

# On Benchmarking the Capability of Symbolic Execution Tools with Logic Bombs

Hui Xu, Zirui Zhao, Yangfan Zhou, and Michael R. Lyu, *Fellow, IEEE*

**Abstract**—Symbolic execution now becomes an indispensable technique for software testing and program analysis. There are several symbolic execution tools available off-the-shelf, and we need a practical benchmark approach to learn their capabilities. Therefore, this paper introduces a novel approach to benchmark symbolic execution tools in a fine-grained and efficient manner. In particular, our approach evaluates the performance of such tools against the known challenges faced by general symbolic execution techniques, such as floating-point numbers and symbolic memories. To this end, we first survey related papers and systematize the challenges of symbolic execution. We extract 12 distinct challenges from the literature and categorize them into two categories: *symbolic-reasoning challenges* and *path-explosion challenges*. Then, we develop a dataset of logic bombs and a framework to benchmark symbolic execution tools automatically. For each challenge, our dataset contains several logic bombs, each of which is guarded by a specific challenging problem. If a symbolic execution tool can find test cases to trigger logic bombs, it indicates that the tool can handle the corresponding problems. We have conducted real-world experiments with three popular symbolic execution tools: KLEE, Angr, and Triton. Experimental results show that our approach can reveal their capabilities and limitations in handling particular issues accurately and efficiently. The benchmark process generally takes only dozens of minutes to evaluate a tool. We release our dataset on GitHub as open source, with an aim to better facilitate the community to conduct future work on benchmarking symbolic execution tools.

**Index Terms**—symbolic execution.

## 1 INTRODUCTION

Symbolic execution is a popular technique for software testing and program analysis [1]. It has achieved rapid development in the last decade with several open-source symbolic execution tools available, such as KLEE [2] and Angr [3]. Current methods for evaluating symbolic execution tools generally rely on the achieved code coverage or the number of bugs detected in real-world programs (e.g., [2], [4]). Nevertheless, such metrics may fluctuate with different types of programs being analyzed, and they cannot manifest the detailed capabilities of a symbolic execution tool. This paper, therefore, aims to propose a fine-grained benchmark approach for symbolic execution tools which is less sensitive to targeting programs.

In particular, we observe that there are some common challenging problems which symbolic execution tools may not handle well, such as floating-point numbers [5] and loops [6]. They are the determinant factors of the code coverage that a symbolic execution tool can achieve. Therefore, we propose to develop a benchmark approach based on these known challenges. To elaborate, we can employ each challenge as an evaluation metric. In this way, our approach can give more meaningful information concerning the capability of a symbolic execution tool. Moreover, the benchmark result is unbiased as it does not depend on

particular programs for analysis.

To this end, we first conduct a systematic survey on the challenges of symbolic execution. This step is essential for our benchmark approach to embrace as many distinct challenges as possible. We categorize existing challenges into two categories: *symbolic-reasoning challenges* and *path-explosion challenges*. Symbolic-reasoning challenges attack the core symbolic reasoning process, where they may pose problems to a symbolic execution tool in generating incorrect test cases for particular control flows. These challenges include symbolic variable declarations, covert propagations, parallel executions, symbolic memories, contextual symbolic values, symbolic jumps, floating-point numbers, buffer overflows, and arithmetic overflows. Path-explosion challenges introduce too many possible control flows to analyze, which may cause a symbolic execution tool starving the computational resources or spending very long time on exploring the paths. Not only large-sized programs but also small-sized programs can cause path-explosion issues, as long as they include complex routines, such as external function calls, loops, and crypto functions. With our approach, consequently, all existing challenges discussed in the literature can be well categorized.

Next, we develop an accurate and efficient approach to benchmark the capability of symbolic execution tools with respect to each of the challenges. To this end, we cannot employ real-world programs for testing because they are too complicated and any challenges could cause a failure. Moreover, symbolic execution itself is not very efficient, and benchmark with real-world programs generally takes a long time. We tackle the problem by designing small programs embedded with logic bombs. A logic bomb is an artificial code block that can only be executed when certain

- H. Xu, and M. R. Lyu are with Shenzhen Research Institute and Dept. of Computer Science & Engineering, The Chinese University of Hong Kong. Email: hxiu, lyu@cse.cuhk.edu.hk
- Z. Zhao is with The Chinese University of Hong Kong and The University of Science and Technology of China
- Y. Zhou<sup>✉</sup> is with School of Computer Science, Fudan University. Email: zyf@fdu.edu.cn

Manuscript submitted to IEEE Transactions on Dependable and Secure Computing.

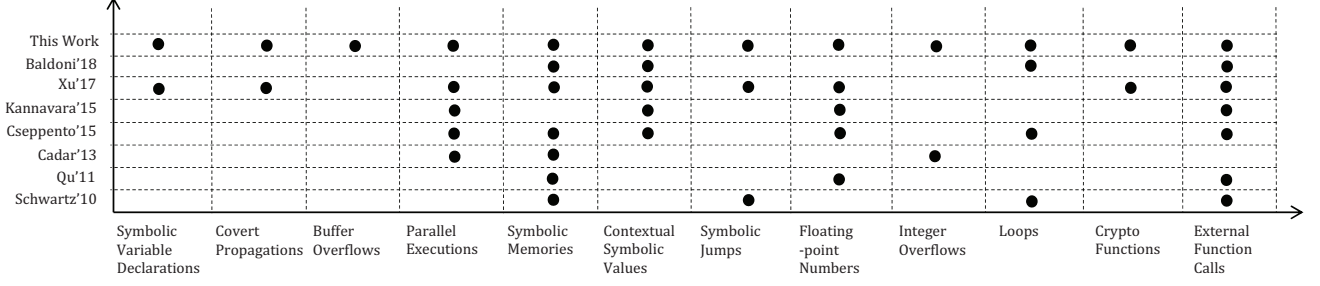


Fig. 1. The challenges of symbolic execution discussed in the literature. The detailed paper references are Schwartz'10 [7], Qu'11 [8], Cadar'13 [6], Cseppento'15 [9], Kannavara'15 [10], Quan'16 [5], Xu'17 [11], and Baldoni'18 [12].

conditions are met. We create such logic bombs that can only be triggered when a challenging problem is solved. The benefits are two folds. By keeping each logic bomb as small as possible, our evaluation result would be less likely affected by other unexpected issues. Also, employing small programs can shorten the required symbolic execution time and provides efficiency.

Following this method, we have designed a dataset of logic bombs covering all the challenges and a framework to run benchmark experiments automatically. For each challenge, our dataset contains several logic bombs with different problem settings or different levels of hardness, such as a one-leveled arrays or a two-leveled arrays for the symbolic memory challenge. Our framework employs the dataset of logic bombs as evaluation metrics. It firstly parses the logic bombs and compiles them to object codes or binaries; then, it directs a symbolic execution tool to perform symbolic execution on the logic bombs in a batch mode; and finally, it verifies the generated test cases and produces reports. We release our dataset and framework tools on GitHub.

We have conducted real-world experiments to benchmark three prevalent symbolic execution tools, KLEE, Triton, and Angr. Although these tools adopt different implementation techniques, our framework can adapt to them with only a little customization. The benchmark process for each tool generally took dozens of minutes. Experimental results show that our benchmark approach can reveal their capabilities and limitations efficiently and accurately. In summary, Angr has achieved the best performance with 21 cases solved, which is roughly one third of the total logic bombs; KLEE solved nine cases; and Triton only solved three cases. We manually checked the reported solutions and confirmed that they are all correct and nontrivial. Besides, the results also demonstrate some interesting findings about these tools. For example, Angr only supports one-leveled arrays but not two-leveled arrays; Triton does not even support the `atoi` function. Most of our findings are unavailable from the tool websites or existing papers, which further justifies the value of our benchmark approach.

The rest of the paper is organized as follows. We first discuss the related work in Section 2 and introduce the preliminary knowledge of symbolic execution in Section 3. Then, we demonstrate our systematic study about the challenges of symbolic execution in Section 4, and we introduce our benchmark methodology in Section 5. We present our experiments and results in Section 6. Finally, we conclude this paper in Section 7.

## 2 RELATED WORK

This section compares our work with present papers that either systematize the challenges of symbolic execution or employ the challenges to evaluate symbolic execution tools. Note that although symbolic execution has received extensive attention for decades, only a few papers include a systematic discussion about the challenges of the technique. Existing work in this area mainly focuses on employing the technique to carry out specific software analysis tasks (e.g., [13], [14], [15]), or proposing new approaches to improve the technology concerning particular challenges (e.g., [16], [17], [18]).

The papers that focus on systematizing the challenges of symbolic execution tools include [8], [9], [10], [19]. Kannavara *et al.* [10] enumerated several challenges that may hinder the adoption of symbolic execution in industrial fields. Qu and Robinson [8] conducted a case study on the limitations of symbolic testing tools and examined their prevalence in real-world programs. However, none of the two papers provides a method to evaluate symbolic execution tools. Cseppento and Micskei [9] proposed several metrics to evaluate source-code-based symbolic execution tools. But their metrics are based on specific program syntax of object-oriented codes rather than the language-independent challenges. These metrics are not very general for symbolic execution. Banescu *et al.* [19] also design several small programs for evaluation. But their purpose is to evaluate the resilience of code obfuscation transformations against symbolic execution-based attacks. They do not investigate the capability of symbolic execution but trust KLEE as a state-of-the-art symbolic executor. Besides, there are several survey papers (e.g., [6], [7], [12]) which also include some discussion about the challenges. But our work provides a more complete list of challenges as shown in Figure 1.

In our previous conference paper [11], we have conducted an empirical study with some of the challenges. This paper notably extends our previous paper with a formal benchmark methodology. It serves as a pilot study on systematically benchmarking symbolic execution tools in handling particular challenges. We consequently design a novel benchmark framework based on logic bombs, which can facilitate the automation of the benchmark process. We further provide a benchmark toolset that can be easily deployed by ordinary users.



TABLE 1  
A list of the challenges faced by symbolic execution, and the symbolic reasoning stages they attack.

Challenge		Idea	Stage of Error		
			$S_{var}$	$S_{inst} \& S_{sem}$	$S_{model}$
Symbolic -reasoning Challenges	Sym. Var. Declaration	Contextual variables besides program arguments	✓	✓	✓
	Covert Propagations	Propagating symbolic values in covert ways	-	✓	✓
	Buffer Overflows	Writing symbolic values without proper boundary check	-	✓	✓
	Parallel Executions	Processing symbolic values with parallel codes	-	✓	✓
	Symbolic Memories	Symbolic values as the offset of memory	-	✓	✓
	Contextual Symbolic Values	Retrieving contextual values with symbolic values	-	✓	✓
	Symbolic Jumps	Sym. values as the addresses of unconditional jump	-	-	✓
	Floating-point Numbers	Symbolic values in float/double type	-	-	✓
	Arithmetic Overflows	Integers outside the scope of an integer type	-	-	✓
Path-explosion Challenges	Loops	Change symbolic values within loops	-	-	-
	Crypto Functions	Processing symbolic values with crypto functions	-	-	-
	External Function Calls	Processing sym. values with some external functions	-	-	-

instructions along with a path on the program control-flow graph (CFG). Dynamic symbolic execution is also known as concolic (concrete and symbolic) execution. It collects instructions which are actually executed. In each round, the concolic execution engine executes the program with concrete values to generate instructions.

We may also classify symbolic execution tools into source-code-based symbolic execution and binary-code-based symbolic execution. In general, we do not perform symbolic reasoning on source codes or binaries directly. A prior step is to interpret the semantics of the program with an intermediate language (IL). For source codes, we can translate the code directly with a compiler frontend. For binaries, we have to lift the assembly codes into IL, which is more difficult. Their main difference lies in the translation process.

## 4 CHALLENGES OF SYMBOLIC EXECUTION

Based on whether a challenge attacks the symbolic reasoning process, we categorize the challenges of symbolic execution into *symbolic-reasoning challenges* and *path-explosion challenges*. A symbolic-reasoning challenge attacks the symbolic reasoning process and leads to incorrect test cases generated. A path-explosion challenge happens when there are too many paths to analyze. It does not attack a single symbolic reasoning process, but it may starve the computational resources or requires a very long time for symbolic execution.

Table 1 demonstrates the challenges that we have investigated in this work. We collect such challenges via a careful survey of existing papers. The survey scope covers several survey papers about symbolic execution techniques (e.g., [6], [7], [12]), several investigations that focus on systemizing the challenges of symbolic execution (e.g., [9], [10]), and other important papers related to symbolic execution (e.g., [15], [16], [22], [23], [24], [25], [26], [27]).

### 4.1 Symbolic-reasoning Challenges

Next, we discuss nine challenges that may incur errors to the symbolic reasoning process.

#### 4.1.1 Symbolic Variable Declarations

Test cases are the solutions of symbolic variables subject to constrain models. Therefore, symbolic variables should be

declared before a symbolic analysis process. For example, in source-code-based symbolic execution tools (e.g., KLEE), users can manually declare symbolic variables in the source codes. Binary-code-based concolic execution tools (e.g., Triton) generally assume a fixed length of program arguments from stdin as the symbolic variable. If some symbolic variables are missing from the declaration, the generated test cases would be insufficient for triggering particular control paths. Since the root cause occurs before symbolic execution, the challenge attacks  $S_{var}$ .

Figure 3(a) is a sample with a symbolic variable declaration problem. It returns a BOMB\_ENDING only when being executed with a particular process id. To explore the path, a symbolic execution tool should treat `pid` as a symbolic variable and then solve the constraint with respect to `pid`. Otherwise, it cannot find test cases that can trigger the path.

To declare symbolic variables precisely, a user should know target programs well. However, the task is impossible when analyzing programs on a large scale, e.g., when performing malware analysis. In an ideal case, a symbolic execution tool may automatically detect such variables which can control program behaviors and reports the solutions accordingly. To our best knowledge, non-existent tools have implemented the ideal feature. Instead, they are generally discussed together with other problems related to the computing environment, such as libraries, kernels, and drivers [28]. In reality, there are several challenges of this work refer to the computing environment, such as contextual symbolic variables, covert propagations, parallel executions, external function calls. We demonstrate that these challenges are different.

#### 4.1.2 Covert Propagations

Some data propagation ways are covert because they cannot be traced easily by data-flow analysis tools. For example, if the symbolic values are propagated via other media (e.g., files) outside of the process memory, the propagation would be untraceable. Such propagation methods are undecidable and can be beyond the capability of pure program analysis. Symbolic execution tools have to handle such cases with ad hoc methods. There are also some propagations challenging only to certain implementations. For example, propagating symbolic values via embedded assembly codes should be a problem for source-code-based symbolic execution tools only. If a symbolic execution tool fails to detect some prop-

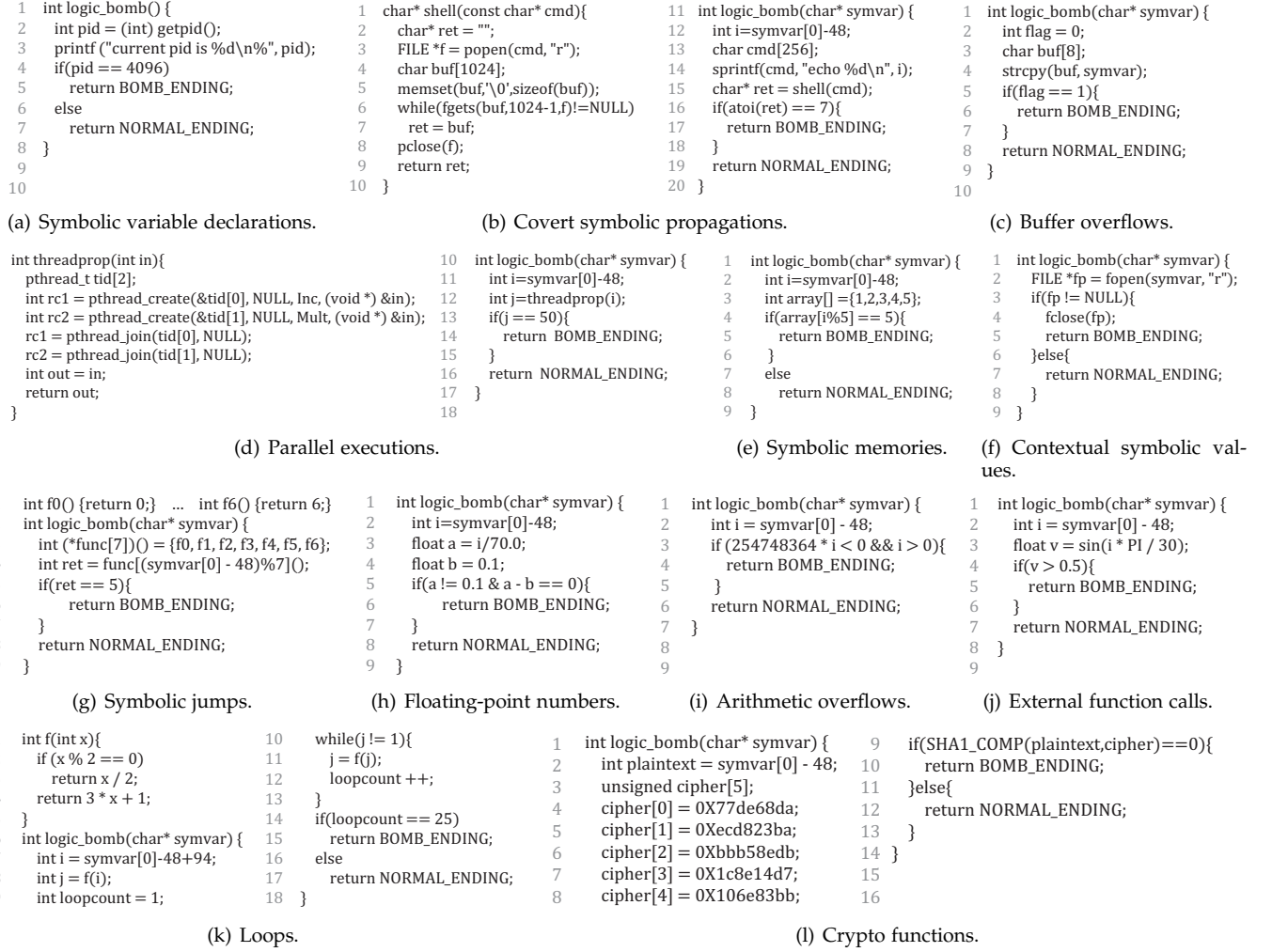


Fig. 3. Logic bomb samples with challenging symbolic execution issues. In each sample, we employ `symvar` to denote a symbolic variable, and `BOMB_ENDING` to denote a macro value indicating a particular program behavior.

agation, the instructions related to the propagated values would be missed from the following analysis. Therefore, the challenge attacks the stages of  $S_{inst}$  and  $S_{sem}$ .

Figure 3(b) shows a covert propagation sample. We define an integer `i` and initiate it with the value of a symbolic variable (`symvar`). So `i` is also a symbolic variable. We then propagate the value of `i` to another variable (`ret`) through a shell command (`echo`), and let `ret` control the return value. To find a test case which can return the `BOMB_ENDING`, a symbolic execution tool should properly track or model the propagation incurred by the shell command.

#### 4.1.3 Buffer Overflows

Buffer overflow is a typical software bug that can bring security issues. Due to insufficient boundary check, the input data may overwrite adjacent memories. Adversaries can employ such bugs to inject data and intentionally tamper the semantics of the original codes. Buffer overflows can happen in either stack or heap regions. If a symbolic execution tool cannot detect the overflow issues, it would fail to track the propagation of symbolic values. Therefore, buffer overflow involves a particular covert propagation issue. Source-code-based symbolic execution tools are prone

to be affected by buffer overflows because the stack layout of a program only exists in assembly codes, which may vary for particular platforms. Therefore, such tools cannot model the stack information with source codes only. In contrast, binary-code-based symbolic execution tools should be more potent in handling buffer overflow issues because there can simulate actual memory layouts. However, even if these tools can precisely track the propagation, they still suffer difficulties in automatically analyzing the unexpected program behaviors caused by overflow. Otherwise, they would be powerful enough to generate exploits for bugs, which is a problem far from being solved [29].

Figure 3(c) demonstrates a buffer overflow example. The program returns a `BOMB_ENDING` if the value of `flag` equals one, which is unlikely because the value is zero and should remain unchanged without explicit modification. However, the program has a buffer overflow bug. It has a buffer (`buf`) of eight bytes and employs no boundary check when copying symbolic values to the buffer with `strcpy`. We can change the value of `flag` to one leveraging the bug, e.g., when `symvar` is `"ANYSTRIN\x01\x00\x00\x00"`.



#### 4.1.4 Parallel Executions

Classic symbolic execution is effective for sequential programs. We can draw an explicit CFG for sequential programs and let a symbolic execution engine traverse the CFG. However, if a program processes symbolic variables in parallel, classic symbolic execution techniques would suffer problems. Parallel programs can be undecidable because the execution order of parallel codes does not only depend on the program but may also depend on the execution context. A parallel program may exhibit different behaviors even with the same test case. This poses a problem for symbolic execution to generate test cases for triggering corresponding control flows. If a symbolic execution tool directly ignores the parallel syntax or addresses the syntax improperly, errors would happen during  $S_{inst}$  and  $S_{sem}$ .

Figure 3(d) demonstrates an example with parallel codes. The symbolic variable `i` is processed by another two additional threads in parallel, and the result is assigned to `j`. Then the value of `j` determines the whether the program should return a `BOMB_ENDING`.

To handle parallel codes, a symbolic execution tool has to interpret the semantics and track parallel executions, *e.g.*, by introducing extra symbolic variables [30]. However, such an approach may not be scalable because the possibility of parallel execution can be a large number. In practice, there are several heuristic approaches to improve the efficiency. For example, we may restrict the exploration time of concurrent regions with a threshold [30]; we may conduct symbolic execution with arbitrary contexts and convert multi-thread programs into equivalent sequential ones [31]; or we can prune unimportant paths leveraging some program codes, such as assertion [32].

#### 4.1.5 Symbolic Memories

Symbolic memory is a situation whereas symbolic variables serve as the offsets or pointers to retrieve values from the memory, such as array indexes. To handle symbolic memories, a symbolic execution engine should take advantage of the memory layout for analysis. For example, we can convert an array selection operation to a `switch/case` clause in which the number of possible cases equals the length of the array. However, the number of possible combinations would grow exponentially when there are several such operations along a control flow. In practice, a symbolic execution tool may directly employ the feature of array operations implemented by some constraint solvers, such as STP [33] and Z3 [34]. It may also analyze the alignment of some pointers in advance, such as CUTE [35]. However, the power of pointer analysis is limited because the problem can be NP-hard or even undecidable for static analysis [36]. If a symbolic execution tool cannot model symbolic memories properly, errors would occur during  $S_{inst}$  and  $S_{sem}$ .

Figure 3(e) demonstrates a sample of symbolic memories. In this example, the symbolic variable `i` serves as an offset to retrieve an element from the array. The retrieved element then determines whether the program returns a `BOMB_ENDING`.

#### 4.1.6 Contextual Symbolic Values

The challenge is similar to symbolic memories but more complicated. Other than retrieving values from the memory

like symbolic memories, symbolic values can also serve as the parameters to retrieve values from the environment, such as loading the contents of a file pointed by symbolic values. By default, this contextual information is unavailable to the program or process, and the analysis is complicated. Moreover, since the contextual information can be changed any time without informing the program, the problem is undecidable. A symbolic tool that does not support such operations would cause errors during  $S_{inst}$  and  $S_{sem}$ .

Figure 3(f) is an example of contextual symbolic values. If `symvar` points to an existed file on the local disk, the program would return a `BOMB_ENDING`.

#### 4.1.7 Symbolic Jumps

In general, symbolic execution only extracts constraint models when encountering conditional jumps, such as `var < 0` in source codes, or `jle 0x400fda` in assembly codes. However, we may also employ unconditional jumps to achieve the same effects as conditional jumps. The idea is to jump to an address controlled by symbolic values. If a symbolic execution engine is not tailored to handle the unconditional jumps, it would fail to extract corresponding constraint models and miss some available control flows. Therefore, the challenge attacks the constraint modeling stage  $S_{model}$ .

Figure 3(g) is an example of symbolic jumps. The program contains an array of function pointers, and each function returns an integer value. The symbolic variable serves as an offset to determine which function should be called during execution. If `f5()` is called, the program would return a `BOMB_ENDING`.

#### 4.1.8 Floating-point Numbers

A floating-point number ( $f \in \mathbb{F}$ ) approximates a real number ( $r \in \mathbb{R}$ ) with a fixed number of digits in the form of  $f = sign * base^{exp}$ . For example, the 32-bit float type compliant to IEEE-754 has 1-bit for *sign*, 23-bit for *base*, and 8-bit for *exp*. The representation is essential for computers, as the memory spaces are limited in comparison with the infinity of  $\mathbb{R}$ . As a tradeoff, floating-point numbers only have limited precision, which makes some unsatisfiable constraints over  $\mathbb{R}$  to be satisfied over  $\mathbb{F}$  with a rounding mode. In order to support reasoning over  $\mathbb{F}$ , a symbolic execution tool should consider such approximations when extracting and solving constraint models. However, recent studies (*e.g.*, [5], [37], [38], [39]) show that there is still no silver bullet for the problem. Floating-point numbers remains a challenge for symbolic execution tools, and the challenge attacks  $S_{model}$ .

Figure 3(h) demonstrates an example with floating-point operations. Because we cannot represent 0.1 with float type precisely, the first predicate `a != 1` is always true. If the second condition `a == b` can be satisfied, the program would return a `BOMB_ENDING`. Therefore, one test case to returning a `BOMB_ENDING` is `symvar` equals '7'.

#### 4.1.9 Arithmetic Overflows

Arithmetic overflow happens when the result of an arithmetic operation is outside the range of an integer type. For example, the range of a 64-bit signed integer is  $[-2^{64}, 2^{64} - 1]$ . In this case, a constraint model (*e.g.*, the result of a positive integer plus another positive integer is negative) may

have no solutions over  $\mathbb{R}$ ; but it can have solutions when we consider arithmetic overflow. Handling such arithmetic overflow issues is not as difficult as previous challenges. However, some preliminary symbolic execution tools may fail to consider these cases and suffer errors when extracting and solving the constraint models.

Figure 3(i) shows a sample with an arithmetic overflow problem. To meet the first condition  $254748364 * i < 0$ ,  $i$  should be a negative value. However, the second condition requires  $i$  to be a positive value. Therefore, it has no solutions in the domain of real numbers. But the conditions can be satisfied when  $254748364 * i$  exceeds the max value that the integer type can represent.

## 4.2 Path-explosion Challenges

Now we discuss three path-explosion challenges existed in small-size programs.

### 4.2.1 External Function Calls

Shared libraries, such as `libc` and `libm` (i.e., a math library), provide some basic function implementations to facilitate software development. An efficient way to employ the functions is via dynamic linkage, which does not pack the function body to the program but only links with the functions dynamically when execution. Therefore, such external functions do not enlarge the size of a program, but they can enlarge the code complexity in nature.

When an external function call is related to the propagation of symbolic values, the control flows within the function body should be analyzed by default. There are two situations. A simple situation is that the external function does not affect the program behaviors after executing it, such as simply printing symbolic values with `printf`. In this case, we may ignore the path alternatives within the function. However, if the function execution affects the follow-up program behaviors, we should not ignore them. Otherwise, the symbolic execution would be based on a wrong assumption that the new test case generated for an alternative path can always trigger the same control flow within the external function. If a small program contains several such function calls, the complexity of external functions may cause path explosion issues. In practice, there are different strategies that symbolic execution tools may adopt with a trade-off between consistency and efficiency [28].

Figure 3(j) demonstrates a sample with an external function call. It computes the sine of a symbolic value via an external function call (i.e., `sin`), and the result is used to determine whether the program should return a `BOMB_ENDING`.

### 4.2.2 Loops

Loop statements, such as `for` and `while`, are widely employed in real-world programs. Even a very small program with loops can include many or even an infinite number of paths. By default, a symbolic execution tool should explore all available paths of a program, which can beyond the capability of the tool if there are too many paths. In practice, a symbolic execution tool may employ a search strategy which favorites the unexplored branches on a program CFG [18], [40], or introduces new symbolic variables as

the counters for each loop [41]. Because loop can incur numerous paths, we can hardly have a perfect solution for this problem.

Figure 3(k) shows a sample with a loop. The loop function is implemented with the Collatz conjecture [42]. No matter what is the initial value of  $i$ , the loop will terminate with  $j$  equals 1.

### 4.2.3 Crypto Functions

Crypto functions generally involve some computationally complex problems to ensure security. For a hash function, the complexity guarantees that adversaries cannot efficiently compute the plaintext of a hash value. For a symmetric encryption function, it promises that one cannot efficiently compute the key when given several pairs of plaintext and ciphertext. Therefore, such programs should also be resistant to symbolic execution attacks. From a program analysis view, the number of possible control paths for the crypto functions can be substantial. For example, the body of the SHA1 algorithm [43] is a loop that iterates 80 rounds, and each round contains several bit-level operations.

Figure 3(l) demonstrates a code snippet which employs a SHA1 function [43]. If the hash result of the symbolic value is equivalent to a predefined value, the program would return a `BOMB_ENDING`. However, it is difficult since SHA1 cannot be reversely calculated.

In general, symbolic execution tools cannot handle such crypto programs. Malware may employ the technique to deter symbolic execution-based program analysis [44]. When analyzing programs with crypto functions, a common way is to avoid exploring the function internals (e.g., [45], [46]). For example, TaintScope [45] first discriminates the symbolic variables corresponding to crypto functions from other variables, and then it employs a fuzzy-based approach to search solutions for such symbolic variables rather than solving the problem via symbolic reasoning.

So far, we have discussed 12 different challenges in total. Note that we do not intend to propose a complete list of challenges for symbolic execution. Instead, we collect all the challenging issues that have been mentioned in the literature and systematically analyze them. This analysis is essential for us to design the dataset of logic bombs in Section 5.2.2.

## 5 BENCHMARK METHODOLOGY

In this section, we introduce our methodology and a framework to benchmark the capability of real-world symbolic execution tools.

### 5.1 Objective and Challenges

Before describing our approach, we first discuss our design goal and the challenges to overcome.

This work aims to design an approach that can benchmark the capabilities of symbolic execution tools. Our purpose is critical and valid in several aspects. As we have discussed, some challenging issues are only engineering issues, such as arithmetic overflows. With enough engineering effort, a symbolic execution tool should be able to handle

**Algorithm 1: Method to design evaluation samples.**


---

```

// Create a function with a symbolic
// variable
Function LogicBomb(symvar)
    // symvar2 is a value computed from a
    // challenging problem related to symvar
    symvar2 ← Challenge(symvar);
    // If symvar2 satisfies a condition
    if Condition(symvar2) then
        // Trigger the bomb
        Bomb();
    end
end

```

---

these issues. On the other hand, some challenges are hard from a theoretical view, such as loops. But some heuristic approaches can tackle certain easy cases. Symbolic execution tools may adopt different heuristics and demonstrate different capabilities in handling them. Therefore, it is worth benchmarking their performances in handling particular challenging issues. Developers generally do not provide much information about the limitations of their tools to users.

A useful benchmark approach should be accurate and efficient. However, it is challenging to benchmark symbolic execution tools accurately and efficiently. Firstly, a real program contains many instructions or lines of codes. When a symbolic execution failure happens, locating the root cause requires much domain knowledge and effort. Since errors may propagate, it is challenging to conjecture whether a symbolic execution tool fails in handling a particular issue. Secondly, the symbolic execution itself is inefficient. Benchmarking a symbolic execution tool generally implies performing several designated symbolic execution tasks, which would be time-consuming. Note that existing symbolic execution papers (e.g., [2], [3], [47], [48]) generally evaluate the performance of their tools by conducting symbolic execution experiments with real programs. The process usually takes several hours or even days. They demonstrate the effectiveness of their work using the achieved code coverage and number of bugs detected, and they do not analyze the root causes of uncovered codes.

## 5.2 Approach based on Logic Bombs

To tackle the challenges of benchmark concerning accuracy and efficiency, we propose an approach based on logic bombs. Below, we discuss our detailed design.

### 5.2.1 Evaluation with Logic Bombs

A logic bomb is a code snippet that can only be executed when certain conditions are met. To evaluate whether a symbolic execution tool can handle a challenge, we can design a logic bomb guarded by a particular issue with the challenge. Then we can perform symbolic execution on the program which embeds the logic bomb. If a symbolic execution tool can generate a test case that can trigger the logic bomb, it indicates the tool can handle the challenging issue, or *vice versa*.

Algorithm 1 demonstrates a general framework to design such logic bombs. It includes four steps: the first step is to create a function with a parameter *symvar* as the

symbolic variable; the second step is to design a challenging problem related to the symbolic variable and save the result to another variable *symvar2*; the third step is to design a condition related to the new variable *symvar2*; the final step is to design a bomb (e.g., return a specific value) which indicates the condition has been satisfied. Note that because the value of *symvar2* is propagated from *symvar*, *symvar2* is also a symbolic variable and should be considered in the symbolic analysis process.

The magic of the logic bomb idea enables us to make the evaluation much precise and efficient. We can create several such small programs; each contains only a challenging issue and a logic bomb that tells the evaluation result. Because the object programs for symbolic execution are small, we can easily avoid unexpected issues that may also cause failures via a careful design. Also because the programs are small, performing symbolic execution on them generally requires a short time. For the programs that unavoidably incur path explosion issues, we can restrict the symbolic execution time either by controlling the problem complexity or by employing a timeout setting.

### 5.2.2 Logic Bomb Dataset

Following Algorithm 1, we have designed a dataset of logic bombs to evaluate the capability of symbolic execution tools. We have already shown several samples in Figure 3. Our full dataset is available on GitHub<sup>1</sup>. The dataset contains over 60 logic bombs for 64-bit Linux platform, which covers all the challenges discussed in Section 4. For each challenge, we implemented several logic bombs. Either each bomb involves a unique challenging issue (e.g., covert propagation via file write/read or via system calls), or introduces a problem with a different complexity setting (e.g., one-leveled arrays or two-leveled arrays).

When designing the logic bombs, we carefully avoid trivial test cases (e.g., `\x00`) that can trigger the bombs. Moreover, we try to employ straightforward implementations, and we hope to ensure that the results would not be affected by other unexpected failures. For example, we avoid using `atoi` to convert `argv[1]` to integers because some tools cannot support `atoi`. However, fully avoiding external function calls is impossible for some logic bombs. For example, we should employ external function calls to create threads when designing parallel codes. Surely that if a symbolic execution tool cannot handle external functions, the result might be affected. To tackle the interference of challenges, we draw a challenge propagation chart among the logic bombs as shown in Figure 4. There are two kinds of challenge propagation relationships: *should* in solid lines, and *may* in dashed lines. A *should* relationship means a logic bomb contains a similar challenging issue in another logic bomb; if a tool cannot solve the precedent logic bomb, it should not be able to solve the later one. For example, the `stackarray_sm_l1` is precedent to `stackarray_sm_l2`. A *may* relationship means a challenge type may be a precedent to other logic bombs, but it is not the determinant one. For example, a parallel program generally involves external function calls. However, although a tool is unable to solve

1. [https://github.com/hxuhack/logic\\_bombs](https://github.com/hxuhack/logic_bombs)



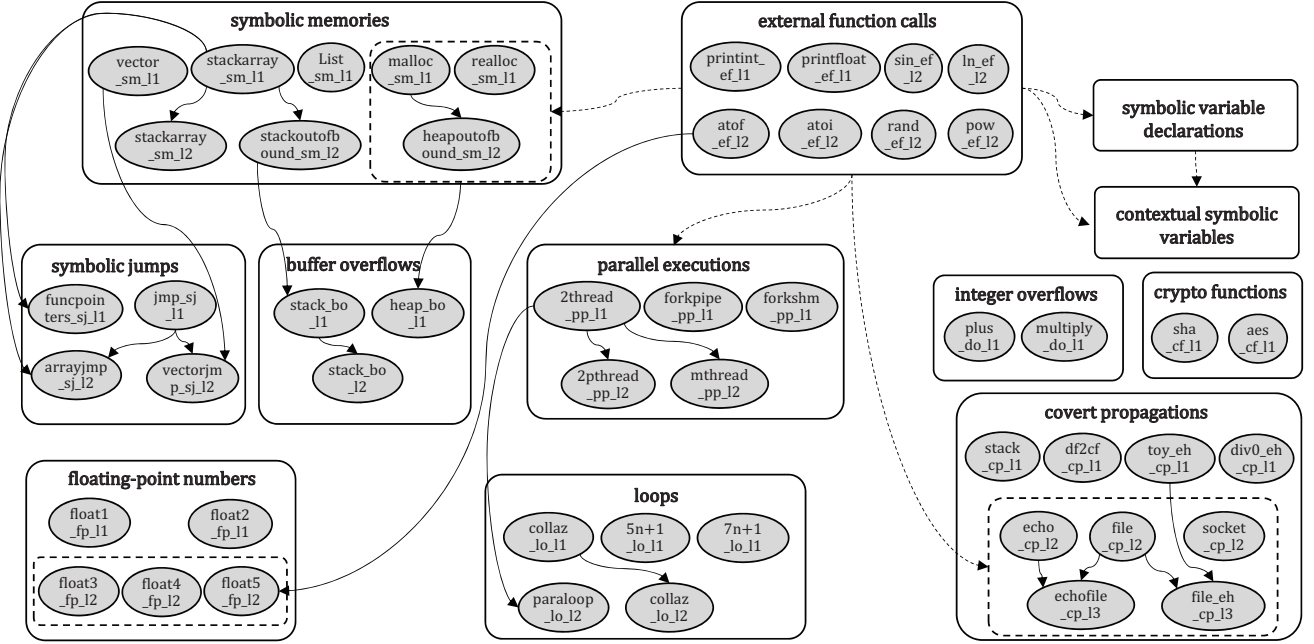


Fig. 4. The challenge propagation relationship among our dataset of logic bombs. A solid line means a logic bomb contains a similar problem defined in another logic bomb; a dashed line means a challenge may affect other logic bombs.

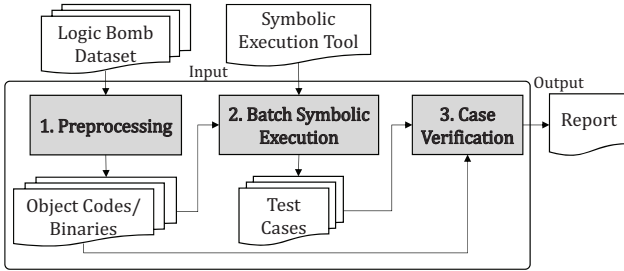


Fig. 5. A framework to benchmark symbolic execution tools.

the external functions well, it might be able to solve some logic bombs with parallel issues as sequential programs.

### 5.3 Automated Benchmark Framework

Base on the evaluation idea with logic bombs, we design a benchmark framework as shown in Figure 5. The framework inputs a dataset of carefully designed logic bombs and outputs the benchmark result for a particular symbolic execution tool. There are three critical steps in the framework: dataset preprocessing, batch symbolic execution, and case verification.

In the preprocessing step, we parse the logic bombs and compile them into object codes or binaries such that a target symbolic execution tool can process them. The parsing process pads each code snippet of a logic bomb with a main function and makes it a self-contained program. By default, we employ `argv[1]` as the symbolic variables. If a target symbolic execution tool requires adding extra instructions to launch tasks, the parser should add such required instructions automatically. For example, we can add symbolic variable declaration codes when benchmarking KLEE. The compilation process compiles the processed

source codes into binaries or other formats that a target symbolic execution tool supports. Symbolic execution is generally performed based on intermediate codes. When benchmarking source-code-based symbolic execution tools such as KLEE, we have to compile the source codes into the supported intermediate codes. When benchmarking binary-code-based symbolic execution tools, we can directly compile them into binaries, and the tool will lift binary codes into intermediate codes automatically.

In the second step, we direct a symbolic execution tool to analyze the compiled logic bombs in a batch mode. This step outputs a set of test cases for each program. Some dynamic symbolic execution tools (e.g., Triton) can directly tell which test case can trigger a logic bomb during runtime. However, other static symbolic execution tools may only output test cases by default, and we need to replay the generated test cases to examine the results further. Besides, some tools may falsely report that a test case can trigger the logic bomb. Therefore, we need a third step to verify the test cases.

In the third step, we replay the test cases with the corresponding programs of logic bombs. If a logic bomb can be triggered, it indicates that the challenging case is solved by the tool. Finally, we can generate a benchmark report based on the case verification results.

## 6 EXPERIMENTAL STUDY

In this section, we conduct an experimental study to demonstrate the effectiveness of our benchmark approach. Below, we discuss the experimental setting and results.

### 6.1 Experimental Setting

We choose three popular symbolic execution tools for benchmark: KLEE [2], Angr [3], and Triton [49]. Because our dataset of logic bombs are written in C/C++, we only

choose symbolic execution tools for C/C++ programs or binaries. The three tools are all released as open source and have a high community impact. Moreover, they adopt different implementation techniques for symbolic execution. By supporting variant tools, we show that our approach is compatible with different symbolic execution implementations.

KLEE [2] is a static symbolic execution tool implemented based on LLVM [50]. It requires program source codes to perform symbolic execution. By default, our benchmark script employs a `klee_make_symbolic` function to declare the symbolic variables of logic bombs in the source-code level. Then, it compiles the source codes into intermediate codes for symbolic execution. The symbolic execution process outputs a set of test cases, and our script finally examines the test cases by replaying them with the binaries. The whole process is automated with our benchmark script. The version of KLEE we benchmark is 1.3.0. Note that because our experiment does not intend to find the best tool for particular challenges, so we do not consider the patches provided by other parties before they are merged into the master branch.

Triton [49] is a dynamic symbolic execution tool based on binaries. It automatically accepts symbolic variables from the standard input. During symbolic execution, Triton firstly runs the programs with concrete values and leverages Intel PinTool [51] to trace related instructions, then it lifts the traced instructions into the SSA (single static assignment) form and performs symbolic analysis. If there are alternative paths found in the trace, Triton generates new test cases via symbolic reasoning and employs them as the concrete values in the following rounds of concrete execution. This symbolic execution process goes on until no alternative path can be found. The version of Triton we adopted is which released on Jul 6, 2017, on GitHub.

Angr [3] is also a tool for binaries but employs different implementations. Before performing any symbolic analysis, Angr firstly lifts the binary program into VEX IR [52]. Then it employs a symbolic analysis engine (*i.e.*, SimuVEX) to analyze the program based on the IR. Angr does not provide ready-to-use symbolic execution script for users but only some APIs. Therefore, we have to implement our own symbolic execution script for Angr. Our script collects all the paths to the CFG leaf nodes and then solves the corresponding path constraints. Angr provides all the critical features via APIs, and we only assemble them. Finally, we check whether the generated test cases can trigger the logic bombs. In our experiment, we employ Angr with version 7.7.9.21.

Note that our benchmark scripts for these tools all follow the framework proposed in Figure 5. During the experiment, we employ our logic bomb dataset for evaluation. A tool can pass a test only if the generated solution can correctly trigger a logic bomb. We finally report which logic bombs can be triggered by the tools.

We conduct our experiments on an Ubuntu 14.04 X86\_64 system with Intel i5 CPU and 8G RAM. Because some symbolic execution tasks may take very long time, our tool allows users to configure a timeout threshold which ensures benchmark efficiency. However, the timeout mechanism may incur some false results if it is too short. To mitigate

the side effects, we adopt two timeout settings (60 seconds and 300 seconds) for each tool. In this way, we can observe the influence of the timeout settings and decide whether we should conduct more experiments with an increased timeout value.

## 6.2 Benchmark Results

### 6.2.1 Result Overview

Table 2 demonstrates our experimental results. We label the results with three options: pass, fail, and timeout. While ‘pass’ and ‘fail’ imply the symbolic execution has finished, ‘timeout’ implies our benchmark script has terminated the symbolic execution process when a timeout threshold is triggered.

From the results, we can observe that Angr has achieved the best performance with 21 cases solved when the timeout is 300 seconds. Comparatively, it only solved 16 cases when the timeout is 60 seconds. KLEE solved nine cases and the result remains the same with different timeout settings. Triton performs much worse with three cases solved. To further verify the correctness of our benchmark results, we compare our experimental results with the previously declared challenge propagation relationships in Figure 4. We find the results are all consistent. It justifies that our dataset can distinguish the capability of different symbolic-execution tools accurately and effectively.

The efficiency of our benchmark approach largely depends on the timeout setting. Note that Table 2 includes some timeout results, and they account for most of our experimental time. Although we try to keep each logic bomb as succinct as possible, our dataset still contains some complex problems or path explosion issues unavoidable. When the timeout value is 60 seconds, our benchmark process for each tool takes only dozens of minutes. When extending the timeout value to 300 seconds, the benchmark takes a bit longer time. However, the benefit is not very obvious, and only Angr can solve 5 more cases. Can the result get further improved with more time? We have tried another group of experiments with 1800 seconds timeout and the results remain unchanged. Therefore, 300 seconds should be a marginal timeout setting for our benchmark experiment. Considering that symbolic execution is computationally expensive, which may take several hours or even several days to test a program, our benchmark process is very efficient. We may further improve the efficiency by employing a parallel mode, such as assigning each process several logic bombs.

### 6.2.2 Case Study

Now we discuss the detailed benchmark results for each challenge. Firstly, there are several challenges that none of the tools can trigger even one logic bomb, including symbolic variable declarations, parallel executions, contextual symbolic values, loops, and crypto functions. For symbolic variable declaration challenge, because all the tools cannot recognize the expected symbolic variables automatically, they fail in modeling the conditions to trigger the logic bombs. The challenges of contextual symbolic values and crypto functions involve tough problems, and it can be expected that all the tools fail in handling them. However, it

TABLE 2

Experimental results on benchmarking three symbolic execution tools (KLEE, Triton, and Angr) in handling our logic bombs. Pass means the tool has successfully triggered the bomb; fail means the tool cannot find test cases to trigger the bomb; timeout means the tool cannot find test cases to trigger the bomb within a given period of time. For each tool, we adopt two timeout settings: 60 seconds and 300 seconds.

Challenge	Case ID	KLEE		Triton		Angr	
		t = 60s	t = 300s	t = 60s	t = 300s	t = 60s	t = 300s
Covert Propagations	df2cf_cp	pass	pass	fail	fail	pass	pass
	echo_cp	fail	fail	timeout	timeout	timeout	timeout
	echofile_cp	fail	fail	fail	fail	timeout	timeout
	file_cp	fail	fail	timeout	timeout	fail	fail
	socket_cp	fail	fail	fail	fail	fail	fail
	stack_cp	fail	fail	pass	pass	pass	pass
	file_eh_cp	fail	fail	fail	fail	timeout	pass
	div0_eh_cp	fail	fail	fail	fail	timeout	pass
Buffer Overflows	file_eh_cp	fail	fail	fail	fail	timeout	fail
	stack_bo_11	fail	fail	fail	fail	pass	pass
	heap_bo_11	fail	fail	fail	fail	fail	fail
Symbolic Memories	stack_bo_12	fail	fail	fail	fail	fail	fail
	malloc_sm_11	pass	pass	timeout	fail	pass	pass
	realloc_sm_11	pass	pass	fail	fail	pass	pass
	stackarray_sm_11	pass	pass	fail	fail	pass	pass
	list_sm_11	fail	fail	fail	fail	timeout	pass
	vector_sm_11	fail	fail	fail	fail	timeout	pass
	stackarray_sm_12	pass	pass	fail	fail	fail	fail
	stackoutofbound_sm_12	pass	pass	fail	fail	pass	pass
Symbolic Jumps	heapoutofbound_sm_12	fail	fail	timeout	fail	pass	pass
	funcpointer_sj_11	pass	pass	fail	fail	fail	fail
	jmp_sj_11	fail	fail	fail	fail	pass	pass
	arrayjmp_sj_12	fail	fail	fail	fail	fail	fail
Floating-point Numbers	vectorjmp_sj_12	fail	fail	fail	fail	timeout	pass
	float1_fp_11	fail	fail	fail	fail	pass	pass
	float2_fp_11	fail	fail	fail	fail	pass	pass
	float3_fp_12	fail	fail	fail	fail	timeout	timeout
	float4_fp_12	fail	fail	fail	fail	timeout	timeout
Arithmetic Overflows	float5_fp_12	fail	fail	fail	fail	timeout	timeout
	plus_do	pass	pass	pass	pass	pass	pass
External Function Calls	multiply_do	pass	pass	fail	fail	pass	pass
	printint_ef_11	fail	fail	pass	pass	pass	pass
	printfloat_ef_11	fail	fail	fail	fail	fail	fail
	atoi_ef_12	fail	fail	fail	fail	pass	pass
	atof_ef_12	fail	fail	fail	fail	timeout	timeout
	ln_ef_12	fail	fail	fail	fail	timeout	fail
	pow_ef_12	fail	fail	fail	fail	pass	pass
	rand_ef_12	fail	fail	timeout	timeout	fail	fail
Symbolic Variable Declarations	sin_ef_12	fail	fail	fail	fail	timeout	timeout
		7 cases, all fail					
Contextual Symbolic Values		4 cases, all fail					
Parallel Executions		5 cases, all fail					
Loops		5 cases, all fail					
Crypto Functions		2 cases, all fail					
pass #	62 cases	9	9	3	3	16	21

is a bit surprising that none of the tools can handle parallel executions and loops.

*Covert Propagations:* Angr have passed four test cases, `df2cf_cp`, `stack_cp`, and two exception handling cases. `df2cf_cp` propagates the symbolic values indirectly by substituting a data assignment operation with equivalent control-flow operations. KLEE also solved the case, but Triton failed. `stack_cp` propagates symbolic values via direct assembly instructions `push` and `pop`. Only KLEE failed the test because it is a source-based analysis tool which does not support assembly codes. Besides, Angr has also passed two test cases that propagate symbolic values via the C++ exception handling mechanism. Not only this exception handling mechanism can be covert for some source-code-based analysis tools, such as KLEE, but also it can be covert for binary-code-based symbolic execution tools, such as Triton. We further break down the details of an exception

handling program in Figure 6. As shown in the box region of Figure 6(b), the mechanism relies on two function calls, which might be the problem that fails Triton. All the tools failed other covert propagation cases that propagate values via file read/write, echo, socket, etc.

*Buffer Overflows:* Only Angr has solved one easy buffer overflow problem `stack_bo_11`. The case has a simple stack overflow issue. Its solution requires modifying the value of the stack that might be illegal. However, Angr cannot solve the heap overflow issue `heap_bo_11`. It also failed on another harder stack overflow issue `stack_bo_12`, which requires composing sophisticated payload, such as employing return-oriented programming methods [53]. We are surprised that Triton failed all the tests because binary-code-based symbolic execution tools should be resilient to buffer overflows in nature.

*Symbolic Memories:* The results show that Triton does not

```

double division(int numerator, int denominator) {
    if( denominator == 0 ) {
        throw "Division by zero condition!";
    }
    return (numerator/denominator);
}

int logic_bomb(char* s) {
    int symvar = s[0] - 48;
    try {
        division(10, symvar-7);
        return NORMAL_ENDING;
    } catch (const char* msg) {
        return BOMB_ENDING;
    }
}

```

(a) Source codes.

```

0x00000000000400aa9 <+41>:    callq 0x400a10 <_Z8divisionii>
0x00000000000400aae <+46>:    movsd %xmm0,0x40(%rbp)
0x00000000000400ab3 <+51>:    jmpq 0x400ab8 <_Z10logic_bombPc+56>
0x00000000000400ab8 <+56>:    movl $0x0,0x4(%rbp)
0x00000000000400abf <+63>:    jmpq 0x400afd <_Z10logic_bombPc+125>
0x00000000000400ac4 <+68>:    mov %edx,%ecx
.....
0x00000000000400ae1 <+97>:    callq 0x4008a0 <_cxa_begin_catch@plt>
0x00000000000400ae6 <+102>:   mov %rax,-0x30(%rbp)
0x00000000000400aea <+106>:   movl $0x1,-0x4(%rbp)
0x00000000000400af1 <+113>:   movl $0x1,-0x34(%rbp)
0x00000000000400af8 <+120>:   callq 0x400890 <_cxa_end_catch@plt>
0x00000000000400afd <+125>:   mov -0x4(%rbp),%eax
0x00000000000400b00 <+128>:   add $0x40,%rsp
0x00000000000400b04 <+132>:   pop %rbp
0x00000000000400b05 <+133>:   retq
0x00000000000400b06 <+134>:   mov -0x20(%rbp),%rdi
0x00000000000400b0a <+138>:   callq 0x4008c0 <_Unwind_Resume@plt>

```

(b) Assembly codes.

Fig. 6. An exemplary program that raises an exception when divided by zero. The assembly codes demonstrates how the `try/catch` mechanism works in low level.

```

int logic_bomb(char* s) {
    int symvar = s[0] - 48;
    int l1_ary[] = {1,2,3,4,5};
    int l2_ary[] = {6,7,8,9,10};

    int x = symvar%5;
}

```

(a) Source codes.

```

(gdb) break *logic_bomb+97
(gdb) run
(gdb) x/20xw $rsp-80
0x7fffffff2d0: 0x00000006 0x00000007 0x00000008 0x00000009
0x7fffffff2e0: 0x0000000a 0x0000000b 0x004003e5 0x0000000c
0x7fffffff2f0: 0x00000001 0x00000002 0x00000003 0x00000004
0x7fffffff300: 0x00000005 0x00007fff 0xf7fe2000 0x00000001
0x7fffffff310: 0xfffffe6db 0x00007fff 0x00000000 0x00000000

```

(b) Memory layout after array initialization.

```

text: 0x00000000004006d <+29>:    mov 0x4007c0,%rdi
text: 0x0000000000400615 <+37>:    mov %rdi,-0x30(%rbp)
text: 0x0000000000400619 <+41>:    mov 0x4007c8,%rdi
text: 0x0000000000400621 <+49>:    mov %rdi,-0x28(%rbp)
text: 0x0000000000400625 <+53>:    mov 0x4007d0,%ecx
text: 0x000000000040062c <+60>:    mov %ecx,-0x20(%rbp)
text: 0x000000000040062f <+63>:    mov 0x4007e0,%rdi
text: 0x0000000000400637 <+71>:    mov %rdi,-0x50(%rbp)
text: 0x000000000040063b <+75>:    mov 0x4007e8,%rdi
text: 0x0000000000400643 <+83>:    mov %rdi,-0x48(%rbp)
text: 0x0000000000400647 <+87>:    mov 0x4007f0,%ecx
text: 0x000000000040064e <+94>:    mov %ecx,-0x40(%rbp)
.....
.rodata:00000000004007C0 dq 200000001h
.rodata:00000000004007C8 dq 400000003h
.rodata:00000000004007D0 dd 5
.rodata:00000000004007D4 align 20h
.rodata:00000000004007E0 dq 700000006h
.rodata:00000000004007E8 dq 9000000008h
.rodata:00000000004007F0 dd 0Ah

```

(c) Assembly codes.

Fig. 7. An exemplary program that demonstrates how the stack works with arrays. There is no information about the size of each array left in assembly codes.

support symbolic memory, but KLEE and Angr provide very good support. Angr has solved seven cases out of eight. It only failed in handling the case (Figure 7(a)) with a two-leveled array `stackarray_sm_12`. Also, it implies that when there are multi-leveled pointers, Angr would fail. Figure 7(c) demonstrates the assembly codes that initialize the arrays, and Figure 7(b) demonstrates the stack layout after initialization. We can observe that the information about array size or boundary does not exist in assembly codes. This justifies why binary-code-based symbolic execution tools do not suffer problems when a challenge requires an out-of-boundary access, e.g., `stackoutofbound_sm_12`. In comparison, KLEE can solve the two-leveled array problem because it is based on STP [33], which is designed for solving such problems related to arrays. However, KLEE does not support C++, so it failed the problems with vectors and lists.

**Symbolic Jumps:** Since symbolic jump demonstrates no explicit conditional branches in the CFG, it should be a hard problem for symbolic execution. However, KLEE and Angr are not likely to be affected much by the trick. KLEE has tackled the problem with an array of function pointers `funcpointer_sj_11`. It failed the other test cases because they employ an assembly instruction `jmp`, which KLEE does not support. Angr successfully handled two cases with

assembly `jmp`, but it failed `funcpointer_sj_11`.

**Floating-point Numbers:** The results show KLEE and Triton do not support floating-point operations, and Angr can support some. During our test, Triton directly reported that it cannot interpret such floating-point instructions. Angr has solved two out of the five designated cases. The two passed cases are easier ones, which only require integer values as the solution. All the failed cases require decimal values as the solution, and they employ the `atof` function to convert `argv[1]` to decimals. Since Angr has also failed the test in handling `atof` in `atof_ef_12`, the failures are likely to be caused by the `atof` function.

**Arithmetic Overflows:** Arithmetic overflow is not a very hard problem, and it only requires symbolic execution tools to handle such cases carefully. In our test, KLEE and Angr have solved all the cases. However, Triton failed in handling the integer overflow case in Figure 3(i). The result shows there is still much room for Triton to improve for this problem.

**External Function Calls:** In this group of logic bombs, each case only contains one external function call. However, the result is very disappointing. Triton only passed a very simple case that print out (with `printf`) a symbolic value of integer type. It does not even support printing out floating-point values. Angr has solved the `printf` cases and two

more complicated cases, `atoi_ef_l2` and `pow_ef_l2`. It cannot support `atof_ef_l2` and other cases. The results show that we should be cautious when designing logic bombs. Even when involving straightforward external function calls, the results could be affected.

## 7 CONCLUSION

This work proposes an approach to benchmark the capability symbolic execution tools in handling particular challenges. To this end, we studied the taxonomy of challenges faced by symbolic execution tools, including nine symbolic-reasoning challenges and three path-explosion challenges. Such a study is essential for us to design the benchmark dataset. Then we proposed a promising benchmark approach based on logic bombs. The idea is to design logic bombs that can only be triggered if a symbolic execution tool solves specific challenging issues. By making the programs of logic bombs as small as possible, we can speed up the benchmark process; and by making them as straightforward as possible, we can avoid unexpected reasons that may affect the benchmark results. In this way, our benchmark approach is both accurate and efficient. Following the idea, we implemented a dataset of logic bombs and a prototype benchmark framework which automates the benchmark process. Then, we conducted real-world experiments on three symbolic execution tools. Experimental results show that the benchmark process for each tool generally takes dozens of minutes. Angr achieved the best benchmark results with 21 cases solved, KLEE solved nine, and Triton only solved three. These results justify the value of a third-party benchmark toolset for symbolic execution tools. Finally, we released our dataset as open source on GitHub for public usage. We hope it would serve as an essential tool for the community to benchmark symbolic execution tools and could facilitate the development of more comprehensive symbolic execution techniques.

## ACKNOWLEDGMENTS

This work was supported by the National Natural Science Foundation of China (Project Nos. 61672164, 61332010), and 2015 Microsoft Research Asia Collaborative Research Program (Project No. FY16-RES-THEME-005).

## REFERENCES

- [1] J. C. King, "Symbolic execution and program testing," *Communications of the ACM*, 1976.
- [2] C. Cadar, D. Dunbar, and D. R. Engler, "Klee: Unassisted and automatic generation of high-coverage tests for complex systems programs," in *Proc. of the 8th USENIX Symposium on Operating Systems Design and Implementation*, 2008.
- [3] Y. Shoshitaishvili and *et al.*, "Sok: (state of) the art of war: Offensive techniques in binary analysis," in *Proc. of the IEEE Symposium on Security and Privacy*, 2016.
- [4] D. A. Ramos and D. R. Engler, "Under-constrained symbolic execution: Correctness checking for real code," in *USENIX Security Symposium*, 2015.
- [5] M. Quan, "Hotspot symbolic execution of floating-point programs," in *Proc. of the 24th ACM SIGSOFT International Symposium on Foundations of Software Engineering (FSE)*, 2016.
- [6] C. Cadar and K. Sen, "Symbolic execution for software testing: three decades later," *Communications of the ACM*, 2013.
- [7] E. J. Schwartz, T. Avgerinos, and D. Brumley, "All you ever wanted to know about dynamic taint analysis and forward symbolic execution (but might have been afraid to ask)," in *Proc. of the 2010 IEEE Symposium on Security and Privacy*, 2010.
- [8] X. Qu and B. Robinson, "A case study of concolic testing tools and their limitations," in *Proc. of the IEEE International Symposium on Empirical Software Engineering and Measurement*, 2011.
- [9] L. Cseppento and Z. Micskei, "Evaluating symbolic execution-based test tools," in *Proc. of the IEEE 8th International Conference on Software Testing, Verification and Validation*, 2015.
- [10] R. Kannavara, C. J. Havlicek, B. Chen, M. R. Tuttle, K. Cong, S. Ray, and F. Xie, "Challenges and opportunities with concolic testing," in *Aerospace and Electronics Conference (NAECON), 2015 National. IEEE*, 2015, pp. 374–378.
- [11] H. Xu, Y. Zhou, Y. Kang, and M. R. Lyu, "Concolic execution on small-size binaries: challenges and empirical study," in *Proc. of the 47th Annual IEEE/IFIP International Conference on Dependable Systems and Networks*, 2017.
- [12] R. Baldoni, E. Coppa, D. C. D'Elia, C. Demetrescu, and I. Finocchi, "A survey of symbolic execution techniques," *ACM Computing Survey*, 2018.
- [13] T. Avgerinos, S. K. Cha, B. L. T. Hao, and D. Brumley, "Aeg: Automatic exploit generation," in *Proc. of the 2011 ACM the Network and Distributed System Security Symposium*, 2011.
- [14] J. Ming, D. Xu, L. Wang, and D. Wu, "Loop: Logic-oriented opaque predicate detection in obfuscated binary code," in *Proc. of the 22nd ACM SIGSAC Conference on Computer and Communications Security*, 2015.
- [15] B. Yadegari and S. Debray, "Symbolic execution of obfuscated code," in *Proc. of the 22nd ACM SIGSAC Conference on Computer and Communications Security*, 2015.
- [16] T. Xie, N. Tillmann, J. de Halleux, and W. Schulte, "Fitness-guided path exploration in dynamic symbolic execution," in *Proc. of the 39th IEEE/IFIP International Conference on Dependable Systems & Networks*, 2009.
- [17] L. Ciortea, C. Zamfir, S. Bucur, V. Chipounov, and G. Candea, "Cloud9: a software testing service," *ACM SIGOPS Operating Systems Review*, vol. 43, no. 4, pp. 5–10, 2010.
- [18] T. Avgerinos, A. Rebert, S. K. Cha, and D. Brumley, "Enhancing symbolic execution with veritesting," in *Proc. of the 36th International Conference on Software Engineering*. ACM, 2014.
- [19] S. Banescu, C. Collberg, V. Ganesh, Z. Newsham, and A. Pretschner, "Code obfuscation against symbolic execution attacks," in *Proceedings of the 32nd Annual Conference on Computer Security Applications (ACSAC)*. ACM, 2016.
- [20] C. A. R. Hoare, "An axiomatic basis for computer programming," *Communications of the ACM*, 1969.
- [21] D. Song, D. Brumley, H. Yin, J. Caballero, I. Jager, M. G. Kang, Z. Liang, J. Newsome, P. Poosankam, and P. Saxena, "Bitblaze: A new approach to computer security via binary analysis," in *Proc. of the International Conference on Information Systems Security*. Springer, 2008.
- [22] P. Godefroid, N. Klarlund, and K. Sen, "Dart: directed automated random testing," in *ACM Sigplan Notices*, 2005.
- [23] D. Brumley, I. Jager, T. Avgerinos, and E. J. Schwartz, "Bap: A binary analysis platform," in *Proc. of the International Conference on Computer Aided Verification*. Springer, 2011.
- [24] N. Razavi, F. Ivančić, V. Kahlon, and A. Gupta, "Concurrent test generation using concolic multi-trace analysis," in *Asian Symposium on Programming Languages and Systems*. Springer, 2012.
- [25] S. K. Cha, T. Avgerinos, A. Rebert, and D. Brumley, "Unleashing mayhem on binary code," in *Proc. of the 2012 IEEE Symposium on Security and Privacy*, 2012.
- [26] D. Davidson, B. Moench, T. Ristenpart, and S. Jha, "Fie on firmware: Finding vulnerabilities in embedded systems using symbolic execution," in *USENIX Security Symposium*, 2013.
- [27] S. Guo, M. Kusano, and C. Wang, "Conc-ise: Incremental symbolic execution of concurrent software," in *Proc. of the 31st IEEE/ACM International Conference on Automated Software Engineering (ASE)*, 2016.
- [28] V. Chipounov, V. Kuznetsov, and G. Candea, "S2e: A platform for in-vivo multi-path analysis of software systems," 2011.
- [29] T. Avgerinos, S. K. Cha, A. Rebert, E. J. Schwartz, M. Woo, and D. Brumley, "Automatic exploit generation," *Communications of the ACM*, 2014.

- [30] A. Farzan, A. Holzer, N. Razavi, and H. Veith, "Con2colic testing," in *Proc. of the 9th Joint Meeting on Foundations of Software Engineering (FSE)*. ACM, 2013.
- [31] T. Bergan, D. Grossman, and L. Ceze, "Symbolic execution of multithreaded programs from arbitrary program contexts," 2014.
- [32] S. Guo, M. Kusano, C. Wang, Z. Yang, and A. Gupta, "Assertion guided symbolic execution of multithreaded programs," in *Proc. of the 2015 10th Joint Meeting on Foundations of Software Engineering*. ACM, 2015.
- [33] V. Ganesh and D. L. Dill, "A decision procedure for bit-vectors and arrays," in *Proc. of the International Conference on Computer Aided Verification*. Springer, 2007.
- [34] L. De Moura and N. Björner, "Z3: An efficient smt solver," in *Proc. of the International Conference on Tools and Algorithms for the Construction and Analysis of Systems*. Springer, 2008.
- [35] K. Sen, D. Marinov, and G. Agha, "Cute: a concolic unit testing engine for c," in *ACM SIGSOFT Software Engineering Notes*, 2005.
- [36] W. Landi and B. G. Ryder, "Pointer-induced aliasing: a problem classification," in *Proc. of the 18th ACM SIGPLAN-SIGACT symposium on Principles of programming languages*, 1991.
- [37] A. Solovyev, C. Jacobsen, Z. Rakamarić, and G. Gopalakrishnan, "Rigorous estimation of floating-point round-off errors with symbolic taylor expansions," in *International Symposium on Formal Methods*. Springer, 2015.
- [38] D. Liew, D. Schemmel, C. Cadar, A. F. Donaldson, R. Zähl, and K. Wehrle, "Floating-point symbolic execution: A case study in n-version programming," in *Proceedings of the 32nd IEEE/ACM International Conference on Automated Software Engineering (ASE)*, 2017.
- [39] D. S. Liew, "Symbolic execution of verification languages and floating-point code," 2018.
- [40] J. Burnim and K. Sen, "Heuristics for scalable dynamic test generation," in *Proc. of the 23rd IEEE/ACM International Conference on Automated Software Engineering*, 2008.
- [41] P. Saxena, P. Poosankam, S. McCamant, and D. Song, "Loop-extended symbolic execution on binary programs," in *Proc. of the eighteenth International Symposium on Software Testing and Analysis*. ACM, 2009.
- [42] J. C. Lagarias, "The  $3x + 1$  problem and its generalizations," *The American Mathematical Monthly*, 1985.
- [43] D. Eastlake 3rd and P. Jones, "Us secure hash algorithm 1 (SHA1)," Tech. Rep., 2001.
- [44] M. I. Sharif, A. Lanzi, J. T. Giffin, and W. Lee, "Impeding malware analysis using conditional code obfuscation." in *NDSS*, 2008.
- [45] T. Wang, T. Wei, G. Gu, and W. Zou, "Taintscope: A checksum-aware directed fuzzing tool for automatic software vulnerability detection," in *IEEE Symposium on Security and Privacy (SP)*, 2010.
- [46] R. Corin and F. A. Manzano, "Efficient symbolic execution for analysing cryptographic protocol implementations." *ESSoS*, 2011.
- [47] V. Kuznetsov, J. Kinder, S. Bucur, and G. Candea, "Efficient state merging in symbolic execution," *ACM Sigplan Notices*, 2012.
- [48] N. Hasabnis and R. Sekar, "Extracting instruction semantics via symbolic execution of code generators," in *Proc. of the 2016 24th ACM SIGSOFT International Symposium on Foundations of Software Engineering*, 2016.
- [49] F. Soudel and J. Salwan, "Triton: a dynamic symbolic execution framework," in *Symposium sur la sécurité des technologies de l'information et des communications, SSTIC, France, Rennes*, 2015.
- [50] C. Lattner and V. Adve, "Llvm: A compilation framework for lifelong program analysis & transformation," in *Proc. of the International Symposium on Code Generation and Optimization: Feedback-Directed and Runtime Optimization*. IEEE Computer Society, 2004.
- [51] C.-K. Luk, R. Cohn, R. Muth, H. Patil, A. Klauser, G. Lowney, S. Wallace, V. J. Reddi, and K. Hazelwood, "Pin: building customized program analysis tools with dynamic instrumentation," in *ACM Sigplan Notices*, vol. 40, no. 6, 2005, pp. 190–200.
- [52] N. Nethercote, "Dynamic binary analysis and instrumentation," Ph.D. dissertation, PhD thesis, University of Cambridge, 2004.
- [53] R. Roemer, E. Buchanan, H. Shacham, and S. Savage, "Return-oriented programming: Systems, languages, and applications," *ACM Transactions on Information and System Security (TISSEC)*, 2012.