

VULSEYE: Detect Smart Contract Vulnerabilities via Stateful Directed Graybox Fuzzing

Ruichao Liang, Jing Chen, Cong Wu, Kun He, Yueming Wu, Ruochen Cao,
Ruiying Du, Yang Liu, Ziming Zhao

Abstract—Smart contracts, the cornerstone of decentralized applications, have become increasingly prominent in revolutionizing the digital landscape. However, vulnerabilities in smart contracts pose great risks to user assets and undermine overall trust in decentralized systems. Fuzzing, a prominent security testing technique, is extensively explored to detect vulnerabilities. But current smart contract fuzzers fall short of expectations in testing efficiency for two primary reasons. Firstly, smart contracts are stateful programs, and existing approaches, primarily coverage-guided, lack effective feedback from the contract state. Consequently, they struggle to effectively explore the contract state space. Secondly, coverage-guided fuzzers, aiming for comprehensive program coverage, may lead to a wastage of testing resources on benign code areas. This wastage worsens in smart contract testing, as the mix of code and state spaces further complicates comprehensive testing.

To address these challenges, we propose VULSEYE, a stateful directed graybox fuzzer for smart contracts guided by vulnerabilities. Different from prior works, VULSEYE achieves stateful directed fuzzing by prioritizing testing resources to code areas and contract states that are more prone to vulnerabilities. We introduce *Code Targets* and *State Targets* into fuzzing loops as the testing targets of VULSEYE. We use static analysis and pattern matching to pinpoint *Code Targets*, and propose a scalable backward analysis algorithm to specify *State Targets*. We design a novel fitness metric that leverages feedback from both the contract code space and state space, directing fuzzing toward these targets. With the guidance of code and state targets, VULSEYE alleviates the wastage of testing resources on benign code areas and achieves effective stateful fuzzing. In comparison with state-of-the-art fuzzers, VULSEYE demonstrated superior effectiveness and efficiency. Notably, it uncovered 4,845 vulnerabilities in 42,738 real-world smart contracts, outperforming existing approaches by up to 9.7×, and identified 11 previously unknown vulnerabilities within the top 50 Ethereum DApps, involving approximately 2,500,000 USD.

Index Terms—Fuzz Testing, Static Analysis, Smart Contract

I. INTRODUCTION

Smart contracts, self-executing programs that automate the requisite actions of agreements or contracts, form a cornerstone of the blockchain and cryptocurrency ecosystem. However, vulnerabilities within these contracts pose considerable risks [1], exemplified by the DAO incident, where a reentrancy vulnerability led to a loss of \$150 million [2], and

the Poly Network breach, resulting in a \$611 million loss [3]. Thus, it is imperative to conduct thorough security testing on smart contracts to detect vulnerabilities and prevent significant financial and operational consequences.

Fuzz testing, an automated technique for program testing, identifies vulnerabilities by generating random inputs and systematically exploring a program’s space [4]. Various efforts have been undertaken to apply fuzz testing techniques to smart contracts [5]–[14]. These initiatives involve generating inputs that mimic real-world transactions for testing smart contracts, with a focus on coverage-guided graybox fuzzing to improve code coverage and facilitate vulnerability discovery. For example, sFUZZ [10] introduces an adaptive strategy for seed selection aimed at reaching hard-to-cover branches for higher overall code coverage. SMARTIAN [12] utilizes data-flow-based feedback for generating high-quality test cases in order to enhance the code coverage. Despite these efforts, there still remain following significant gaps.

i) Insufficient State Space Exploration in Smart Contract Fuzzing. Although existing contract fuzzers have thoroughly studied the use of feedback from code space, such as code coverage, to guide fuzzing, they lack an effective feedback mechanism for contract states, resulting in indiscriminate exploration of the contract state space. As smart contracts are intrinsically stateful [15], with state variables playing a pivotal role in their functionality, accurately identifying and examining specific states is essential for vulnerability exploitation. Existing approaches such as CONFUZZIUS [11] and HARVEY [8] typically generate contract states through combinations of transactions. ITYFUZZ [16] utilizes snapshot for contract state to reduce the time required for transaction re-execution. However, they are insufficient for a targeted and systematic examination of the contract state space, thereby impeding effective vulnerability detection.

ii) Inefficient Code Space Exploration in Smart Contract Fuzzing. Most existing contract fuzzers adopt the coverage-guided fuzzing strategy, aiming for comprehensive program coverage. However, research has shown that vulnerabilities within programs tend to be concentrated in specific code areas, leaving the majority of code benign [17]. Therefore, coverage-guided fuzzers may waste testing resources on these non-vulnerable areas [18]. This issue becomes more pronounced in smart contract testing, where the complex interaction between contract state and code spaces complicates comprehensive testing, resulting in a significant portion of testing efforts being allocated to benign code areas and contract states.

To bridge the gaps, we aim to design a stateful directed fuzzer for smart contracts, inspired by the concept of directed graybox fuzzing (DGF) [18], [19]. DGF, known for

Ruichao Liang, Jing Chen, Kun He, Ruochen Cao, and Ruiying Du are with Key Laboratory of Aerospace Information Security and Trusted Computing, Ministry of Education, School of Cyber Science and Engineering, Wuhan University, Wuhan, 430072, China. Email: liangruichao, chenjing, hekun, crc2019, duraying@whu.edu.cn. Cong Wu, Yang Liu, and Yueming Wu are with Cyber Security Lab, College of Computing and Data Science, Nanyang Technological University, Singapore. Email: cong.wu, yangliu@ntu.edu.sg, wuyueming21@gmail.com. Ziming Zhao is with Department of Computer Science and Engineering, University at Buffalo, New York. Email: zimingzh@buffalo.edu.

its efficiency in allocating testing resources to specific areas of interest while minimizing unnecessary stress on unrelated parts, has shown effectiveness in various contexts [20]–[23]. However, existing DGFs lack capabilities for stateful fuzzing, primarily due to their focus on the code space while overlooking the essential state space in smart contracts. Moreover, these tools depend heavily on external information like bug reports for target identification [18], limiting their autonomy in identifying smart contract vulnerabilities. Different from the existing DGFs, our motivations include: i) achieving stateful fuzzing in the contract state space through the use of feedback from critical states, and ii) enabling directed fuzzing in the contract code space via automated target identification.

Our Solution. In this paper, we propose VULSEYE, a stateful directed graybox fuzzer for smart contracts. It prioritizes resources on vulnerable code areas and contract states, enhancing the efficiency of vulnerability detection. We utilize static analysis and pattern matching to identify code areas vulnerable to exploits, termed as *code targets*. We design a backward analysis algorithm for bytecode to determine the critical contract states required for vulnerability exploitation, termed as *state targets*. Code and state targets provide feedback in both code and state levels during the fuzzing process of VULSEYE. We propose a novel fitness metric integrating the feedback from code and state levels to achieve directed fuzzing across both contract code and state spaces.

We evaluated VULSEYE’s performance on both closed datasets and real-world scenarios. The evaluation reveals that VULSEYE outperforms SOTA tools in vulnerability detection and proves effective in identifying real-world vulnerabilities. Notably, VULSEYE detected 4,845 vulnerabilities among 42,738 real-world smart contracts, outperforming SOTA tools by up to 9.7 \times . Besides, utilizing VULSEYE, we found 11 previously unknown vulnerabilities in the top 50 Ethereum DApps and have reported them to corresponding authorities.

Contributions. This paper makes following contributions.

- We propose VULSEYE, the first stateful directed graybox fuzzer for smart contracts. It autonomously identifies vulnerable code areas and contract states, then prioritizes testing resources on these targets, achieving stateful and directed exploration of both contract code and state spaces.
- We design a backward analysis algorithm for contract bytecode. It efficiently identifies vulnerable states to specify state targets and provides feedback on the contract state space throughout the fuzzing process.
- We contribute a fitness metric that evaluates the proximity of seeds to testing targets at the contract code and state levels, effectively gauging the potential of seeds to trigger vulnerabilities. This algorithm guides seed scheduling in fuzzing, facilitating targeted vulnerability detection within both code and state spaces.
- We comprehensively evaluated VULSEYE and benchmarked it against SOTA tools. The results show that: i) VULSEYE reaches specific test targets up to 8.4 \times faster than the baselines. ii) It outperforms SOTA tools in terms of both efficiency and quantity on a ground truth vulnerability dataset. iii) It identified 4,845 vulnerabilities among 42,738 real-world smart contracts, surpassing SOTA tools in both

TABLE I
COMPARISON OF REPRESENTATIVE RELATED WORKS.

Features	Data Dependency [†]	Directed Seed Mutation	Adaptive Seed Priority	Exploring Code Space	Exploring State Space [‡]
CONTRACTFUZZER [24]	○	○	○	●	○
HARVEY [8]	○	●	○	●	●
SFUZZ [10]	○	○	○	●	○
CONFUZZIUS [11]	●	○	●	●	●
SMARTIAN [12]	●	○	○	●	○
IR-FUZZ [14]	●	○	●	●	○
ITYFUZZ [16]	○	○	●	●	●
VULSEYE	●	●	●	●	●

* ● = true, ○ = partially true, ○ = false.

[†] Support data dependency analysis to generate meaningful transaction sequences.

[‡] Support testing on specific contract states.

quantity (9.7 \times) and true positive rate. iv) Additionally, it uncovered 11 previously unknown vulnerabilities in the top 50 Ethereum DApps, involving about 2,500,000 USD.

II. BACKGROUND

A. Stateful Smart Contract

Smart contracts [25] are programs running on top of blockchain platforms. They are Turing-complete programs used to automate the execution of an agreement and process assets on blockchain [26]. Within smart contracts, there are a kind of variables called state variables that are created at the contract level and stored permanently on the blockchain [27]. The values of state variables can change based on the contract’s execution, provide a direct reflection of the contract’s states, and impact the contract’s behavior and functionality. State variables make smart contracts *stateful* programs by preserving and managing their states throughout their lifecycle.

B. Graybox Fuzzing

Graybox fuzzing is a software testing method that combines aspects of white-box [28] and black-box testing [29]. It involves injecting semi-random input data, known as seeds, into a program to uncover vulnerabilities and potential issues [30]. Based on the feedback information from the execution of the program under test cases, greybox fuzzers generate new inputs and perform seed scheduling to manage and prioritize these inputs to effectively test the program. Most greybox fuzzing tools are coverage-guided, striving to explore as many program paths as possible [18].

C. Limitations of Existing Contract Fuzzers

Firstly, being stateful programs, smart contracts differ from typical programs by not only having their code space but also having an implicit state space. This complexity adds a layer of difficulty to the exploration of the program space to detect vulnerabilities [31]. Some existing fuzzers take into account the stateful nature of smart contracts. For instance, some [8], [11], [12], [14] attempt to manipulate the contract state by leveraging the read-write relationship of variables to create transaction sequences, while some [16] employ snapshot to preserve the contract state for rapid resumption. However, these approaches lack knowledge of the specific state required

```

1
2 contract FancyBank{
3   mapping(address => uint256) private balances;
4   uint256 dueDate = 0, unlock = 0;
5   event WithdrawalFailed(address user, uint256 amount);
6
7   function deposit(uint256 amount) public payable{
8       require(msg.value >= amount)
9       balances[msg.sender] += amount;
10
11  function setState(uint256 time,uint256 State) public{
12      dueDate = time;
13      unlock = State;
14
15  function withdraw(uint256 amount) public {
16      require(balance[msg.sender] >= amount);
17      if(dueDate >30 && dueDate < 40 && unlock == 1){
18          msg.sender.call{value: amount}();
19          balance[msg.sender] -= amount;
20      }else{
21          emit WithdrawalFailed(msg.sender, amount);
22      }
23  }
24  /** Other functions */
25 }

```

Fig. 1. Motivating Example

to trigger vulnerabilities, resulting in blind and inefficient exploration of the contract state space, as shown in Table I.

Secondly, existing tools predominantly rely on coverage-guided fuzzing strategies to achieve greater code coverage. However, this may result in an inefficient allocation of effort on testing non-vulnerable code areas, while overlooking critical code areas. The stateful nature of smart contracts further compounds this issue.

These limitations are exemplified in Section III. As shown in Table I, VULSEYE surpasses existing state-of-the-art (SOTA) fuzzers by facilitating targeted exploration in both the code and state spaces of smart contracts. These advancements collectively lead to more efficient vulnerability detection.

III. MOTIVATING EXAMPLE

This section presents a motivating example: a contract that allows users to deposit funds and withdraw them as needed. The contract's simplified segment is illustrated in Figure 1. It is vulnerable to reentrancy because of the `withdraw()` function. Intended for users to withdraw funds, this function mistakenly transfers funds before updating the balance variable. This allows an attacker to reenter the `withdraw()` function multiple times during the external call in Line 19, ultimately draining the contract's funds. However, it is not easy to exploit this vulnerability in practice, as illustrated below.

A. State Feedback is Necessary

As discussed in Section I, stateful programs typically necessitate reaching specific states before the vulnerability is manifested. In the contract depicted in Figure 1, invoking the `setState()` function to set the contract states to: $\text{dueDate} \in (30, 40)$ and $\text{unlock} \in \{1\}$ is necessary before triggering the reentrancy vulnerability in Line 19.

Now considering a scenario with three seeds, each seed comprises two transactions triggering two functions:

Seed A: ①deposit(10) \rightarrow ②withdraw(10)

Seed B: ①setState(100,0) \rightarrow ②withdraw(10)

Seed C: ①setState(50,1) \rightarrow ②withdraw(10)

While all three seeds ultimately follow the *else* branch (Line 21) and miss the vulnerable code location (Line 19), seed C is closer to exploiting the vulnerability. This is because seed C modifies the critical state variables to be closer to the states that are necessary for triggering the vulnerability (i.e., $\text{dueDate} \in (30, 40)$, $\text{unlock} \in \{1\}$). Consequently, allocating more testing resources to seed C during the fuzzing process can expedite the discovery of the vulnerability. However, current approaches such as CONFUZZIUS [11] and HARVEY [8] rely solely on combining transactions in a Read-After-Write order (a trait present in all three seeds in this example) to manipulate the contract state, but they lack feedback from the contract state for systematic stateful exploration. As a result, they fail to recognize `dueDate` and `unlock` as key state variables for triggering the vulnerability, thus hindering their ability to differentiate between the relative effectiveness of these three seeds. This observation underscores the significance of using feedback from the state space in the fuzzing of stateful programs like smart contracts.

B. Pursuing Comprehensive Coverage May Reduce Efficiency

The vulnerability can be exploited when an attacker manages to pass the *if* condition in Line 18, and execute the external call in Line 19. However, exploiting this vulnerability in the *if* branch is time-consuming due to the strict state conditions it requires. Test cases are more likely to take the *else* branch. In such situations, the current approaches with a coverage-guided fuzzing strategy would typically shift focus to other non-critical functions to enhance overall code coverage, since intensifying code coverage through testing the `withdraw()` function has presented significant challenges. This observation underscores the inefficiency in blindly chasing code coverage, as it often results in disproportionate testing resources being spent on non-vulnerable code sections. A more effective strategy for uncovering vulnerabilities would be to strategically focus testing efforts on code segments with a higher susceptibility to vulnerabilities.

IV. OVERVIEW

We brief threat model and give an overview of VULSEYE.

A. Threat Model

Smart contracts, with their nature of blockchain-based execution and capability of handling digital assets, introduce a set of unique vulnerabilities distinct from traditional platforms. These vulnerability threats often arise from developer oversights or inherent blockchain characteristics, and exploiting them can cause significant and irreversible financial damage. They can be categorized into different types including *Ether Leaking*, *Reentrancy*, *Controlled Delegatecall*, *Suicidal*, etc., which have been well illustrated in previous studies [13], [24]. The attacker has the ability to access data on the public blockchain, analyze smart contracts to identify vulnerabilities, and knows how to exploit them by sending transactions or deploying new contracts.

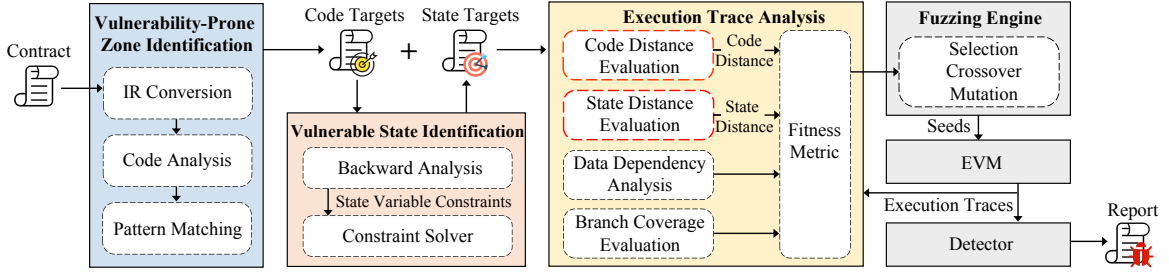


Fig. 2. Overview of VULSEYE.

Algorithm 1: Code Target Identification.

Input : Contract C
Output: CodeTargets

```

1:  $irCode, cfg \leftarrow IRConversion(C)$ 
2:  $res \leftarrow CodeAnalysis(irCode, cfg)$ 
3: for  $node$  in  $cfg.nodes$  do
4:   for  $ir$  in  $node.irs$  do
5:      $offset \leftarrow PatternMatching(ir, cfg, res)$ 
6:      $CodeTargets \leftarrow Locating(offset, cfg)$ 
7:   end
8: end

```

VULSEYE aims to provide an efficient stateful directed fuzzing solution for vulnerability detection, ensuring the safety of smart contracts before deployment and enabling defenders to timely implement preventive measures for on-chain vulnerable contracts. This will mitigate the threats that contract vulnerabilities pose to blockchain and decentralized ecosystems.

B. Overview of Vulseye

To implement stateful directed fuzzing in both contract code and state spaces to hunt vulnerabilities, we address three main challenges: *locating vulnerability-prone code* in smart contracts, *identifying critical contract states* required to exploit vulnerabilities in these code areas, and *guiding fuzzing to test specific states and code* for accelerated vulnerability detection.

As shown in Figure 2, to locate vulnerable code areas, a smart contract is first compiled and converted into an Intermediate Representation (IR), which serves as the foundation for subsequent static code analysis. Based on the result of code analysis, VULSEYE conducts pattern matching to identify code areas that could potentially harbor vulnerabilities, referred to as code targets. To identify vulnerable contract states necessary for exploiting vulnerabilities, VULSEYE conducts backward analysis on the control flow graph (CFG) starting from code targets, gathering the constraints associated with state variables. Then, a constraint solver is applied to resolve the state constraints and specify the state targets. During the fuzzing loop, VULSEYE analyzes execution traces, evaluating both the code and state distances of seeds to estimate their propensity to trigger vulnerabilities. Seeds with a greater potential to trigger vulnerabilities are given higher priority in the seed scheduling process, and values from state targets are incorporated into the mutation pool to generate vulnerable contract states. Vulnerability detection is carried out by a detector that evaluates the execution trace based on test oracles.

TABLE II
HAZARDOUS BEHAVIORS THAT MAY LEAD TO VULNERABILITIES.

Hazardous Behaviors Description	Potential Bugs
Has <i>Payable</i> function.	
No <i>Suicide</i> or <i>Selfdestruct</i> operation.	Lock Ether
No high-level or low-level <i>calls</i> that send ether.	
<i>Delegatecall</i> using <i>msg.data</i> as destination.	Controlled Delegatecall
<i>Delegatecall</i> using <i>msg.data</i> as arguments.	Dangerous Delegatecall
Use <i>block.data</i> in <i>Require</i> or <i>If</i> statement.	Block Dependency
Send ether following this statement.	
<i>Read</i> a state variable.	
External <i>call</i> after the <i>read</i> .	Reentrancy
<i>Write</i> the same state variable after the <i>call</i> .	
No use of <i>msg.sender</i> as index.	
Send ether with no dependency on <i>msg.value</i> .	Arbitrary Send Ether
Use <i>msg.sender</i> or <i>msg.data</i> as receiver.	
Has function that is not protected.	
Set this function as <i>public</i> or <i>external</i> visible.	Suicidal
Has <i>Suicide</i> or <i>Selfdestruct</i> in this function.	

V. METHODOLOGY

In this section, we detail the design of VULSEYE.

A. Target Identification

To implement directed gray box fuzzing, the first step is to specify potential vulnerable locations as test targets, which, in our solution, involves identifying code and state targets.

Vulnerability-prone Zone Identification. We statically analyze the contracts to identify potentially vulnerable code areas prior to conducting fuzzing, as illustrated in Algorithm 1. We utilize SlithIR [32], an intermediate representation, to represent Solidity code of the contracts under test. A series of code analyses are carried out based on the intermediate representation, such as function property analysis and data dependency analysis (Line 2). Utilizing insights gained from these analyses, we navigate through the contract's CFG, scrutinizing semantic contract behaviors through pattern matching to identify hazardous behaviors that could potentially lead to vulnerabilities (Lines 3-5). The rules of pattern matching are tailored to specific vulnerable behaviors, as outlined in Table II. For instance, a sequence of basic blocks where a state variable is ① initially read before an external call, and ② subsequently written after the external call, may suggest a potential *Reentrancy* vulnerability. To enhance the detection of vulnerable code areas and reduce the risk of false negatives, we traverse the contract code space and conduct pattern matching with an over-approximation principle. To pinpoint

Algorithm 2: State Target Identification.

Input : CFG \mathcal{C} , Target \mathcal{T}
Output: stateTarget

```

1: stateCons  $\leftarrow \emptyset$ 
2: paths  $\leftarrow \text{FindPath}(\mathcal{C}, \mathcal{T})$ 
3: for  $p$  in paths do
4:   pathCons  $\leftarrow \emptyset$ 
5:   branchPoints  $\leftarrow p.\text{branchPoints}$ 
6:   for  $bp$  in branchPoints do
7:     branchCons  $\leftarrow \text{BackwardAnalysis}(p, bp)$ 
8:     pathCons  $\leftarrow \text{pathCons} \cup \text{branchCons}$ 
9:   end
10:  stateCons  $\leftarrow \text{stateCons} \cup \text{pathCons}$ 
11: end
12: stateTarget  $\leftarrow \text{Solve}(\text{stateCons})$ 

```

Algorithm 3: Contract Bytecode Backward Analysis.

Input : path \mathcal{P} , branchPoint \mathcal{B}
Output: branchCons

```

1: branchCons  $\leftarrow \text{InitCons}(\mathcal{P}, \mathcal{B})$ 
2:  $t \leftarrow \text{InitTraceBack}(\mathcal{P}, \mathcal{B})$ 
3: stack  $\leftarrow []$ 
4: reversedIns  $\leftarrow \mathcal{P}.\text{BackwardFrom}(\mathcal{B})$ 
5: while  $\text{NotEmpty}(t) \ \&\& \ \text{NotEmpty}(\text{reversedIns})$  do
6:    $\text{ins}, pc \leftarrow \text{reversedIns}.\text{Pop}()$ 
7:    $\text{stack}[pc] \leftarrow \text{StackReconstruction}(\text{ins}, pc, \text{stack})$ 
8:    $\text{cons}, t \leftarrow \text{ConsCollection}(\text{ins}, pc, \text{stack}, t)$ 
9:    $\text{branchCons} \leftarrow \text{branchCons} \cup \text{cons}$ 
10: end

```

these hazardous behaviors, we locate them within basic blocks of the contract’s CFG as the *code targets* (Line 6).

It is important to note that, the *code targets* identified in this phase do not directly report vulnerabilities; rather, they highlight potential vulnerable code areas (which may not necessarily result in vulnerabilities). These flagged areas undergo further validation through subsequent fuzzing tests, thereby mitigating the risk of high false positive rates inherent in static analysis [33].

Vulnerable State Identification. Smart contracts are stateful programs, and vulnerabilities necessitate specific states for exploitation. While we’ve identified vulnerable code targets and plan to utilize them as guidance for fuzzing, blindly altering the contract’s state can still be inefficient in reaching those code targets and exploiting vulnerabilities. To address this challenge, we introduce a backward analysis algorithm designed to identify vulnerable contract states, facilitating the stateful exploration in the contract state space. Unlike symbolic execution [34]–[36] which is slow and suffers from poor scalability due to the path explosion problem, we focus solely on the code targets-related paths and states, ensuring the efficiency and scalability of our tool. Algorithm 2 reveals how we perform backward analysis on the contract’s CFG and define the corresponding state targets. Algorithm 3 further details the process by which backward analysis systematically gathers all state constraints for a given path.

Algorithm 2 takes as input the contract’s CFG and the code target obtained from the vulnerability-prone zone identification process, and outputs a set of state targets. We handle loops in the CFG using loop unrolling with an upper bound of 20. This compromise balances efficiency and soundness in our analysis, as smart contracts typically avoid extensive loops

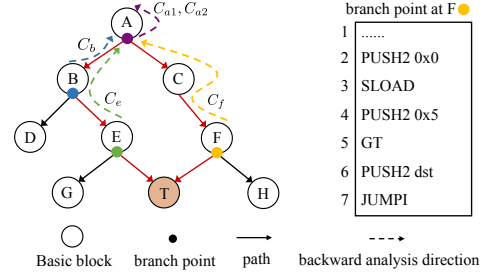


Fig. 3. Example for State Identification.

due to gas limits. We initialize the set of state constraints (stateCons) and use NetworkX [37], an efficient tool for manipulating graph structures, to identify all paths leading to the target basic blocks (Lines 1-2). For each path, we extract the involved branch jumps (typically the JUMPI instruction) in Line 5 and conduct backward analysis at these branch points, gathering constraints on state variables in Line 7. Next, we intersect these constraints in Line 8 to derive state requirements for this path to reach the code target. The union of these path-specific constraints (Line 10) defines the overall state requirements for all paths. Finally, we employ a constraint solver to resolve these constraints, determining the specific contract states required to reach the code target, referred to as the *state targets*. It is worth noting that for state variables of mapping type that take a symbolic address as the key, we do not regard them as state targets because the key cannot be concretely determined in static analysis.

Algorithm 3 further details the implementation of backward analysis at each branch point. In Line 1, we generate an initial constraint based on the jump direction at the current branch point. To be specific, the jump direction is determined by the second parameter of the JUMPI instruction. If the second parameter is non-zero, the program jumps to the *if* branch; otherwise, it proceeds to the *else* branch. As the concrete values of the instruction parameters cannot be known statically, we use symbolic variables to represent them. To transform this initial constraint into a state-related constraint, we designate the symbolic variables in this constraint as the trace-back target (Line 2) and collect their relationships with state variables during the backward traversal of instructions (Lines 5-10). In Line 7, we reconstruct the EVM stack by analyzing instructions in reverse order. Utilizing the recovered stack information, we generate new constraints derived from the initial constraint, subsequently refreshing the trace-back targets in a last-in-first-out (LIFO) manner in Line 8. This iterative process continues until all traced targets manifest as either state variables or state-independent constants, yielding a set of state-related constraints.

Example for State Identification. In Figure 3, the left side illustrates a simplified contract CFG with a determined code target, denoted as node T. As described in Algorithm 2, there are two paths reaching the target, i.e., A-B-E-T and A-C-F-T, with branch points at nodes A, B, E, and F. Backward analysis on these branch points yields constraints C_{a1} , C_{a2} , C_b , C_e , and C_f . The state constraints for the code target are then expressed as $(C_{a1} \cap C_b \cap C_e) \cup (C_{a2} \cap C_f)$. The right

PC	Instruction	Stack (before instruction)	Description	Constraints	Trace Back Targets
7	JUMPI	top →	Jump to dst	$C_f \leftarrow (A \neq 0)$	A
6	PUSH2 dst	top →	Push jump destination	$C_f \leftarrow C_f$	A
5	GT	top →	Greater-than comparison, i.e., $A = GT(B, C)$	$C_f \leftarrow C_f \cap (A = 1 \cap B - C > 0 \cup A = 0 \cap B - C \leq 0)$	$A \begin{cases} \nearrow B \\ \searrow C \end{cases}$
4	PUSH2 0x5	top →	Push 0x5, i.e., $B = PUSH2(0x5)$	$C_f \leftarrow C_f \cap B = 5$	$A \begin{cases} \nearrow B \\ \searrow C \end{cases}$
3	SLOAD	top →	Load from Storage, i.e., $C = SLOAD(D)$	$C_f \leftarrow C_f \cap C = \text{Storage}(D)$	$A \begin{cases} \nearrow B \\ \searrow C \end{cases} \rightarrow D$
2	PUSH2 0x0	top →	Push 0x0, i.e., $D = PUSH2(0x0)$	$C_f \leftarrow C_f \cap D = 0$	End

Fig. 4. Example for Backward Analysis at the branch point F.

side of Figure 3 is a simplified opcode segment of the branch point at node F, and Figure 4 demonstrates the backward analysis process at this branch point. In Figure 4, columns one and two display the program counter and corresponding instructions, while the third column represents the EVM stack information restored based on these instructions. The fourth column describes the current instruction's impact on stack elements, and the fifth column presents constraints generated during backward analysis. The sixth column displays elements in the trace-back list. Backward analysis initiates at the JUMPI instruction (pc=7), where the top of the EVM stack contains two elements, i.e., the jump destination `dst` and value `A`. The JUMPI instruction specifies that if $A \neq 0$, pc jumps to `dst`; otherwise, it jumps to pc+1. Consequently, the first constraint, $A \neq 0$, is derived, based on the jump direction, and the value `A` is added to the trace-back list. The PUSH2 instruction (pc=6) pushes `dst` onto the stack, which does not affect the tracing target. Analyzing the GT instruction (pc=5) which compares the sizes of `B` and `C` and then assigns the result to `A`, produces a new constraint: $(A = 1 \cap B - C > 0) \cup (A = 0 \cap B - C \leq 0)$. Consequently, the element in the trace-back list transitions from `A` to `B` and `C`. The subsequent PUSH2 instruction (pc=4) reveals `B` as 5, and the SLOAD instruction (pc=3) indicates `C` is obtained from the current contract states at slot `D`. The backward analysis continues until the trace-back list is empty. These derived constraints constitute the state constraints for this branch point, denoted as C_f in Figure 3

B. Execution Trace Analysis

We analyze the execution trace of each seed during the fuzzing loop and employ a fitness metric algorithm to assess the proximity of the seeds to both code and state targets. This evaluation serves to measure the likelihood of a seed triggering a vulnerability. By utilizing this approach, we can strategically allocate testing resources to prioritize better seeds, thereby accelerating the testing process.

Code Distance Evaluation. The code target is essentially a program basic block in the contract's CFG. In this section,

we refer to it as the target block. For a basic block BB , we define its distance to a target block TB as

$$d_{TB}(BB) = \begin{cases} d(BB, TB) & \text{if } BB \rightsquigarrow TB, \\ \infty & \text{Otherwise.} \end{cases} \quad (1)$$

where $BB \rightsquigarrow TB$ represents that TB is reachable from BB , and $d(BB, TB)$ represents the number of edges covered by the shortest path from BB to TB . Based on the above definition, for the block BB , we define its block distance $D(BB)$, as the harmonic mean of its distances to all target blocks, denoted as

$$D(BB) = \begin{cases} 0 & \text{if } BB \in \mathcal{Q}, \\ \text{undefined} & \text{if } (TB \in \mathcal{Q} \mid d_{TB}(BB) \equiv \infty), \\ |\mathcal{Q}| \times [\sum_{TB \in \mathcal{Q}} d_{TB}(BB)^{-1}]^{-1} & \text{Otherwise.} \end{cases} \quad (2)$$

where \mathcal{Q} represents the set of target blocks and $|\mathcal{Q}|$ represents the size of \mathcal{Q} . Compared to the arithmetic mean, the harmonic mean distinguishes basic blocks that are closer to specific targets, and the more targets a block can reach, the smaller its distance will be, which facilitates covering as many targets as possible. After a seed S is executed, we analyze its execution trace and compute its code distance. We define the seed's code distance as the average of the smallest n block distances along the seed's execution trace, denoted as

$$\text{CodeDistance}(S) = \frac{1}{n} \times \sum_{BB \in \mathcal{N}} D(BB) \quad (3)$$

where \mathcal{N} represents the set of n basic blocks with the smallest block distances along the execution trace. The value of n is contingent on the number of basic blocks. Notably, we avoid using the average of all block distances as the seed's code distance. Because the global optimal strategy can lead to discrepancy when there are multiple targets, where seeds reaching the target may exhibit larger distances than those that do not reach any target [18]. We also avoid using the smallest block distance as the seed's code distance, as it fails to effectively differentiate between different seeds.

State Distance Evaluation. Each code target corresponds to a specific state target, which is a collection of value ranges

for state variables. For a seed S , we define its distance to a state target ST as

$$d_{ST}(S) = \begin{cases} 0 & \text{if } ST == \emptyset, \\ \sum_{SV \in \mathcal{M}} R(SV, ST) & \text{Otherwise.} \end{cases} \quad (4)$$

where SV represents the state variable, \mathcal{M} represents all state variables involved in ST . $R(SV, ST)$ represents the distance of the value SV to the range ST throughout the entire execution trace. If SV falls within the range ST , $R(SV, ST)$ is considered to be 0. Based on the above definition, for seed S , we define its state distance as the harmonic mean of its distances to all state targets, denoted as

$$StateDistance(S) = \begin{cases} 0 & \text{if } (\exists ST \in \mathcal{P} \mid d_{ST}(S) == 0), \\ |\mathcal{P}| \times [\sum_{ST \in \mathcal{P}} d_{ST}(S)^{-1}]^{-1} & \text{Otherwise.} \end{cases} \quad (5)$$

where \mathcal{P} represents the set of state targets, and $|\mathcal{P}|$ represents the size of \mathcal{P} .

Fitness Metric. We evaluate the potential of a seed to trigger vulnerabilities during execution by assigning a score that integrates both its code distance and state distance. It is defined as

$$D = \alpha \times \mathcal{F}(CodeDistance) + (1 - \alpha) \times \mathcal{F}(StateDistance) \quad (6)$$

$$SC_{bug} = (D + \beta)^{-1} \quad (7)$$

where α and β are constants less than one, and \mathcal{F} is a normalization function to eliminate the difference in magnitude between the code distance and state distance. Specifically, we assign α a value of 0.5 and β a value of 0.1. If a seed consists of multiple transactions, the score for the seed is determined by the transaction with the highest score. The fitness metric also incorporates feedback from branch coverage and data dependencies which are commonly used in graybox fuzzing, resulting in the following representation

$$Fitness = \gamma \times SC_{bug} + (1 - \gamma) \times (SC_{branch} + SC_{dep}) \quad (8)$$

where γ is a constant less than one, SC_{branch} is computed based on the count of newly covered branch edges, and SC_{dep} is computed based on the frequency of state variable writes.

C. Fuzzing Loop

During the fuzzing loop stage, we use a genetic algorithm [38] guided by our proposed fitness metric for seed scheduling. We execute contracts using an independent EVM, which we instrument to manage the persistent storage of contract states and provide execution traces for subsequent analysis.

Seed Initialization. The testing seeds mainly consist of a *function selector*, *calldata*, and *value*. The *function selector* determines the function to be invoked and is computed using the contract ABI. At the beginning of the test, we initialize two seeds for each function of the contract. *Calldata* refers to the function parameters. We first determine the parameter type according to the ABI. For fixed data types, such as `uint256`, we randomly select values within the valid input range. For non-fixed data types, like `string`, we determine a positive number as the data length and generate an input of that length.

Value represents the funds transferred into the contract by the transaction. If the function is marked as payable, we attach an appropriate value to the transaction; otherwise, the *value* is 0.

Selection. We employ the aforementioned fitness metrics to select valuable seeds and allocate more testing resources to them, steering the fuzzing process toward vulnerabilities. Specifically, it adjusts the seed selection strategy by manipulating the probability of a seed being chosen. Seeds with higher fitness are more likely to be selected for subsequent crossover and mutation. In a fuzzing generation \mathcal{G} , we define $P(S)$ as the probability of selecting a seed S . It is computed as

$$P(S) = Fitness(S) / \sum_{G \in \mathcal{G}} Fitness(G). \quad (9)$$

Crossover. We use crossover to generate transaction sequences. If two transaction sequences have a Read-After-Write (RAW) dependency, we combine them in the order of RAW. Otherwise, the two transaction sequences are combined with a certain probability to generate two longer transaction sequences.

Mutation. The calldata of the seeds is mutated either randomly or from a mutation pool with a certain probability P . The mutation pool consists of state targets we identified and values that appeared in previous transactions, and it is continuously expanded during execution. This heuristic method increases the likelihood of passing narrow conditional branches.

Vulnerability Detection. We define a vulnerability as being discovered by analyzing the execution trace to see if it violates the test oracles. Our test oracles are created following those used in previous works such as CONTRACTFUZZER [24] and SMARTIAN [12].

VI. IMPLEMENTATION.

VULSEYE is implemented in Python, and we use py-evm [39], the official Python implementation of EVM, for the execution of smart contracts. We implement the static analysis phase of VULSEYE in an over-approximate way, so that more targets can be tested during the fuzzing phase, reducing the false negatives. Due to the uncertainty of the interactive objects when contracts initiate external calls, ensuring the success of every external call during testing is challenging [40]. To address this, we instrument py-evm to ensure that external calls within contracts consistently return appropriate values, thus preventing execution failures and facilitating cross-contract testing. We modularize our test oracle to facilitate easy extension for supporting more vulnerability detections.

VII. EXPERIMENTS

Research Questions. We conduct experiments to answer the following four questions. Our testing environment is comprised of a server with a 16-core Intel(R)-Xeon(R)-Gold-5218 CPU @2.30 GHz, 340GB of RAM, and the Ubuntu 18.04 LTS operating system.

- **RQ1:** How effective is VULSEYE in guiding testing toward specific targets, and how significantly do the code and state targets contribute to its performance?

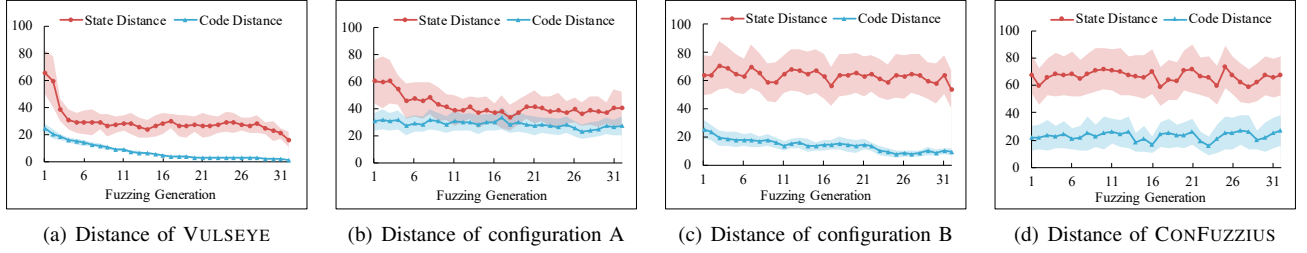


Fig. 5. Tendency of code/state distance on different ablation configurations.

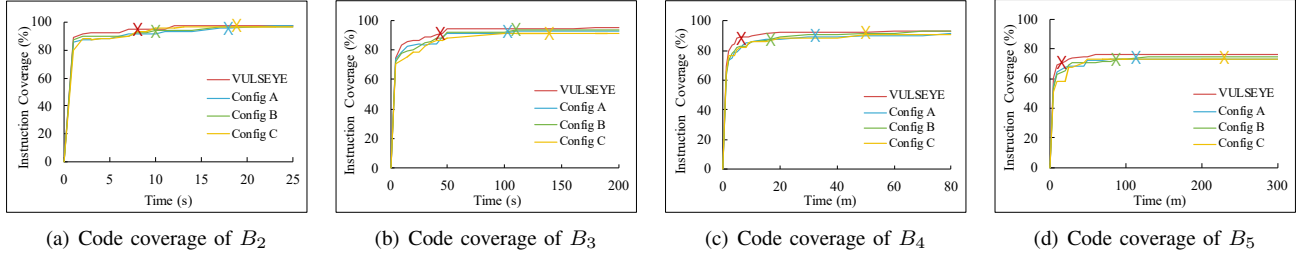


Fig. 6. Tendency of code coverage on different benchmarks. The X indicates when and at what coverage the test target is reached.

TABLE III
COMPARISON OF TIME TO TARGETS.

	T_V	T_A	T_B	T_C	$\frac{T_A}{T_V}$	$\frac{T_B}{T_V}$	$\frac{T_C}{T_V}$
B₁	2.30	4.35	2.71	4.14	1.89	1.18	1.80
B₂	7.70	17.79	9.32	18.25	2.31	1.21	2.37
B₃	42.90	98.67	110.25	138.57	2.30	2.57	3.23
B₄	494.42	1888.68	1112.45	2882.47	3.82	2.25	5.83
B₅	1618.34	6667.56	5502.36	13691.16	4.12	3.40	8.46

* B₁-B₅ are five benchmarks. T_V is the time (seconds) VULSEYE takes to reach targets, and T_A , T_B , T_C are the time of configuration A, B, and C.

- **RQ2:** How effective is VULSEYE in vulnerability detection compared to state-of-the-art approaches.
- **RQ3:** How does VULSEYE perform on code coverage?
- **RQ4:** How does VULSEYE perform on testing real-world smart contracts?

A. RQ1: Testing Specific Targets

We validate the efficacy of VULSEYE in testing specific code areas and create three ablations to evaluate the impact of both code and state targets that we introduced.

Dataset. We craft five benchmarks for RQ1 as existing datasets don't meet our needs and cannot contribute to a fair evaluation. We require: i) each contract to contain a vulnerability in a suitable branch to serve as a test target; ii) sequentially increasing difficulty in reaching these targets; iii) sufficiently challenging branch conditions to extend fuzzing time for observation. To achieve this, each contract is instrumented with hard-to-trigger branches by intentionally imposing restrictions on contract states. From these branches within each contract, we select one to inject a suicidal vulnerability as the testing target. The difficulty of reaching these targets escalates sequentially across these five benchmarks, from B₁ to B₅. Each benchmark consists of 10 contracts, and their average results are used to mitigate variability.

Ablation Study. We test VULSEYE on these five benchmarks, calculating the average time consumed to reach these

targets. We create three ablations, excluding the guidance of either code targets, state targets, or both (designated as configurations A, B, and C respectively). We conducted this experiment five times, and the results are presented in Table III. The first column denotes the five benchmarks, while the second to fifth columns display the average time taken by VULSEYE, configurations A, B, and C, respectively, to reach the targets. Columns six to eight present the speed-up factor achieved by VULSEYE. As depicted in Table III, VULSEYE achieves the shortest time to reach the targets in all five benchmarks compared to the ablations. Notably, the speed-up factor is more pronounced in the tests involving more intricate benchmarks, showcasing a more substantial efficiency improvement, reaching up to $8.46\times$ when compared with configuration C. Moreover, configurations A and B take more time than VULSEYE (1.18 to $4.12\times$), suggesting that both code target and state target contribute to VULSEYE's ability in testing specific code areas.

To visually illustrate how VULSEYE guides testing toward these targets, we depict the tendency of code and state distance during the test of a contract within B₄ in Figure 5, with values averaged over all seeds within a fuzzing generation. In Figure 5a, we note a consistent decrease in the state and code distances of VULSEYE with the increase in fuzzing generations, exhibiting minimal standard deviations. In contrast, configurations A and B exhibit larger fluctuations and standard deviations in code distance and state distance, respectively. Configuration C, lacking guidance in both code and state space, exhibits random changes in both state and code distance in Figure 5d. Additionally, we depict the tendency of average code coverage across different benchmarks in Figure 6. It can be seen that VULSEYE outperforms the ablations in terms of coverage while reaching the targets faster.

Answer to RQ1: With the guidance of code and state targets we introduced, VULSEYE directs testing toward specific code areas effectively. It reaches challenging targets up to $8.4\times$ faster than the baselines.

B. RQ2: Effectiveness of Vulnerability Detection

To answer RQ2, we conduct comparative experiments between VULSEYE and SOTA approaches on a ground-truth dataset consisting of 42 previously exploited projects.

Dataset. Our second dataset is derived from previously exploited projects curated by Durieux et al [41] and Torres et al [11]. However, these contracts lack the desired quantity and complexity for certain vulnerability types. To address this limitation, we have expanded and enriched these contracts to augment both their quantity and complexity. As a result, our dataset encompasses vulnerabilities of various types, including ether leaking (EL), block dependency (BD), reentrancy (RE), controlled delegatecall (CD), dangerous delegatecall (DD), suicidal (SD), and locking ether (LE).

Baselines Selection. In the experiment of RQ2, we compare VULSEYE to SOTA graybox fuzzers that are publicly available, including CONFUZZIUS [11], SFUZZ [10], SMARTIAN [12], IR-FUZZ [14], and ITYFUZZ [16]. We choose CONFUZZIUS because it is the first hybrid fuzzer designed for smart contracts. We choose SFUZZ because its fuzzing strategy is inspired by AFL [42], a highly effective and widely-used fuzzer for C programs. We choose ITYFUZZ because, to our knowledge, it is the first tool to employ snapshot technology in smart contract fuzzing. We choose SMARTIAN and IR-FUZZ since they represent the latest advancements in research and assert superior performance compared to prior works.

Vulnerability Detection. We applied VULSEYE and SOTA tools on this ground-truth dataset. We performed five runs for each contract, employing a fuzzing timeout of 15 minutes for contracts with a code size below 200 lines and extending it to 30 minutes for larger contracts. The experimental results are presented in Table IV, where the first column indicates the contract, the second column specifies the vulnerability type within the contract, and the third column identifies the SWC ID according to the Smart Contract Weakness Classification [43]. The average time to detect vulnerabilities is recorded in columns four through nine.

As shown in Table IV, VULSEYE successfully detected 42 out of 42 vulnerabilities, CONFUZZIUS detected 38 vulnerabilities, and SMARTIAN detected 29 vulnerabilities. SFUZZ and IR-FUZZ lacked test oracles for suicidal and ether leak, and exhibited poor performance in detecting reentrancy. As a result, they only detected 16 and 14 out of 42 vulnerabilities, respectively. ITYFUZZ lacked test oracles for delegatecall types and lock Ether, identifying 17 vulnerabilities. In terms of the time consumed to discover vulnerabilities, VULSEYE outperforms the other tools significantly. Particularly, our tool achieves the shortest time for the detection of 34 out of the 42 vulnerabilities (account for **81%**). Notably, this efficiency improvement is more significant especially for contracts larger than 200 lines (13 out of 14 vulnerabilities, account for **93%**).

TABLE IV
COMPARISON OF TIME TO VULNERABILITIES.

CID	BID	SWC-ID	VE	CF	SF	IF	IT	ST
1	EL	SWC-105	16.87	3.56	-	-	1.02	7.70
2	EL	SWC-105	6.63	26.10	-	-	10.30	900.00
3	EL	SWC-105	13.87	2.24	-	-	2.25	4.96
4	EL	SWC-105	2.86	15.39	-	-	3.63	900.00
5	EL	SWC-105	0.15	4.02	-	-	1.11	30.12
6	EL	SWC-105	23.8	185.60	-	-	32.95	1800.00
7	BD	SWC-120	0.64	4.21	10.38	55.50	0.90	6.57
8	BD	SWC-120	1.01	2.33	900.00	900.00	900.00	4.79
9	BD	SWC-120	1.99	64.50	5.50	5.63	900.00	18.94
10	BD	SWC-120	0.11	3.23	15.37	77.68	900.00	13.61
11	BD	SWC-120	9.15	32.11	5.89	369.22	4.56	68.42
12	BD	SWC-120	1.06	37.24	355.14	488.45	1800.00	539.67
13	RE	SWC-107	2.05	2.14	21.33	9.12	4.55	200.50
14	RE	SWC-107	6.27	15.02	900.00	900.00	9.03	87.43
15	RE	SWC-107	16.89	26.33	900.00	900.00	900.00	109.11
16	RE	SWC-107	3.06	8.12	900.00	900.00	19.48	900.00
17	RE	SWC-107	8.73	23.45	1800.00	1800.00	1800.00	55.34
18	RE	SWC-107	211.29	377.72	1800.00	1800.00	1800.00	1800.00
19	CD	SWC-112	5.45	3.30	5.29	24.06	-	2.23
20	CD	SWC-112	4.12	2.50	5.60	19.23	-	5.28
21	CD	SWC-112	42.21	224.34	900.00	476.45	-	48.33
22	CD	SWC-112	99.09	7.49	5.12	64.22	-	6.69
23	CD	SWC-112	3.06	69.68	10.88	343.78	-	11.61
24	CD	SWC-112	5.45	29.97	35.43	1800.00	-	26.78
25	DD	SWC-112	3.92	6.89	5.67	900.00	-	900.00
26	DD	SWC-112	0.15	2.02	5.30	13.10	-	2.29
27	DD	SWC-112	126.59	900.00	900.00	900.00	-	900.00
28	DD	SWC-112	24.53	34.23	5.98	4.54	-	5.32
29	DD	SWC-112	1.25	53.30	5.70	3.20	-	2.71
30	DD	SWC-112	12.13	31.54	30.31	1800.00	-	25.42
31	SD	SWC-106	1.01	2.21	-	-	1.33	1.88
32	SD	SWC-106	1.42	2.09	-	-	0.89	1.34
33	SD	SWC-106	0.49	4.18	-	-	1.33	6.61
34	SD	SWC-106	4.47	14.71	-	-	5.78	5.92
35	SD	SWC-106	10.54	19.66	-	-	31.26	226.45
36	SD	SWC-106	107.29	382.64	-	-	198.14	159.03
37	LE	SWC-132	✓	✗	✗	✗	✗	✗
38	LE	SWC-132	✓	✓	✗	✗	✗	✗
39	LE	SWC-132	✓	✗	✗	✗	✗	✗
40	LE	SWC-132	✓	✓	✗	✗	✗	✗
41	LE	SWC-132	✓	✓	✗	✗	✗	✗
42	LE	SWC-132	✓	✗	✗	✗	✗	✗

* VE represents VULSEYE, CF represents CONFUZZIUS, SF represents SFUZZ, IF represents IR-FUZZ, IT represents ITYFUZZ, and ST represents SMARTIAN. Since locking ether (LE) vulnerabilities are only reported upon the fuzzing process is terminated, ✓ represents successful detections and ✗ represents failed detections within the specified time limit.

This is due to VULSEYE's prioritization of testing resources in vulnerable areas.

Answer to RQ2: VULSEYE outperforms SOTA tools in smart contract vulnerability detection. It exploited all vulnerabilities in the 42 projects and exhibited the fastest detection time for 81% of them; this figure rises to 93% when the project is larger.

C. RQ3: Code Coverage

We categorized the contracts in RQ2 according to their size into three groups: small, medium, and large, and recorded the average instruction coverage (a form of code coverage) achieved by the six tools. As shown in Figure 7, VULSEYE achieved the final instruction coverage of 96.9%, 95.5%, and 89.9%, surpassing all other tools. While achieving notable results in final instruction coverage, VULSEYE does not solely prioritize maximizing this metric. The light shading in Figure 7 represents the difference between the final instruction coverage and the coverage upon vulnerability discovery. A larger light shading area indicates a lower instruction coverage requirement for vulnerability detection. Comparing the light shading areas, they exhibit the largest proportion in VULSEYE (16.1%, 18.9%, 27.4%) within all three groups. This suggests that with

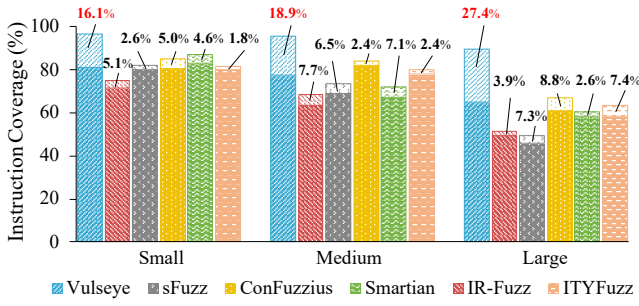


Fig. 7. Instruction coverage upon vulnerability discovery (dark shading) and upon fuzzing is terminated (dark shading and light shading).

the fuzzing strategy we employed, VULSEYE doesn't require achieving such a high code coverage to expose vulnerabilities, leading to a significant improvement in vulnerability discovery efficiency. It is worth noting that, although the code coverage at vulnerability discovery (dark shading) of comparison tools is lower than that of VULSEYE in Figure 7, it's due to their limitation in the notably lower final code coverage. However, their proportion of the dark shading areas remains higher than that of VULSEYE (the lower the better).

Answer to RQ3: VULSEYE achieves the highest code coverage among the SOTA tools. It doesn't necessitate as high code coverage as SOTA tools to detect vulnerabilities, highlighting its significant advantage of prioritizing testing resources on vulnerable code and states.

D. RQ4: Performance on Real-World Contracts

To answer RQ4, we employed VULSEYE on 42,738 real-world smart contracts for bug detection and compared it with other tools. Additionally, we collected top 50 DApps from Ethereum for testing and found previously unknown bugs.

Finding Bugs in 42,738 Real-World Smart Contracts. We obtained the dataset from the study of Durieux et al [41], consisting of 47,587 real-world smart contracts extracted from Ethereum. After excluding contracts that couldn't be compiled by our compiler, this dataset comprised 42,738 real-world contracts. The source code for all contracts in this dataset is available on Etherscan [44].

We applied VULSEYE and other comparison tools on these contracts. We set the maximum number of test cases for each smart contract as 2,000, which we consider sufficient for most exploits. To further evaluate the performance of VULSEYE, we examine the identified vulnerable contracts manually to see whether they are true positives or not. Due to the large number of reported vulnerabilities, we are unable to check all of them. Instead, for each tool, we randomly sample 350 vulnerable contracts for manual verification (50 contracts for each vulnerability type). In cases where the reported results involve less than 50 contracts, all of them will be subject to manual verification. Table V reveals that VULSEYE identified 4,845 vulnerabilities, surpassing other tools by up to $9.7\times$. Additionally, the overall true positive rate for VULSEYE in manual checking is the highest. Specifically,

TABLE V
COMPARISON OF DETECTED VULNERABILITIES AND TRUE POSITIVE RATE.

BID	VULSEYE		CONFUZZIUS		SFUZZ		IR-FUZZ		SMARTIAN		ITYFUZZ	
	#	tp	#	tp	#	tp	#	tp	#	tp	#	tp
EL	408	88%	351	86%	0	N/A	0	N/A	387	86%	838	78%
BD	3720	100%	2032	78%	2737	100%	343	84%	1825	80%	455	84%
RE	426	92%	397	82%	278	78%	99	86%	25	68%	132	64%
CD	12	100%	5	100%	36	33%	10	70%	1	0%	0	N/A
DD	26	100%	7	100%	66	42%	23	48%	8	100%	0	N/A
SD	87	100%	131	80%	0	N/A	0	N/A	94	92%	22	86%
LE	166	64%	113	84%	92	58%	25	60%	0	N/A	0	N/A
Total	4,845	91%	3036	83%	3209	64%	500	75%	2340	84%	1447	77%

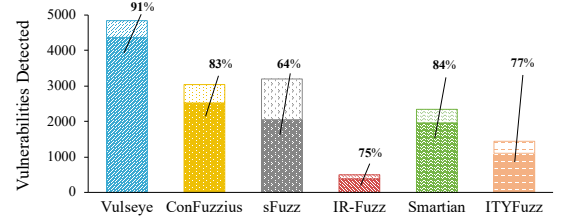


Fig. 8. The total number of identified vulnerabilities reported by each method (dark shading and light shading) and the true positive rates obtained from manual inspection (dark shading).

VULSEYE detected more vulnerabilities in 4 out of 7 vulnerability categories and demonstrated the highest true positive rate among the 6 vulnerability categories. Figure 8 provides an overview of the experimental results, clearly demonstrating that VULSEYE outperforms other tools in both vulnerability discovery quantity and true positive rate.

Case Study. To provide further insight on how the fuzzing strategies implemented in VULSEYE enhance bug detection and surpass existing fuzzers, we present a case study of a contract named *THE_BANK*¹, which contains a reentrancy vulnerability. This vulnerability was detected by VULSEYE within the aforementioned real-world dataset but remained undiscovered by other tools. As shown in Figure 9, to trigger the reentrancy in Line 2, specific values of the state variables `acc.balance` and `acc.unlockTime` are required, which are influenced by other functions. During testing of *THE_BANK* by VULSEYE and other fuzzers, we recorded the average code and state distances to the reentrancy vulnerability in each fuzzing generation, as depicted in Figure 10. VULSEYE showed a consistent decrease in both code and state distances, uncovering the vulnerability by the 37th test generation. It demonstrated that VULSEYE successfully identified this vulnerable location, determined the necessary state conditions, and set corresponding targets to guide the fuzzing. In contrast, other fuzzers exhibited erratic changes in distance and failed to detect the vulnerability in the limited time. Their inefficient allocation of testing resources hinders their ability to set the contract state appropriately and exploit the vulnerability.

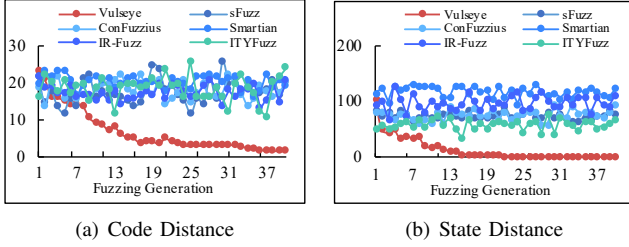
False Positives. During manual verification, we summarized several reasons for false positives reported by VULSEYE. For ether leaking (EL), the false positives attribute to contracts designed for lotteries, where the outward sending of ether is the original intention of the contract designer and should not be considered a vulnerability. For reentrancy (RE), the false positives result from instances where contracts are immune

¹0xCB6fe98097fE7D6E00415Bb6623D5fC3EFFA4E83

```

1 if(acc.balance>=MinSum && acc.balance>=_am && now>acc.
  unlockTime){
2   if(msg.sender.call.value(_am)()){acc.balance-=_am;} }

```

Fig. 9. Code Snippet of *THE_BANK*Fig. 10. Tendency of code/state distance during testing of *THE_BANK*.

to reentrancy attacks due to inherent reentrancy protection or restrictions on external calls. In fact, it is non-trivial to automatically conduct reentrancy attacks on smart contracts. False positives in locking ether (LE) stem from contracts that cannot directly send funds but instead employ `delegatecall` to engage other contracts in fund transfers.

Finding 0-Day Bugs in Top DApps. We collected the top 50 popular DApps listed on BitDegree [45], a famous Web3 platform. These DApps primarily include projects in DeFi, gaming, and gambling, involving millions of transactions on the blockchain, and we utilized VULSEYE to test them. After five runs for each DApp, VULSEYE reported a total of 16 warnings involving 6 projects. After manual verification, we identified 11 previously unknown bugs related to reentrancy, controlled `delegatecall`, and block dependencies within 4 projects. The estimated balance of these four projects on BitDegree is about 2,500,000 USD, and these vulnerabilities put these funds at risk. We have responsibly reported these vulnerabilities to the corresponding authorities.

Answer to RQ4: VULSEYE is effective at identifying vulnerabilities in real-world scenarios, outperforming existing tools by reporting more vulnerabilities and achieving the highest true positive rate.

E. Threats to Validity

First, the performance of VULSEYE varies with the choice of the initial seeds. To ensure fairness, we used the same initial seeds in our comparative experiments with SOTA works, but other initial seeds or benchmarks may produce different results. We have open-sourced VULSEYE to facilitate further evaluation in other studies. Second, the false positive rate in RQ4 may not be entirely accurate. Due to the large number of contracts, we randomly selected 350 contracts from the positive results and manually checked them to estimate the tool’s false positive rate for the entire dataset. Another threat to validity relates to the type of vulnerabilities reported by ITYFUZZ in the comparative experiment. Since ITYFUZZ’s outputs do not directly indicate the vulnerability type, we determine it by manual inspection.

VIII. RELATED WORK

We briefly review related works on smart contract vulnerability discovery and techniques of directed-graybox fuzzing.

A. Smart Contracts Vulnerability Discovery

Smart Contracts Fuzzing. VULSEYE is closely related to work on smart contracts fuzzing [8], [10], [11], [13], [24], [46]. CONTRACTFUZZER [24] is the first black-box fuzzer for Ethereum smart contracts. CONFUZZIUS [11] leverages constraint solving to generate inputs that satisfy complex conditions and generates meaningful sequences of inputs at runtime using dynamic data dependency analysis. Most of these existing approaches are coverage-guided fuzzers that may waste testing resources on benign code areas and lack guidance in exploring the contract’s state space. By contrast, VULSEYE is the first directed graybox fuzzer for smart contracts that concentrates testing resources on vulnerability-prone areas and performs directed fuzzing in both contract code and state spaces.

Smart Contracts Analysis. VULSEYE is related to work on smart contracts static analysis. Numerous static analysis frameworks for smart contracts have been proposed [32], [47]–[50]. ZEUS [48] is a symbolic model checking framework for verification of smart contract correctness and fairness policies. SECURIFY [49] analyzes smart contracts using a dependency graph to check compliance and violation patterns that capture sufficient conditions for proving if a property holds or not. SLITHER [32] is a widely used tool that integrates IR transfer and flow analysis. In our work, we use SlithIR as the foundation of static analysis, to identify code targets prior to conducting fuzz testing.

B. Directed-Graybox Fuzzing Techniques

VULSEYE is closely related to work on directed-graybox techniques [20], [21], [51], [52]. AFLGO [20] is one of the first directed-graybox systems. It utilizes a simulated annealing-based power schedule to drive the seed toward target sites. HAWKEYE [21] boosts the speed for fuzzer to the target sites by utilizing power scheduling, adaptive mutation, and seed prioritization. Existing directed-graybox fuzzers cannot be used to detect vulnerabilities in smart contracts. Firstly, these tools lack the capability to autonomously identify vulnerable areas within smart contracts, relying instead on manual or external information to specify testing targets. Additionally, their exclusive focus on the code space makes them less effective in testing stateful programs such as smart contracts. VULSEYE addresses these challenges and achieves effective stateful directed fuzzing.

IX. CONCLUSION

In this paper, we propose VULSEYE, a stateful directed fuzzer for smart contracts. The key idea is to enhance fuzzing efficiency by exploring critical contract states through feedback from the contract state space and prioritizing testing resources to vulnerable code areas. To achieve this, we introduce code targets and state targets to flag the vulnerable contract code areas and states. We design a novel fitness metric algorithm to steer the fuzzing toward these identified targets. Experimental results show that VULSEYE outperforms existing fuzzers and is effective in finding bugs in real-world scenarios.

REFERENCES

- [1] Z. Zhang, B. Zhang, W. Xu, and Z. Lin, “Demystifying exploitable bugs in smart contracts,” in *Proceedings of International Conference on Software Engineering (ICSE)*, 2023.
- [2] C. Staff, “What was the dao?” 2023. [Online]. Available: <https://www.gemini.com/cryptopedia/the-dao-hack-makerdao>
- [3] REKT, “Poly network - rekt,” 2021. [Online]. Available: <https://www.gemini.com/cryptopedia/the-dao-hack-makerdao>
- [4] X. Zhu, S. Wen, S. Camtepe, and Y. Xiang, “Fuzzing: A survey for roadmap,” *ACM Comput. Surv.*, vol. 54, 2022.
- [5] S. Wu, Z. Li, L. Yan, W. Chen, M. Jiang, C. Wang, X. Luo, and H. Zhou, “Are we there yet? unraveling the state-of-the-art smart contract fuzzers,” in *Proceedings of International Conference on Software Engineering (ICSE)*, 2024.
- [6] J. He, M. Balunović, N. Ambroladze, P. Tsankov, and M. Vechev, “Learning to fuzz from symbolic execution with application to smart contracts,” in *Proceedings of Conference on Computer and Communications Security (CCS)*, 2019.
- [7] G. Grieco, W. Song, A. Cygan, J. Feist, and A. Groce, “Echidna: Effective, usable, and fast fuzzing for smart contracts,” in *Proceedings of International Symposium on Software Testing and Analysis (ISSTA)*, 2020.
- [8] V. Wüstholtz and M. Christakis, “Harvey: A greybox fuzzer for smart contracts,” in *Proceedings of Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE)*, 2020.
- [9] V. Wüstholtz and M. Christakis, “Targeted greybox fuzzing with static lookahead analysis,” in *Proceedings of International Conference on Software Engineering (ICSE)*, 2020.
- [10] T. D. Nguyen, L. H. Pham, J. Sun, Y. Lin, and Q. T. Minh, “Sfuzz: An efficient adaptive fuzzer for solidity smart contracts,” in *Proceedings of International Conference on Software Engineering (ICSE)*, 2020.
- [11] C. F. Torres, A. K. Iannillo, A. Gervais, and R. State, “Confuzzius: A data dependency-aware hybrid fuzzer for smart contracts,” in *Proceedings of European Symposium on Security and Privacy (EuroS&P)*, 2021.
- [12] J. Choi, D. Kim, S. Kim, G. Grieco, A. Groce, and S. K. Cha, “Smartian: Enhancing smart contract fuzzing with static and dynamic data-flow analyses,” in *Proceedings of International Conference on Automated Software Engineering (ASE)*, 2021.
- [13] J. Su, H.-N. Dai, L. Zhao, Z. Zheng, and X. Luo, “Effectively generating vulnerable transaction sequences in smart contracts with reinforcement learning-guided fuzzing,” in *Proceedings of International Conference on Automated Software Engineering (ASE)*, 2023.
- [14] Z. Liu, P. Qian, J. Yang, L. Liu, X. Xu, Q. He, and X. Zhang, “Rethinking smart contract fuzzing: Fuzzing with invocation ordering and important branch revisiting,” *IEEE Transactions on Information Forensics and Security (TIFS)*, vol. 18, 2023.
- [15] S. So, S. Hong, and H. Oh, “SmarTest: Effectively hunting vulnerable transaction sequences in smart contracts through language Model-Guided symbolic execution,” in *Proceedings of USENIX Security Symposium (USENIX Security)*, 2021.
- [16] C. Shou, S. Tan, and K. Sen, “Ityfuzz: Snapshot-based fuzzer for smart contract,” in *Proceedings of International Symposium on Software Testing and Analysis (ISSTA)*, 2023.
- [17] M. Böhme, L. Szekeres, and J. Metzman, “On the reliability of coverage-based fuzzer benchmarking,” in *Proceedings of International Conference on Software Engineering (ICSE)*, 2022.
- [18] P. Wang and X. Zhou, “Sok: The progress, challenges, and perspectives of directed greybox fuzzing,” *CoRR*, vol. abs/2005.11907, 2020. [Online]. Available: <https://arxiv.org/abs/2005.11907>
- [19] M. Böhme, V.-T. Pham, M.-D. Nguyen, and A. Roychoudhury, “Directed greybox fuzzing,” in *Proceedings of Conference on Computer and Communications Security (CCS)*, 2017.
- [20] —, “Directed greybox fuzzing,” in *Proceedings of Conference on Computer and Communications Security (CCS)*, 2017.
- [21] H. Chen, Y. Xue, Y. Li, B. Chen, X. Xie, X. Wu, and Y. Liu, “Hawkeye: Towards a desired directed grey-box fuzzer,” in *Proceedings of Conference on Computer and Communications Security (CCS)*, 2018.
- [22] H. Zheng, J. Zhang, Y. Huang, Z. Ren, H. Wang, C. Cao, Y. Zhang, F. Toffalini, and M. Payer, “FISHFUZZ: Catch deeper bugs by throwing larger nets,” in *Proceedings of USENIX Security Symposium (USENIX Security)*, 2023.
- [23] C. Luo, W. Meng, and P. Li, “Selectfuzz: Efficient directed fuzzing with selective path exploration,” in *Proceedings of IEEE Symposium on Security and Privacy (S&P)*, 2023.
- [24] B. Jiang, Y. Liu, and W. K. Chan, “Contractfuzzer: Fuzzing smart contracts for vulnerability detection,” in *Proceedings of International Conference on Automated Software Engineering (ASE)*, 2018.
- [25] C. D. Clack, V. A. Bakshi, and L. Braine, “Smart contract templates: foundations, design landscape and research directions,” *arXiv preprint arXiv:1608.00771*, 2016.
- [26] G. Wood, “Ethereum: A secure decentralised generalised transaction ledger berlin version,” 2022. [Online]. Available: <https://ethereum.github.io/yellowpaper/paper.pdf>
- [27] P. Bose, D. Das, Y. Chen, Y. Feng, C. Kruegel, and G. Vigna, “Sailfish: Vetting smart contract state-inconsistency bugs in seconds,” in *Proceedings of IEEE Symposium on Security and Privacy (S&P)*, 2022.
- [28] P. Godefroid, M. Y. Levin, D. A. Molnar et al., “Automated whitebox fuzz testing,” in *Proceedings of Network and Distributed System Security Symposium (NDSS)*, 2008.
- [29] K. Kim, S. Kim, K. R. B. Butler, A. Bianchi, R. Kennell, and D. Tian, “Fuzz the power: Dual-role state guided black-box fuzzing for USB power delivery,” in *Proceedings of USENIX Security Symposium (USENIX Security)*, 2023.
- [30] H. Liang, X. Pei, X. Jia, W. Shen, and J. Zhang, “Fuzzing: State of the art,” *IEEE Transactions on Reliability*, vol. 67, 2018.
- [31] J. Ba, M. Böhme, Z. Mirzamomen, and A. Roychoudhury, “Stateful greybox fuzzing,” in *Proceedings of USENIX Security Symposium (USENIX Security)*, 2022.
- [32] J. Feist, G. Grieco, and A. Groce, “Slither: a static analysis framework for smart contracts,” in *IEEE/ACM International Workshop on Emerging Trends in Software Engineering for Blockchain (WETSEB)*, 2019.
- [33] M. Nachtigall, L. Nguyen Quang Do, and E. Bodden, “Explaining static analysis - a perspective,” in *Proceedings of International Conference on Automated Software Engineering Workshop (ASEW)*, 2019.
- [34] L. Luu, D.-H. Chu, H. Olickel, P. Saxena, and A. Hobor, “Making smart contracts smarter,” in *Proceedings of Conference on Computer and Communications Security (CCS)*, 2016.
- [35] C. F. Torres, J. Schütte, and R. State, “Osiris: Hunting for integer bugs in ethereum smart contracts,” in *Proceedings of Annual Computer Security Applications Conference (ACSAC)*, 2018.
- [36] I. Nikolic, A. Kolluri, I. Sergey, P. Saxena, and A. Hobor, “Finding the greedy, prodigal, and suicidal contracts at scale,” in *Proceedings of Annual Computer Security Applications Conference (ACSAC)*, 2018.
- [37] NetworkX, “networkx,” 2023. [Online]. Available: <https://networkx.org/>
- [38] V. Chahar, S. Katoch, and S. Chauhan, “A review on genetic algorithm: Past, present, and future,” *Multimedia Tools and Applications*, vol. 80, 2021.
- [39] ethereum, “py-evm,” 2023. [Online]. Available: <https://github.com/ethereum/py-evm>
- [40] Y. Xue, J. Ye, W. Zhang, J. Sun, L. Ma, H. Wang, and J. Zhao, “xfuzz: Machine learning guided cross-contract fuzzing,” *IEEE Transactions on Dependable and Secure Computing (TDSC)*, 2022.
- [41] T. Durieux, J. F. Ferreira, R. Abreu, and P. Cruz, “Empirical review of automated analysis tools on 47,587 ethereum smart contracts,” in *Proceedings of International conference on software engineering (ICSE)*, 2020.
- [42] google, “american fuzzy lop,” 2022. [Online]. Available: <https://github.com/google/AFL>
- [43] SmartContractSecurity, “Smart contract weakness classification,” 2020. [Online]. Available: <https://swcregistry.io/>
- [44] Etherscan, “Etherscan.io,” 2023. [Online]. Available: <https://etherscan.io/>
- [45] BitDegree, “bitdegree,” 2023. [Online]. Available: <https://cn.bitdegree.org/>
- [46] M. Ye, Y. Nan, Z. Zheng, D. Wu, and H. Li, “Detecting state inconsistency bugs in dapps via on-chain transaction replay and fuzzing,” in *Proceedings of International Symposium on Software Testing and Analysis (ISSTA)*, 2023.
- [47] N. Grech, M. Kong, A. Jurisevic, L. Brent, B. Scholz, and Y. Smaragdakis, “Madmax: Surviving out-of-gas conditions in ethereum smart contracts,” *Proceedings of the ACM on Programming Languages*, vol. 2, no. OOPSLA, 2018.
- [48] S. Kalra, S. Goel, M. Dhawan, and S. Sharma, “ZEUS: analyzing safety of smart contracts,” in *Proceedings of Network and Distributed System Security Symposium (NDSS)*, 2018.
- [49] P. Tsankov, A. M. Dan, D. Drachler-Cohen, A. Gervais, F. Bünzli, and M. T. Vechev, “Securify: Practical security analysis of smart contracts,” in *Proceedings of Conference on Computer and Communications Security (CCS)*, 2018.
- [50] Y. Sun, D. Wu, Y. Xue, H. Liu, H. Wang, Z. Xu, X. Xie, and Y. Liu, “Gptscan: Detecting logic vulnerabilities in smart contracts by

combining gpt with program analysis,” in *Proceedings of International Conference on Software Engineering (ICSE)*, 2024.

- [51] V. Wüstholz and M. Christakis, “Targeted greybox fuzzing with static lookahead analysis,” in *Proceedings of International Conference on Software Engineering (ICSE)*, 2020.
- [52] H. Huang, Y. Guo, Q. Shi, P. Yao, R. Wu, and C. Zhang, “Beacon: Directed grey-box fuzzing with provable path pruning,” in *Proceedings of IEEE Symposium on Security and Privacy (S&P)*, 2022.



Ruichao Liang received the B.E. degree in cyberspace security from Wuhan University, Wuhan, China, in 2020. He is currently working toward the Ph.D. degree with the School of Cyber Science and Engineering, Wuhan University, China. His research interests include web3 security and vulnerability analysis.



IEEE TPDS, IEEE TSC, etc. He was twice runner-up for the best paper at INFOCOM 2018 and INFOCOM 2021. His research interests include the areas of network security, cloud security, and mobile security.



Cong Wu received the Ph.D. degree from the School of Cyber Science and Engineering, Wuhan University, in 2022. He is currently a Research Fellow with the Cyber Security Laboratory, Nanyang Technological University, Singapore. His leading research outcomes have appeared in USENIX Security, ACM CCS, IEEE TRANSACTIONS ON DEPENDABLE AND SECURE COMPUTING, and IEEE TRANSACTIONS ON MOBILE COMPUTING. His research interests include biometric security, system security, mobile security, and web3 security.



Kun He received his Ph.D. from Wuhan University, Wuhan, China. He is currently an associate professor with Wuhan University. His research interests include cryptography and data security. He has published more than 30 research papers in various journals and conferences, such as TIFS, TDSC, TMC, USENIX Security, CCS, and INFOCOM.



Yueming Wu received the B.E. degree in Computer Science and Technology at Southwest Jiaotong University, Chengdu, China, in 2016 and the Ph.D. degree in School of Cyber Science and Engineering at Huazhong University of Science and Technology, Wuhan, China, in 2021. He is currently a research fellow in the School of Computer Science and Engineering at Nanyang Technological University. His primary research interests lie in malware analysis and vulnerability analysis.



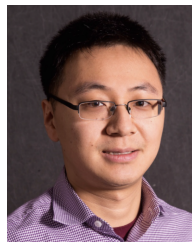
Ruochen Cao received the B.E. degree in cyberspace security from Wuhan University, Wuhan, China, in 2023. He is currently working toward the master's degree with the School of Cyber Science and Engineering, Wuhan University, China.



Ruiying Du received the BS, MS, PH. D degrees in computer science in 1987, 1994 and 2008, from Wuhan University, Wuhan, China. She is a professor at School of Cyber Science and Engineering, Wuhan University. Her research interests include network security, wireless network, cloud computing and mobile computing. She has published more than 80 research papers in many international journals and conferences, such as TPDS, USENIX Security, CCS, INFOCOM, SECON, TrustCom, NSS.



Yang Liu (Senior Member, IEEE) is currently a Full Professor and the Director of the Cyber Security Laboratory, Nanyang Technological University, Singapore. He specializes in software security, verification, software engineering, and artificial intelligence. His research has bridged the gap between the theory and practical usage of formal methods and program analysis to evaluate the design and implementation of software for high assurance and security. His work led to the development of state-of-the-art model checker and process analysis toolkit (PAT). He has more than 200 publications and six best paper awards in top tier conferences and journals. With more than 50 million Singapore dollar funding support, he is leading a large research team working on the state-of-the-art software engineering and cyber security problems.



Ziming Zhao received the Ph.D. degree in computer science from Arizona State University, Tempe, AZ, USA, in 2014. He is an Associate Professor with the Khoury College of Computer Sciences and the Director of the CyberspAcE security and forensics Laboratory (CactiLab), Northeastern University, Boston, USA. His research has been supported by the U.S. National Science Foundation (NSF), the U.S. Department of Defense, the U.S. Air Force Office of Scientific Research, and the U.S. National Centers of Academic Excellence in Cybersecurity (part of the National Security Agency). His research outcomes have appeared in IEEE SECURITY AND PRIVACY, USENIX Security, ACM CCS, NDSS, ACM MobiSys, ACM/IEEE DAC, IEEE RTAS, ACM TISSEC/TOPS, IEEE TRANSACTIONS ON DEPENDABLE AND SECURE COMPUTING, and IEEE TRANSACTIONS ON INFORMATION FORENSICS AND SECURITY. His current research interests include systems and software security, network and web security, and human-centric security. He was a recipient of the NSF CAREER Award and the NSF CRII Award. He was also a recipient of the Test-of-Time Paper Award at ACM SACMAT 2024. Additionally, he has received Best/Distinguished Paper Awards from several prestigious conferences, including USENIX Security 2019, ACM AsiaCCS 2022, ACM CODASPY 2014, and ITU Kaleidoscope 2016.