)
Load := load Type = Value
Update := update Index[Fields] = Index
Assert := assert Rule
File := file write | file read Index
ArgList := Index | Index, ArgList
Rule := non_null(Index) | eq(Index, Index)
Index := <numeric literal>
Name := function name
Type := type name of value
Value := serialized value
Fields := path to locate a value inside struct or array

```

Figure 6: Grammar of HOPPER’s DSL.

```

415 struct vpx_codec_ctx *ctx_hEP_1; // = &ctx_hEP_0;
416
417 // Dependence family #3 Definition
418 dep_3 = (struct vpx_codec_ctx *)ctx_hEP_1;
419 // initializing argument 'cfg_Ywn'
420 struct vpx_codec_dec_cfg cfg_Ywn_0;
421
422 *(uint32_t*)((uint64_t)&cfg_Ywn_0 + 0) = (E.eat1() & 0x3f)
    + 1; /* UNKNOWN */
423 *(uint32_t*)((uint64_t)&cfg_Ywn_0 + 4) = 0; /* UNKNOWN */
424 *(uint32_t*)((uint64_t)&cfg_Ywn_0 + 8) = 0; /* UNKNOWN */
425 struct vpx_codec_dec_cfg *cfg_Ywn_1 = &cfg_Ywn_0;
426
427 if (vpx_codec_dec_init_ver(dep_3, dep_6, cfg_Ywn_1, 0, 12))
    { /* vertex #4 */
428     return 0;
429 }

```

Figure 7: An example of misuse of consumer code in *libvpx*’s fuzz driver [1] generated by FuzzGen.

pointer, which is then used directly as an argument in the call on Line 427. In addition, upon examining the fuzz driver for *libaom* that was released by the authors of FuzzGen, we found an incompatibility with the newest version of *libaom*. Specifically, the initialization of `aom_codec_dec_cfg` objects in the driver overwrites 17-20 bytes despite the size of the object being reduced to only 16 bytes.

The schema [2] for *sqlite3* is written by the authors of GraphFuzz. When the program fails to invoke `sqlite3_prepare_v2` on line 72, it calls `close_all` to release all `sqlite3*` pointers, but other types of allocated resources, such as `sqlite3_str*` pointers may still be leaked. Moreover, the schema does not verify whether `sqlite3*` is initialized successfully or not in `new_database`.