

# Precise Concolic Unit Testing of C Programs with Alarm Filtering Using Symbolic Calling Contexts

Yunho Kim

KAIST

Daejeon, South Korea  
yunho.kim03@gmail.com

Yunja Choi

Kyungpook National University

Daegu, South Korea  
yuchoi76@knu.ac.kr

Moonzoo Kim

KAIST

Daejeon, South Korea  
moonzoo@cs.kaist.ac.kr

## ABSTRACT

Automated unit testing techniques can reduce manual effort to write unit test drivers/stubs and generate unit test inputs. However, automatically generated unit test drivers/stubs may raise many false alarms because they often over-approximate real context of a target function  $f$  and allow infeasible unit executions of  $f$ . To alleviate this limitation, we have developed a concolic unit testing technique CONBRIO. CONBRIO generates a symbolic unit test driver of  $f$  using *closely relevant* functions to  $f$ . Also, CONBRIO filters out a false alarm by checking feasibility of a corresponding symbolic execution path with regard to  $f$ 's *symbolic calling contexts* obtained by combining symbolic execution paths of  $f$ 's closely related predecessor functions.

In the experiments targeting the crash bugs of the 15 real-world C programs, CONBRIO shows both high bug detection ability (i.e. 91.0% of the target bugs detected) and precision (i.e. a true to false alarm ratio is 1:4.5). Furthermore, CONBRIO detects 14 new crash bugs in the latest versions of the 9 target C programs studied in other papers on crash bug detection techniques.

### ACM Reference Format:

Yunho Kim, Yunja Choi, and Moonzoo Kim. 2017. Precise Concolic Unit Testing of C Programs with Alarm Filtering Using Symbolic Calling Contexts. In *Proceedings of ACM Conference, Washington, DC, USA, July 2017 (Conference'17)*, 12 pages.  
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

## 1 INTRODUCTION

Although unit testing is effective to detect software bugs, field engineers often have burden of manually creating test drivers/stubs and test inputs for each target unit. To reduce manual effort to generate test inputs, automated test generation techniques have been developed (e.g. concolic testing have been applied to detect bugs in open source programs [1–3, 26, 30, 40] and industrial projects [5, 18, 24, 32, 45] through system-level testing). Also, to reduce manual effort to generate test drivers/stubs of target units, automated unit testing techniques have been developed to detect bugs in open source programs [9, 14, 36] and large industrial software [27].

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

Conference'17, July 2017, Washington, DC, USA

© 2017 Association for Computing Machinery.  
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00  
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

A main drawback of the automated unit testing techniques is a large number of *false alarms* raised by *infeasible unit executions* (i.e. unit executions that are infeasible at system-level). Infeasible unit executions are generated due to inaccurate unit test drivers/stubs that over-approximate real contexts of a target unit (Sect. 2.4). This false alarm problem is a serious obstacle to apply automated unit testing in practice since field engineers would not like to spend time to manually filter out many false alarms.

To overcome this limitation, we have developed an automated concolic unit testing framework CONBRIO (CONcolic unit testing with symBolic alaRm filtering using symbolic calling cOntexts) which operates in the following two stages:

1. CONBRIO generates a *symbolic unit test driver* for a target function  $f$ . To provide realistic context to  $f$ , a unit test driver of  $f$  contains *closely relevant* functions to  $f$ . The relevance of a function  $g$  to  $f$  is measured by the degree of dependency of  $f$  on  $g$  (Sect.3.2). Then, CONBRIO performs concolic execution of a unit test driver of  $f$ .
2. To filter out false alarms by checking feasibility of a corresponding symbolic unit execution of  $f$ , CONBRIO generates *symbolic calling context* of  $f$  by combining symbolic paths of *closely relevant* predecessor functions of  $f$  in a static call graph.

As a result, CONBRIO detects bugs effectively and precisely because it enforces various and realistic executions of  $f$  through concolic execution of  $f$  with  $f$ 's realistic contexts (i.e. with the functions closely relevant to  $f$ ) and accurately filters out false alarms using  $f$ 's symbolic calling contexts.<sup>1</sup>

Note that it is important to construct a unit test driver and calling context formulas of  $f$  to contain *only* functions *closely relevant* to  $f$  since including more functions will enlarge symbolic search space and degrade unit testing effectiveness and efficiency. For example, at one extreme end, a unit test driver may contain all successor functions of  $f$  and fail to detect bugs due to too large symbolic search space to explore. Also, symbolic calling context of  $f$  may contain symbolic execution paths of all predecessor functions of  $f$  up to main and fail to detect bugs. This is because the context become too strong to satisfy with a symbolic unit execution of  $f$  that raises an alarm (Sect. 3.5).<sup>2</sup>

<sup>1</sup> CONBRIO may not be applicable at the very early stage of SW development when  $f$ 's relevant functions are not yet available. However, considering long period of SW development cycles including many revisions and releases, this limitation seems not critical.

<sup>2</sup> This is because symbolic calling context may *not* represent all feasible calling context of  $f$  due to the limitation of symbolic execution for complex C code (e.g. pointer operations, loops and/or external binary libraries). Thus, if a symbolic calling context becomes more strong by adding (i.e. via logical  $\wedge$ ) symbolic execution paths of more predecessor functions of  $f$ , the context easily becomes unsatisfiable. See Sect. 3.5 for more detail.

We have applied CONBRIO to the 15 real-world C programs in SIR [11] and SPEC2006 [42] benchmarks and CONBRIO shows both high bug detection ability (i.e. 91.0% of all target bugs detected) and high bug detection precision (i.e. a true to false alarm ratio is 1:4.5) which is more precise than the latest concolic unit testing techniques for C programs (e.g. 1:5.7 by UC-KLEE [36]). Also, CONBRIO detects 14 *new* bugs in the latest versions of the nine target C programs studied in other papers on crash bug detection techniques.

The contributions of this paper are as follows:

- CONBRIO achieves both high bug detection ability (91.0% of the target bugs detected) and high precision (false alarm ratio is 1:4.5) based on the two core ideas: 1) building and utilizing contexts of a target function explicitly based on relevance of functions measured by a function dependency metric, 2) a new alarm filtering strategy that constructs symbolic calling contexts compositionally and utilizes them to check feasibility of a violating unit execution.
- The extensive empirical evaluation on both bug detection ability and precision of CONBRIO and the other concolic unit testing techniques on the 15 real-world C programs supports researchers and practitioners to learn the pros and cons of the related techniques (Sect. 4–5).
- By applying CONBRIO, we have detected and reported 14 new crash bugs in the latest versions of the 9 target programs that were studied in other papers on crash bug detection techniques (Sect. 5.5).
- We have made the real-world crash bug data of the C benchmark programs publicly available, which were collected and organized after examining the bug reports of the last 12–24 years (<https://sites.google.com/view/conbrio-icse2018/>), so that researchers can use them for various testing research purposes (Sect. 4.2.1).

The remainder of the paper is as follows. Section 2 explains the background of automated concolic unit testing. Section 3 describes the detail of CONBRIO. Section 4 explains the experiment setup to evaluate CONBRIO compared to other techniques. Section 5 reports the experiment results. Section 6 discusses related work and Section 7 concludes the paper with future work.

## 2 BACKGROUND

### 2.1 Preliminary

Unit testing uses *drivers* and *stubs* (or mock objects) to test a target function in isolation (i.e. without the rest of a target program). Suppose that a target function under test  $f$  takes  $n$  arguments  $a_1, \dots, a_n$  and accesses  $m$  global variables  $v_1, \dots, v_m$ , and directly calls  $l$  other functions  $g_1, \dots, g_l$ . To enforce diverse test executions of  $f$ , a tester develops various unit test drivers  $drv_i^f$ s each of which generates argument values  $a_1^i, \dots, a_n^i$ , global variable values  $v_1^i, \dots, v_m^i$  and finally invokes  $f$  with these input values. Also, a tester builds stub functions  $s_{g_1}^i, \dots, s_{g_l}^i$  to replace  $g_1, \dots, g_l$ . Furthermore, test drivers/stubs should satisfy the constraints on the interface between  $f$  and the rest of a target program to avoid infeasible unit test executions of  $f$ .

### 2.2 Concolic Unit Test Driver/Stub Generation

For each target function  $f$ , a concolic unit testing technique automatically generates *symbolic* stubs and a *symbolic* unit test driver. Symbolic stubs simply return symbolic values (without updating global variables and output parameters for simplicity) and a symbolic driver invokes  $f$  after assigning symbolic values to the input variables of  $f$  according to their types as follows:<sup>3</sup>

- *primitive types*: primitive variables are directly assigned with primitive symbolic values of the corresponding types.
- *array types*: each array element is assigned with a symbolic variable according to the type of the array element (for a large array, only the first  $n$  elements are assigned with symbolic values where  $n$  is given by a user).
- *pointer types*: for a pointer variable  $ptr$  pointing to a variable of a type  $T$ , a driver allocates memory whose size is equal to the size of  $T$  and assigns the address of the allocated memory to  $ptr$  (i.e.  $ptr = \text{malloc}(\text{sizeof}(T))$ ). Then, a driver assigns  $*ptr$  with a symbolic value of type  $T$ . If the size of  $T$  is not known (e.g. `FILE` in standard C library), `NULL` is assigned to  $ptr$ . If there exists a pointer variable  $ptr2$  pointing to a symbolic variable of the same type  $T$ , a driver assigns  $ptr2$  to  $ptr$ .
- *structure types*: a unit test driver specifies all fields of struct variable  $s$  as symbolic variables recursively (i.e. if  $s$  contains struct variable  $t$ , a unit test driver specifies the fields of  $t$  as symbolic too).

A limitation of this approach is that the drivers and stubs often over-approximate the real environment of  $f$  and allow *infeasible unit executions* (i.e. executions of  $f$  which are *not* feasible at system-level) that may raise *false alarms*.

### 2.3 Insertion of Assertions Targeting Crash Bugs

Concolic unit testing techniques aim to detect crashes/run-time failures such as null-pointer dereference (NPD), array index out-of-bounds (OOB), and divide-by-zero (DBZ) as well as violations of user-given assertions. They often focus on crashes because user-given assertions are usually not available.

Concolic unit testing techniques insert `assert( $exp$ )` into  $f$  where  $exp$  specifies a condition to avoid crashes (e.g. *denominator*  $\neq 0$  to avoid DBZ). Because of `assert( $exp$ )` in  $f$ , concolic testing tries to generate a test input with which  $f$  makes  $exp$  false and increases a chance to detect crash bugs.

### 2.4 Example of False Alarm

Figure 1 shows a target program with a target function  $f$  under test (lines 10–16). `main` calls `a1` if the first parameter  $x$  of `main` is greater than 0 or calls `a2`, otherwise (line 3). `a1` and `a2` call `b` at line 5 and line 6, respectively, and `b` calls  $f$  at line 7.  $f$  takes an integer parameter  $x$  and calls  $g(x)$  (line 12) (a sanity check function for accessing array through an index  $x$ ) and  $h(x)$  (line 15). A concolic unit testing technique generates a unit test driver `driver_f` and symbolic stubs `stub_g` and `stub_h` for  $f$ . Also, it modifies  $f$  to call `stub_g` and `stub_h` instead of  $g$  and  $h$  respectively (see the

<sup>3</sup>This subsection is excerpted from [27].

```

01:// x and y are inputs of a target program
02:int main(int x,int y){
03: return (x>0) ? a1(x,y) : a2(y);}
04:
05:int a1(int x, int y){if(y>0) return b(x); else return 0;}
06:int a2(int x){if(x>0) return b(x);}
07:int b(int x){if(x>0) return f(x);else return 0;}
08:
09:// Target function under test
10:int f(int x){
11: int array[5] = {1,3,5,7,9}, result;
12: if (g(x) != 0){ //=> if (stub_g(x) != 0) {
13: // => assert(0<=x && x<5);
14: result = array[x];
15: }else result=h(x);//=> else result=stub_h(x);
16: return result;}
17:
18:int g(int x){ return (x<5)? 1:0;}
19:
20:int h(int x){ return x + 2;}

```

**Figure 1: Target program with a target function f**

```

01: int driver_f(){
02: int arg1 = SYM_int();
03: f(arg1);}
04:
05: int stub_g(int x){
06: int ret = SYM_int();
07: return ret;}
08:
09: int stub_h(int x){
10: int ret = SYM_int();
11: return ret;}

```

**Figure 2: Generated unit test driver and stubs for f**

comments at line 12 and line 15) and inserts an OOB assertion at line 13.

Figure 2 shows a unit test driver and stubs for  $f$ .  $driver\_f$  invokes  $f$  with a symbolic argument  $arg1$  (lines 2–3) where  $int\ arg1 = SYM\_int()$  sets  $arg1$  as a symbolic integer value (line 2).  $stub\_g$  and  $stub\_h$  return symbolic integer values as  $g$  and  $h$  return integer values (lines 5–7 and lines 9–11 respectively). Concolic execution of  $driver\_f$  violates the OOB assertion at line 13 of  $f$  if a unit test execution satisfies the following two conditions:

- a symbolic argument  $arg1$  to  $f$  (line 3 of  $driver\_f$ ) is larger than or equal to the size of array (e.g.  $arg1$  is 5)
- $stub\_g$  returns a non-zero value (e.g. 1)

However, an alarm raised in such unit test execution is a *false alarm* because such unit test execution of  $f$  is *infeasible* with the real target program where  $g$  (not  $stub\_g$ ) is invoked ( $g$  returns 0 if  $arg1 \geq 5$  (line 18 of Figure 1) unlike  $stub\_g$ ). In other words, a concolic unit testing technique can raise a false alarm if it generates unit test drivers/stubs different from real environment of  $f$  which consist of  $main$ ,  $a1$ ,  $a2$ ,  $b$ ,  $g$ , and  $h$ .

### 3 CONBRIO TECHNIQUE

Figure 3 shows the overall process of CONBRIO in the following steps:

1. CONBRIO receives source code of a target program, a list of target functions to test, and system test cases of the target program as

inputs. CONBRIO obtains function call profiles from the system test executions (Section 3.1).

2. It checks function relevance by calculating *dependency* of a target function  $f$  on other function  $g$  using conditional probability  $p(g|f)$  based on the observed function call profiles (Section 3.2). With a given dependency threshold  $\tau$ , we consider  $f$  has a *high dependency* on  $g$  if  $p(g|f) \geq \tau$ .
3. Based on the calculated dependency of  $f$  on other functions,
  - It generates a unit test driver that contains  $f$ ,  $f$ 's successor functions in a static function call graph on which  $f$  has high dependency, and symbolic stubs.
  - It identifies *calling contexts* of  $f$  each of which is a maximal call path  $a^1 \rightarrow a^2 \rightarrow \dots \rightarrow f$  in a static function call graph such that  $f$  has high dependency on all  $a^i$ s.
4. CONBRIO applies concolic testing to a unit test driver to explore diverse and realistic target unit test executions. During concolic execution, it records a symbolic path formula  $\sigma_{fv_i}$  that violates a given assertion  $v_i$  in  $f$  (Section 3.4).
5. It filters out an alarm raised at  $v_i$  by checking the feasibility of  $\sigma_{fv_i}$  with regard to  $f$ 's *calling contexts* (see Step 3). For this purpose, CONBRIO constructs  $f$ 's *calling context formulas*  $\Sigma_{ctx(f,k)}$  and uses a SMT solver to check satisfiability of  $\sigma_{fv_i}$  (see Step 4) conjuncted with  $\Sigma_{ctx(f,k)}$  (Section 3.5).

If the result is UNSAT for all calling contexts (i.e. there exists *no* feasible execution in any calling context of  $f$  to make  $\sigma_{fv_i}$  feasible), a target alarm is considered as false and ignored. Otherwise (i.e. the result is SAT with at least one calling context), a corresponding alarm is reported as a violation of  $v_i$  in  $f$ .

#### 3.1 Obtaining Function Call Profile from System Test Executions

CONBRIO executes a target program with given system test cases and obtains function call profiles. For example, suppose that a target program in Figure 1 has three system test cases  $(x,y)$  to  $main$ :  $(-1,1)$ ,  $(1,1)$ ,  $(5,1)$ . Then, the function call profiles are obtained as follows:  $\{main \rightarrow a2, a2 \rightarrow b, b \rightarrow f, f \rightarrow g\}$  with  $(-1,1)$ ,  $\{main \rightarrow a1, a1 \rightarrow b, b \rightarrow f, f \rightarrow g\}$  with  $(1,1)$ , and  $\{main \rightarrow a1, a1 \rightarrow b, b \rightarrow f, f \rightarrow g, f \rightarrow h\}$  with  $(5,1)$ .

#### 3.2 Computing Dependency of a Target Function on Other Functions

Suppose that a program has a target function  $f$  and other function  $g$  and it has  $n_f$  system test executions that invokes  $f$ . Based on the observed function call profiles in the system executions, we compute *dependency* of  $f$  on  $g$  as *conditional probability*  $p(g|f)$ . Given a static call graph  $G(V, E)$  (see Definition 1) and system test executions, we compute  $p(g|f)$  as follows:

- Case 1: for  $g$  which is a *predecessor* of  $f$  in  $G(V, E)$ ,  $p(g|f)$  is calculated as  $\frac{n_1}{n_f}$  where  $n_1$  is a number of system executions where  $g$  calls  $f$  directly or transitively
- Case 2: for  $g$  which is a *successor* of  $f$  in  $G(V, E)$ ,  $p(g|f)$  is calculated as  $\frac{n_2}{n_f}$  where  $n_2$  is a number of system executions where  $f$  calls  $g$  directly or transitively

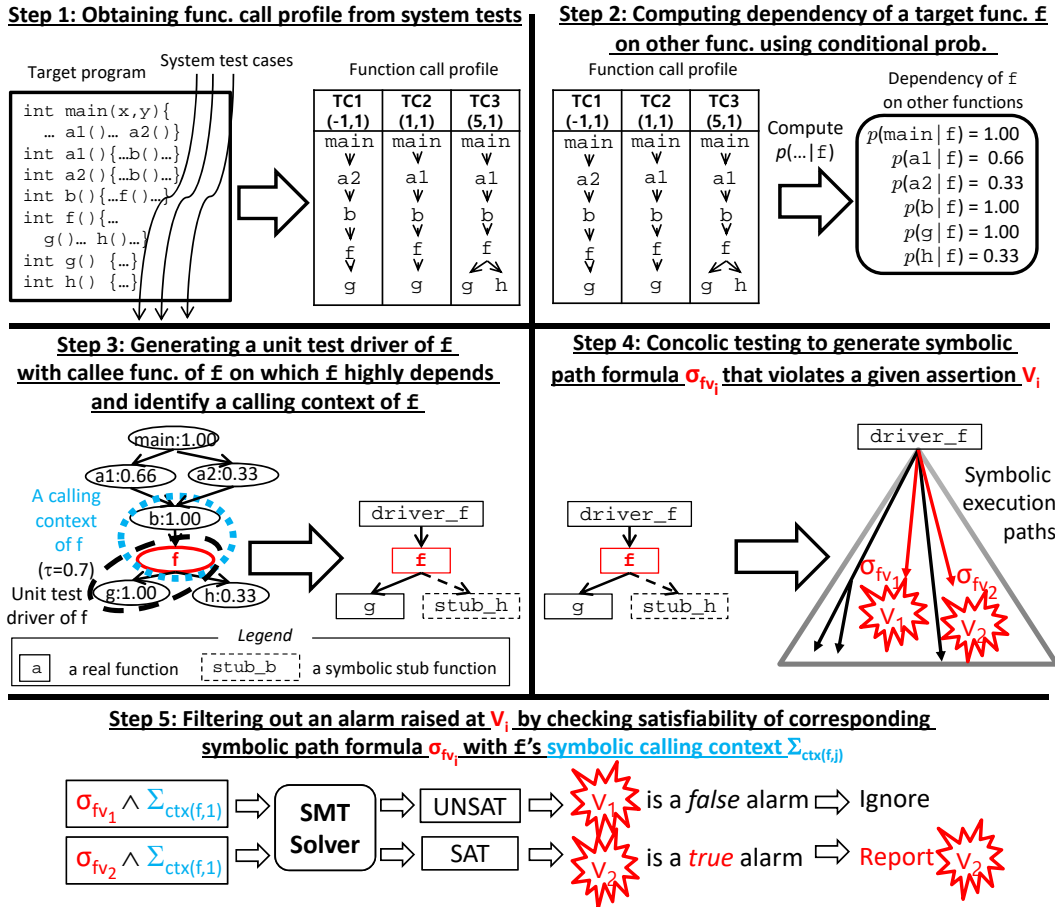


Figure 3: Overall process of CONBRIO

- Case 3: for  $g$  which is a successor and predecessor of  $f$  in  $G(V, E)$  (i.e. there exists a recursive call cycle between  $f$  and  $g$ ),  $p(g|f)$  is calculated as  $\frac{n_3}{n_f}$  where  $n_3$  is a number of system executions where  $f$  calls  $g$  or  $g$  calls  $f$  directly or transitively

For example, Step 1 of Figure 3 shows three test cases (-1,1), (1,1), and (5,1) and their corresponding function call profiles for the program in Figure 1. Based on the profiles, we calculated dependency of  $f$  on other functions as follows (where  $n_f$  is three since we have three test cases):

- $p(\text{main}|f) = 1.00 (= \frac{n_1}{n_f} = \frac{3}{3})$
- $p(a1|f) = 0.66 (= \frac{n_1}{n_f} = \frac{2}{3})$
- $p(a2|f) = 0.33 (= \frac{n_1}{n_f} = \frac{1}{3})$
- $p(b|f) = 1.00 (= \frac{n_1}{n_f} = \frac{3}{3})$
- $p(g|f) = 1.00 (= \frac{n_2}{n_f} = \frac{3}{3})$
- $p(h|f) = 0.33 (= \frac{n_2}{n_f} = \frac{1}{3})$

### 3.3 Generating Unit Test Drivers and Calling Contexts

Given a static call graph  $G(V, E)$  of a target program (Definition 1), a target function  $f$ 's dependency on other functions (i.e.  $p(g|f)$ ),

and a dependency threshold  $\tau$ , CONBRIO generates a unit test driver and calling contexts of  $f$ .

DEFINITION 1. A static call graph  $G(V, E)$  is a directed graph where  $V$  is a set of nodes representing functions in a program and  $E$  is a relation  $V \times V$ . Each edge  $(a, b) \in E$  indicates that  $a$  directly calls  $b$ . We call a node  $p$  as a predecessor of  $f$  if there exists a path from  $p$  to  $f$ . We call a node  $s$  as a successor of  $f$  if there exists a path from  $f$  to  $s$ .

For example, Step 3 of Figure 3 shows how CONBRIO generates a unit test driver and calling context of  $f$  for a program in Figure 1 which contains  $\text{main}$ ,  $a1$ ,  $a2$ ,  $b$ , a target function  $f$ ,  $g$ , and  $h$ . Given a static call graph whose nodes are labelled with dependency of  $f$ , CONBRIO generates a unit test driver of  $f$  (i.e.  $\text{driver}_f$ ) that contains  $f$  and  $g$ , and  $\text{stub}_h$  (instead of  $h$ ) since  $f$  has high dependency on  $g$ , but not  $h$  (i.e.  $p(g|f) = 1.00 \geq \tau$  but  $p(h|f) = 0.33 < \tau$  where  $\tau = 0.7$ ). Finally,  $\text{driver}_f$  invokes  $f$  with symbolic inputs. Note that CONBRIO does not raise a false alarm in this example unlike concolic unit testing in Sect. 2.4 because CONBRIO generates a unit test driver providing realistic environment to  $f$  by using  $g$  which is closely relevant to  $f$ . In addition, CONBRIO determines a calling context of  $f$  as  $b \rightarrow f$  since  $f$  has high dependency on  $b$ , but not  $a1$  nor  $a2$  (i.e.  $p(b|f) = 1.00 \geq \tau$  but  $p(a1|f), p(a2|f) < \tau$ ).

**3.3.1 Unit Test Driver.** For each target function  $f$ , CONBRIO generates a unit test driver that contains  $f$ ,  $f$ 's successor functions  $g$  such that  $f$  has high dependency on all function nodes in a call path from  $f$  to  $g$  in a static function call graph (i.e. for all nodes  $n_i$  between  $f$  and  $g$ ,  $p(n_i|f) \geq \tau$ ), and symbolic stubs for the other reachable functions from  $f$  (i.e. functions on which  $f$  has low dependency). The generated unit test driver sets all arguments of  $f$  and all global variables accessed by the unit test driver as symbolic inputs as described in Sect. 2 and invokes  $f$ .

For example, Figure 4 shows a static call graph whose nodes are labeled with dependency of  $f$ . Figure 4 shows that a *unit test driver* of a target function  $f$  (marked with black dashed line at the bottom of the figure) contains code of  $n12$ ,  $n13$ , and  $n14$  functions on which  $f$  has high dependency (i.e.  $p(n12|f), p(n13|f), p(n14|f) \geq \tau$ ). However, the unit test driver does not contain code of  $n15$  and  $n16$  functions, but symbolic stubs for those functions because  $p(n15|f), p(n16|f) < \tau$ .

In addition, as a false alarm reduction heuristic, CONBRIO adds  $\text{SYM\_assume}(expr)$ <sup>4</sup> at the beginning of  $f$  in a unit test driver where  $expr$  represents possible value ranges of symbolic input variables of  $f$  (which are obtained by applying static value range analyzer [35] to an entire target program code). If an input value is not in the estimated range, a current test execution immediately terminates without raising any alarms and CONBRIO continues a next test execution. As another heuristic, CONBRIO generates unit test drivers to keep consistency between a pointer input variable to dynamically allocated memory and its size variable by figuring out such relation between input variables based on the variable names.

**3.3.2 Calling Contexts.** In a static call graph  $G(V, E)$  labelled with dependency of  $f$ , we define a *calling context* of  $f$  as a maximal call path from a predecessor node of  $f$  (saying  $n_1$ ) to  $f$  as follows.

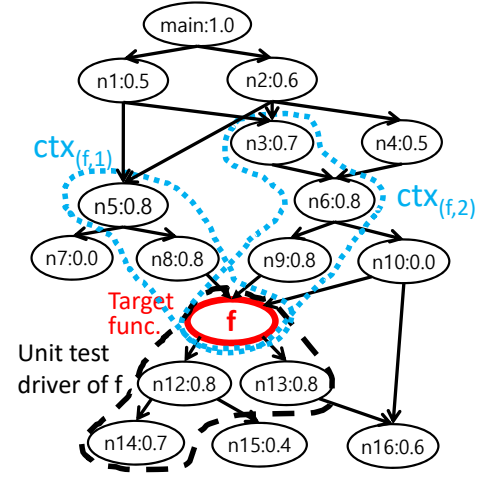
**DEFINITION 2.** A *calling context* of  $f \in V$  (saying  $ctx_{(f,i)}$ ) is a maximal call path  $a^1 \rightarrow a^2 \rightarrow \dots \rightarrow f$  in a static call graph  $G(V, E)$  satisfying the following conditions:

- $a^1$  is a predecessor of  $f$
- for all  $a^j$  in  $ctx_{(f,i)}$ ,  $p(a^j|f) \geq \tau$
- there exists no other calling context of  $f$  that contains  $ctx_{(f,i)}$  as its sub path (i.e.  $ctx_{(f,i)}$  is maximal).

CONBRIO generates a calling context by traversing a static call graph from  $f$  in a reverse direction until it reaches a node labelled with low dependency of  $f$ . For example, Figure 4 shows two calling contexts of  $f$ :  $ctx_{(f,1)}$  and  $ctx_{(f,2)}$ .  $ctx_{(f,1)}$  is a call path from  $n5$  to  $f$  (see the blue dotted line in the right part of the figure) where  $p(n5|f) = p(n8|f) = 0.8 > \tau$  and  $p(n1|f) = 0.5, p(n2|f) = 0.6$ . Thus,  $ctx_{(f,1)} = n5 \rightarrow n8 \rightarrow f$ . Similarly,  $ctx_{(f,2)} = n3 \rightarrow n6 \rightarrow n9 \rightarrow f$ .

### 3.4 Concolic Testing to Generate Violating Symbolic Path Formulas

CONBRIO applies concolic testing to a unit test driver to explore diverse and realistic target unit test executions. During concolic execution, it obtains a set of symbolic execution path formulas



**Figure 4: Static call graph showing a unit test driver and two calling contexts of  $f$  ( $ctx_{(f,1)}$  and  $ctx_{(f,2)}$ ) with  $\tau = 0.7$**

$SE_f$ <sup>5</sup> and records a symbolic path formula  $\sigma_{f(v_i,j)}$  that violates an assertion  $v_i$  in  $f$  ( $j$  is an index to a symbolic path formula violating  $v_i$  since there can be multiple such symbolic path formulas). We use  $\sigma_{f v_i}$  to denote  $\bigvee_j \sigma_{f(v_i,j)}$ .

To focus on exploring  $f$ 's diverse symbolic execution paths, CONBRIO modifies DFS search strategies by using two queues to contain branch conditions in symbolic execution paths: a priority queue for branch conditions of a target function  $f$  and a normal queue for those of the other functions in a unit test driver of  $f$  (e.g.  $g$  in Figure 1). CONBRIO explores various behaviors of  $f$  first by negating the branch conditions in the priority queue first (the branch conditions in the normal queue are negated when the priority queue becomes empty).

### 3.5 Alarm Filtering by Checking Satisfiability of $f$ 's Violating Symbolic Path Formula $\sigma_{f v_i}$ with $f$ 's Calling Context Formula

To filter out false alarms raised at  $v_i$  in  $f$ , CONBRIO checks the feasibility of  $\sigma_{f v_i}$  with regard to  $f$ 's calling contexts (see Sect. 3.3.2). For this purpose, CONBRIO constructs  $\Sigma_{ctx_{(f,k)}}$  which serves as a *symbolic calling context* of  $f$  and checks satisfiability of  $\sigma_{f v_i} \wedge \Sigma_{ctx_{(f,k)}}$  using a SMT solver.

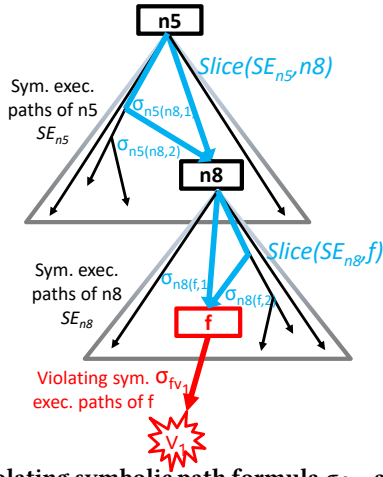
A symbolic calling context formula  $\Sigma_{ctx_{(f,k)}}$  of  $f$  is constructed as follows (see Figure 5):

1. For each function  $a^i$  in a calling context of  $f$  (i.e.  $ctx_{(f,k)}$ ), CONBRIO obtains  $SE_{a^i}$  which is a set of symbolic execution path formulas of  $a^i$  with its successor functions on which  $a^i$  has high dependency.

Note that this task is the same task in Step 4 (Sect. 3.4). If  $a^i$  was already tested as a target function and  $SE_{a^i}$  was generated in Step 4, this alarm filtering step (Step 5) reuses  $SE_{a^i}$  or vice versa. In other words, CONBRIO performs concolic execution of a function  $f$  (with its successor functions on which  $f$  has high dependency) in a target program *only once*.

<sup>4</sup> $\text{SYM\_assume}(expr)$  is a macro of  $\text{if}(!expr) \text{exit}(0);$ .

<sup>5</sup>  $SE_f$  is reused to filter alarms of other target function  $f'$  if a calling context of  $f'$  contains  $f$ . See Section 3.5 for more detail.



**Figure 5: Violating symbolic path formula  $\sigma_{fv_i}$  and a symbolic calling context of  $f$  (i.e.  $\Sigma_{ctx(f,1)}$ ) with  $ctx(f,1)$  in Figure 4**

- From  $SE_{a^i}$ , CONBRIO obtains a slice of  $SE_{a^i}$  with regard to  $a^{i+1}$  (saying  $Slice(SE_{a^i}, a^{i+1})$ ) as follows:

$Slice(SE_{a^i}, a^{i+1}) \stackrel{def}{=} \{\sigma' \mid \sigma' \text{ is a prefix of } \sigma \in SE_{a^i} \text{ such that } \sigma \text{ contains invocation of } a^{i+1} \text{ and } \sigma' \text{ does not contain a suffix of } \sigma \text{ after the invocation}\}.$

- For example, for  $ctx(f,1)$  in Figure 4, Figure 5 shows that  $Slice(SE_{n5}, n8)$  has two symbolic path formulas that call  $n8$ :  $\sigma_{n5(n8,1)}$  and  $\sigma_{n5(n8,2)}$  and  $Slice(SE_{n8}, f)$  has two symbolic path formulas that call  $f$ :  $\sigma_{n8(f,1)}$  and  $\sigma_{n8(f,2)}$  (shown as thick blue arrows in the figure).
- CONBRIO obtains *symbolic calling context formula* of  $f$  with  $ctx(f,k)$  (i.e.  $\Sigma_{ctx(f,k)}$ ) by combining sets of sliced symbolic execution path formulas of  $a^1$  (i.e.  $Slice(SE_{a^1}, a^2)$ ),  $a^2$  (i.e.  $Slice(SE_{a^2}, a^3)$ ), ... of  $ctx(f,k)$  until reaching fusing logical conjunction. Thus,  $\Sigma_{ctx(f,k)}$  with  $ctx(f,k) = a^1 \rightarrow a^2 \rightarrow \dots \rightarrow f$  is defined as follows:

$$\Sigma_{ctx(f,k)} \stackrel{def}{=} \bigwedge_{a^i \in ctx(f,k) - \{f\}} (\bigvee_{\sigma_j \in Slice(SE_{a^i}, a^{i+1})} \sigma_j)$$

For example, Figure 5 shows  $\Sigma_{ctx(f,1)}$  for  $f$  with  $ctx_f = n5 \rightarrow n8 \rightarrow f$  in Figure 4 as follows (see thick blue arrows representing  $\sigma_{n5(n8,1)}$ ,  $\sigma_{n5(n8,2)}$ ,  $\sigma_{n8(f,1)}$  and  $\sigma_{n8(f,2)}$ ):

$$\Sigma_{ctx(f,1)} = \Sigma_{n5 \rightarrow n8 \rightarrow f} = (\sigma_{n5(n8,1)} \vee \sigma_{n5(n8,2)}) \wedge (\sigma_{n8(f,1)} \vee \sigma_{n8(f,2)})$$

Finally, CONBRIO applies a SMT solver to  $\sigma_{fv_i} \wedge \Sigma_{ctx(f,k)}$  for every symbolic calling context of  $f$ . If the result is UNSAT for all calling contexts (i.e. there exists *no* feasible execution in any calling contexts of  $f$  to make  $\sigma_{fv_i}$  feasible), a target alarm is considered as false and ignored. Otherwise (i.e. the result is SAT with at least one calling context), a corresponding alarm is reported as a violation of  $v_i$  in  $f$ .

Although the alarm filtering strategy is effective (Sect. 5.4), the true/false decision on alarms can be still inaccurate. This is because  $\sigma_{fv_i}$  and/or  $\Sigma_{ctx(f,k)}$  might be inaccurate due to the limitation of concolic execution (i.e.  $\sigma_{fv_i}$  and  $\Sigma_{ctx(f,k)}$  may over-approximate and under-approximate real execution space). Thus, CONBRIO tries to improve accuracy of  $\sigma_{fv_i}$  by performing concolic execution on not only a target function but also with its closely relevant successor functions (the same idea applies to  $\Sigma_{ctx(f,k)}$ ).

**Table 1: Target programs and bugs for RQ1 to RQ4**

Target programs and versions	Lines	# of functions	# of sys. test cases	Branch coverage (%)	# of target bugs
Bash-2.0	32714	1214	1100	46.2%	6
Flex-2.4.3	7471	147	567	45.7%	2
Grep-2.0	5956	132	809	50.3%	5
Gzip-1.0.7	3054	82	214	55.8%	2
Make-3.75	28715	555	1043	64.5%	3
Sed-1.17	4085	73	360	47.3%	2
Vim-5.0	66209	1749	975	35.8%	6
Perl-5.8.7	79873	2240	1201	52.3%	6
Bzip2-1.0.3	4737	114	6	67.4%	2
Gcc-3.2	342561	5553	9	43.7%	15
Gobmk-3.3.14	154583	2682	1354	65.2%	5
Hmmer-2.0.42	35992	539	4	75.6%	3
Sjeng-11.2	10146	144	3	77.9%	2
Libquantum-0.2.4	2255	101	3	68.5%	3
H264ref-9.3	51578	590	6	63.6%	5
Sum	829929	15915	7654	N/A	67
Average	55328.6	1061.0	510.3	57.3%	4.5

## 4 EXPERIMENT SETUP

We have designed the five research questions (Sect. 4.1) to evaluate bug detection ability and precision of CONBRIO and compare CONBRIO with other concolic unit testing techniques (Sect. 4.3) on the 15 real-world C programs (Sect. 4.2). Note that it is important to evaluate both bug detection ability and precision together because there is a trade-off between them (i.e. a technique may improve bug detection ability at the cost of low precision or vice versa). Furthermore, we have applied CONBRIO to the latest versions of the nine C programs studied in other papers on crash bug detection techniques. Section 4.4 describes what we measured to answer the research questions and Section 4.5 describes the testbed setup. Section 4.6 discusses threats to validity of the experiment.

### 4.1 Research Questions

**RQ1. Bug Detection Ability:** How many crash bugs among the target crash bugs does CONBRIO detect compared to the other concolic unit testing techniques (i.e. measuring recall) ?

**RQ2. Bug Detection Precision:** How much is a false alarm ratio of CONBRIO compared to the other techniques?

**RQ3. Effectiveness of the Alarm Filtering Strategy:** How much does the alarm filtering strategy using symbolic calling contexts affect the number of target bugs detected and a false alarm ratio?

**RQ4. Effect of the Function Selection Strategy on Bug Detection Ability and Precision:** How much does the function selection strategy that includes closely relevant functions in unit test drivers and calling contexts affect a number of target bugs detected and a false alarm ratio?

**RQ5. Effectiveness of Detecting New Crash Bugs:** How many new crash bugs does CONBRIO detect?

### 4.2 Target Bugs and Programs

We target crash bugs described in Section 2.3 by inserting corresponding crash assertions in target programs because crash bugs are serious problems and CONBRIO can automatically generate

**Table 2: Target programs for RQ5**

Target programs and versions	Lines	# of functions	# of sys. test cases	Branch coverage (%)
abcm2ps-8.13.9	36595	499	12	74.0
autotrace-0.31.1	18495	343	5	69.3
bib2xml-5.11	77216	1032	24	73.2
catdvi-0.14	12693	187	7	53.0
eog-3.14.1	43463	605	42	73.3
gif2png-2.5.11	4058	76	2	60.1
jpegtran-1.3.1	51828	817	33	72.0
mp3gain-1.5.2	5786	100	3	53.7
xpdf-3.03	22309	381	13	54.9
Sum	272443	4040	141	N/A
Average	30271.4	448.9	15.7	64.8

and insert such assertions without user-given test oracles, which are rarely available in target programs. We use two benchmarks: *known crash bug benchmark* for RQ1 to RQ4 and *unknown crash bug benchmark* for RQ5.

**4.2.1 Known Crash Bug Benchmark.** The known crash bug benchmark consists of all C programs in SIR [11] (except Siemens programs and space which do not have available bug-fix histories) and SPEC2006 integer benchmarks (except mcf-1.2 which has only one system test case). We target the crash bugs of the benchmark programs that satisfy both of the following conditions:

- crash bugs that exist in a target program version and have been confirmed by the original developers through the bug-fix commits since the release of the target program version (e.g. Dec 1996 for bash-2.0) until April 2017
- crash bugs that can be detected by unit testing (i.e. both the buggy statement(s) reported in a bug-fix commit and the violated assertions are located in the same target function)

Table 1 describes the 15 target programs including their sizes (in LoC including comments and empty lines), a number of functions to test, a number of system test cases used, branch coverage achieved by the system test cases, and a number of the target crash bugs. For all target programs, we used all system test cases provided in the benchmarks. Each target program has two to fifteen target crash bugs (4.5 on average). Note that no system test case detects a target bug.

For example, we have reviewed 28 crash bug-fix commits reported since the release of vim-5.0 (i.e. Feb 1998) until April 2017. 11 among them report crash bugs existing in vim-5.0. Among the 11 crash bugs, unit testing can detect six bugs of them, which we target for vim-5.0 (see the eighth row of the table).

**4.2.2 Unknown Crash Bug Benchmark.** The unknown crash bug benchmark programs were selected from the literature on crash bug detection techniques. This is because SIR and SPEC benchmark programs do not satisfy the following criteria: we selected target programs whose sizes are 1,000 to 100,000 LoC and which have more than three crash bug fixes in the last three years (i.e. between April 2014 to April 2017). We excluded very large programs due to huge manual effort required to check validity of alarms. We also excluded programs with three or less crash bug fixes in the last three years because such programs may not have a crash bug.

To obtain new crash bug benchmark programs, we surveyed papers on crash bug detection techniques published in major SE (ICSE, FSE, ASE, ISSTA), PL (PLDI, POPL, SPLASH), and security conferences (IEEE S&P, ACM CCS, USENIX Security) in the last three years and obtained the nine relevant papers [2, 7, 19, 29, 36, 37, 46–48]. Then, we applied the above criteria to the latest versions of the target programs studied in these papers and obtained the nine target programs in Table 2. Again, for all target programs, we used all system test cases provided in the target program versions and no system test case violated the crash assertions.

### 4.3 Concolic Unit Testing Techniques to Compare

We have compared CONBRIO with the following concolic unit testing techniques:

- *Symbolic unit testing (SUT):* It generates a symbolic unit testing driver with symbolic arguments to a target function  $f$  and symbolic global variables without any constraints on the symbolic values, as described in Section 2.2. Also, SUT uses symbolic stubs to replace all functions called by  $f$ .
- *Static call-graph distance techniques:* It generates a symbolic unit test driver to include all successor functions of  $f$  within a certain distance bound from  $f$  in a static function call graph. Also, a calling context of  $f$  contains predecessor functions of  $f$  within a certain distance bound from  $f$ . We use distance bounds 3, 6 and 9. SUT corresponds to a static call-graph distance technique with a distance bound 0.

SUT uses DFS as a concolic search strategy. Call-graph distance techniques and CONBRIO use the modified DFS (Sect. 3.4).<sup>6</sup> We implemented CONBRIO and the other concolic unit testing techniques in 6,100 lines of C++ code using Clang/LLVM-3.4 [28] and CREST-BV [26]. CONBRIO uses Z3 [10] as a SMT solver and the LLVM-based static variable range analyzer [35] to compute the possible ranges of variables.

### 4.4 Measurement

We consider that a target bug is detected if a unit test execution that violates an assertion covers one of the buggy statements in a target unit. To identify the buggy statements, we have manually analyzed all crash bug-fix commits of the all subsequent releases of the target program versions included in SIR and SPEC2006 benchmark programs. We consider that a statement  $s$  of a target program is a buggy statement if  $s$  corresponds to the changed/fixed statements in a crash bug-fix commit.

We analyze alarms reported by the alarm filtering strategy (Sect. 3.5). For true alarms, we count a number of violated assert statements which satisfy all of the following conditions:

- There exists a unit test execution that covers a buggy statement and violates an assert statement in a target function.
- We can confirm that the unit test execution is feasible at system level by manually creating a system-level test that includes the unit test execution and violates the assert statement.

<sup>6</sup>We report the experiment using only DFS and modified DFS since the experiments using other search strategies such as random negation and CFG heuristic show only negligible difference.

We consider all the other alarms as false ones.

#### 4.5 Testbed Setting

For SUT, call-graph distance techniques, and CONBRIO, we set the timeout of concolic testing (Step 4 in Figure 3) as 180 seconds per a target function.<sup>7</sup> After test case generation terminates, call-graph distance techniques and CONBRIO performs the false alarm filtering task (Step 5 in Figure 3). We set a function dependency threshold  $\tau$  as 0.7.

Since the experiment scale is large (i.e. targeting 15,915 functions for the known crash bugs and 4,040 functions for the unknown crash bugs), the experiments were performed on 100 machines each of which is equipped with Intel quad-core i5 4670K and 8GB ram, running Ubuntu 14.04.2 64 bit version. We run four concolic unit test runs on a machine in parallel.

#### 4.6 Threats to Validity

A threat to external validity is the representativeness of our target programs. But we expect that this threat is limited since the target programs are widely used real-world ones and tested by many other researchers. Another threat to external validity is the possible bias of the system tests we used to obtain dependency between functions. We tried to reduce this threat by utilizing all available system test cases in the benchmarks.

A threat to internal validity is possible faults in the implementation of the concolic unit testing techniques we studied. To address this threat, we extensively tested our implementation. A threat to construct validity is the use of the crash bugs that were fixed by the bug-fix commits reported so far (i.e. the target programs may have unknown/unreported crash bugs which we do not count). We target crash bugs confirmed by the developers through the bug-fix commits because it would require too much effort to manually validate numerous alarms without confirmed reports in this large scale experiment. However, this threat seems limited because all target programs are well-maintained so that these programs may not have many new bugs.

## 5 EXPERIMENT RESULT

For all comparison in the experiments in this section, we applied Wilcoxon test with a significance level 0.05 to show the statistical significance. All comparison results in this section are statistically significant unless mentioned otherwise. The experiment data missing in the paper due to page limit are available at <https://sites.google.com/view/conbrrio-icse2018/>

### 5.1 Experiment Data

**5.1.1 Data on Unit Test Drivers and Calling Contexts.** For the 15 known crash bug benchmark programs, each unit test driver generated by CONBRIO contains 6.2 functions on average. CONBRIO generated 3.0 calling contexts per target function where each calling context has 6.6 functions on average. Call-graph distance techniques with bound 3, 6, and 9 generate a unit test driver that contains 5.8, 13.8, and 22.5 functions on average, respectively. Also,

<sup>7</sup>We selected timeout as 180 seconds because exploratory study with timeout 60, 180, 300, and 600 seconds suggested that timeout beyond 180 seconds had negligible effect on the overall experiment results of CONBRIO and the other techniques.

**Table 3: Numbers of the target bugs detected by the concolic unit testing techniques**

Target program	No. of target bugs	Bound of static call graph distance tech.				CONBRIO
		0(SUT)	3	6	9	
Bash-2.0	6	5	3	3	3	5
Flex-2.4.3	2	2	1	1	1	1
Grep-2.0	5	3	4	2	2	4
Gzip-1.0.7	2	2	1	1	1	2
Make-3.75	3	3	3	2	2	3
Sed-1.17	2	2	2	2	2	2
Vim-5.0	6	5	4	2	2	5
Perl-5.8.7	6	6	5	4	3	6
Bzip2-1.0.3	2	2	2	2	2	2
Gcc-3.2	15	14	12	9	8	14
Gobmk-3.3.14	5	4	3	3	3	5
Hmmer-2.0.42	3	3	3	3	3	3
Sjeng-11.2	2	2	2	2	2	2
Libquantum-0.2.4	3	3	2	2	2	3
H264ref-9.3	5	5	4	3	3	4
Sum	67	61	51	41	39	61

they generate 5.9, 11.1, and 24.3 calling contexts per target function on average, respectively.

**5.1.2 Data on Unit Tests Generated and Alarm Filtering.** For the 15 known crash bug benchmark programs, CONBRIO spent 1.8 hours to generate 7,979,781 unit test cases for 15,915 target functions and 2.3 hours to filter out false alarms using Z3 on 100 quad-core machines. Z3 reports that a symbolic calling context formula with a violating symbolic unit execution consists of 1.5 million clauses on 0.1 million Boolean variables on average and its maximum memory usage is around 7.6 GB. Call-graph distance techniques with a distance bound 0, 3, 6, and 9 spent the almost same 1.8 hours for unit test generation (i.e. most target functions reach the timeout) and 0, 2.6, 3.9, and 6.3 hours for the alarm filtering respectively.

CONBRIO covered 69.8% to 88.0% of the branches of a target program (82.5% on average) with the unit tests cases and the given system test cases (i.e. the unit test cases increase the branch coverage 25.2%p more on average (= 82.5% - 57.3% where 57.3% is the average branch coverage achieved by the system test cases (see the last row of Table 1)).

### 5.2 RQ1: Bug Detection Ability

Table 3 describes a number of the target bugs detected by the concolic unit testing techniques and shows that CONBRIO has high bug detection ability. CONBRIO and static call-graph distance technique with bound zero (i.e. SUT) achieve the highest bug detection ability (i.e. 91.0% (=61/67)) (but SUT achieves this at the cost of many false alarms (see Sect. 5.3)). Note that the given system tests do not detect any of the target bugs. In addition, we applied concolic testing at system level using distributed concolic testing tool [25] with the same amount of total time on 100 machines and found that no target bug is detected.

As a distance bound of the call-graph distance techniques increases to 3, 6, and 9, the number of detected bugs severely decreases



**Table 4: Numbers of false alarms and numbers of false alarms per true alarm of the concolic unit testing techniques**

Target programs	Static call-graph distance techniques								CONBRIO	
	0 (SUT)		3		6		9		# of false alarms	F/T ratio
	# of alarms	F/T ratio	# of false alarms	F/T ratio	# of false alarms	F/T ratio	# of false alarms	F/T ratio		
Bash-2.0	484	96.8	137	45.7	69	23.0	54	18.0	18	3.6
Flex-2.4.3	142	71.0	25	25.0	12	12.0	12	12.0	6	6.0
Grep-2.0	120	40.0	34	8.5	18	9.0	18	9.0	13	3.3
Gzip-1.0.7	33	16.5	7	7.0	3	3.0	3	3.0	5	2.5
Make-3.75	664	221.3	106	35.3	59	29.5	46	23.0	9	3.0
Sed-1.17	31	15.5	9	4.5	4	2.0	4	2.0	5	2.5
Vim-5.0	906	181.2	207	51.8	123	61.5	72	36.0	25	5.0
Perl-5.8.7	392	65.3	187	37.4	64	16.0	44	14.7	57	9.5
Bzip2-1.0.3	34	17.0	12	6.0	7	3.5	7	3.5	10	5.0
Gcc-3.2	2026	144.7	503	41.9	195	21.7	147	18.4	79	5.6
Gobmk-3.3.14	791	197.8	133	44.3	62	20.7	45	15.0	39	7.8
Hmmer-2.0.42	162	54.0	48	16.0	22	7.3	22	7.3	12	4.0
Sjeng-11.2	108	54.0	13	6.5	7	3.5	7	3.5	8	4.0
Libquantum-0.2.4	55	18.3	9	4.5	4	2.0	4	2.0	5	1.7
H264ref-9.3	232	46.4	34	8.5	15	5.0	15	5.0	17	4.3
Average	412.0	82.7	97.6	22.9	44.3	14.6	33.3	11.5	20.5	4.5

to 51, 41, and 39, respectively because larger symbolic search space should be explored within the timeout.

Among the undetected six target bugs (=67-61), three target bugs in bash, grep, and gcc were missed because concolic execution did not cover corresponding buggy statements within the timeout, two bugs in flex and h264ref were missed because of the alarm filtering strategy, and one in vim was missed because a unit execution covered the corresponding buggy statement and an assert statement but did not violate the assert statement.

### 5.3 RQ2: Bug Detection Precision

Table 4 describes a number of false alarms and a ratio of false alarms per true alarm of the techniques and shows that CONBRIO achieves high bug detection precision. Among the techniques, CONBRIO raises the lowest number of false alarms (i.e. 20.5 false alarms per target program on average) and the lowest false alarms per true alarms ratio (i.e. 4.5 false alarms per true alarm on average).<sup>8</sup>

The static call-graph distance technique with distance 0 (i.e. SUT) suffers the largest number of false alarms (412.0 false alarms per target program on average). CONBRIO raises only 5.0% (=20.5/412.0), 21.0%, and 46.4% and 61.6% of the false alarms raised by the static call-graph distance techniques with distance bounds 0, 3, 6, and 9 on average respectively (see the last row of the table).

### 5.4 RQ3. Effectiveness of the Alarm Filtering Strategy

The comparison of the experiment results of CONBRIO and CONBRIO without the alarm filtering strategy using symbolic calling context formulas (Sect. 3.5) demonstrates that the alarm filtering strategy improves bug detection precision significantly. In other words, CONBRIO without the alarm filtering strategy detects two more target bugs (i.e. 63 bugs) in all target programs but with five times higher false alarm ratio (i.e. 20.3 false alarms per true alarm on average). Although the alarm filtering strategy spent more time

<sup>8</sup>The static alarm reduction heuristics of CONBRIO decrease the number of false alarms (23.5 to 20.5 on average) and the number of false alarms per true alarm (5.2 to 4.5 on average) without decreasing the bug detection ability (i.e. CONBRIO without the static alarm reduction heuristics detects the same 61 bugs and raises 23.5 false alarms per target program on average).

(2.3 hours) than the unit test generation (1.8 hours), this strategy is worthwhile to apply to improve bug detection precision. Detailed experiment data is available at <https://sites.google.com/view/conbr-io-icse2018/>.

### 5.5 RQ4. Effect of the Function Selection Strategy on Bug Detection Ability and Precision

The comparison on the experiment results of CONBRIO and the call-graph distance techniques with bounds confirm that the idea of including only *closely relevant* functions to a target function in unit test drivers and calling contexts and the proposed dependency metric to measure function relevance (Sect. 3.2) are effective.

For example, CONBRIO and the call-graph distance technique with bound 3 generate a unit test driver of similar size (i.e. 6.2 vs. 5.8 functions on average) and the amount of generated calling contexts are also comparable (3.0 calling contexts each of which has 6.6 functions vs. 5.9 calling contexts each of which has 2.8 functions on average) (see Sect. 5.1.1). The time taken to generate unit test execution is almost same 1.8 hours and the time taken to apply the alarm filtering strategy is also similar (2.3 vs 2.6 hours).

However, CONBRIO achieves much higher bug detection ability and precision than the call-graph distance technique with bound 3 (i.e. 91.0% vs 76.1% (=51/67) for bug detection ability and 4.5 vs. 22.9 false alarms per true alarm on average). With larger distance bounds 6 and 9, a number of the detected bugs drops to 41 and 39 and the false alarm ratio decreases to 14.6 and 11.5 respectively, which is still three to two times less precise than CONBRIO.

### 5.6 RQ5. Effectiveness of Detecting New Crash Bugs

CONBRIO detects 14 new crash bugs in the seven target programs. CONBRIO detects five new crash bugs in autotrace, two bugs in each of abcm2ps, gif2png, and mp3gain, one bug in each of bib2xml, eog, and jpegtran, and no bug in catdvi and xpdf.<sup>9</sup>

Note that we have confirmed the 14 new crash bugs by manually creating system-level test cases that crash a target program due to the bugs detected by CONBRIO. CONBRIO raises 71 false alarms over the all target programs and its true to false alarm ratio for each program ranges from 1:3.0 to 1:6.0 (1:4.3 on average except catdvi and xpdf). We have reported these 14 new crash bugs to the original developers and been waiting the responses from them (detailed example and explanation of the newly detected bugs is available at <https://sites.google.com/view/conbr-io-icse2018/>).

## 6 RELATED WORK

### 6.1 Concolic Unit Testing Techniques

There exist concolic unit testing techniques (e.g. [6, 34, 41, 43]) which require a user to build symbolic unit test drivers and stubs. DART [17] generates symbolic unit test drivers (but not symbolic

<sup>9</sup> CONBRIO generated 3.3 calling contexts per target function each of which has 4.3 functions on average. It spent 11.2 minutes to generate 725,584 unit test cases for 4,040 target functions and 20.7 minutes to apply alarm filtering strategy on 100 machines and covered 80.9% of the branches on average.

stubs) like SUT (Section 2.2) and test inputs for C programs. CONBOL [27] generates symbolic unit test drivers/stubs and test inputs targeting large-scale embedded C programs. DART and CONBOL generate symbolic unit test drivers without utilizing contexts of a target function  $f$  and may suffer many false alarms.<sup>10</sup>

Chakrabarti and Godefroid [8] developed a unit testing technique which statically partitions a static call graph using topological information and tests each partition as a unit through symbolic execution. This technique may suffer many false alarms because the obtained partitions may not represent groups of relevant functions due to insufficient information to generate partitions (i.e. using only topological information of a static call graph without semantic or dynamic information). Their tool is not publicly available and the paper does not report bug detection ability nor precision [8]. Tomb. et al [44] reported that interprocedural program analysis with deeper call depth bound raise fewer false alarms. However, they did not report how to set a proper call depth bound.

Recently, UC-KLEE [36] directly starts symbolic execution from a target function using lazy initialization [23]. Through the manual analysis of the thousands of alarms, the authors of UC-KLEE detected 67 new bugs in BIND, openssl, Linux kernel and its true to false alarm ratio is 1:5.7 on average. We could not directly compare CONBRIO with UC-KLEE because UC-KLEE is not publicly available and all of BIND, openssl, and Linux kernel are too large (i.e. million lines of code).

## 6.2 Random Method Sequences Generation Techniques for Object-Oriented Programs

Randoop [31] invokes a random sequence of public methods including constructors of a target method's class. EvoSuite [13, 14] tests Java methods using search-based strategies with symbolic execution. TestFul [4] combines genetic algorithm and a local search to improve the speed of Java unit test generation. Garg et al. [16] improves Randoop by generating input test cases of the generated method sequence using concolic testing for C++ programs. These techniques may also suffer false alarms due to infeasible test inputs/method sequences generated. For example, Gross et al. [20] reported that Randoop raised 181 alarms without detecting any bug (i.e. all alarms were false ones) on five Java programs although the authors of Randoop reported that Randoop's true to false alarm ratio is 1:0.67 on 8 Java libraries and 6 .NET libraries on average [31]. Fraser et al. [14] reported that the statistically estimated true to false alarm ratios range from 1:0.6 to 1:4.2 in their experiments on randomly selected 100 projects hosted on sourceforge.net. Garg et al. [16] does not report detected bugs or false alarm ratios but branch coverage obtained using the proposed technique on eight programs (except gnuccness on which the authors reported nine new bugs and that a true to false alarm ratio was 1:1.0). In spite of the lack of explicit context information like class/object information in C programs, CONBRIO detects bugs precisely in C programs (i.e. a true to false alarm ratio is 1:4.5 on average) while keeping

high bug detection ability (i.e. 91.0% of the target bug detected on average).<sup>11</sup>

## 6.3 Automated Unit Testing Techniques based on System Tests

Elbaum et al. [12] proposed a technique to generate unit tests from the system tests; the technique captures program states before and after an invocation of a target function  $f$  to generate unit test inputs and oracles for  $f$ . OCAT [22] captures object instances during system executions and generates unit tests using Randoop with the captured object and the mutated object instances as seed objects. GenUtest [33] automatically generates unit tests and mock objects using captured method sequences during system testing. A limitation of these techniques is that the executions of the generated unit tests just replay the same behaviors [12, 33] (or similar behaviors [22]) of the target unit in already performed system testing (i.e. they are applicable to only regression testing of evolving software, not to a single version of software). Also, the aforementioned papers do not report bug detection ability nor precision.

## 6.4 Automatic Generation of Mock Objects/Testing Stubs

Dsc+Mock [21] generates mock objects for testing the interfaces of Java programs. Dsc+Mock collects type constraints on the interfaces during symbolic execution and generates mock objects using the solution of the type constraints. Galler et al. [15] proposed a technique to generate mock objects using the design-by-contract specification. Saff et al. [38, 39] proposed a technique to generate mocking objects from system test executions (i.e. mock objects are generated based on the interactions between the target code and its environment captured during system executions). Since these techniques generate only concrete mock objects, not symbolic mocks, they have a limitation in providing various unit contexts.

## 7 CONCLUSION AND FUTURE WORK

We have presented an automated concolic unit testing technique CONBRIO which generates unit test drivers to closely mimic the real contexts of a target function  $f$  and filters out false alarms using symbolic calling context formulas of  $f$  using relevant functions to  $f$ . Through the experiments, CONBRIO demonstrates both high bug detection ability (91.0% of all target bugs detected) and high bug detection precision (i.e. a true to false alarm ratio is 1:4.5). Furthermore, CONBRIO detects 14 new crash bugs in the latest versions of the target C programs studied in other papers on crash bug detection techniques.

As future work, to improve the precision of automated unit testing further, we plan to refine the function dependency metric by analyzing more semantic characteristic of target program executions. Also, we will improve bug detection precision by synthesizing a common-likely symbolic calling context based on multiple calling contexts of a target function  $f$ .

<sup>10</sup> DART [17] bypasses the false alarm issue by targeting public API functions of libraries which should work with all possible inputs (i.e. no false alarms caused by infeasible unit executions).

<sup>11</sup> The aforementioned papers report only bug detection precision (RQ2), not bug detection ability (RQ1), which makes fair comparison between these techniques and CONBRIO difficult. This is because these techniques may improve a true to false alarm ratio at the cost of missing bugs. We report both bug detection ability and precision due to the trade-off between precision and recall of bug detection.

## REFERENCES

- [1] Shay Artzi, Adam Kiezun, Julian Dolby, Frank Tip, Danny Dig, Amit Paradkar, and Michael D. Ernst. 2008. Finding Bugs in Dynamic Web Applications. In *Proceedings of the 2008 International Symposium on Software Testing and Analysis (ISSTA '08)*. ACM, New York, NY, USA, 261–272.
- [2] Thanassis Avgerinos, Alexandre Rebert, Sang Kil Cha, and David Brumley. 2014. Enhancing Symbolic Execution with Veritest. In *Proceedings of the 36th International Conference on Software Engineering (ICSE 2014)*. ACM, New York, NY, USA, 1083–1094.
- [3] Radu Banabic, George Candea, and Rachid Guerraoui. 2014. Finding Trojan Message Vulnerabilities in Distributed Systems. In *Proceedings of the 19th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS '14)*. ACM, New York, NY, USA, 113–126.
- [4] Luciano Baresi, Pier Luca Lanzi, and Matteo Miraz. 2010. TestFul: An Evolutionary Test Approach for Java. In *Proceedings of the 2010 Third International Conference on Software Testing, Verification and Validation (ICST '10)*. IEEE Computer Society, Washington, DC, USA, 185–194.
- [5] Ella Bounimova, Patrice Godefroid, and David Molnar. 2013. Billions and Billions of Constraints: Whitebox Fuzz Testing in Production. In *Proceedings of the 2013 International Conference on Software Engineering (ICSE '13)*. IEEE Press, Piscataway, NJ, USA, 122–131.
- [6] Cristian Cadar, Daniel Dunbar, and Dawson Engler. 2008. KLEE: Unassisted and Automatic Generation of High-coverage Tests for Complex Systems Programs. In *Proceedings of the 8th USENIX Conference on Operating Systems Design and Implementation (OSDI'08)*. USENIX Association, Berkeley, CA, USA, 209–224.
- [7] S. K. Cha, M. Woo, and D. Brumley. 2015. Program-Adaptive Mutational Fuzzing. In *2015 IEEE Symposium on Security and Privacy*. 725–741. <https://doi.org/10.1109/SP.2015.50>
- [8] Arindam Chakrabarti and Patrice Godefroid. 2006. Software Partitioning for Effective Automated Unit Testing. In *Proceedings of the 6th International Conference on Embedded Software (EMSOFT '06)*. ACM, New York, NY, USA, 262–271.
- [9] Christoph Csallner and Yannis Smaragdakis. 2004. JCrasher: An Automatic Robustness Tester for Java. *Software Practical Experience* 34, 11 (Sept. 2004), 1025–1050.
- [10] Leonardo De Moura and Nikolaj Bjørner. 2008. Z3: An Efficient SMT Solver. In *Proceedings of the Theory and Practice of Software, 14th International Conference on Tools and Algorithms for the Construction and Analysis of Systems (TACAS'08/ETAPS'08)*. Springer-Verlag, Berlin, Heidelberg, 337–340.
- [11] Hyunsook Do, Sebastian Elbaum, and Gregg Rothermel. 2005. Supporting Controlled Experimentation with Testing Techniques: An Infrastructure and Its Potential Impact. *Empirical Software Engineering* 10, 4 (Oct. 2005), 405–435.
- [12] S. Elbaum, H. Chin, M. Dwyer, and M. Jorde. 2009. Carving and Replaying Differential Unit Test Cases from System Test Cases. *IEEE Transactions on Software Engineering (TSE)* 35, 1 (Jan 2009), 29–45.
- [13] Gordon Fraser and Andrea Arcuri. 2011. EvoSuite: Automatic Test Suite Generation for Object-oriented Software. In *Proceedings of the 19th ACM SIGSOFT Symposium and the 13th European Conference on Foundations of Software Engineering (ESEC/FSE '11)*. ACM, New York, NY, USA, 416–419.
- [14] Gordon Fraser and Andrea Arcuri. 2015. 1600 Faults in 100 Projects: Automatically Finding Faults While Achieving High Coverage with EvoSuite. *Empirical Software Engineering* 20, 3 (June 2015), 611–639.
- [15] Stefan J. Galler, Andreas Maller, and Franz Wotawa. 2010. Automatically Extracting Mock Object Behavior from Design by Contract&Trade; Specification for Test Data Generation. In *Proceedings of the 5th Workshop on Automation of Software Test (AST '10)*. ACM, New York, NY, USA, 43–50.
- [16] Pranav Garg, Franjo Ivancic, Gogul Balakrishnan, Naoto Maeda, and Aarti Gupta. 2013. Feedback-directed Unit Test Generation for C/C++ Using Concolic Execution. In *Proceedings of the 2013 International Conference on Software Engineering (ICSE '13)*. IEEE Press, Piscataway, NJ, USA, 132–141.
- [17] Patrice Godefroid, Nils Klarlund, and Koushik Sen. 2005. DART: Directed Automated Random Testing. In *Proceedings of the 2005 ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '05)*. ACM, New York, NY, USA, 213–223.
- [18] Patrice Godefroid, Michael Y Levin, and David A Molnar. 2008. Automated Whitebox Fuzz Testing. In *Proceedings of the 2008 Network and Distributed System Symposium*, Vol. 8. 151–166.
- [19] Denis Gopan, Evan Driscoll, Ducson Nguyen, Dimitri Naydich, Alexey Loginov, and David Melski. 2015. Data-delineation in Software Binaries and Its Application to Buffer-overflow Discovery. In *Proceedings of the 37th International Conference on Software Engineering - Volume 1 (ICSE '15)*. IEEE Press, Piscataway, NJ, USA, 145–155. <http://dl.acm.org/citation.cfm?id=2818754.2818775>
- [20] F. Gross, G. Fraser, and A. Zeller. 2012. Search-based System Testing: High Coverage, No False Alarms. In *International Symposium on Software Testing and Analysis (ISSTA)*.
- [21] Mainul Islam and Christoph Csallner. 2010. Dsc+Mock: A Test Case + Mock Class Generator in Support of Coding Against Interfaces. In *Proceedings of the Eighth International Workshop on Dynamic Analysis (WODA '10)*. ACM, New York, NY, USA, 26–31.
- [22] Hojun Jaygarl, Sunghun Kim, Tao Xie, and Carl K. Chang. 2010. OCAT: Object Capture-based Automated Testing. In *Proceedings of the 19th International Symposium on Software Testing and Analysis (ISSTA '10)*. ACM, New York, NY, USA, 159–170.
- [23] Sarfaraz Khurshid, Corina S. Păsăreanu, and Willem Visser. 2003. Generalized Symbolic Execution for Model Checking and Testing. In *Proceedings of the 9th International Conference on Tools and Algorithms for the Construction and Analysis of Systems (TACAS'03)*. Springer-Verlag, Berlin, Heidelberg, 553–568.
- [24] Moonzoo Kim, Yunho Kim, and Yoonkyu Jang. 2012. Industrial Application of Concolic Testing on Embedded Software: Case Studies. In *Proceedings of the 2012 IEEE Fifth International Conference on Software Testing, Verification and Validation (ICST '12)*. IEEE Computer Society, Washington, DC, USA, 390–399.
- [25] M. Kim, Y. Kim, and G. Rothermel. 2012. A Scalable Distributed Concolic Testing Approach: An Empirical Evaluation. In *International Conference on Software Testing, Verification and Validation (ICST)*.
- [26] Yunho Kim, Moonzoo Kim, YoungJoo Kim, and Yoonkyu Jang. 2012. Industrial Application of Concolic Testing Approach: A Case Study on Libexif by Using CREST-BV and KLEE. In *Proceedings of the 34th International Conference on Software Engineering (ICSE '12)*. IEEE Press, Piscataway, NJ, USA, 1143–1152.
- [27] Y. Kim, Y. Kim, T. Kim, G. Lee, Y. Jang, and M. Kim. 2013. Automated unit testing of large industrial embedded software using concolic testing. In *Proceedings of the 2013 IEEE/ACM 28th International Conference on Automated Software Engineering*. 519–528.
- [28] Chris Lattner and Vikram Adve. 2004. LLVM: A Compilation Framework for Lifelong Program Analysis & Transformation. In *Proceedings of the International Symposium on Code Generation and Optimization: Feedback-directed and Runtime Optimization (CGO '04)*. IEEE Computer Society, Washington, DC, USA, 75–.
- [29] Wei Le and Shannon D. Pattison. 2014. Patch Verification via Multiversion Interprocedural Control Flow Graphs. In *Proceedings of the 36th International Conference on Software Engineering (ICSE '14)*. ACM, New York, NY, USA, 1047–1058. <https://doi.org/10.1145/2568225.2568304>
- [30] David Molnar, Xue Cong Li, and David A. Wagner. 2009. Dynamic Test Generation to Find Integer Bugs in x86 Binary Linux Programs. In *Proceedings of the 18th Conference on USENIX Security Symposium (SSYM'09)*. USENIX Association, Berkeley, CA, USA, 67–82.
- [31] Carlos Pacheco, Shuvendu K. Lahiri, Michael D. Ernst, and Thomas Ball. 2007. Feedback-Directed Random Test Generation. In *Proceedings of the 29th International Conference on Software Engineering (ICSE '07)*. IEEE Computer Society, Washington, DC, USA, 75–84.
- [32] Yongbae Park, Shin Hong, Moonzoo Kim, Dongju Lee, and Junhee Cho. 2015. Systematic Testing of Reactive Software with Non-deterministic Events: A Case Study on LG Electric Oven. In *Proceedings of the 37th International Conference on Software Engineering - Volume 2 (ICSE '15)*. IEEE Press, Piscataway, NJ, USA, 29–38.
- [33] B. Pasternak, S. Syszberowicz, and A. Yehudai. 2009. GenUTest: A Unit TEst and Mock Aspect Generation Tool. *Software Tools for Technology Transfer* 11, 4 (2009), 273–290.
- [34] Corina S. Păsăreanu, Peter C. Mehrlitz, David H. Bushnell, Karen Gundy-Burlet, Michael Lowry, Suzette Person, and Mark Pape. 2008. Combining Unit-level Symbolic Execution and System-level Concrete Execution for Testing NASA Software. In *Proceedings of the 2008 International Symposium on Software Testing and Analysis (ISSTA '08)*. ACM, New York, NY, USA, 15–26.
- [35] Fernando Magno Quintao Pereira, Raphael Ernani Rodrigues, and Victor Hugo Sperle Campos. 2013. A Fast and Low-overhead Technique to Secure Programs Against Integer Overflows. In *Proceedings of the 2013 IEEE/ACM International Symposium on Code Generation and Optimization (CGO) (CGO '13)*. IEEE Computer Society, Washington, DC, USA, 1–11.
- [36] David A. Ramos and Dawson Engler. 2015. Under-constrained Symbolic Execution: Correctness Checking for Real Code. In *Proceedings of the 24th USENIX Conference on Security Symposium (SEC'15)*. USENIX Association, Berkeley, CA, USA, 49–64.
- [37] Alexandre Rebert, Sang Kil Cha, Thanassis Avgerinos, Jonathan Foote, David Warren, Gustavo Grieco, and David Brumley. 2014. Optimizing Seed Selection for Fuzzing. In *23rd USENIX Security Symposium (USENIX Security 14)*. USENIX Association, San Diego, CA, 861–875. <https://www.usenix.org/conference/usenixsecurity14/technical-sessions/presentation/rebert>
- [38] David Saff, Shay Artzi, Jeff H. Perkins, and Michael D. Ernst. 2005. Automatic Test Factoring for Java. In *Proceedings of the 20th IEEE/ACM International Conference on Automated Software Engineering (ASE '05)*. ACM, New York, NY, USA, 114–123.
- [39] David Saff and Michael D. Ernst. 2004. Mock Object Creation for Test Factoring. In *Proceedings of the 5th ACM SIGPLAN-SIGSOFT Workshop on Program Analysis for Software Tools and Engineering (PASTE '04)*. ACM, New York, NY, USA, 49–51.
- [40] Raimondas Sasnauskas, Olaf Landsiedel, Muhammad Hamad Alizai, Carsten Weise, Stefan Kowalewski, and Klaus Wehrle. 2010. KleeNet: Discovering Insidious Interaction Bugs in Wireless Sensor Networks Before Deployment. In *Proceedings of the 9th ACM/IEEE International Conference on Information Processing in Sensor Networks (IPSN '10)*. ACM, New York, NY, USA, 186–196.

- [41] Koushik Sen, Darko Marinov, and Gul Agha. 2005. CUTE: A Concolic Unit Testing Engine for C. In *Proceedings of the 10th European Software Engineering Conference Held Jointly with 13th ACM SIGSOFT International Symposium on Foundations of Software Engineering (ESEC/FSE-13)*. ACM, New York, NY, USA, 263–272.
- [42] spec2006 [n. d.]. The SPEC CPU 2006 Benchmark Suite. ([n. d.]). <https://www.spec.org/cpu2006/>.
- [43] Nikolai Tillmann and Jonathan De Halleux. 2008. Pex: White Box Test Generation for .NET. In *Proceedings of the 2Nd International Conference on Tests and Proofs (TAP'08)*. Springer-Verlag, Berlin, Heidelberg, 134–153.
- [44] Aaron Tomb, Guillaume Brat, and Willem Visser. 2007. Variably Interprocedural Program Analysis for Runtime Error Detection. In *Proceedings of the 2007 International Symposium on Software Testing and Analysis (ISSTA '07)*. ACM, New York, NY, USA, 97–107. <https://doi.org/10.1145/1273463.1273478>
- [45] Willem Visser, Corina S. Păsăreanu, and Sarfraz Khurshid. 2004. Test Input Generation with Java PathFinder. In *Proceedings of the 2004 ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA '04)*. ACM, New York, NY, USA, 97–107.
- [46] J. Wagner, V. Kuznetsov, G. Candea, and J. Kinder. 2015. High System-Code Security with Low Overhead. In *2015 IEEE Symposium on Security and Privacy*. 866–879. <https://doi.org/10.1109/SP.2015.58>
- [47] Xiaofei Xie, Yang Liu, Wei Le, Xiaohong Li, and Hongxu Chen. 2015. S-looper: Automatic Summarization for Multipath String Loops. In *Proceedings of the 2015 International Symposium on Software Testing and Analysis (ISSTA 2015)*. ACM, New York, NY, USA, 188–198. <https://doi.org/10.1145/2771783.2771815>
- [48] F. Yamaguchi, A. Maier, H. Gascon, and K. Rieck. 2015. Automatic Inference of Search Patterns for Taint-Style Vulnerabilities. In *2015 IEEE Symposium on Security and Privacy*. 797–812. <https://doi.org/10.1109/SP.2015.54>