

Targeted Greybox Fuzzing with Static Lookahead Analysis

Valentin Wüstholz
ConsenSys Diligence, Germany
valentin.wustholz@consensys.net

Maria Christakis
MPI-SWS, Germany
maria@mpi-sws.org

Abstract—Automatic test generation typically aims to generate inputs that explore new paths in the program under test in order to find bugs. Existing work has, therefore, focused on guiding the exploration toward program parts that are more likely to contain bugs by using an offline static analysis.

In this paper, we introduce a novel technique for targeted greybox fuzzing using an *online* static analysis that guides the fuzzer toward a set of *target locations*, for instance, located in recently modified parts of the program. This is achieved by first *semantically* analyzing each program path that is explored by an input in the fuzzer’s test suite. The results of this analysis are then used to control the fuzzer’s specialized power schedule, which determines how often to fuzz inputs from the test suite. We implemented our technique by extending a state-of-the-art, industrial fuzzer for Ethereum smart contracts and evaluate its effectiveness on 27 real-world benchmarks. Using an online analysis is particularly suitable for the domain of smart contracts since it does not require any code instrumentation—adding instrumentation to contracts changes their semantics. Our experiments show that targeted fuzzing significantly outperforms standard greybox fuzzing for reaching 83% of the challenging target locations (up to 14x of median speed-up).

I. INTRODUCTION

Automatic test generation is known to help find bugs and security vulnerabilities, and therefore, improve software quality. As a result, there has emerged a wide variety of test-generation tools that implement techniques such as random testing [1], [2], [3] and blackbox fuzzing [4], [5], greybox fuzzing [6], [7] as well as dynamic symbolic execution [8], [9] and whitebox fuzzing [10], [11], [12].

These techniques differ from each other in how much of the program structure they take into account. In general, the more structure a testing tool may leverage, the more effective it becomes in discovering new paths, but the less efficient it is in generating new inputs. For example, greybox fuzzing lies in the middle of this spectrum between performance and effectiveness in increasing coverage. In particular, it uses lightweight runtime monitoring that makes it possible to distinguish different paths, but it may not access any additional information about the program under test.

What these techniques have in common is that, just like any (static or dynamic) path-based program analysis, they can usually only explore a subset of all feasible paths in a program under test; for instance, in the presence of input-dependent loops. For this reason, path-based program analyses are typically not able to prove the absence of errors in a program, only their existence.

To make bug detection more effective, existing work has focused on guiding the exploration toward warnings reported by a static analysis (e.g., [13], [14], [15]), unverified program executions (e.g., [16], [17]), or sets of dangerous program locations (e.g., [18]). The motivation behind these approaches is to identify safe program paths at compile time and avoid them at runtime. This is often achieved with an offline static analysis whose results are recorded and used to prune parts of the search space that is then explored by test generation.

The offline static analysis may be semantic, e.g., based on abstract interpretation, or not, e.g., based on the program text or its control-flow graph. A semantic analysis must consider all possible program inputs and states in which a piece of code may be executed. As a result, the analysis can quickly become imprecise, thus impeding its purpose of pruning as much of the search space as possible. For better results, one could resort to a more precise analysis, which would be less efficient, or to a more unsound analysis. The latter would limit the number of considered execution states in order to increase precision, but may also prune paths that are unsoundly verified [19].

Our approach. In this paper, we present a technique that *semantically* guides greybox fuzzing toward *target locations*, for instance, locations reported by another analysis or located in recently modified parts of the program. This is achieved with an *online* static analysis. In particular, the fuzzer invokes this online analysis right before adding a new input to its test suite. For the program path π that the new input explores (see bold path in Fig. 1), the goal of the analysis is to determine a path prefix π_{pre} for which all suffix paths are unable to reach a target location (e.g., T_x and T_y in Fig. 1). This additional information allows the fuzzer to allocate its resources more strategically such that more effort is spent on exercising program paths that might reach the target locations, thereby enabling *targeted fuzzing*. More precisely, this information feeds into a specialized power schedule of the fuzzer that determines how often to fuzz an input from the test suite.

We refer to our online static analysis as a *lookahead analysis* since, given a path prefix π_{pre} , it looks for reachable target locations along all suffix paths (sub-tree rooted at P_i in Fig. 1). We call the last program location of prefix π_{pre} a *split point* (P_i in Fig. 1). Unlike a traditional static analysis, the lookahead analysis does not consider all possible execution states at the split point when analyzing all suffix paths—only the ones that are feasible along π_{pre} . In other words, the

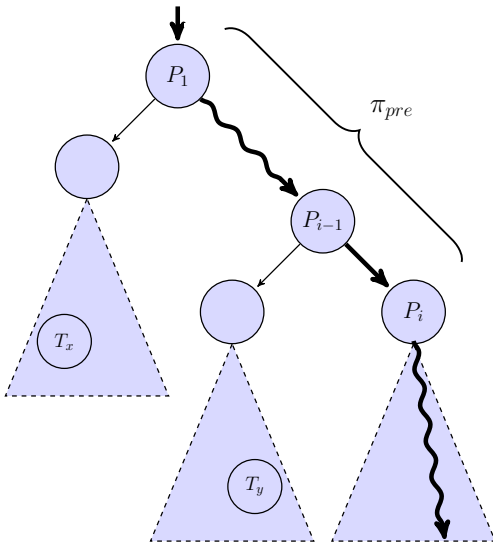


Fig. 1: Execution tree of a program containing target locations T_x and T_y . The lookahead analysis analyzes a path π (bold) to identify a prefix π_{pre} such that no suffix paths reach a target location.

lookahead analysis combines the precision of a path-sensitive analysis along a feasible path prefix with the scalability of a path-insensitive suffix analysis. Intuitively, for a given path π , the precision of the lookahead analysis is determined by the number of suffix paths that are proved not to reach any target locations. Therefore, to optimize precision, the analysis tries to identify the *first* split point (P_i in Fig. 1) along π such that all targets are unreachable. Note that the lookahead analysis may consider any program location along π as a split point.

When combining greybox fuzzing with an online lookahead analysis, we faced four main challenges, which we address in this paper. In particular, we provide answers to the following questions: (1) How can the lookahead analysis effectively communicate its results to the fuzzer? (2) How lightweight can the analysis be to improve the effectiveness of the fuzzer in reaching target locations without having a negative impact on its performance? (3) How can the analysis be invoked from a certain split point along a path? (4) What are suitable split points for invoking the analysis to check all suffix paths?

Our implementation uses HARVEY, a state-of-the-art, industrial greybox fuzzer for Ethereum smart contracts, which are programs managing crypto-currency accounts on a blockchain. We extended HARVEY to incorporate BRAN, a new static-analysis framework for smart contracts. A main reason for targeting the domain of smart contracts is that adding code instrumentation to contracts changes their semantics, and all existing techniques that use an offline static analysis require instrumentation of the program under test. Our experiments on 27 benchmarks show that targeted fuzzing significantly outperforms standard greybox fuzzing for reaching 83% of the challenging target locations (up to 14x of median speed-up).

Contributions. We make the following contributions:

- We introduce a greybox-fuzzing algorithm that uses a

lightweight, online static analysis and a specialized power schedule to guide the exploration toward target locations.

- We implement this fuzzing algorithm by extending the HARVEY greybox fuzzer with BRAN, a static analysis for smart contracts.
- We evaluate our technique on 27 real-world benchmarks and demonstrate that our lookahead analysis and power schedule significantly increase the effectiveness of greybox fuzzing in reaching target locations.

Outline. The next section provides background on greybox fuzzing and smart contracts. In Sect. III, we give an overview of our technique through an example. Sect. IV explains the technical details, and Sect. V describes our implementation. We present our experimental evaluation in Sect. VI, discuss related work in Sect. VII, and conclude in Sect. VIII.

II. BACKGROUND

In this section, we review background on greybox fuzzing and smart contracts.

A. Greybox Fuzzing

Greybox fuzzing [6], [7] is a practical test-generation technique that has been shown to be very effective in detecting bugs and security vulnerabilities (e.g., [20]). Alg. 1 shows exactly how it works. (The grey boxes should be ignored.)

A greybox fuzzer takes as input the program under test *prog* and a set of seed inputs S . The fuzzer runs the program with the seeds (line 1) and associates each input with the unique identifier of the path it exercises, or *PID*. The *PIDs* data structure, therefore, represents a map from a *PID* to the corresponding input. Note that a path identifier is computed with lightweight runtime monitoring that allows the fuzzer to distinguish different program paths.

Next, the fuzzer selects an input from *PIDs* for mutation (line 3), which is typically performed randomly. This input is assigned an “energy” value, which indicates how long it should be fuzzed (line 5). The input is then mutated (line 8), and the program is run again with this new input (line 9). If the new input exercises a path that has not been seen before, it is added to *PIDs* with the corresponding path identifier (lines 10, 12).

This process terminates when a bound is reached, such as a timeout or a number of generated inputs (line 2). When that happens, the fuzzer returns a test suite comprising all inputs in *PIDs*, each exercising a different path in the program.

B. Smart Contracts

Ethereum [21] is one of the most well known blockchain-based [22], [23] computing platforms. Like a bank, Ethereum supports accounts that store a balance (in digital assets) and are owned by a user. More specifically, there is support for two account types, namely user and contract accounts.

Contract accounts are not managed by a user, but instead by a program. The program associated with a certain contract account describes an agreement between the account and any users that interact with it. For example, such a program could encode the rules of a gambling game. To store information,

such as bets from various users, a contract account also comes with persistent state that the program may access and modify.

A contract account together with its managing program and persistent state is called a *smart contract*. However, the term may also refer to the code alone. Ethereum smart contracts can be developed in several high-level languages, such as Solidity and Vyper, which compile to Ethereum Virtual Machine (EVM) [24] bytecode.

Users interact with a smart contract, for instance to place a bet, by issuing a transaction with the contract. The transaction simply calls one of the contract functions, but in order to be carried out, users need to provide a fee. This fee is called *gas* and is approximately proportional to how much code needs to run. Any transaction that runs out of gas is aborted.

III. OVERVIEW

We now give an overview of our approach through the example of Fig. 2.

Example. The figure shows a constructed function `Bar`, which is written in Solidity and contained in a smart contract. (The comments should be ignored for now.) There are three assertions in this function, on lines 14, 19, and 22. A compiler will typically introduce a conditional jump for each assertion, where one branch leads to a location that fails. Let us assume that we select the failing locations (t_{14} , t_{19} , and t_{22}) of the three assertions as our target locations. Note that any target locations could be (automatically) selected based on various strategies, e.g., recently modified code, assertions, etc. Out of the above locations, t_{14} and t_{19} are unreachable, whereas t_{22} is reachable when input parameter `a` has value 42.

Generating a test input that reaches location t_{22} is difficult for a greybox fuzzer for two reasons. First, the probability of generating value 42 for parameter `a` is tiny, namely 1 out of 2^{256} . This means that, for the fuzzer to increase the chances of reaching t_{22} , it would need to fuzz certain “promising” inputs with a large amount of energy. However, standard greybox fuzzers are agnostic to what constitutes a promising input that is more likely to reach a target location when mutated.

Second, there are more than 100’000 program paths in function `Bar`. In fact, the then-branch of the first if-statement (line 5) contains two input-dependent loops (lines 11 and 16), whose number of iterations depends on parameters `w` and `z`, respectively. Recall that a greybox fuzzer generates new inputs by mutating existing ones from the test suite. Therefore, the larger the size of the test suite, the larger the space of possible mutations, and the lower the chances of generating an input that reaches the target location.

Existing work. As discussed earlier, there is existing work that leverages the results of an offline static analysis to guide automatic test generation toward unverified executions (e.g., [13], [14], [15], [16], [17]). To apply such a technique on the example of Fig. 2, let us assume a very lightweight static analysis that is based on abstract interpretation [25], [26] and uses the simple constant-propagation domain [27]. Note that, for each program variable, the constant-propagation domain can only infer a single constant value. When run offline, this

```

1 function Bar(uint256 w, uint256 x, uint256 y,
2   uint256 z, uint256 a) returns (uint256)
3 {
4   uint256 ret = 0;
5   if (x % 2 == 0) { // if (x % 1000 != 42) {
6     ret = 256;
7     if (y % 2 == 0) {
8       ret = 257;
9     }
10    w = w % ret;
11    while (w != 0) {
12      w--;
13    }
14    assert(w == 0); // drop this line
15    z = z % ret;
16    while (ret != z) {
17      z++;
18    }
19    assert(ret == z); // assert(x != 42 - w*z);
20  } else {
21    ret = 3*a*a + 7*a + 101;
22    assert(ret != 5687);
23  }
24  return ret;
25 }

```

Fig. 2: The running example.

analysis is able to prove that target location t_{14} is unreachable. This is because, after the loop on line 11, the analysis assumes the negation of the loop condition (that is, `w == 0`), which is equivalent to the asserted condition.

However, the analysis cannot prove that location t_{19} is also unreachable. This is because, after the if-statement on line 7, variable `ret` has abstract value \top . In other words, the analysis finds `ret` to be unconstrained since the constant-propagation domain is not able to express that its value is either 256 or 257. Given that `ret` is \top , `z` also becomes \top (line 15). It is, therefore, not possible for the analysis to determine whether these two variables always have the same value on line 19 and verify the assertion. As a result, automatic test generation needs to explore function `Bar` as if no static analysis had previously run. To check whether the assertion on line 19 always holds, a testing tool would have to generate inputs for all paths leading to it, thus including each iteration of the loop on line 11.

On the other hand, an existing technique for directed greybox fuzzing [18] performs lightweight instrumentation of the program under test to extract a distance metric for each input, which is then used as feedback for the fuzzer. So, the instrumentation encodes a static metric that measures the distance between the instrumented and the target locations in the control-flow graph. In our example, such metrics are less effective since all instructions are relatively close to the target locations, and the control-flow graph alone is not precise enough to determine more semantic reachability conditions. In addition, when directly fuzzing bytecode or assembly, a control-flow graph might not be easily recoverable, for instance due to indirect jumps.

Lookahead analysis. In contrast, our approach alleviates the imprecision of a static analysis by running it online and

does not require a control-flow graph. Our greybox fuzzer invokes the lookahead analysis for each input that is added to the test suite. Starting from split points (e.g., P_1 , P_{i-1} , and P_i in Fig. 1) along an explored program path, the analysis computes a path prefix (π_{pre}) for which all suffix paths do not reach any target location (e.g., T_x and T_y). We refer to such a path prefix as a *no-target-ahead prefix* (see Def. 2 for more details). As we explain below, the lookahead analysis aims to identify short no-target-ahead prefixes.

As an example, let us consider the constant-propagation analysis and an input for function `Bar` with an even value for x (thus making execution take the then-branch of the first if-statement on line 5). Along the path exercised by this input, the analysis fails to show that both target locations t_{14} and t_{19} are unreachable for the suffix paths starting from line 7. In fact, the analysis is as imprecise as when run offline on the entire function. However, it does verify the unreachability of the target locations for all suffix paths from line 9 by propagating forward the constant value of variable `ret` (either 256 or 257, depending on the value of y). Out of the many paths with an even value for x , the two no-target-ahead prefixes until line 9 (through the then- and else-branches of the if-statement on line 7) are actually the shortest ones for which the lookahead analysis proves that target locations t_{14} and t_{19} are unreachable.

Power schedule. The no-target-ahead prefixes computed by the lookahead analysis are used to control the fuzzer’s power schedule [28], which assigns more energy to certain inputs according to two criteria.

First, it assigns more energy to inputs that exercise a rare (i.e., rarely explored) no-target-ahead prefix. The intuition is to fuzz these inputs more in order to increase the chances of flipping a branch along the rare prefix, and thereby, reaching a target location. Note that flipping a branch in a suffix path can never lead to a target location. For this reason, our power schedule no longer distinguishes inputs based on the program path they exercise, but rather based on their no-target-ahead prefix. To maximize the chances of discovering a target location with fuzzing, the lookahead analysis tries to identify the shortest no-target-ahead prefixes, which are shared by the most suffix paths.

For the example of Fig. 2, consider the two no-target-ahead prefixes (until line 9) that we discussed above. Consider also the no-target-ahead prefix until the successful branch of the assertion on line 22. The inputs that exercise these prefixes are dynamically assigned roughly the same energy by our schedule—if one of them is exercised more rarely than the others, it is given more energy. This makes reaching target location t_{22} significantly more likely than with standard power schedules based on path identifiers, which assign roughly the same energy to each input exercising one of the thousands of paths in `Bar`.

Second, our power schedule also assigns more energy to inputs exercising rare split points in a no-target-ahead prefix, similarly to how existing work assigns more energy to rare branches [29]. The intuition is the following. Any newly

discovered no-target-ahead prefix is by definition rare—it has not been fuzzed before. Since it is rare, the power schedule will assign more energy to it, as discussed above. However, there are programs where new no-target-ahead prefixes can be easily discovered, for instance due to an input-dependent loop. In such cases, a power schedule only focusing on rare prefixes would prioritize these new prefixes at the expense of older ones that explore rare program locations, such as split points. For this reason, when a split point in a no-target-ahead prefix becomes rare, the power schedule tries to explore it more often.

As an example, consider the code in Fig. 2 while taking the comments into account, that is, replace lines 5 and 19 with the comments and drop line 14. The assertion on line 19 holds, but the constant-propagation analysis is too weak to prove it. As a result, for any path through this assertion, its no-target-ahead prefix has to include line 19. However, new no-target-ahead prefixes are very easily discovered; for instance, by exploring a different number of iterations in any of the two loops. So, even if at some point the fuzzer discovers the path that successfully exercises the assertion on line 22, its no-target-ahead prefix will quickly become less rare than any new prefixes going through the loops. The corresponding input will, therefore, be fuzzed less often even though it is very close to revealing the assertion violation. By prioritizing rare split points, for instance line 21, our power schedule will assign more energy to that input. This increases the chances of mutating the value of a to be 42 and reaching target t_{22} .

Both of these criteria effectively guide the fuzzer toward the target locations. For Fig. 2, our technique generates a test that reaches t_{22} in 27s on average (between 13 and 48s in 5 runs). Standard greybox fuzzing does not reach t_{22} in 4 out of 5 runs, with a timeout of 300s. The target location is reached in 113s during a fifth run, so in 263s on average. For this example, our technique achieves at least a 10x speed-up.

Why smart contracts. While our approach could be applied to regular programs, it is particularly useful in the context of smart contracts. One reason is that, in this setting, combining an offline static analysis with test generation using code instrumentation would change the program semantics. Recall that a transaction with a smart contract is carried out when users provide enough gas, which is roughly proportional to how much code is run. Since instrumentation consumes gas at execution time, it could cause a testing tool to report spurious out-of-gas errors. Another reason is that most deployed smart contracts are only available as bytecode, and recovering the control-flow graph from the bytecode is challenging.

IV. TECHNIQUE

In this section, we describe our technique in detail by first formally defining a lookahead analysis (Sect. IV-A). We then discuss how to integrate such an analysis with greybox fuzzing to enable a more targeted exploration of the search space (Sect. IV-B). Lastly, we present a concrete algorithm for a lookahead analysis based on abstract interpretation.

A. Lookahead Analysis

Let us first define a prefix and a no-target-ahead prefix of a given path.

Definition 1 (Prefix). Given a program P and a path π in P , we say that π_{pre} is a *prefix* of π iff there exists a suffix ρ such that $\pi = \text{concat}(\pi_{pre}, \rho)$.

Note that, in the above definition, ρ may be empty, in which case $\pi = \pi_{pre}$.

Definition 2 (No-target-ahead prefix). Given a program P , target locations T , and a prefix π_{pre} of a path in P , we say that π_{pre} is a *no-target-ahead prefix* iff the suffix ρ of every path $\pi = \text{concat}(\pi_{pre}, \rho)$ in P does not contain a target location $\tau \in T$.

Note that any path π in a program P is trivially a no-target-ahead prefix since there cannot be any target locations after reaching the end of its execution.

For a given no-target-ahead prefix, the analysis computes a *lookahead identifier* (LID) that will later be used to guide the greybox fuzzer.

Definition 3 (Lookahead identifier). Given a no-target-ahead prefix π_{pre} , the *lookahead identifier* λ is a cryptographic hash $\text{hash}(\pi_{pre})$.

The above definition ensures that it is very unlikely that two different no-target-ahead prefixes map to the same lookahead identifier.

Unlike a path identifier (PID) in standard greybox fuzzing, which is computed purely syntactically, a LID captures a no-target-ahead prefix, which is computed by semantically analyzing a program path. As a result, two program paths with different $PIDs$ may share the same LID . In other words, lookahead identifiers define equivalence classes of paths that share the same no-target-ahead prefix.

Definition 4 (Lookahead analysis). Given a program P , an input I , and a set of target locations T , a *lookahead analysis* computes a lookahead identifier λ for the corresponding no-target-ahead prefix π_{pre} (of path π exercised by input I) and a set of split points SPs along π_{pre} .

Any analysis that satisfies the above definition is a sound lookahead analysis. For instance, one that simply returns the hash of path π exercised by input I and all locations along π is trivially sound. For a given input, the precision of the analysis is determined by the length of the no-target-ahead prefix, and thereby, the number of suffix paths that are proved not to contain any target locations. In other words, the shorter the no-target-ahead prefix for a given input, the more precise the lookahead analysis.

B. Fuzzing with Lookahead Analysis

The integration of greybox fuzzing with a lookahead analysis builds on the following core idea. For each input in the test suite, the lookahead analysis determines a set of split points, that is, program locations along the explored path. It

Algorithm 1: Greybox fuzzing with lookahead analysis.

Input: Program $prog$, Seeds S , Target locations T

```

1  $PIDs \leftarrow \text{RUNSEEDS}(S, prog)$ 
2 while  $\neg \text{INTERRUPTED}()$  do
3    $input \leftarrow \text{PICKINPUT}(PIDs)$ 
4    $energy \leftarrow 0$ 
5    $maxEnergy \leftarrow \text{ASSIGNENERGY}(input)$ 
6    $maxEnergy \leftarrow \text{LOOKAHEADASSIGNENERGY}(input)$ 
7   while  $energy < maxEnergy$  do
8      $input' \leftarrow \text{FUZZINPUT}(input)$ 
9      $PID' \leftarrow \text{RUN}(input', prog)$ 
10    if  $\text{ISNEW}(PID', PIDs)$  then
11       $LID, SPs \leftarrow \text{LOOKAHEADANALYZE}(prog, input', T)$ 
12       $PIDs \leftarrow \text{ADD}(PID', input', LID, SPs, PIDs)$ 
13     $energy \leftarrow energy + 1$ 

```

Output: Test suite $\text{INPUTS}(PIDs)$

then computes a no-target-ahead prefix, which spans until one of these split points and is identified by a lookahead identifier. The fuzzer uses the rarity of the lookahead identifier as well as of the split points that are located along the no-target-ahead prefix to assign energy to the corresponding input.

The grey boxes in Alg. 1 highlight the key extensions we made to standard greybox fuzzing. For one, our algorithm invokes the lookahead analysis on line 11. This is done for every new input that is added to the test suite and computes the LID of the no-target-ahead prefix as well as the split points SPs along the prefix. Both are stored in the $PIDs$ data structure for efficient lookups (e.g., when assigning energy).

We also replace the existing power schedule on line 5 with a specialized one given by $\text{LOOKAHEADASSIGNENERGY}$ (line 6). As discussed in Sect. III, our power schedule assigns more energy to inputs that exercise either a *rare LID* or a *rare split point* along a no-target-ahead prefix. We define the new power schedule in the following.

Definition 5 (Rare LID). Given a test suite with $LIDs$ Λ , a LID λ is rare iff

$$\text{fuzz}(\lambda) < \text{rarity_cutoff},$$

where $\text{fuzz}(\lambda)$ measures the number of fuzzed inputs that exercised λ so far and $\text{rarity_cutoff} = 2^i$ such that

$$2^{i-1} < \min_{\lambda' \in \Lambda} \text{fuzz}(\lambda') \leq 2^i.$$

For example, if the LID with the fewest fuzzed inputs has been explored 42 times, then any LID that has been explored less than 2^6 times is rare.

The above definition is inspired by an existing power schedule for targeting rare branches [29] that introduced such a dynamically adjusted rarity_cutoff . Their experience shows that this metric performs better than simply considering the n $LIDs$ with the lowest number of fuzzed inputs as rare.

Definition 6 (Rare split point). Given a test suite with split points SPs along the no-target-ahead prefixes, a split point p is rare iff

$$\text{fuzz}(p) < \text{rarity_cutoff},$$

Algorithm 2: Lookahead algorithm.

Input: Program $prog$, Input $input$, Target locations T

```
1  $\pi \leftarrow \text{RUN}(input, prog)$ 
2  $i \leftarrow 0$ 
3  $SPs \leftarrow \emptyset$ 
4 while  $i < |\pi|$  do
5   if  $\text{ISSPLITPOINT}(i, \pi)$  then
6      $\pi_{pre} \leftarrow \pi[0..i + 1]$ 
7      $SPs \leftarrow SPs \cup \{\pi[i]\}$ 
8      $\phi, loc \leftarrow \text{PREFIXINFERENCE}(\pi_{pre})$ 
9     if  $\text{ARETARGETSUNREACHABLE}(prog, loc, \phi, T)$  then
10      return  $\text{COMPUTEHASH}(\pi_{pre}, SPs)$ 
11    $i \leftarrow i + 1$ 
12 return  $\text{COMPUTEHASH}(\pi), SPs$ 
```

Output: Lookahead identifier λ , Split points SPs

where $fuzz(p)$ measures the number of fuzzed inputs that exercised p so far and $rarity_cutoff = 2^i$ such that

$$2^{i-1} < \min_{p' \in SPs} fuzz(p') \leq 2^i.$$

Power schedule. Our power schedule is defined as follows for an input I with LID λ and split points SPs along the no-target-ahead prefix:

$$\begin{cases} \min(2^{selected(I)}, K), & \text{if } \lambda \text{ is rare } \vee \exists p \in SPs \cdot p \text{ is rare} \\ 1, & \text{otherwise.} \end{cases}$$

In the above definition, $selected(I)$ denotes the number of times that I was selected for fuzzing (line 3 in Alg. 1), and K is a constant (1024 in our implementation). Intuitively, our power schedule assigns little energy to inputs whose LID is not rare and whose no-target-ahead prefix does not contain any rare split points. Otherwise, it assigns much more energy, the amount of which depends on how often the input has been selected for fuzzing before. The energy grows exponentially up to some bound K , similarly to the cut-off-exponential schedule in AFLFast [28].

C. Lookahead Algorithm

Alg. 2 shows the algorithm for the lookahead analysis, which is implemented in function `LOOKAHEADANALYZE` from Alg. 1 and uses abstract interpretation [25], [26].

First, the lookahead analysis executes the program input concretely to collect the exercised path π (line 1 in Alg. 2). Given path π , it searches for the shortest no-target-ahead prefix π_{pre} by iterating over possible split points p (lines 4–11). Let us explain these lines in detail.

On line 5, the algorithm calls a predicate `ISSPLITPOINT`, which is parametric in which locations constitute split points. All locations along π could be split points, but to narrow down the search, the implementation may consider only a subset of them, for instance, at conditional jumps.

At each split point, the analysis performs two separate steps: (1) *prefix inference* and (2) *suffix checking*. The prefix inference (line 8) statically analyzes the prefix π_{pre} using abstract interpretation to infer its postcondition ϕ . This step

essentially executes the prefix in the abstract for all possible inputs that exercise this path.

Given condition ϕ , the analysis then performs the suffix checking to determine if all target locations are unreachable (line 9). This analysis performs standard, forward abstract interpretation by computing a fixed-point. If all target locations are unreachable, the analysis terminates and returns a non-empty LID by computing a hash over the program locations along the path prefix π_{pre} (line 10). This ensures that the analysis returns as soon as it reaches the first split point for which all targets are unreachable. In addition, it returns the set of all split points along prefix π_{pre} .

Even though off-the-shelf abstract interpreters are not designed to perform prefix inference and suffix checking, it is relatively straightforward to extend them. Essentially, when invoking a standard abstract interpreter on a program, the path prefix is always empty, whereas our lookahead analysis is partially path-sensitive (i.e., for the prefix, but not the suffix). Due to this partial path-sensitivity, even an inexpensive abstract domain (e.g., constant propagation or intervals) might be able to prove unreachability of a certain target location, which would otherwise require a more precise domain (for an empty prefix).

Split points. In practice, it is important to choose split points with care since too many split points will have a negative impact on the performance of the lookahead analysis. In our implementation, we only consider split points when entering a basic block for the first time along a given path. The intuition is that the lookahead analysis should run every time “new code” is discovered. Our experiments show that this design decision results in negligible overhead.

Calls. To keep the lookahead analysis lightweight, the suffix-checking step is modular. More specifically, any calls to other contracts are conservatively treated as potentially leading to target locations. (Note that inter-contract calls are used very sparingly in smart contracts and that intra-contract calls are simply jumps.) In contrast, during the prefix-inference step, we compute a summary LID for the callee context by recursively invoking the lookahead algorithm on the callee. This requires separating the parts of path π (from Alg. 2) that belong to the caller and the callee. It is also necessary to conservatively model the effect of a call on the caller context (e.g., by havocking return values).

V. IMPLEMENTATION

Our implementation extends HARVEY [30], [31], an existing greybox fuzzer for Ethereum smart contracts. It is actively used at ConsenSys Diligence, one of the largest blockchain-security consulting companies, and is one of the tools that power the MythX analysis platform. For our purposes, we integrated HARVEY with BRAN, our new abstract-interpretation framework for EVM bytecode, which is open source¹.

BRAN is designed to be scalable by performing a very lightweight, modular analysis that checks functional-correctness properties. Unlike other static analyzers for EVM

¹<https://github.com/Practical-Formal-Methods/bran>

bytecode (e.g., Securify [32] and MadMax [33]), BRAN runs directly on the bytecode without having to reconstruct the control-flow graph or decompile to an intermediate language. BRAN is equipped with a constant-propagation domain [27], which is commonly used in compiler optimizations. It handles all opcodes and integrates the go-ethereum virtual machine to concretely execute any opcodes with all-constant arguments.

Prefix length. During our preliminary experiments with the integration of HARVEY and BRAN, we observed that the prefix length may become quite large, for instance in the presence of input-dependent loops. However, the running time of the lookahead analysis is proportional to the prefix length, and our goal is to keep the analysis as lightweight as possible. For this reason, our implementation ignores any split points after the first 8’192 bytecode locations of the prefix. Note that this design decision does not affect the soundness of the lookahead analysis; it only reduces the search space of prefixes and might result in considering the entire path as the no-target-ahead prefix.

VI. EXPERIMENTAL EVALUATION

We now evaluate our technique on real-world Ethereum smart contracts. First, we discuss the benchmark selection (Sect. VI-A) and describe our experimental setup (Sect. VI-B). We then evaluate the effectiveness of the static lookahead analysis in greybox fuzzing (Sect. VI-C) and identify potential threats to the validity of our experiments (Sect. VI-D).

A. Benchmark Selection

We evaluated our technique on a total of 27 smart contracts, which originate from 17 GitHub repositories. Tab. I gives an overview. The first column lists a benchmark identifier for each smart contract under test, while the second and last columns provide the name and description of the containing project. Note that a repository may contain more than one contract, for instance including libraries; from each repository, we selected one or more contracts for our evaluation. The third and fourth columns of the table show the number of public functions and lines of Solidity code in each benchmark. (We provide links to all repositories as well as the changesets used for our experiments in the appendix.)

It is important to note that the majority of smart contracts are under 1’000 lines of code. Still, contracts of this size are complex programs, and each of them might take several weeks to audit. However, as it becomes clear from the example of Fig. 2, code size is not necessarily proportional to the number of feasible program paths or the difficulty to reach a particular target location with greybox fuzzing.

The repositories were selected with the goal of ensuring a diverse set of benchmarks. In particular, they include popular projects, such as the ENS domain name auction, the ConsenSys multisig wallet, and the MicroRaiden payment service. In addition to being widely known in the Ethereum community, these projects are highly starred on GitHub (4’857 stars in total on 2019-05-07, median 132), have been independently audited, and regularly transfer large amounts of assets. Moreover, our

BIDs	Name	Functions	LoSC	Description
1	ENS	24	1205	ENS domain name auction
2–3	CMSW	49	503	ConsenSys multisig wallet
4–5	GMSW	49	704	Gnosis multisig wallet
6	BAT	23	191	BAT token (advertising)
7	CT	12	200	ConsenSys token library
8	ERCF	19	747	ERC Fund (investment fund)
9	FBT	34	385	FirstBlood token (e-sports)
10–13	HPN	173	3065	Haven payment network
14	MR	25	1053	MicroRaiden payment service
15	MT	38	437	MOD token (supply-chain)
16	PC	7	69	Payment channel
17–18	RNTS	49	749	Request Network token sale
19	DAO	23	783	The DAO organization
20	VT	18	242	Valid token (personal data)
21	USCC1	4	57	USCC’17 entry
22	USCC2	14	89	USCC’17 (honorable mention)
23	USCC3	21	535	USCC’17 (3rd place)
24	USCC4	7	164	USCC’17 (1st place)
25	USCC5	10	188	USCC’17 (2nd place)
26	PW	19	549	Parity multisig wallet
27	BNK	44	649	Bankera token
Total		662	12564	

Table I: Overview of benchmarks. The third and fourth columns provide the number of public functions and lines of source code in each benchmark, respectively.

selection includes contracts from various application domains (like auctions, wallets, and tokens), attacked contracts (namely, The DAO and Parity wallet) as well as contracts submitted to the first Underhanded Solidity Coding Contest (USCC) [34]. Entries in this contest aim to conceal subtle vulnerabilities.

For selecting these repositories, we followed guidelines on how to evaluate fuzzers [35]. We do not randomly collect smart contracts from the Ethereum blockchain since this would likely contaminate our benchmarks with duplicates or bad-quality contracts—that is, contracts without users, assets, or dependencies, for instance, on libraries or other contracts.

B. Experimental Setup

Our experiments compare the integration of HARVEY and BRAN (incl. three variants) with HARVEY alone to evaluate the effectiveness of targeted fuzzing. The comparison focuses on the time it takes for each configuration to cover a set of target locations.

Targets. We randomly selected up to four target locations for each benchmark. In particular, we picked contract locations of varying difficulty to reach, based on when they were first discovered during a 1h standard greybox-fuzzing run. So, we randomly picked at most one newly discovered location, if one existed, from each of the following time brackets in this order: 30–60m, 15–30m, 7.5–15m, 3.75–7.5m, and 1.875–3.75m.

Runs. We performed 24 runs of each configuration on the 27 benchmarks of Tab. I. For each run, we used a different randomness seed, the same seed input, and a time limit of 1h (i.e., 3’600s). In our results, we report medians and use Wilcoxon-Mann-Whitney U tests to determine if differences in medians between configurations are statistically significant.

Machine. We used an Intel® Xeon® CPU @ 2.67GHz 24-core machine with 50GB of memory running Debian 9.5.

C. Results

We now evaluate the effectiveness of our technique by investigating five research questions.

RQ1: Effectiveness of targeted fuzzing. Tab. II compares our baseline configuration A, which does not enable the static lookahead analysis, with configuration B, which does. Note that configuration A uses the cut-off-exponential power schedule of AFLFast [28], whereas B uses our specialized schedule. The first two columns of the table indicate the benchmark and target IDs. Columns 3 and 4 show the median time (in seconds) required to discover the first input that reaches the target location (time-to-target) for both configurations, and column 5 shows the speed-up factor. Column 6 shows the p -value, which indicates the level of statistical significance; here, we use $p < 0.05$ for “significant” differences. The last two columns show Vargha-Delaney A12 effect sizes [36]. Intuitively, these measure the probability that configuration A is faster than B and vice versa.

For 32 (out of 60) target locations, we observe significant differences in time (i.e., $p < 0.05$), marked in bold in the table. *Configuration B significantly outperforms A for 31 (out of 32) of these target locations, with a median speed-up of up to 14x for one of the targets in benchmark 26.* In general, the results suggest that targeted fuzzing is very effective, and unsurprisingly, its impact is most significant for difficult targets (i.e., with high time-to-target for configuration A). Specifically, *for the 24 targets with $T_A \geq 900$ or $T_B \geq 900$, configuration B is significantly faster for 20, with insignificant differences between A and B for the remaining 4 targets.*

Note that running the static analysis with an empty prefix (resembling an offline analysis) on these benchmarks is not able to guide the fuzzer at all. Since all our target locations are reachable by construction, the analysis soundly reports them as reachable. Therefore, the fuzzer still needs to explore the entire contract to see if they indeed are.

RQ2: Effectiveness of lookahead analysis. To measure the effect of the lookahead analysis, we created configuration C, which is identical to configuration B except that the analysis is maximally imprecise and inexpensive. Specifically, ARETARGETSUNREACHABLE from Alg. 2 simply returns false, and consequently, the computed *LIDs* capture entire program paths, similarly to *PIDs*.

As shown in Tab. III, there are significant differences between configurations B and C for 21 target locations. *Configuration B is significantly faster than C for 17 out of 21 targets, and they are equally fast for 2 of the remaining 4 target locations.*

Interestingly, configuration C is faster than A (for all 12 target locations with significant differences). This suggests that our power schedule regarding rare split points is effective independently of the lookahead analysis.

RQ3: Effectiveness of power schedule. To measure the effect of targeting rare *LIDs* and rare split points in our power schedule, we created configuration D. It is identical to configuration B except that it uses a variant of AFLFast’s cut-off-exponential power schedule [28]. The original power

BID	Target ID	T_A	T_B	T_A/T_B	p	A12 _A	A12 _B
1	79145a51:35ec	324.15	90.25	3.59	0.049	0.33	0.67
1	79145a51:bd4	32.69	69.53	0.47	0.130	0.63	0.37
2	060a46c9:d03	3385.55	706.71	4.79	0.000	0.20	0.80
2	060a46c9:e29	161.66	106.57	1.52	0.197	0.39	0.61
2	060a46c9:16a5	701.39	339.86	2.06	0.008	0.27	0.73
2	060a46c9:1f11	346.06	63.14	5.48	0.000	0.11	0.89
3	708721b5:1485	396.11	394.54	1.00	0.477	0.44	0.56
3	708721b5:4ac	2292.00	775.93	2.95	0.000	0.19	0.81
3	708721b5:1ca0	1248.59	817.76	1.53	0.005	0.26	0.74
3	708721b5:1132	413.00	216.72	1.91	0.003	0.24	0.76
4	9b8e6b2a:d08	3600.00	867.65	4.15	0.000	0.15	0.85
4	9b8e6b2a:18f0	1657.33	432.50	3.83	0.002	0.24	0.76
4	9b8e6b2a:1fee	143.96	47.13	3.05	0.062	0.34	0.66
4	9b8e6b2a:553	3600.00	833.70	4.32	0.001	0.22	0.78
5	5a3e5a7f:c09	3600.00	1282.42	2.81	0.000	0.08	0.92
5	5a3e5a7f:23f	900.53	466.99	1.93	0.017	0.30	0.70
5	5a3e5a7f:1da8	1355.07	646.41	2.10	0.000	0.16	0.84
5	5a3e5a7f:1d67	1497.96	524.08	2.86	0.000	0.15	0.85
6	387bdf82:da7	61.66	22.70	2.72	0.089	0.36	0.64
8	e2aedada:15a7	2592.56	1135.37	2.28	0.002	0.24	0.76
8	e2aedada:17bb	1783.03	612.39	2.91	0.001	0.22	0.78
8	e2aedada:d71	73.93	47.89	1.54	0.307	0.41	0.59
8	e2aedada:13a8	258.14	74.87	3.45	0.035	0.32	0.68
9	dada6ee2:1693	334.82	49.38	6.78	0.000	0.13	0.87
9	dada6ee2:bee	225.12	72.14	3.12	0.000	0.19	0.81
9	dada6ee2:90e	84.62	50.39	1.68	0.338	0.42	0.58
10	d98d1d6b:1f10	1124.84	281.45	4.00	0.004	0.26	0.74
10	d98d1d6b:401a	164.12	153.95	1.07	0.861	0.48	0.52
10	d98d1d6b:3cdd	1669.91	1817.05	0.92	0.729	0.53	0.47
10	d98d1d6b:3ce8	3600.00	3600.00	1.00	0.713	0.47	0.53
11	3ae06fbc:34db	3600.00	3600.00	1.00	0.105	0.38	0.62
11	3ae06fbc:3de2	150.22	81.77	1.84	0.557	0.45	0.55
11	3ae06fbc:3ef3	284.34	395.15	0.72	0.703	0.47	0.53
11	3ae06fbc:10b2	238.35	142.03	1.68	0.228	0.40	0.60
12	0203d94d:713	76.82	60.27	1.27	0.910	0.49	0.51
14	b8c706d1:125e	3600.00	3600.00	1.00	0.085	0.39	0.61
14	b8c706d1:3479	290.73	299.26	0.97	0.861	0.52	0.48
14	b8c706d1:2023	34.65	43.72	0.79	0.992	0.50	0.50
15	06ef1a9c:27ce	3365.87	467.90	7.19	0.000	0.10	0.90
15	06ef1a9c:b41	100.00	73.83	1.35	0.877	0.49	0.51
15	06ef1a9c:a16	71.00	39.46	1.80	0.106	0.36	0.64
17	1c57401c:ef1	186.24	218.20	0.85	0.101	0.64	0.36
17	1c57401c:558	45.72	111.38	0.41	0.130	0.63	0.37
18	ac0bf5ee:15e4	1827.66	321.36	5.69	0.000	0.12	0.88
18	ac0bf5ee:171b	176.36	48.04	3.67	0.000	0.16	0.84
18	ac0bf5ee:15e0	133.84	27.80	4.81	0.001	0.22	0.78
18	ac0bf5ee:70c	24.87	61.47	0.40	0.036	0.68	0.32
20	54142e12:1555	29.57	15.42	1.92	0.298	0.41	0.59
23	d047b56e:5fb	42.01	20.70	2.03	0.279	0.41	0.59
24	b9ebdb99:40c	980.79	139.78	7.02	0.000	0.13	0.87
24	b9ebdb99:3d1	282.28	57.21	4.93	0.000	0.18	0.82
25	f1e90f8f:9fd	316.48	24.61	12.86	0.000	0.09	0.91
26	a788e7af:1f07	1778.07	130.34	13.64	0.000	0.07	0.93
26	a788e7af:1e29	2005.67	336.04	5.97	0.000	0.12	0.88
26	a788e7af:544	395.22	47.84	8.26	0.140	0.38	0.62
26	a788e7af:32b	44.67	45.92	0.97	0.813	0.48	0.52
27	9473c978:1541	2445.87	324.46	7.54	0.020	0.31	0.69
27	9473c978:e33	1493.03	637.16	2.34	0.023	0.31	0.69
27	9473c978:150e	178.11	97.60	1.82	0.120	0.37	0.63
27	9473c978:8e8	102.29	150.72	0.68	0.236	0.60	0.40

Table II: Comparing time-to-target between configuration A (w/o lookahead analysis) and B (w/ lookahead analysis).

schedule assigns energy to an input I based on how often its *PID* has been exercised. In contrast, our variant is based on how often its *LID* has been exercised and corresponds to using the results of the lookahead analysis with a standard power schedule.

However, as shown in Tab. IV, *configuration B is faster than configuration D for 28 of 30 targets (with significant differences)*. This indicates that our power schedule significantly reduces the time-to-target, thus effectively guiding the fuzzer.

Nonetheless, configuration D is faster than A for all 6 targets with significant differences. This shows the effectiveness of the lookahead analysis independently of the power schedule.

RQ4: Scalability of lookahead analysis. One key design decision for using an *online* static analysis as part of a

diverse application domains. Moreover, in the appendix, we provide the versions of all contracts used in our experiments so that others can also test them. The results may also not generalize to other target locations, but we alleviate this threat by selecting them at random and with varying difficulty to reach.

Internal validity. Internal validity [37] is compromised when systematic errors are introduced in the experimental setup. A common pitfall in evaluating randomized approaches, such as fuzzing, is the potentially biased selection of seeds. During our experiments, when comparing the different configurations of our technique, we consistently used the same seed inputs for HARVEY.

Construct validity. Construct validity ensures that any improvements, for instance in effectiveness or performance, achieved by a particular technique are due to that technique alone, and not due to other factors, such as better engineering. In our experiments, we compare different configurations of the same greybox fuzzer, and consequently, any effect on the results is exclusively caused by their differences.

VII. RELATED WORK

Our technique for targeted greybox fuzzing leverages an online static analysis to semantically analyze each new path that is added to the fuzzer’s test suite. The feedback collected by the static analysis is used to guide the fuzzer toward a set of target locations using a novel power schedule that takes inspiration from two existing ones [28], [29].

In contrast, the most closely related work [18] performs an offline instrumentation of the program under test encoding a static distance metric between the instrumented and the target locations in the control-flow graph. When running a given input, the instrumentation is used to obtain a dynamic (aggregated) distance. This distance subsequently guides the fuzzer toward the target locations.

Since a control-flow graph cannot always be easily recovered from EVM bytecode (e.g., due to indirect jumps), our lookahead analysis directly analyzes the bytecode using abstract interpretation [25], [26]. Our implementation uses the constant-propagation domain [27] to track the current state of the EVM (for instance, to resolve jump targets that are pushed to the execution stack). Unlike traditional static analyses, it aims to improve precision by performing a partially path-sensitive analysis—that is, path-sensitive for a prefix of a feasible path recorded at runtime by the fuzzer, and path-insensitive for all suffix paths.

Guiding greybox fuzzers. Besides directed greybox fuzzing [18], there is a number of greybox fuzzers that target specific program locations [38], rare branches [29], uncovered branches [39], [40], or suspected vulnerabilities [12], [41], [42], [43]. While several of these fuzzers use an offline static analysis to guide the exploration, none of them leverages an online analysis.

Guiding other program analyzers. There is a large body of work on guiding analyzers toward specific target locations [44], [45] or potential failures [13], [46], [47], [48], [14],

[15], [49], [16], [17], [50] by combining static and dynamic analysis. These combinations typically perform an offline static analysis first and use it to improve the effectiveness of a subsequent dynamic analysis; for instance, by pruning parts of the program. For example, Check ‘n’ Crash [13] integrates the ESC/Java static checker [51] with the JCrasher test-generation tool [2]. Similarly, DyTa [14] combines the .NET static analyzer Clousot [52] with the dynamic symbolic execution engine Pex [53]. YOGI [47], [48] constantly refines its over- and under-approximations in the style of counterexample-driven refinement [54]. In contrast, our lookahead analysis is online and constitutes a core component of our targeted greybox fuzzer.

Hybrid concolic testing [55] combines random testing with concolic testing [8], [9], [56]. Even though the technique significantly differs from ours, it shares an interesting similarity: it uses online concolic testing during a concrete program execution to discover uncovered code on-the-fly. When successful, the inputs for covering the code are used to resume the concrete program execution.

Symbolic execution. In the context of symbolic execution [57], there have emerged numerous search strategies for guiding the exploration; for instance, to target deeper paths (in depth-first search), uncovered statements [58], or “less-traveled paths” [59]. Our technique resembles a search strategy in that it prioritizes exploration of certain inputs over others.

Compositional symbolic execution [60], [61] has been shown to be effective in merging different program paths by means of summaries in order to alleviate path explosion. Dynamic state merging [62] and veritesting [63] can also be seen as forms of summarization. Similarly, our technique merges different paths that share the same lookahead identifier for the purpose of assigning energy. The more precise the lookahead analysis, the shorter the no-target-ahead prefixes, and thus, the more effective the merging.

Program analysis for smart contracts. There is a growing number of program analyzers for smart contracts, ranging from random test generation frameworks to static analyzers and verifiers [64], [65], [66], [67], [68], [69], [70], [71], [72], [33], [73], [74], [75], [32], [76], [77], [78]. In contrast, we present a targeted greybox fuzzer for smart contracts, the first analyzer for contracts that incorporates static and dynamic analysis.

VIII. CONCLUSION

We have presented a novel technique for targeted fuzzing using static lookahead analysis. The key idea is to enable a symbiotic collaboration between the greybox fuzzer and an online static analysis. On one hand, dynamic information (i.e., feasible program paths) are used to boost the precision of the static analysis. On the other hand, static information about reachable target locations—more specifically, lookahead identifiers and split points—is used to guide the greybox fuzzer toward target locations. Our experiments on 27 real-world benchmarks show that targeted fuzzing significantly outperforms standard greybox fuzzing for reaching 83% of the challenging target locations (up to 14x of median speed-up).

In future work, we plan to investigate other combinations of dynamic and online static analysis; for instance, to guide dynamic symbolic execution.

REFERENCES

- [1] K. Claessen and J. Hughes, “QuickCheck: A lightweight tool for random testing of Haskell programs,” in *ICFP*. ACM, 2000, pp. 268–279.
- [2] C. Csallner and Y. Smaragdakis, “JCrasher: An automatic robustness tester for Java,” *SPE*, vol. 34, pp. 1025–1050, 2004.
- [3] C. Pacheco, S. K. Lahiri, M. D. Ernst, and T. Ball, “Feedback-directed random test generation,” in *ICSE*. IEEE Computer Society, 2007, pp. 75–84.
- [4] “Peach Fuzzer Platform,” <https://www.peach.tech/products/peach-fuzzer/peach-platform/>.
- [5] “zzuf—Multi-Purpose Fuzzer,” <http://caca.zoy.org/wiki/zzuf>.
- [6] “Technical “whitepaper” for AFL,” http://lcamtuf.coredump.cx/afl/technical_details.txt.
- [7] “Libfuzzer—A library for coverage-guided fuzz testing,” <https://lvm.org/docs/LibFuzzer.html>.
- [8] P. Godefroid, N. Klarlund, and K. Sen, “DART: Directed automated random testing,” in *PLDI*. ACM, 2005, pp. 213–223.
- [9] C. Cadar and D. R. Engler, “Execution generated test cases: How to make systems code crash itself,” in *SPIN*, ser. LNCS, vol. 3639. Springer, 2005, pp. 2–23.
- [10] P. Godefroid, M. Y. Levin, and D. A. Molnar, “Automated whitebox fuzz testing,” in *NDSS*. The Internet Society, 2008, pp. 151–166.
- [11] C. Cadar, D. Dunbar, and D. R. Engler, “KLEE: Unassisted and automatic generation of high-coverage tests for complex systems programs,” in *OSDI*. USENIX, 2008, pp. 209–224.
- [12] V. Ganesh, T. Leek, and M. C. Rinard, “Taint-based directed whitebox fuzzing,” in *ICSE*. IEEE Computer Society, 2009, pp. 474–484.
- [13] C. Csallner and Y. Smaragdakis, “Check ‘n’ Crash: Combining static checking and testing,” in *ICSE*. ACM, 2005, pp. 422–431.
- [14] X. Ge, K. Taneja, T. Xie, and N. Tillmann, “DyTa: Dynamic symbolic execution guided with static verification results,” in *ICSE*. ACM, 2011, pp. 992–994.
- [15] M. Czech, M.-C. Jakobs, and H. Wehrheim, “Just test what you cannot verify!” in *FASE*, ser. LNCS, vol. 9033. Springer, 2015, pp. 100–114.
- [16] M. Christakis, P. Müller, and V. Wüstholtz, “Guiding dynamic symbolic execution toward unverified program executions,” in *ICSE*. ACM, 2016, pp. 144–155.
- [17] K. Ferles, V. Wüstholtz, M. Christakis, and I. Dillig, “Failure-directed program trimming,” in *ESEC/FSE*. ACM, 2017, pp. 174–185.
- [18] M. Böhme, V. Pham, M. Nguyen, and A. Roychoudhury, “Directed greybox fuzzing,” in *CCS*. ACM, 2017, pp. 2329–2344.
- [19] B. Livshits, M. Sridharan, Y. Smaragdakis, O. Lhoták, J. N. Amaral, B.-Y. E. Chang, S. Z. Guyer, U. P. Khedker, A. Møller, and D. Vardoulakis, “In defense of soundness: A manifesto,” *CACM*, vol. 58, pp. 44–46, 2015.
- [20] “The AFL vulnerability trophy case,” <http://lcamtuf.coredump.cx/afl/#bugs>.
- [21] “Ethereum white paper,” 2014, <https://github.com/ethereum/wiki/wiki/White-Paper>.
- [22] M. Swan, *Blockchain: Blueprint for a New Economy*. O’Reilly Media, 2015.
- [23] S. Raval, *Decentralized Applications: Harnessing Bitcoin’s Blockchain Technology*. O’Reilly Media, 2016.
- [24] G. Wood, “Ethereum: A secure decentralised generalised transaction ledger,” 2014, <http://gavwood.com/paper.pdf>.
- [25] P. Cousot and R. Cousot, “Abstract interpretation: A unified lattice model for static analysis of programs by construction or approximation of fixpoints,” in *POPL*. ACM, 1977, pp. 238–252.
- [26] P. Cousot and R. Cousot, “Systematic design of program analysis frameworks,” in *POPL*. ACM, 1979, pp. 269–282.
- [27] G. A. Kildall, “A unified approach to global program optimization,” in *POPL*. ACM, 1973, pp. 194–206.
- [28] M. Böhme, V. Pham, and A. Roychoudhury, “Coverage-based greybox fuzzing as Markov chain,” in *CCS*. ACM, 2016, pp. 1032–1043.
- [29] C. Lemieux and K. Sen, “FairFuzz: A targeted mutation strategy for increasing greybox fuzz testing coverage,” in *ASE*. ACM, 2018, pp. 475–485.
- [30] V. Wüstholtz and M. Christakis, “Learning inputs in greybox fuzzing,” *CoRR*, vol. abs/1807.07875, 2018.
- [31] V. Wüstholtz and M. Christakis, “Harvey: A greybox fuzzer for smart contracts,” *CoRR*, vol. abs/1905.06944, 2019.

- [32] P. Tsankov, A. M. Dan, D. Drachler-Cohen, A. Gervais, F. Bünzli, and M. T. Vechev, “Securify: Practical security analysis of smart contracts,” in *CCS*. ACM, 2018, pp. 67–82.
- [33] N. Grech, M. Kong, A. Jurisevic, L. Brent, B. Scholz, and Y. Smaragdakis, “MadMax: Surviving out-of-gas conditions in Ethereum smart contracts,” *PACMPL*, vol. 2, pp. 116:1–116:27, 2018.
- [34] “Underhanded solidity coding contest,” <http://u.solidity.cc/>.
- [35] G. Klees, A. Ruef, B. Cooper, S. Wei, and M. Hicks, “Evaluating fuzz testing,” in *CCS*. ACM, 2018, pp. 2123–2138.
- [36] A. Vargha and H. D. Delaney, “A critique and improvement of the CL common language effect size statistics of McGraw and Wong,” *JEBBS*, vol. 25, pp. 101–132, 2000.
- [37] J. Siegmund, N. Siegmund, and S. Apel, “Views on internal and external validity in empirical software engineering,” in *ICSE*. IEEE Computer Society, 2015, pp. 9–19.
- [38] H. Chen, Y. Xue, Y. Li, B. Chen, X. Xie, X. Wu, and Y. Liu, “Hawkeye: Towards a desired directed grey-box fuzzer,” in *CCS*. ACM, 2018, pp. 2095–2108.
- [39] Y. Li, B. Chen, M. Chandramohan, S. Lin, Y. Liu, and A. Tiu, “Steelix: Program-state based binary fuzzing,” in *ESEC/FSE*. ACM, 2017, pp. 627–637.
- [40] M. Wang, J. Liang, Y. Chen, Y. Jiang, X. Jiao, H. Liu, X. Zhao, and J. Sun, “SAFL: Increasing and accelerating testing coverage with symbolic execution and guided fuzzing,” in *ICSE Companion*. ACM, 2018, pp. 61–64.
- [41] I. Haller, A. Slowinska, M. Neugschwandtner, and H. Bos, “Dowsing for overflows: A guided fuzzer to find buffer boundary violations,” in *Security*. USENIX, 2013, pp. 49–64.
- [42] Y. Li, S. Ji, C. Lv, Y. Chen, J. Chen, Q. Gu, and C. Wu, “V-Fuzz: Vulnerability-oriented evolutionary fuzzing,” *CoRR*, vol. abs/1901.01142, 2019.
- [43] A. B. Chowdhury, R. K. Medicherla, and R. Venkatesh, “VeriFuzz: Program aware fuzzing—(competition contribution),” in *TACAS*, ser. LNCS, vol. 11429. Springer, 2019, pp. 244–249.
- [44] K.-K. Ma, Y. P. Khoo, J. S. Foster, and M. Hicks, “Directed symbolic execution,” in *SAS*, ser. LNCS, vol. 6887. Springer, 2011, pp. 95–111.
- [45] P. D. Marinescu and C. Cadar, “KATCH: High-coverage testing of software patches,” in *ESEC/FSE*. ACM, 2013, pp. 235–245.
- [46] M. B. Dwyer and R. Purandare, “Residual dynamic typestate analysis exploiting static analysis: Results to reformulate and reduce the cost of dynamic analysis,” in *ASE*. ACM, 2007, pp. 124–133.
- [47] A. V. Nori, S. K. Rajamani, S. Tetali, and A. V. Thakur, “The YOGI project: Software property checking via static analysis and testing,” in *TACAS*, ser. LNCS, vol. 5505. Springer, 2009, pp. 178–181.
- [48] P. Godefroid, A. V. Nori, S. K. Rajamani, and S. Tetali, “Compositional may-must program analysis: Unleashing the power of alternation,” in *POPL*. ACM, 2010, pp. 43–56.
- [49] L. Ma, C. Artho, C. Zhang, H. Sato, J. Gmeiner, and R. Ramler, “GRT: Program-analysis-guided random testing,” in *ASE*. IEEE Computer Society, 2015, pp. 212–223.
- [50] D. Devescary, P. M. Chen, J. Flinn, and S. Narayanasamy, “Optimistic hybrid analysis: Accelerating dynamic analysis through predicated static analysis,” in *ASPLOS*. ACM, 2018, pp. 348–362.
- [51] C. Flanagan, K. R. M. Leino, M. Lillibridge, G. Nelson, J. B. Saxe, and R. Stata, “Extended static checking for Java,” in *PLDI*. ACM, 2002, pp. 234–245.
- [52] M. Fähndrich and F. Logozzo, “Static contract checking with abstract interpretation,” in *FoVeOOS*, ser. LNCS, vol. 6528. Springer, 2010, pp. 10–30.
- [53] N. Tillmann and J. de Halleux, “Pex—White box test generation for .NET,” in *TAP*, ser. LNCS, vol. 4966. Springer, 2008, pp. 134–153.
- [54] E. M. Clarke, O. Grumberg, S. Jha, Y. Lu, and H. Veith, “Counterexample-guided abstraction refinement,” in *CAV*, ser. LNCS, vol. 1855. Springer, 2000, pp. 154–169.
- [55] R. Majumdar and K. Sen, “Hybrid concolic testing,” in *ICSE*. IEEE Computer Society, 2007, pp. 416–426.
- [56] K. Sen and G. Agha, “CUTE and jCUTE: Concolic unit testing and explicit path model-checking tools,” in *CAV*, ser. LNCS, vol. 4144. Springer, 2006, pp. 419–423.
- [57] J. C. King, “Symbolic execution and program testing,” *CACM*, vol. 19, pp. 385–394, 1976.
- [58] S. Park, B. M. M. Hossain, I. Hussain, C. Csallner, M. Grechanik, K. Taneja, C. Fu, and Q. Xie, “CarFast: Achieving higher statement coverage faster,” in *FSE*. ACM, 2012, p. 35.
- [59] Y. Li, Z. Su, L. Wang, and X. Li, “Steering symbolic execution to less traveled paths,” in *OOPSLA*. ACM, 2013, pp. 19–32.
- [60] P. Godefroid, “Compositional dynamic test generation,” in *POPL*. ACM, 2007, pp. 47–54.
- [61] S. Anand, P. Godefroid, and N. Tillmann, “Demand-driven compositional symbolic execution,” in *TACAS*, ser. LNCS, vol. 4963. Springer, 2008, pp. 367–381.
- [62] V. Kuznetsov, J. Kinder, S. Bucur, and G. Candea, “Efficient state merging in symbolic execution,” in *PLDI*. ACM, 2012, pp. 193–204.
- [63] T. Avgerinos, A. Rebert, S. K. Cha, and D. Brumley, “Enhancing symbolic execution with veritesting,” in *ICSE*. ACM, 2014, pp. 1083–1094.
- [64] L. Luu, D. Chu, H. Olickel, P. Saxena, and A. Hobor, “Making smart contracts smarter,” in *CCS*. ACM, 2016, pp. 254–269.
- [65] K. Bhargavan, A. Delignat-Lavaud, C. Fournet, A. Gollamudi, G. Gonthier, N. Kobeissi, N. Kulatova, A. Rastogi, T. Sibut-Pinote, N. Swamy, and S. Zanella-Béguélin, “Formal verification of smart contracts: Short paper,” in *PLAS*. ACM, 2016, pp. 91–96.
- [66] N. Atzei, M. Bartoletti, and T. Cimoli, “A survey of attacks on Ethereum smart contracts,” in *POST*, ser. LNCS, vol. 10204. Springer, 2017, pp. 164–186.
- [67] T. Chen, X. Li, X. Luo, and X. Zhang, “Under-optimized smart contracts devour your money,” in *SANER*. IEEE Computer Society, 2017, pp. 442–446.
- [68] I. Sergey and A. Hobor, “A concurrent perspective on smart contracts,” in *FC*, ser. LNCS, vol. 10323. Springer, 2017, pp. 478–493.
- [69] B. Jiang, Y. Liu, and W. K. Chan, “ContractFuzzer: Fuzzing smart contracts for vulnerability detection,” in *ASE*. ACM, 2018, pp. 259–269.
- [70] K. Chatterjee, A. K. Goharshady, and Y. Velnér, “Quantitative analysis of smart contracts,” in *ESOP*, ser. LNCS, vol. 10801. Springer, 2018, pp. 739–767.
- [71] S. Amani, M. Bégué, M. Bortin, and M. Staples, “Towards verifying Ethereum smart contract bytecode in Isabelle/HOL,” in *CPP*. ACM, 2018, pp. 66–77.
- [72] L. Brent, A. Jurisevic, M. Kong, E. Liu, F. Gauthier, V. Gramoli, R. Holz, and B. Scholz, “Vandal: A scalable security analysis framework for smart contracts,” *CoRR*, vol. abs/1809.03981, 2018.
- [73] S. Grossman, I. Abraham, G. Golan-Gueta, Y. Michalevsky, N. Rinetzky, M. Sagiv, and Y. Zohar, “Online detection of effectively callback free objects with applications to smart contracts,” *PACMPL*, vol. 2, pp. 48:1–48:28, 2018.
- [74] S. Kalra, S. Goel, M. Dhawan, and S. Sharma, “ZEUS: Analyzing safety of smart contracts,” in *NDSS*. The Internet Society, 2018.
- [75] I. Nikolic, A. Kolluri, I. Sergey, P. Saxena, and A. Hobor, “Finding the greedy, prodigal, and suicidal contracts at scale,” pp. 653–663, 2018.
- [76] “Echidna,” <https://github.com/trailofbits/echidna>.
- [77] “Manticore,” <https://github.com/trailofbits/manticore>.
- [78] “Mythril,” <https://github.com/ConsenSys/mythril-classic>.

APPENDIX

All tested smart contracts are open source. Tab. V provides the changeset IDs and links to their repositories.

BIDs	Name	Changeset ID	Repository
1	ENS	5108f51d656f201dc0054e55f5fd000d00ef9ef3	https://github.com/ethereum/ens
2-3	CMSW	2582787a14dd861b51df6f815fab122ff51fb574	https://github.com/ConsenSys/MultiSigWallet
4-5	GMSW	8ac8ba7effe6c3845719e480defb5f2ecafd2fd4	https://github.com/gnosis/MultiSigWallet
6	BAT	15bebd0642dac614d56709477c7c31d5c993ae1	https://github.com/brave-intl/basic-attention-token-crowdsale
7	CT	1f62e1ba3bf32dc22fe2de94a9ee486d667edef2	https://github.com/ConsenSys/Tokens
8	ERCF	c7d025220a1388326b926d8983e47184e249d979	https://github.com/ScJa/ercfund
9	FBT	ae71053e0656b0ceba7e229e1d67c09f271191dc	https://github.com/Firstbloodio/token
10-13	HPN	540006e0e2e5ef729482ad8bebcf7eafcd5198c2	https://github.com/Havven/havven
14	MR	527eb90c614ff4178b269d48ea063eb49ee0f254	https://github.com/raiden-network/microraiden
15	MT	7009cc95affa5a2a41a013b85903b14602c25b4f	https://github.com/modum-io/tokenapp-smartcontract
16	PC	515c1b935ac43afc6bf54caff68cf8521595b0b	https://github.com/mattdf/payment-channel
17-18	RNTS	6c39082eff65b2d3035a89a3f3dd94bde6cca60f	https://github.com/RequestNetwork/RequestTokenSale
19	DAO	f347c0e177edcf99d64fe589d236754fa375658	https://github.com/slockit/DAO
20	VT	30ede971bb682f245e5be11f544e305ef033a765	https://github.com/valid-global/token
21	USCC1	3b26643a85d182a9b8f0b6fe8c1153f3bd510a96	https://github.com/Arachnid/uscc
22	USCC2	3b26643a85d182a9b8f0b6fe8c1153f3bd510a96	https://github.com/Arachnid/uscc
23	USCC3	3b26643a85d182a9b8f0b6fe8c1153f3bd510a96	https://github.com/Arachnid/uscc
24	USCC4	3b26643a85d182a9b8f0b6fe8c1153f3bd510a96	https://github.com/Arachnid/uscc
25	USCC5	3b26643a85d182a9b8f0b6fe8c1153f3bd510a96	https://github.com/Arachnid/uscc
26	PW	657da22245dcfe0fe1cccc58ee8cd86924d65cdd	https://github.com/paritytech/contracts
27	BNK	97f1c3195b6cf4d8b3393016ecf106b42a2b1d97	https://github.com/Bankera-token/BNK-ETH-Contract

Table V: Smart contract repositories.