

# SCRUTINEER: Detecting Logic-Level Usage Violations of Reusable Components in Smart Contracts

Xingshuang Lin\*, Binbin Zhao\*, Jinwen Wang\*, Qinge Xie<sup>†</sup>, Xibin Zhao<sup>‡</sup>, Shouling Ji\*

<sup>\*</sup>Zhejiang University, <sup>†</sup>Georgia Institute of Technology, <sup>‡</sup>Tsinghua University

E-mails: cs.xslin@zju.edu.cn, binbinz@zju.edu.cn, jinwen.22@intl.zju.edu.cn, qxie47@gatech.edu, zxb@tsinghua.edu.cn, sji@zju.edu.cn

**Abstract**—Smart Contract Reusable Components (SCRs) play a vital role in accelerating the development of business-specific contracts by promoting modularity and code reuse. However, the risks associated with SCR usage violations have become a growing concern. One particular type of SCR usage violation, known as a logic-level usage violation, is becoming especially harmful. This violation occurs when the SCR adheres to its specified usage rules but fails to align with the specific business logic of the current context, leading to significant vulnerabilities. Detecting such violations necessitates a deep semantic understanding of the contract’s business logic, including the ability to extract implicit usage patterns and analyze fine-grained logical behaviors.

To address these challenges, we propose SCRUTINEER, the first automated and practical system for detecting logic-level usage violations of SCRs. First, we design a composite feature extraction approach that produces three complementary feature representations, supporting subsequent analysis. We then introduce a Large Language Model (LLM)-powered knowledge construction framework, which leverages comprehension-oriented prompts and domain-specific tools to extract logic-level usage and build the SCR knowledge base. Next, we develop a Retrieval-Augmented Generation (RAG)-driven inspector, which combines a rapid retrieval strategy with both comprehensive and targeted analysis to identify potentially insecure logic-level usages. Finally, we implement a logic-level usage violation analysis engine that integrates a similarity-based checker and a snapshot-based inference conflict checker to enable accurate and robust detection. We evaluate SCRUTINEER from multiple perspectives on 3 ground-truth datasets. The results show that SCRUTINEER achieves a precision of 80.77%, a recall of 82.35%, and an F1-score of 81.55% in detecting logic-level usage violations of SCRs. Moreover, we conduct a real-world analysis on 1,181 on-chain contracts, successfully identifying 13 previously undocumented logic-level usage violations of SCRs, which can result in vulnerabilities. As of now, 9 of these vulnerabilities have been assigned CVE IDs. Notably, the zero-day vulnerabilities identified in one of the most prominent DeFi platforms, with a market capitalization exceeding \$1 billion and a 24-hour trading volume surpassing \$150 million, have been confirmed by their security teams.

## 1. Introduction

Smart contracts are self-executing digital agreements deployed on blockchain platforms, with predefined rules that facilitate trustless and automated transactions without intermediaries. Smart contracts are predominantly implemented in Solidity, the most widely used programming language for Ethereum and other EVM-compatible blockchains. They serve as the foundational infrastructure for decentralized applications across various blockchain ecosystems. In modern smart contract development, Smart Contract Reusable Components (SCRs) significantly streamline the construction of business-specific contracts. SCRs are modular code units, such as interfaces, abstract contracts, libraries, and base contracts, that encapsulate general-purpose logic and can be imported, inherited, or cloned by business contracts [1].

A recent study [2] reveals that approximately 70% of smart contracts incorporate SCRs, highlighting their widespread adoption in practice. However, the risks associated with SCR usage violations have become a growing concern. One particular type of SCR usage violation, known as logic-level usage violation, is becoming increasingly detrimental. This violation exhibits distinct characteristics compared to those in traditional software libraries. Specifically, usage violations in traditional libraries typically stem from breaches of explicitly defined syntactic rules, referred to as syntax-level usage violations. In contrast, most usage violations involving SCRs often relate to logic-level usage violations, which arise from implicit usage constraints that are not explicitly documented in the SCR specifications. These implicit constraints are closely tied to the underlying business logic of smart contracts and must be analyzed across various scenarios. Logic-level usage violations of SCRs can lead to significant vulnerabilities, potentially resulting in substantial financial losses.

Currently, many works have focused on smart contract security, which can be categorized into two main directions based on their analytical perspectives. The first line of work focuses on the detection of contract-level vulnerabilities. These approaches have demonstrated strong effectiveness in identifying common vulnerability types such as reentrancy, improper permission validation, and others [3], [4], [5]. The

second direction centers on the security analysis of contract reuse. Prior efforts in this area have made progress in detecting the misuse of reusable components [2], identifying vulnerabilities in component variants [6], and analyzing security risks in cross-chain component reuse [7]. Despite these advances, logic-level usage violations of SCRs remain largely unaddressed. For instance, in 2024, the ATM token on the BNB Chain suffered an exploit due to logic-level usage violations of SCRs, resulting in a financial loss of approximately \$180,000 [8]. The core limitation of existing detection tools lies in their focus on either syntax-level usage violations or general contract vulnerabilities, without the ability to detect logic-level usage violations. In practice, identifying such violations requires a deeper semantic understanding of the contract source code and the ability to extract implicit usage patterns of SCRs, capabilities that are beyond the scope of current methodologies. This gap highlights the pressing need for a practical and automated system capable of detecting logic-level usage violations of SCRs, thereby addressing an important yet underexplored aspect of smart contract security.

## 1.1. Challenges

To detect logic-level usage violations of SCRs automatically, we have the following key challenges.

**Challenge I: Inference of Relevant Logic-Level Usage Knowledge of SCRs.** Detecting logic-level usage violations of SCRs requires a comprehensive understanding of their intended logic-level usage. The first challenge lies in the automated inference of relevant logic-level usage knowledge, which is a foundational step for enabling accurate violation analysis. This task is particularly difficult because most SCRs are distributed as source code with limited or no accompanying documentation. Even for well-documented components, the implicit nature of logic-level usage constraints poses a substantial challenge for automated inference.

**Challenge II: Identification of Potentially Insecure Logic-Level Usages of SCRs.** The logic-level usage knowledge inferred from SCRs is often abstract and cannot be directly applied to analysis scenarios. The second challenge, therefore, is to bridge the gap between inferred usage knowledge and effective analysis by designing an effective inspector. This process entails interpreting abstract logic-level usage knowledge, translating it into actionable criteria, and analyzing contract scenarios to identify potentially insecure logic-level usages of SCRs. The task is nontrivial, as it involves multi-step reasoning over the contextual constraints of SCRs and the business logic of smart contracts. Furthermore, the inspector must analyze multi-dimensional contract features, making both the reasoning strategy and the construction of an effective reasoning pipeline significantly more complex.

**Challenge III: Accurate Detection of Logic-Level Usage Violation of SCRs.** Even after identifying potentially insecure logic-level usages of SCRs, the analysis may still suffer from false positives due to the coarse granularity of

multi-dimensional reasoning. As such, effective detection of logic-level usage violations of SCRs requires support from fine-grained, context-aware analysis mechanisms. The third challenge, therefore, is to design a robust and practical detection method capable of validating suspected violations with higher precision. This necessitates the integration of fine-grained structural insights, logic-level usage features, and mitigation of false inferences propagated from earlier phases of reasoning.

## 1.2. Methodology

In this paper, we aim to address these challenges to detect logic-level usage violations of SCRs. To this end, we propose SCRUTINEER, an automated and practical system to conduct SCR logic-level usage violation detection in smart contracts. Our design philosophy is as follows.

**First**, to support automated SCR usage analysis, we design a composite feature extraction approach that captures critical features of SCR usages. This approach is built upon a tri-component compression strategy that produces three complementary representations: a composite signature, a logical sequence, and a latent structural embedding. These representations are dynamically adaptable and can be selectively applied based on the requirements of downstream analysis. **Next**, to address **Challenge I**, we propose a Large Language Model (LLM)-powered framework for constructing the SCR knowledge base. This framework operates in two stages. In the first stage, we implement a targeted crawler that collects SCRs from diverse sources, prioritizing those with high usage frequency and a history of associated violation incidents. In the second stage, we deploy an SCR analysis agent, UKEAGENT, which is equipped with a task-planning middleware, comprehension-oriented prompts, and domain-specific tools. UKEAGENT facilitates accurate and robust extraction of logic-level usage knowledge of SCRs. **Then**, to address **Challenge II**, we develop a Retrieval-Augmented Generation (RAG)-driven inspector, LUVAGENT, which integrates two core strategies: a retrieval strategy that efficiently locates relevant reference data from the knowledge base, and an augmented generation strategy that performs both comprehensive and targeted analyses. Together, these strategies enable the detection of potentially insecure logic-level usages. **Finally**, to address **Challenge III**, we implement a logic-level usage violation analysis engine that integrates two complementary techniques: a similarity-based checker and a snapshot-based inference conflict checker. This multi-layered design enables effective and reliable detection of logic-level usage violations of SCRs. By integrating the above components into a unified system, we achieve high precision, recall, and F1-score in detecting logic-level usage violations of SCRs.

## 1.3. Contribution

We summarize our main contributions as follows:

- To the best of our knowledge, we propose SCRUTINEER, the first automated and practical system for detecting

logic-level usage violations of SCRs, bridging a critical gap in existing smart contract security research. Furthermore, we construct the first ground-truth dataset of verified logic-level usage violation cases, encompassing 382 smart contracts collected from the real-world attack incidents. To support future research and facilitate reproducibility, we will release the full implementation of SCRUTINEER at <https://anonymous.4open.science/r/SCRUTINEER-97CF> upon paper acceptance.

- We design an SCR analysis agent, UKEAGENT, which is equipped with a task-planning middleware, comprehension-oriented prompts, and domain-specific tools to automatically construct a reliable SCR usage knowledge base. Our RAG-driven inspector, LUVAGENT, integrates retrieved knowledge with both comprehensive and targeted analyses to identify potentially insecure logic-level usages. Moreover, the detection engine employs two tailored analysis strategies that ensure effective and reliable detection.

- We implement and evaluate SCRUTINEER on ground-truth datasets. The results show that SCRUTINEER achieves a precision of 80.77%, a recall of 82.35%, and an F1-score of 81.55% in detecting logic-level usage violations of SCRs, significantly outperforming existing approaches.

- Our real-world analysis of 1,181 on-chain contracts demonstrates the effectiveness of SCRUTINEER in detecting previously overlooked logic-level usage violations of SCRs. SCRUTINEER identifies 13 zero-day vulnerabilities, 9 of which have been assigned CVE IDs. Notably, the zero-day vulnerabilities identified in one of the most prominent DeFi platforms, with a market capitalization exceeding \$1 billion and a 24-hour trading volume surpassing \$150 million, have been confirmed by their security team.

## 2. Background

In this section, we provide a brief introduction to retrieval-augmented generation, usage violations, and a motivating example.

### 2.1. Retrieval-Augmented Generation

To analyze logic-level SCR usages, SCRUTINEER often requires context about how an SCR is intended to be used and aligned with its logic constraints. Retrieval-Augmented Generation (RAG) enables SCRUTINEER to obtain such information by retrieving relevant SCR instances from the knowledge base and grounding its inference.

RAG is a hybrid framework that supplements LLMs with external knowledge retrieved from a vector database, thereby enhancing the in-context reasoning capability of LLMs [9] [10]. Several mainstream agent development frameworks, such as LangChain [11], RAGflow [12], and LlamaIndex [13], have integrated support for building RAG systems. A typical RAG system follows a two-stage pipeline. In the first stage, the input query is encoded into a dense vector representation and used to retrieve semantically relevant documents from a pre-constructed vector database. In the second stage, the retrieved content is incorporated

into the prompt context of the LLM, enabling it to generate more accurate and informed responses.

RAG systems have demonstrated substantial progress in a wide range of domains, including code generation [14], automated program repair [15], and vulnerability detection [16]. By incorporating external domain-specific knowledge into the reasoning process [17], RAG enhances LLM performance along three key dimensions. First, RAG enables the generation of responses grounded in up-to-date domain knowledge, significantly improving both accuracy and robustness [18]. Second, hallucinations, which are confident yet incorrect outputs, are a well-known limitation of LLMs. RAG mitigates this issue by grounding responses in verifiable external knowledge retrieved from curated sources [19]. Third, RAG enhances resource efficiency. Instead of retraining or fine-tuning large models as datasets evolve, developers can simply expand the underlying vector store and rebuild the retrieval index. This lightweight update mechanism reduces computational overhead and lowers associated operational costs [20] [21].

### 2.2. Usage Violations of SCRs

As SCRs provide standardized building blocks for smart contract development, their correct usage is essential for preventing risky contract behaviors. In practice, securely using SCRs requires attention to two types of usage violations: syntax-level usage violations and logic-level usage violations.

- **Syntax-level Usage Violations.** In smart contracts, a syntax-level usage violation refers to a violation arising from the misinterpretation or misunderstanding of the syntactic rules in SCR specifications. Specifically, such violations may involve passing incorrect out-of-scope arguments to SCRs [22], misusing their return values [23], invoking inappropriate SCRs or overestimating their functional capabilities [2].
- **Logic-Level Usage Violations.** In smart contracts, a logic-level usage violation refers to a violation arising from the neglect of implicit usage constraints of SCRs that are not explicitly documented in their specifications. These implicit usage constraints are often context-dependent, closely tied to the business logic of smart contracts, and require dynamic reasoning across different scenarios. In this paper, we primarily focus on addressing logic-level usage violations of SCRs, as mentioned in Section 1.

### 2.3. Motivating Example

We provide an SCR logic-level usage violation discovered on DeFiHackLabs [24], the case of misaligned parameter logic resulting from an implicit usage constraint, as a motivating example, as shown in Figure 1.

In the `SimpleContract`, the function `swapTokensForCurrency` swaps tokens by calling the SCR `_swapRouter.swap(...)` derived from `Uniswap` [25]. The logic-level usage violation arises from setting the second

```

1 contract SimpleContract{
2   IUniswapV2Router02 public _swapRouter;
3   TokenDistributor public _tokenDistributor;
4   address public currency;
5   function swapTokensForCurrency(uint256 tokenAmount)
6     private {
7     address[] memory path = new address[](2);
8     path[0] = address(this);
9     path[1] = currency;
10    _approve(address(this), address(_swapRouter),
11      tokenAmount);
12    try
13      _swapRouter.swap(tokenAmount,0,path,address(
14        _tokenDistributor),block.timestamp)
15    {} catch { emit Failed_swap(1); }
16    uint256 currencyBal = IERC20(currency).balanceOf(
17      address(_tokenDistributor));
18    if (currencyBal != 0) {
19      IERC20(currency).transferFrom(address(
20        _tokenDistributor), address(this),
21        currencyBal); } } }

```

Figure 1: The violation code snippet in the DeFiHackLabs case (line 11).

parameter of `_swapRouter.swap(...)` to a fixed value of zero, which denotes the minimum acceptable output amount. Despite being syntactically correct, the second parameter must be logically coupled with the first parameter, `tokenAmount` (the input amount). In real token-swap business logic, the minimum output must scale with the input to protect against slippage and price manipulation. By breaking this implicit constraint, the contract unintentionally allows an attacker to execute a sandwich attack, manipulating prices so the victim receives significantly fewer tokens.

While prior studies [2] have investigated the security of SCRs, they fail to detect such issues. The core limitation lies in their inability to capture and reason about implicit logic-level usage patterns, as well as their lack of understanding of contract-specific business semantics.

To reveal such implicit violations, SCRUTINEER employs a unified analysis pipeline that integrates contract feature extraction, knowledge base construction, retrieval-augmented inspector, and a violation detection engine. Through these modules, SCRUTINEER captures hidden logic constraints like the one violated in this example and reliably detects the logic-level usage violation of SCRs.

### 3. SCRUTINEER Design

In this section, we present the detailed design of SCRUTINEER, which aims to automatically detect logic-level usage violations of SCRs within smart contracts. As shown in Figure 2, SCRUTINEER comprises the following four key modules:

- **Contract Feature Extraction:** SCRUTINEER first utilizes the contract feature extraction module to analyze the input smart contract to construct multi-level feature representations that serve as the foundation for subsequent analysis. These representations include a composite signature, a logical sequence, and a latent structural embedding, each capturing the contract’s underlying logic and behavioral structure at different levels of abstraction.

- **SCR Knowledge Base Construction:** Next, the SCR knowledge base construction module is used to collect SCR-related content from diverse sources and constructs a structured knowledge base. This process is guided by an SCR analysis agent, UKEAGENT, equipped with a task-planning middleware, comprehension-oriented prompts, and domain-specific tools, ensuring accurate parsing, normalization, and integration of SCR information.
- **RAG-Driven Inspector:** the RAG-driven inspector module, LUVAGENT, enhances the LLM’s reasoning capabilities by querying the SCR knowledge base to retrieve relevant reference data. It employs a RAG strategy, integrated with both comprehensive and targeted analyses, to identify potential insecure logic-level usages of SCRs.
- **Logic-Level Usage Violation Detection Engine:** the logic-level usage violation detection engine conducts a fine-grained analysis of the potential cases. It employs two complementary techniques: a similarity-based checker and a snapshot-based inference conflict checker. These techniques work in concert to validate potential logic-level usage violations of SCRs within smart contracts.

#### 3.1. Contract Feature Extraction

The purpose of smart contract feature extraction is to identify both explicit and implicit features from the contract and abstract them into a structured representation that supports subsequent analysis. Detecting logic-level usage violations of SCRs necessitates leveraging multi-dimensional program features, as such issues cannot be accurately identified from a single perspective alone. Nevertheless, it is challenging to construct a precise and informative feature representation. The main problem lies in balancing conciseness with expressiveness, ensuring the representation is compact while still capturing the unique structural and semantic properties of the contract. To address the above challenges, we propose a composite analysis approach for contract feature extraction, incorporating the following designs.

**Multi-Level Contract Feature Representation.** Performing in-depth analysis directly at the source-code level often yields noisy or incomplete information, as smart contract source codes lack the structural annotations required for reliable program analysis. To obtain a more comprehensive representation, we construct a composite feature graph for each smart contract. We first generate the Abstract Syntax Tree (AST) of the smart contract by compiling its source code. To ensure compatibility with different Solidity versions, we design and implement an adaptive contract compiler capable of automatically identifying the Solidity version specified in the source code and configuring the corresponding compiler, thereby enabling accurate and efficient AST construction. Based on the generated AST, we extract three fundamental program graphs, including the Program Call Graph (PCG), the Data Flow Graph (DFG), and the Control Flow Graph (CFG). We then integrate these graphs into a unified multi-level composite feature graph that holistically captures the contract’s structural and behavioral characteristics.



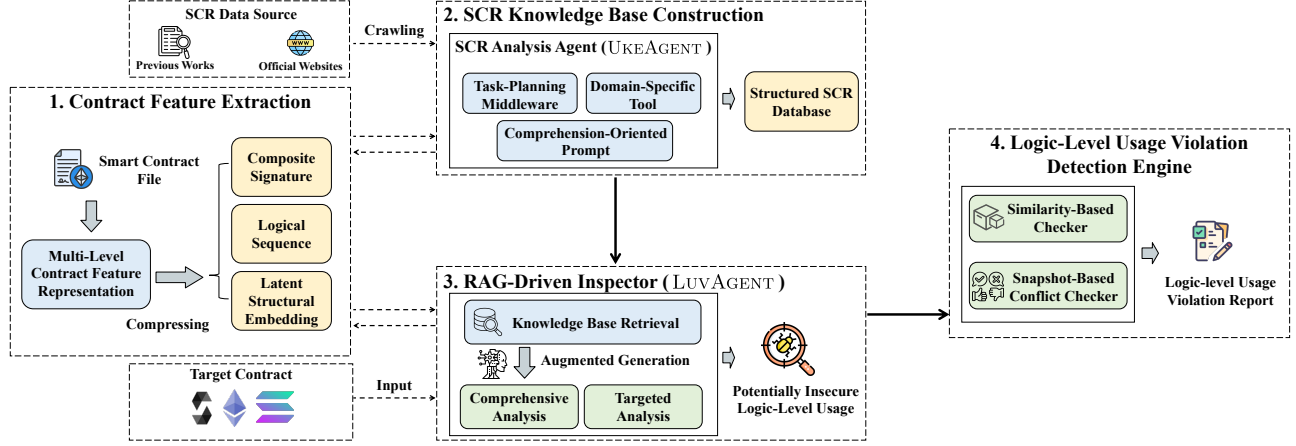


Figure 2: Framework of SCRUTINEER.

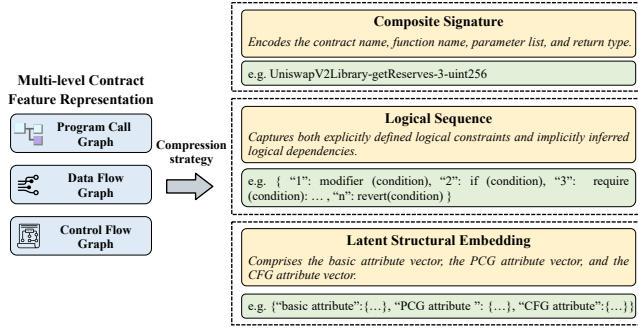


Figure 3: Feature compression process.

**SCR Feature Extraction.** Given the composite graph representation constructed for each contract, it becomes feasible to extract and analyze the SCR usage within smart contracts. However, directly utilizing the complete composite graph introduces substantial computational overhead, which can degrade system performance. Therefore, an effective graph compression strategy is critical. In designing this strategy, the compressed representation must strike a balance between efficiency and informativeness. It should be compact enough to enable efficient matching and retrieval within the SCR knowledge base, while still preserving the essential logic-level usage features and internal structural features. Therefore, we propose a tri-component compression strategy that extracts the composite signature, the logical sequence, and the latent structural embedding. These representations can be used individually or in combination to support the subsequent logic-level usage violation analysis with both flexibility and efficiency.

**1) Composite Signature for Efficient Matching.** Since the composite signature is primarily utilized in scenarios requiring rapid matching, it is intentionally designed to retain only essential information related to SCR usage. Specifically, it encapsulates the contract name, function name, parameter list, and return types. For instance, as shown in Figure 3, we en-

code the function `getReserves(address factory, address tokenA, address tokenB)` in the contract `UniswapV2Library` with a return type `uint256` to `UniswapV2Library-getReserves-3-uint256`. This compact representation enables fast retrieval and facilitates efficient message passing between system modules.

**2) Logical Sequence for Capturing Implicit Constraints.** The logic-level usage feature associated with SCR usage is often entangled with various structural and semantic elements of the smart contract. In our approach, we define the logical sequence to represent the logic-level usage feature. The logical sequence consists of an SCR usage along with its surrounding logical constraints. The algorithm traverses each function node related to SCR usage, examining its structural and semantic context, including both component calls, inheritance, and overriding. For SCR calls, we examine control-flow constructs such as `modifier`, `if-else`, `require`, `revert`, and other relevant nodes to infer the explicit logical constraints that govern the invocation of the SCR. For SCR inheritance and overriding, we analyze both inherited contracts and overridden functions to identify behavioral differences. These differences are then used to infer implicit logical constraints that affect SCR usage in the context of polymorphic behavior. For instance, as shown in Figure 3, the logical sequence is represented in the following format: `{ "1": "modifier (condition)", "2": "if (condition)", "3": "require (condition)", ..., "n": "revert (condition)" }`.

**3) Latent Structural Embedding for Internal Analysis.** The internal structural features of an SCR are also crucial for understanding how its behavior emerges from underlying structural patterns and analyzing the usage of SCRs. In our approach, we categorize SCR structural attributes into three categories: basic attributes, PCG attributes, and CFG attributes. The basic attributes capture essential contract-level information, including the number of nodes, the number of parameters, and the return type. The PCG attributes describe inter-component relationships, such as the number of internal function calls and external calls.

The CFG attributes characterize the fine-grained internal control flow structure, including metrics such as the number of entry-point nodes, variable nodes, expression nodes, conditional nodes, loop nodes, and return nodes. For instance, as shown in Figure 3, the latent structural embedding of `getReserves(...)` is represented in the following format: `{"basic attribute": ..., "PCG attribute": ..., "CFG attribute": ... }`. This structured attribute-based compression enables efficient representation and comparison of SCR internal logic.

### 3.2. SCR Knowledge Base Construction

The purpose of SCR knowledge base construction is to collect SCR code from diverse sources and build a structured reference database to support subsequent violation analysis. Detecting logic-level usage violations of SCRs requires a comprehensive understanding of their intended logic-level usage specifications. Nevertheless, it is challenging to design an automated approach to infer the concealed logic-level usage knowledge of SCRs. The main problem lies in two aspects. First, the sheer number and diversity of SCRs used across smart contracts make it difficult to devise an effective collection strategy. Second, most SCRs are distributed as source code with limited or no accompanying documentation. Even when documentation exists, it tends to emphasize syntax-level usage correctness rather than deeper logic-level usage specifications. To address these challenges, we propose an LLM-powered construction framework that integrates a tailored crawling mechanism with an SCR analysis agent. The framework consists of two core stages: SCR Collection and SCR Analysis.

**SCR Collection.** In this stage, we collect SCR source codes and documents to construct a raw dataset. Nevertheless, it is challenging to download SCRs directly since there is no unified access channel available. To solve the problem of dispersed data resources, we implement a crawler equipped with two strategies to collect SCRs. Specifically, we determine the first data source based on usage frequency. Previous research has analyzed the usage of libraries in real-world Ethereum smart contracts [2]. They find 3,864 types of libraries used in 156,761 smart contracts, and rank all libraries based on their usage frequency. We take the Top 40 libraries among them as one of our data sources [2]. Additionally, we determine the second data source based on historical logic-level usage violation incidents.

**SCR Analysis.** In this stage, we aim to extract precise logic-level usage knowledge of SCRs from the corresponding source code. The knowledge is subsequently converted into structured representations and utilized by both the RAG-driven inspector and the logic-level usage violation detection engine. Nevertheless, accurately extracting SCR logic-level usage knowledge presents significant challenges, as it typically requires deep code understanding, which cannot be effectively achieved through simple keyword matching or template-based heuristics. To address these challenges, we propose an LLM-powered agent, UKEAGENT, for SCR

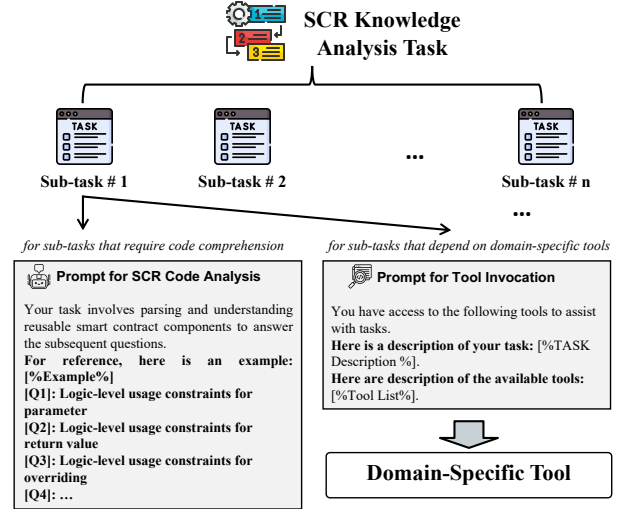


Figure 4: The workflow of UKEAGENT.

analysis, which integrates an LLM with program analysis tools to extract reliable and semantically rich logic-level usage knowledge of SCRs. Designing UKEAGENT for SCR analysis involves overcoming three key challenges: First, the LLMs integrated into UKEAGENT lack inherent capabilities for autonomous task planning and step-wise reasoning, both of which are essential for managing complex extraction workflows. Second, there is a lack of prior research or reference cases to guide the design of prompts specifically tailored for code comprehension and logic-level usage knowledge extraction, which is a critical aspect of UKEAGENT. Third, relying solely on generative LLMs for extracting low-level contract features proves unreliable due to their limited precision in program analysis. We elaborate on our solutions to these challenges as follows.

**1) Task-Planning Middleware.** To address the challenge of task planning, we design a middleware component that orchestrates the overall analysis workflow. As shown in Figure 4, the middleware decomposes each analysis task into a sequence of sub-tasks using a set of predefined rules. Specifically, for sub-tasks that require code comprehension, the middleware directs UKEAGENT to perform comprehension-oriented analysis. For sub-tasks that depend on domain-specific tools, the middleware instructs UKEAGENT to resolve them by invoking the appropriate external tools. This modular task delegation mechanism ensures that each task is handled by the most suitable execution strategy, thereby improving both accuracy and efficiency in SCR analysis.

**2) Comprehension-Oriented Prompt.** To address the challenge of knowledge extraction, we systematically analyze the potential causes of different types of logic-level usage violations of SCRs and formulate comprehension-oriented prompts to guide the LLM. For instance, common sources of logic-level usage violations of SCRs include the absence of proper validation for input parameters passed to SCRs, and the failure to verify the return values of such calls

TABLE 1: Description for each plugin, along with the associated plugin type.

Plugin Type	Plugin Name	Description	Parameter	Return Type
Tool-invoking Plugin	<code>complee_solidity_contract</code>	Complie the specified smart contract	The solidity file path	Class
	<code>read_specified_range</code>	Read a range of lines in a solidity file	The array containing line numbers	List
	<code>get_all_contracts</code>	Get all contract instances in a solidity file	The intermediate representation of solidity	List
	<code>judge_interface</code>	Assess whether a contract has been implemented	The instance of contract class	Bool
	<code>get_all_functions_by_contract</code>	Get all function instances in the specified contract	The instance of contract class	List
	<code>extract_calls_by_function</code>	Get all calls in the specified function	The instance of function class	List
	<code>extract_CFG_by_function</code>	Get the control flow graph in the specified function	The instance of function class	JSON
Response-processing Plugin	<code>json_data_parsing</code>	Parse JSON format data generated by the tools using the specified strategy	1) JSON format data 2) Strategy configuration file	Map
	<code>list_data_parsing</code>	Parse List format data generated by the tools using the specified strategy	1) List format data 2) Strategy configuration file	Map
	<code>exception_parsing</code>	Parse the exceptions generated by the tools and generate the prompt using prompt template	1) Exception data 2) Prompt template ID	String

for correctness or expected behavior. Therefore, we design prompts that instruct UKEAGENT to identify all implicit operations necessary for compliant usages of the given SCR.

Figure 4 illustrates the prompt template designed for UKEAGENT. The questions included in the prompt are carefully crafted to support the extraction of logic-level usage knowledge of SCRs, and are effective for two primary reasons. First, the questions are derived from a comprehensive analysis of hundreds of real-world cases of logic-level usage violations of SCRs, along with an in-depth examination of a large corpus of SCR source code. As such, they are highly relevant to common root causes of usage-related violation issues and collectively cover the vast majority of logic-level usage violation patterns. Second, each question is formulated to elicit clear, concise, and structured responses from the LLM. This design ensures the outputs are machine-readable and directly usable in subsequent analysis stages.

**3) Domain-Specific Tool.** To address the challenge of analyzing low-level features in smart contracts, we integrate a suite of specialized program analysis tools into UKEAGENT. These tools enhance UKEAGENT’s capabilities by enabling it to perform complex, domain-specific analysis tasks that are otherwise beyond the scope of LLMs alone. One such tool is *Slither* [26], a static analysis framework widely used for examining smart contracts. It supports a variety of program analysis techniques and provides rich semantic information. However, integrating *Slither* directly into an LLM-based workflow presents practical obstacles. The LLM is not capable of directly invoking *Slither*, and it also lacks the precision to interpret and filter the tool’s raw output effectively. To bridge this gap, we develop a set of dedicated plugins that mediate between UKEAGENT and the program analysis tools. These plugins are categorized into two functional types: tool-invoking plugins and response-processing plugins. The tool-invoking plugins encapsulate commonly used program analysis functionalities and expose them in a format that can be invoked by UKEAGENT. The response-processing plugins reformat and restructure the outputs returned by these tools into a representation that is compatible with the UKEAGENT’s reasoning process. As summarized in Table 1, each plugin is designed to ensure seamless integration and efficient interaction between pro-

gram analysis tools and UKEAGENT. Figure 4 illustrates the prompt designed for UKEAGENT to invoke these plugins. The prompt begins by presenting a list of tools from which UKEAGENT can make a selection. Once a tool is selected, it is executed via a tool-invoking plugin, and the results of tools are returned to UKEAGENT through a response-processing plugin.

### 3.3. RAG-Driven Inspector

While we have constructed a comprehensive SCR knowledge base, a practical and effective method for leveraging this knowledge in identifying potential insecure logic-level usage of SCRs remains absent. The challenge lies in the fact that SCRs often encapsulate highly business-specific semantics and exhibit complex usage patterns. Therefore, their correct usage requires deep semantic understanding and contextual reasoning, which cannot be achieved through simple rule-based matching. To address the above challenges, we propose a RAG-driven inspector for identifying potential insecure logic-level usages of SCRs, which leverages the in-context learning capabilities of LLMs to perform nuanced behavior analysis. Nevertheless, implementing this approach requires the design of a well-structured and carefully organized workflow that facilitates seamless integration between retrieved knowledge and LLM-based reasoning. In the following, we present our solutions to construct this RAG-driven inspector, LUVAGENT, as shown in Figure 6.

**Retrieval Strategy.** As described in Section 3.2, our knowledge base serves as a reference source to support the LLM during logic-level usage violation analysis. To effectively integrate this knowledge into LUVAGENT, we design a retrieval strategy comprising the following steps. First, we generate composite signatures for all SCRs within the knowledge base. These signatures serve as the indexing keys for similarity-based retrieval. Next, we compute the similarity score between the composite signature of the target SCR and those stored in the knowledge base. As shown in Figure 5, the edit-distance similarity of each item in the composite signature is computed individually, and the overall similarity score is obtained as the weighted sum of the individual item similarities. The weighting coefficients

for each item are determined through distribution-based statistical analysis, with  $\omega^s = (0.079, 0.421, 0.313, 0.187)$ , as detailed in Appendix A. Finally, the SCR with the highest similarity score is selected as the reference candidate, providing contextual guidance for the LLM to assess logic-level usage violations of SCRs.

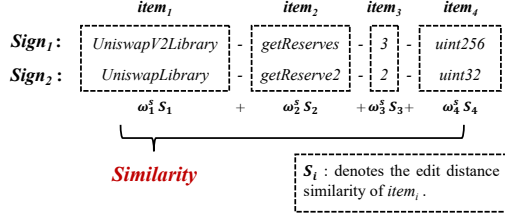


Figure 5: Principle of similarity computation between the composite signature of the target SCR and those archived in the knowledge base.

**Augmented Violation Analysis.** In real-world auditing scenarios, auditors typically follow a systematic process when identifying logic-level usage violations of SCRs. They decompose the analysis of potentially insecure logic-level usages into two phases: comprehensive analysis and targeted analysis. The comprehensive analysis focuses on the security risks introduced by SCR overriding and inheritance across the contract structure. During this phase, auditors assess whether the contract has removed security logic or introduced new logic that compromises security. The targeted analysis, on the other hand, focuses on the security risks associated with SCR calls within the contract. In this phase, auditors begin by identifying potentially insecure SCR usages and then verifying their correctness and security relevance. Inspired by this practical workflow, we design our prompt strategy to emulate the reasoning patterns commonly adopted by human auditors, thereby enabling the LLM to perform analogous stepwise analyses. Furthermore, we integrate a RAG mechanism into our prompts to provide relevant contextual information from the SCR knowledge base, thereby enhancing the model’s precision in assessing logic-level usage violations. This design also mirrors real-world auditing practices, where human auditors routinely consult extensive reference materials to make informed and accurate judgments.

Figure 6 illustrates the prompt strategy designed for LUVAGENT, which emulates the reasoning patterns commonly adopted by human auditors. The strategy consists of two phases: comprehensive analysis and targeted analysis. It first instructs the LLM to assume the role of a code auditor and then decomposes the reasoning process into six sequential stages. In the comprehensive analysis phase (Stages C.1–C.3), the prompts guide the LLM to reason about SCR overriding and inheritance through a series of structured questions. The input for these stages includes the relevant smart contract code and reference information retrieved from the SCR knowledge base. Stage C.1 focuses on the newly added logic introduced through overriding and assesses whether such changes create new security risks.

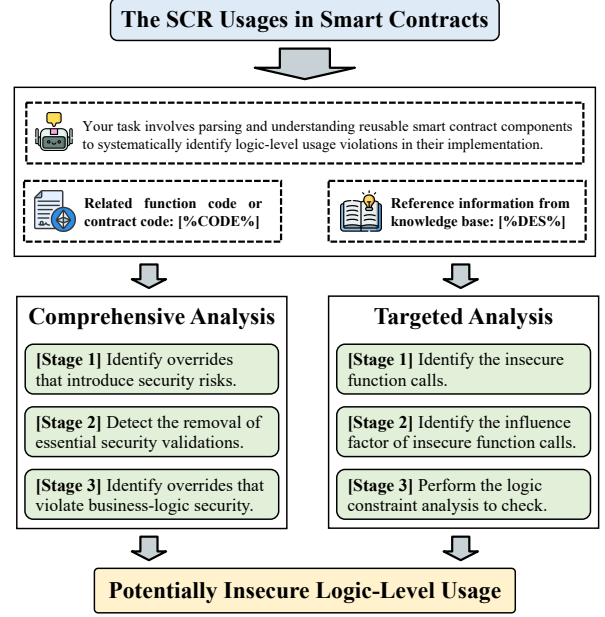


Figure 6: The workflow of LUVAGENT.

Stage C.2 examines the logic removed during overriding to determine whether essential protective measures have been eliminated. Stage C.3 inspects the modified business logic to evaluate whether the functional changes introduce potential security vulnerabilities. In the targeted analysis phase (Stages T.1–T.3), the prompts narrow the focus to the security implications of each SCR call within the contract. Stage T.1 takes as input the contract code, a list of extracted SCR calls, and a preliminary query to identify usages that may pose security risks. Stage T.2 introduces a follow-up prompt that guides the LLM to reason about how each identified SCR usage affects the overall contract security and to articulate its corresponding threat vectors. Stage T.3 enhances the model’s reasoning and verification by retrieving relevant reference information from the SCR knowledge base, informed by the previously identified usage patterns and threat vectors. This reference data supports a more accurate and explainable violation assessment through secondary validation. Moreover, to ensure coherent reasoning and seamless interaction among stages and subsequent system modules, we design specific prompt rules for each stage. These rules constrain the model’s responses to structured formats, such as Yes/No, key–value pairs, or JSON, thereby facilitating consistent interpretation and integration throughout the analysis pipeline.

### 3.4. Logic-Level Usage Violation Detection Engine

Though SCRUTINEER is capable of identifying potential insecure logic-level usages of SCRs through LUVAGENT, we cannot fully rely on the LLM’s output due to its inherent randomness. In particular, even secure SCR usages may occasionally trigger false alerts. This issue arises primarily



from the LLM’s hallucinations and sensitivity to similar SCR usage patterns, which can influence its judgment and lead to inconsistent assessments. To address these challenges, we implement a logic-level usage violation detection engine that performs rigorous validation. This engine complements the LLM’s reasoning with additional layers of structural and context-sensitive verification, ensuring greater accuracy and robustness in logic-level usage violation detection. We craft a similarity-based checker and a snapshot-based inference conflict checker with the following design.

**Similarity-Based Checker.** In the analysis of SCR usages, both the internal structure and the associated logic-level usage features of an SCR are closely related to its compliance with logic-level usage requirements. To capture these factors, we design a similarity-based checker that computes the composite similarity between a target SCR usage instance and a reference instance retrieved from the knowledge base. The similarity-based checker evaluates two key dimensions: the logical sequence and the latent structural embedding. In the logical sequence dimension, we measure the distance-based similarity of relevant logical constraints, denoted as  $S_l$ . In the latent structural embedding dimension, we compute the cosine similarity of numerical and control features, denoted as  $S_n$  and  $S_c$ , respectively. The numerical features capture basic and PCG attributes, while the control features are derived from CFG attributes. All relevant features are extracted using the multi-level feature extraction approach described in Section 3.1. Finally, the principle of the similarity-based checker is defined in Equation 1, which identifies whether potentially insecure logic-level usages ( $u$ ) satisfy the checking criteria. For potentially insecure logic-level usages identified through targeted analysis ( $u \in G_T$ ),  $S_n$ ,  $S_c$ , and  $S_l$  are considered. For those identified through comprehensive analysis ( $u \in G_O$ ), only  $S_n$  and  $S_c$  are considered, since SCR overriding and inheritance do not involve external logical constraints as SCR calls do. A lower composite similarity score indicates a larger deviation of the SCR usage from the compliant usage patterns in the knowledge base. Specifically, lower values of  $S_n$  and  $S_c$  imply that the internal structure of the SCR has been inconsistently modified during its usage, whereas a lower  $S_l$  suggests that fewer logical security constraints are enforced when invoking the SCR.

$$S(u) = \begin{cases} \{u \mid \omega^o \cdot (S_n(u), S_c(u)) < \tau_o\}, & u \in G_O, \\ \{u \mid \omega^t \cdot (S_n(u), S_c(u)) < \tau_t, \\ \quad S_l(u) < \tau_l\}, & u \in G_T. \end{cases} \quad (1)$$

where  $\omega^o$  is (0.27, 0.73),  $\omega^t$  is (0.42, 0.58),  $\tau_o$  is 0.92,  $\tau_t$  is 0.68, and  $\tau_l$  is 0.90. The weights and thresholds are empirically determined, and their soundness is verified through sensitivity analyses, as detailed in Appendix B.

**Snapshot-Based Inference Conflict Checker.** LLMs exhibit strong capabilities in code comprehension and step-by-step security reasoning. Nevertheless, their propensity to generate hallucinated outputs, which are articulate yet incorrect responses, undermines the reliability of the system.

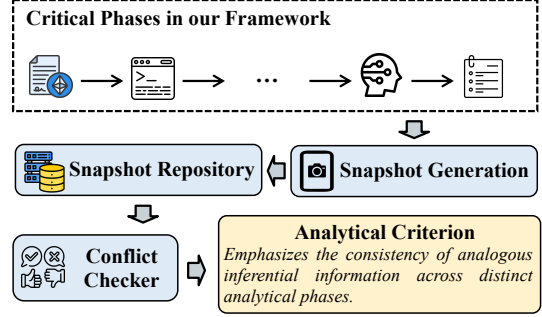


Figure 7: The principles of snapshot-based inference conflict checker

TABLE 2: Snapshot content across critical phases in SCRUTINEER.

Snapshot Content	Feature Extraction (3.1)	LUVAGENT (3.3)			
		Stage 1	Stage 2	Stage 3	Stage 4
1) Signature of SCR	✓	✓	✓	✓	✓
2) Security of SCR	✗	✓	✗	✓	✗
3) Definition of SCR	✓	✓	✗	✓	✗
4) Parameters of SCR	✓	✓	✗	✗	✓
5) Return Type of SCR	✓	✓	✗	✗	✓
6) Parent Contract of SCR	✓	✓	✗	✗	✗
7) Overridden Function of SCR	✓	✓	✗	✗	✗
8) List of Related SCR calls	✓	✗	✓	✗	✓

**Note:** ✓ indicates the information is captured at the corresponding phase, ✗ indicates absence.

Specifically, when inference tasks exceed the LLM’s capacity, it often generates plausible-sounding answers regardless of correctness, leading to a high rate of false positives. To mitigate this issue, we propose a snapshot-based inference conflict checker that validates the consistency of intermediate reasoning steps throughout the analysis pipeline of SCRUTINEER.

Figure 7 illustrates the core principles of the snapshot-based inference conflict checker. The criterion of this checker is to ensure the consistency of analogous inferential information across distinct reasoning phases, thereby enhancing the overall trustworthiness of the analysis. We capture snapshots at critical phases of the SCRUTINEER workflow, each representing a summary of inferential outputs. Through empirical observation, we identify a set of analogous inference points across these phases, as summarized in Table 2. For instance, if our feature extraction identifies an SCR usage with no parameters, but LUVAGENT incorrectly interprets it as an insecure parameter passing, this constitutes a conflict. We will treat such instances as LLM hallucinations and discard them. These snapshots are stored in the snapshot repository, enabling systematic cross-phase consistency checks. The consistency analysis evaluates whether analogous inferential information, such as SCR signature, SCR security, or SCR definition, remains coherent across different analytical phases. If significant inconsistencies are detected, we infer that the task likely exceeds the LLM’s reliable inference capacity and discard the result to avoid propagating hallucinated outputs.

## 4. System Evaluation

In this section, we evaluate the performance and capabilities of SCRUTINEER from multiple perspectives. Our evaluation is structured around the four key research questions:

- **RQ1 (Knowledge Base Construction):** How accurately can SCRUTINEER analyze SCR data and construct the associated logic-level usage knowledge base? (Section 4.2)
- **RQ2 (Violation Detection):** How effective is SCRUTINEER in detecting logic-level usage violations of SCRs? (Section 4.3)
- **RQ3 (Ablation Study):** How do the critical components of SCRUTINEER contribute to its overall detection performance? (Section 4.4)
- **RQ4 (Real-world Impact):** How well does SCRUTINEER identify previously overlooked logic-level usage violations of SCRs in real-world on-chain smart contracts? (Section 4.5)

Additionally, **DeepSeek-V3** serves as the foundational model for SCRUTINEER’s evaluation.

### 4.1. Datasets

To enable our evaluation of SCRUTINEER’s performance, we systematically collect 1,563 smart contracts and 246 SCRs from various sources. We create 3 datasets. The detailed description of these datasets is provided as follows. **SCR Dataset ( $D_{SCR}$ ).** The  $D_{SCR}$  dataset consists of 246 SCRs with 7,422 usage instances, collected from two sources as described in Section 3.2. Specifically, among the 246 SCRs, 40 were selected based on usage frequency and 196 were associated with historical logic-level usage violation incidents. These SCRs are chosen for knowledge base construction for two primary reasons: 1) they are widely adopted in real-world smart contracts, as demonstrated by prior research [2]; and 2) they exhibit a higher likelihood of inducing security risks, as evidenced by our analysis of historical attack statistics [24].

**Historical Violation Dataset ( $D_{HV}$ ).** The  $D_{HV}$  dataset consists of smart contracts that have experienced attacks due to logic-level usage violations of SCRs. In particular,  $D_{HV}$  includes 382 smart contracts associated with 51 logic-level usage violation cases, reported by BlockSec [27], DefiHack-Labs [24]. These real-world attack instances serve as high-fidelity benchmarks to evaluate the practical effectiveness of SCRUTINEER in detecting logic-level usage violations of SCRs.

**Real-World On-Chain Dataset ( $D_{RO}$ ).** The  $D_{RO}$  dataset comprises smart contracts deployed on prominent blockchains and sourced from leading audit platforms, such as Code4rena [28]. It includes 1,181 contracts collected from these sources. These contracts represent modern development trends, incorporating contemporary business models and programming paradigms, thus providing a representative testbed for assessing the generalizability of SCRUTINEER.

TABLE 3: Precision of UKEAGENT in SCR analysis.

Strategy	TP	FP	Precision
Complexity-Based Sampling	192	18	91.43%
Frequency-Based Sampling	202	10	95.28%
Random-Based Sampling	209	13	94.14%
<b>Overall</b>	<b>603</b>	<b>41</b>	<b>93.63%</b>

### 4.2. Evaluation of Knowledge Base Construction

In this step, we evaluate the construction quality of the SCR knowledge base using precision as the primary metric. Within SCRUTINEER, we design an SCR analysis agent, UKEAGENT, to facilitate automated and accurate knowledge extraction for SCR logic-level usage. We begin by assessing the precision of the SCR knowledge extracted by UKEAGENT, which is foundational for downstream logic-level usage violation detection. To further understand the design efficacy of UKEAGENT, we conduct an ablation study. The critical design principle of UKEAGENT lies in its integration of a suite of specialized program analysis tools with LLM, thereby enhancing UKEAGENT’s analytical precision and contextual understanding. For the ablation, we compare two versions of the agent: the full version, UKEAGENT, and a variant without integrated tools, UKEAGENT *w/o Tools*. We conduct the above evaluation on  $D_{SCR}$ .

To evaluate the performance of UKEAGENT on  $D_{SCR}$ , we obtain a total of 7,422 analysis results. Given the difficulty of manually verifying each result at scale, we conduct a sampling-based evaluation to assess correctness and reliability. To enhance the credibility of our evaluation, we adopt three complementary sampling strategies from different perspectives: complexity-based, frequency-based, and random-based sampling. **First**, to assess UKEAGENT’s capability in analyzing complex SCRs, we employ a complexity-based strategy. SCRs in  $D_{SCR}$  exhibit varying levels of complexity, ranging from simple components comprising only 1–2 lines of code, to more complex components involving multiple statements, external calls, or nested invocations. We select SCRs that contain more than four function calls, resulting in a total of 210 SCR usage instances. **Second**, to evaluate UKEAGENT’s capability in handling SCR with similar signatures but variations in internal structures, we adopt a frequency-based sampling strategy. Specifically, we group SCRs by function name and select those groups that contain more than six distinct usage variations, resulting in a total of 212 SCR usage instances. **Third**, we apply a random-based strategy by randomly selecting 222 SCR usage instances from  $D_{SCR}$ . This ensures an unbiased and diverse coverage across different usage scenarios, enhancing the generalizability of our evaluation.

**Precision Evaluation.** As shown in Table 3, UKEAGENT achieves high precision in analyzing SCR usage, with an overall precision of 93.63%. Specifically, it attains a precision of 91.43% for SCR usage instances selected via the complexity-based strategy, 95.28% via the frequency-based strategy, and 94.14% via the random-based strategy. Further

TABLE 4: Comparison of UKEAGENT and UKEAGENT *w/o Tools*.

Strategy	Precision	
	UKEAGENT	UKEAGENT <i>w/o Tools</i>
Complexity-Based Sampling	91.43%	82.86%
Frequency-Based Sampling	95.28%	90.57%
Random-Based Sampling	94.14%	88.29%
<b>Overall</b>	<b>93.63%</b>	<b>87.27%</b>

analysis of the false positives reveals that UKEAGENT occasionally extracts incorrect logic-level usage patterns for certain SCRs. For instance, when the input parameter involves an Ethereum address, UKEAGENT tends to output “non-zero address”, which represents a syntax-level rule rather than a logic-level usage constraint. In contrast, for input parameters involving token amounts, UKEAGENT generates recommended value ranges and their logic relationships with user token balances, which appropriately reflect logic-level usage.

**Tool Integration Impact.** As shown in Table 4, UKEAGENT demonstrates substantial performance improvements over UKEAGENT *w/o Tools*. Specifically, UKEAGENT achieves a precision of 93.63%, whereas UKEAGENT *w/o Tools* yields a lower precision of 87.27%, due to both a reduction in true positives and an increase in false positives. The integration of program analysis tools enables UKEAGENT to significantly enhance the precision of SCR usage extraction by improving the identification of valid patterns and reducing incorrect inferences. We further analyze the false positives by UKEAGENT *w/o Tools*. The primary issue arises from its difficulty in correctly interpreting low-level features of SCR code. For instance, when instructed to infer the logic-level secure range of input parameters, some false positives from UKEAGENT *w/o Tools* can be attributed to its misinterpretation of low-level code features. Specifically, in some SCRs with no or very few input parameters, the agent misidentifies locally defined variables as input parameters. This behavior stems from its limited ability to distinguish between actual input parameters and internal variables. In contrast, the tool-integrated version of UKEAGENT accurately identifies these low-level features, thereby avoiding such misinterpretations.

**RQ1:** UKEAGENT in SCRUTINEER exhibits strong performance in SCR analysis, achieving a precision of 93.63%. Moreover, our tool integration strategy enhances the UKEAGENT’s precision by 6.36%.

### 4.3. Evaluation of Violation Detection

This subsection addresses RQ2. We evaluate the effectiveness of SCRUTINEER in detecting logic-level usage violations of SCRs using three standard metrics: Precision, Recall, and F1-score. To the best of our knowledge, we are the first to address logic-level usage violations of SCRs. Though previous smart contract vulnerability detectors are not specifically designed for this purpose, some vulnerabilities reported by these detectors stem from the logic-level

TABLE 5: Logic-level usage violation detection accuracy.

Tool	$D_{HV}$					
	TP	FP	FN	Precision	Recall	F1-Score
SCRUTINEER	42	10	9	80.77%	82.35%	81.55%
<i>Smartinv</i>	19	7	32	73.08%	37.25%	49.35%
<i>Smartian</i>	3	0	48	100.00%	5.88%	11.11%
<i>Mythril</i>	2	1	49	66.67%	3.92%	7.41%
<i>Falcon</i>	18	59	33	23.38%	35.29%	28.13%
<i>Slither</i>	16	65	35	19.75%	31.38%	24.24%

usage violations of SCRs. Therefore, we compare SCRUTINEER against representative vulnerability detectors. Specifically, we select five widely used detectors: *Smartinv* [5], *Smartian* [29], *Mythril* [30], *Falcon* [31], and *Slither* [26]. *Smartinv* is an LLM-powered verification framework that leverages LLMs to identify the vulnerability and infer security invariants. Comparing SCRUTINEER with LLM-based approaches such as *Smartinv* not only allows for evaluating its performance advantages but also helps mitigate potential risks of LLM-induced data leakage. *Smartian* is a grey-box fuzzer that utilizes data-flow analysis to guide input generation. *Mythril* is a symbolic execution engine tailored for smart contract vulnerability detection. Both *Falcon* and *Slither* are static analysis tools that detect vulnerabilities based on a large corpus of hand-crafted rules. For a fair comparison, we only consider the subset of vulnerabilities reported by these tools that are attributable to logic-level usage violations of SCRs. We evaluate all tools on  $D_{HV}$ .

As shown in Table 5, SCRUTINEER outperforms all baselines in detecting logic-level usage violations of SCRs. It achieves 42 true positives, 10 false positives, and 9 false negatives, corresponding to a precision of 80.77%, a recall of 82.35%, and an F1-score of 81.55%. We further analyze the causes of false positives and false negatives in SCRUTINEER. First, imprecise SCR knowledge extracted by UKEAGENT constitutes the primary cause of both false positives and false negatives, as perfect extraction of logic-level usage semantics remains inherently challenging. Second, inaccuracies in the analysis performed by UKEAGENT may also contribute to certain false positives and false negatives.

In addition, we further investigate the false positives and false negatives produced by the baselines. For *Smartinv*, its primary limitation lies in the inability to generate accurate and verifiable invariants necessary for effective detection. *Smartian* and *Mythril* struggle to simulate the execution paths that lead to logic-level usage violations, and their predefined detection patterns are often insufficient for capturing such issues. *Falcon* and *Slither* employ a large set of detection rules derived from historical vulnerability patterns. While this rule-based approach enables them to identify a greater number of true positives compared to other baselines, it also leads to a high incidence of false positives due to the limited contextual understanding and coarse granularity of these predefined rules.



TABLE 6: Comparison of SCRUTINEER and three variants.

Variant	$D_{HV}$		
	Precision	Recall	F1-Score
SCRUTINEER	<b>80.77%</b>	<b>82.35%</b>	<b>81.55%</b>
SCRUTINEER w/o UKEAGENT	55.77%	56.86%	56.31%
SCRUTINEER w/o LUVAGENT	30.07%	84.31%	44.33%
SCRUTINEER w/o <i>Sim Checker</i>	47.25%	84.31%	60.56%

**RQ2:** SCRUTINEER demonstrates strong performance in detecting logic-level usage violations of SCRs, achieving a precision of 80.77%, a recall of 82.35%, and an F1-score of 81.55%.

#### 4.4. Ablation Study

This subsection addresses RQ3. We conduct the ablation study from two perspectives. **First**, to evaluate the contribution of critical modules to violation detection accuracy, we evaluate three variants of SCRUTINEER: 1) SCRUTINEER w/o UKEAGENT, 2) SCRUTINEER w/o LUVAGENT, and 3) SCRUTINEER w/o *Sim Checker*. Specifically, SCRUTINEER w/o UKEAGENT disables the SCR knowledge base constructed by UKEAGENT (Section 3.2). SCRUTINEER w/o LUVAGENT disables the RAG-driven inspector powered by LUVAGENT (Section 3.3). And SCRUTINEER w/o *Sim Checker* disables the similarity-based checker (Section 3.4). We perform the above variants in the dataset  $D_{HV}$ . **Second**, to evaluate the role of the snapshot-based inference conflict analysis in mitigating LLM hallucinations, we perform an in-depth analysis of the outputs generated by LUVAGENT.

**Impact on Detection Accuracy.** As shown in Table 6, SCRUTINEER consistently outperforms all ablated variants. In particular, SCRUTINEER w/o UKEAGENT achieves a precision of 55.77%, a recall of 56.86%, and an F1-score of 56.31%, indicating that the absence of a structured knowledge base significantly degrades performance. SCRUTINEER w/o LUVAGENT achieves a precision of 30.07%, a recall of 84.31%, and an F1-score of 44.33%, underscoring the importance of the RAG-driven inspector in balancing precision and recall. SCRUTINEER w/o *Sim Checker* achieves a precision of 47.25%, a recall of 84.31%, and an F1-score of 60.56%, demonstrating the necessity of reducing false positives.

**Effectiveness in Mitigating Hallucinated Outputs.** As shown in Table 7, LUVAGENT initially produces 2,684 candidate outputs. The snapshot-based inference conflict checker identifies 1,957 outputs (72.91%) as hallucinated due to the inconsistencies among their snapshots. Consequently, only 727 outputs are retained for further analysis. This high filtering rate confirms the checker’s effectiveness in suppressing unreliable LUVAGENT-generated information, thereby enhancing the overall robustness and reliability of SCRUTINEER.

TABLE 7: Effectiveness of snapshot-based inference conflict checker in reducing LUVAGENT-generated hallucinations.

Processing Stage	Count	Percentage
Initial Generation	2,684	100.00%
Conflicts Identified	1,957	72.91%
Final Validated Entries	727	27.09%

**RQ3:** The integration of UKEAGENT, LUVAGENT, and a similarity-based checker significantly enhances the violation detection performance of SCRUTINEER. Furthermore, the snapshot-based inference conflict checker effectively mitigates the LLM’s hallucinations, substantially enhancing the robustness of SCRUTINEER.

#### 4.5. Real-world Impact

This subsection addresses RQ4. We deploy SCRUTINEER on the real-world on-chain dataset  $D_{RO}$  to assess its effectiveness in detecting logic-level usage violations of SCRs, including the identification of zero-day vulnerabilities resulting from logic-level usage violations in real-world on-chain smart contracts.

SCRUTINEER successfully identified 13 zero-day vulnerabilities within the above smart contracts. As of now, 9 of these zero vulnerabilities have been assigned CVE IDs. These findings have been responsibly disclosed to the respective vendors. Notably, one of the most prominent DeFi vendors, anonymized as *AnonymDeFi*, has acknowledged our report. According to their security team, the identified vulnerability poses a severe threat to their platform’s financial security. *AnonymDeFi* exhibits high liquidity, with a **24-hour Trading Volume-to-Market Capitalization (Vol/Mkt Cap)** ratio ranging from 0.13 to 0.18. Specifically, its **market capitalization exceeds \$1 billion**, and its **24-hour Trading Volume surpasses \$150 million**. Therefore, any vulnerabilities within its contracts could lead to significant economic consequences. Due to the sensitive nature of these zero-day vulnerabilities and the fact that they remain unpatched at the time of writing, we anonymize all identifying details. A representative case study, with critical information redacted for confidentiality, is provided below.

As shown in Figure 8, the code exhibits a logic-level usage violation of SCRs, allowing the sender to mint tokens without enforcing the intended supply cap. Specifically, the `SimpleToken` contract inherits from `ERC20Capped`, an SCR designed to impose a maximum token supply of `1e27` at line 2. Subsequently, the contract invokes `ERC20._mint(msg.sender, _initialSupply)` at line 3 to mint the initial token supply. However, this invocation ultimately calls `ERC20._update` at line 7 to update the total supply, where no check against the `cap` is performed. As a result, the supply cap enforcement is silently bypassed, allowing the sender to mint tokens arbitrarily, which poses significant financial risks to other users. To address logic-level usage violations of `ERC20Capped`, the correct usage requires overriding the `_mint` function in `SimpleToken` to delegate



```

1 contract SimpleToken is ERC20Capped, Ownable {
2   constructor(uint256 _initialSupply, address
      initialOwner) ERC20("SimpleToken", "Simple")
      ERC20Capped(10000000000 * 10 ** 18) Ownable(
      initialOwner)
3   { ERC20._mint(msg.sender, _initialSupply); } }
4 contract ERC20 {
5   function _mint(address account, uint256 value)
      internal {
6     if (account == address(0)) { revert
      ERC20InvalidReceiver(address(0)); }
7     _update(address(0), account, value); }
8   function _update(address from, address to, uint256
      value) internal virtual { ... } }
9 contract ERC20Capped is ERC20 {
10  uint256 private immutable _cap;
11  function _update(address from, address to, uint256
      value) internal virtual override {
12    super._update(from, to, value);
13    if (from == address(0)) {
14      uint256 maxSupply = cap();
15      uint256 supply = totalSupply();
16      if (supply > maxSupply) {
17        revert ERC20ExceededCap(supply, maxSupply);
18      } } } }

```

Figure 8: Bypassing the supply cap enables arbitrary token minting (lines 2-3).

minting to `ERC20Capped.update`, thereby enforcing the required supply cap verification (lines 14–17).

**RQ4:** SCRUTINEER successfully detects 13 zero-day vulnerabilities in real-world on-chain smart contracts, 9 of which have been assigned CVE IDs. Notably, the zero-day vulnerabilities identified in one of the most prominent DeFi platforms, with a market capitalization exceeding \$1 billion and a 24-hour trading volume surpassing \$150 million, have been confirmed by their security team.

## 5. Discussion

Although SCRUTINEER demonstrates strong performance in detecting logic-level usage violations of SCRs, it still has several limitations. First, the presence of false positives and false negatives is primarily attributed to imprecise SCR knowledge extraction and occasional failures in the RAG-driven analysis. While UKEAGENT achieves high precision in extracting SCR knowledge, it struggles with particularly complex cases. Similarly, LUVAGENT currently focuses on analyzing parameter and return-value constraints, and reasoning about the logical differences arising from inheritance and overriding, which may still overlook certain exceptional cases. Future work will focus on enhancing both UKEAGENT and LUVAGENT by integrating additional program analysis tools and developing more sophisticated reasoning strategies. Second, the inherent randomness and hallucination tendencies of LLMs can undermine the overall reliability of SCRUTINEER. Although the snapshot-based inference conflict checker effectively mitigates these issues, it does not fully eliminate them. As a future direction, we plan to fine-tune a dedicated LLM tailored to logic-level usage violation detection, aiming to further improve reasoning accuracy and system robustness.

## 6. Related work

**Smart Contract Analysis.** Many works have been conducted on smart contract security, primarily through static and dynamic analysis techniques. Static analysis methods examine source code or bytecode to capture deep structural or semantic features for vulnerability detection. Representative tools include *Slither* [26], *AutoAR* [32], *Securify* [33], *Zeus* [34], and *MadMax* [35]. Dynamic analysis approaches assess contract behavior through runtime testing. Notable tools include *Smartian* [29], *ItyFuzz* [4], *Harvey* [36], *ContractFuzzer* [37], and *Echidna* [38]. Currently, LLMs have also been applied to smart contract analysis, such as *PROMFUZZ* [39], *PropertyGPT* [3], *SmartInv* [5], and *GPTScan* [40]. For instance, *PropertyGPT* leverages LLMs to automatically generate verification properties for unknown code, enabling comprehensive formal analysis. Nevertheless, applying these methods in our task is not trivial since they are not designed to detect logic-level usage violations of SCRs.

**Reusable Components Analysis.** Numerous studies have been devoted to detecting security issues arising from component reuse across a wide range of software systems. These approaches can be broadly categorized into two groups based on their analytical perspectives. The first category aims to assess the internal security of reusable components themselves, such as *Vulture* [41], *VAScanner* [42], *ATVHUNTER* [43], *LibID* [44], and *LibScout* [45]. The second category focuses on detecting the misuse of reusable components during integration, such as *APP-Miner* [46], *GP-TAid* [47], *ARBITRAR* [48], *Advance* [49], *MutApi* [50], *APIScan* [51] and *APEX* [52]. However, these detection approaches primarily target misuse patterns and internal security issues within their respective software environments. Even when adapted to smart contracts, they remain ineffective at identifying logic-level usage violations of SCRs. In addition, several studies have investigated the security of reusable components in smart contracts, such as *RHT* [2], *EquivGuard* [7], and *ZepScope* [6]. For example, *RHT* introduces eight patterns of library misuse to assist developers in avoiding common pitfalls and mitigating potential financial losses. *ZepScope* presents a systematic analysis of the security posture of the widely used *OpenZeppelin* library in real-world contracts. *EquivGuard* leverages a combination of static taint analysis and symbolic execution to identify six categories of EVM-inequivalence issues in cross-chain contract reuse. Although these methods offer valuable insights into the SCR security, they overlook the analysis of the logic-level usage feature of such components, thereby failing to detect logic-level usage violations of SCRs.

## 7. Conclusion

In this paper, we introduce SCRUTINEER, the first automated and practical system for detecting logic-level usage violations of SCRs, which bridges a critical gap in existing smart contract research. SCRUTINEER achieves a precision of 80.77%, a recall of 82.35%, and an F1-score of 81.55% in

detecting logic-level usage violations of SCRs, representing a significant performance improvement over state-of-the-art vulnerability detectors. Additionally, we have created the first ground-truth dataset for logic-level usage violations of SCRs, including 382 contracts from various DeFi projects. Furthermore, our analysis of 1,181 real-world contracts has identified 13 zero-day vulnerabilities, with 9 of which have been assigned CVE IDs. The identified vulnerabilities have been confirmed by one of the most prominent DeFi platforms with a market capitalization exceeding \$1 billion and a 24-hour trading volume surpassing \$150 million.

## Ethics Considerations

We pay careful attention to potential ethical concerns associated with this work. First, all smart contracts analyzed in our study are sourced from publicly available blockchain platforms or reputable auditing platforms, ensuring transparency and legitimacy. Second, all proof-of-concept attacks are conducted exclusively in a controlled local environment, without interacting with deployed contracts, thereby avoiding any real-world consequences. Finally, we have responsibly disclosed all identified bugs to the corresponding vendors to support the improvement of their security practices.

## LLM Usage Considerations

In this work, we strictly follow the policy on the responsible use of LLMs, as detailed below.

**Originality.** LLMs are used to improve the clarity and fluency of writing. All research content, experimental design, and conclusions are authored, reviewed, and verified entirely by the authors. We therefore take full responsibility for the originality of this paper.

**Transparency and Responsibility.** LLMs are employed in the SCR Knowledge Base Construction and RAG-Driven Inspector modules to assist with information extraction and reasoning. To ensure transparency, reproducibility, and responsible usage, we summarize our practices as follows. First, we adopt the open-source model DeepSeek-V3 without performing any fine-tuning, thereby ensuring the transparency and integrity of LLM usage. Second, the raw data used for SCR knowledge construction are collected entirely from publicly accessible blockchain platforms, ensuring that our analysis is based on open and verifiable sources. Third, all processes and configurations involving LLM usage are explicitly documented in the paper to ensure experimental reproducibility.

## References

- [1] K. Sun, Z. Xu, C. Liu, K. Li, and Y. Liu, “Demystifying the composition and code reuse in solidity smart contracts,” in *ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering, ESEC/FSE*, 2023.

- [2] M. Huang, J. Chen, Z. Jiang, and Z. Zheng, “Revealing hidden threats: An empirical study of library misuse in smart contracts,” in *IEEE/ACM International Conference on Software Engineering, ICSE*, 2024.
- [3] Y. Liu, Y. Xue, D. Wu, Y. Sun, Y. Li, M. Shi, and Y. Liu, “Propertytpt: Llm-driven formal verification of smart contracts through retrieval-augmented property generation,” in *Network and Distributed System Security Symposium, NDSS*, 2025.
- [4] C. Shou, S. Tan, and K. Sen, “Ityfuzz: Snapshot-based fuzzer for smart contract,” in *ACM SIGSOFT International Symposium on Software Testing and Analysis, ISSA*, 2023.
- [5] S. J. Wang, K. Pei, and J. Yang, “Smartinv: Multimodal learning for smart contract invariant inference,” in *IEEE Symposium on Security and Privacy, SP*, 2024.
- [6] H. Liu, D. Wu, Y. Sun, H. Wang, K. Li, Y. Liu, and Y. Chen, “Using my functions should follow my checks: Understanding and detecting insecure openzeppelin code in smart contracts,” in *USENIX Security Symposium*, 2024.
- [7] Z. Wang, J. Chen, T. Zhang, Y. Zhang, W. Zhang, Y. Feng, and Z. Zheng, “Copy-and-paste? identifying EVM-Inequivalent code smells in multi-chain reuse contracts,” in *ACM SIGSOFT International Symposium on Software Testing and Analysis, ISSA*, 2025.
- [8] “ATM Token Incident,” <https://app.blocksec.com/explorer/tx/bsc/0xee10553c26742bec9a4761fd717642d19012bab1704cbced048425070ee21a8a>, 2024.
- [9] Z. Jiang, F. F. Xu, L. Gao, Z. Sun, Q. Liu, J. Dwivedi-Yu, Y. Yang, J. Callan, and G. Neubig, “Active retrieval augmented generation,” in *Conference on Empirical Methods in Natural Language Processing, EMNLP*, 2023.
- [10] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W. Yih, T. Rocktäschel, S. Riedel, and D. Kiela, “Retrieval-augmented generation for knowledge-intensive NLP tasks,” in *Conference on Neural Information Processing Systems, NeurIPS*, 2020.
- [11] “LangChain,” <https://github.com/langchain-ai/langchain>, 2025.
- [12] “Ragflow,” <https://github.com/infiniflow/ragflow>, 2025.
- [13] “LlamaIndex,” [https://github.com/run-llama/llama\\_index](https://github.com/run-llama/llama_index), 2024.
- [14] X. Gao, Y. Xiong, D. Wang, Z. Guan, Z. Shi, H. Wang, and S. Li, “Preference-guided refactored tuning for retrieval augmented code generation,” in *IEEE/ACM International Conference on Automated Software Engineering, ASE*, 2024.
- [15] S. Ouyang, J. M. Zhang, Z. Sun, and A. Meroño-Peñuela, “Knowledge-enhanced program repair for data science code,” in *IEEE/ACM International Conference on Software Engineering, ICSE*, 2025.
- [16] Z. Li, X. Li, W. Li, and X. Wang, “SCALM: detecting bad practices in smart contracts through llms,” in *AAAI Conference on Artificial Intelligence, AAAI*, 2025.
- [17] J. R. Chowdhury, Y. Zhuang, and S. Wang, “Novelty controlled paraphrase generation with retrieval augmented conditional prompt tuning,” in *AAAI Conference on Artificial Intelligence, AAAI*, 2022.
- [18] J. Yu, S. Wu, J. Chen, and W. Zhou, “Llms as collaborator: Demands-guided collaborative retrieval-augmented generation for common-sense knowledge-grounded open-domain dialogue systems,” in *Findings of the Association for Computational Linguistics: EMNLP*, 2024.
- [19] J. Chen, H. Lin, X. Han, and L. Sun, “Benchmarking large language models in retrieval-augmented generation,” in *AAAI Conference on Artificial Intelligence, AAAI*, 2024.
- [20] S. Xu, L. Pang, M. Yu, F. Meng, H. Shen, X. Cheng, and J. Zhou, “Unsupervised information refinement training of large language models for retrieval-augmented generation,” in *Annual Meeting of the Association for Computational Linguistics, ACL*, 2024.

- [21] Z. Wang, S. X. Teo, J. Ouyang, Y. Xu, and W. Shi, "M-RAG: reinforcing large language model performance through retrieval-augmented generation with multiple partitions," in *Annual Meeting of the Association for Computational Linguistics, ACL*, 2024.
- [22] X. Li, J. Jiang, S. Benton, Y. Xiong, and L. Zhang, "A large-scale study on API misuses in the wild," in *IEEE Conference on Software Testing, Verification and Validation, ICST*, 2021.
- [23] J. Chen, M. Huang, Z. Lin, P. Zheng, and Z. Zheng, "To healthier ethereum: A comprehensive and iterative smart contract weakness enumeration," *CoRR*, vol. abs/2308.10227, 2023.
- [24] "SunWeb3Sec/DeFiHackLabs," <https://github.com/SunWeb3Sec/DeFiHackLabs>, 2025.
- [25] "Uniswap," <https://github.com/Uniswap>, 2025.
- [26] J. Feist, G. Grieco, and A. Groce, "Slither: a static analysis framework for smart contracts," in *International Workshop on Emerging Trends in Software Engineering for Blockchain, WETSEB@ICSE*, 2019.
- [27] "BlockSec," <https://blocksec.com/blog>, 2025.
- [28] "Code4rena," <https://code4rena.com>, 2024.
- [29] 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 *IEEE/ACM International Conference on Automated Software Engineering, ASE*, 2021.
- [30] "Mythril," <https://github.com/ConsenSysDiligence/mythril>, 2024.
- [31] "falcon," <https://github.com/MetaTrustLabs/falcon-metatruster>, 2024.
- [32] Q. Song, H. Huang, X. Jia, Y. Xie, and J. Cao, "Silence false alarms: Identifying anti-reentrancy patterns on ethereum to refine smart contract reentrancy detection," in *Network and Distributed System Security Symposium, NDSS*, 2025.
- [33] P. Tsankov, A. M. Dan, D. Drachler-Cohen, A. Gervais, F. Bünzli, and M. T. Vechev, "Securify: Practical security analysis of smart contracts," in *ACM SIGSAC Conference on Computer and Communications Security, CCS*, 2018.
- [34] S. Kalra, S. Goel, M. Dhawan, and S. Sharma, "ZEUS: analyzing safety of smart contracts," in *Network and Distributed System Security Symposium, NDSS*, 2018.
- [35] 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, pp. 116:1–116:27, 2018.
- [36] V. Wüstholtz and M. Christakis, "Harvey: a greybox fuzzer for smart contracts," in *ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering, ESEC/FSE*, 2020.
- [37] B. Jiang, Y. Liu, and W. K. Chan, "Contractfuzzer: fuzzing smart contracts for vulnerability detection," in *ACM/IEEE International Conference on Automated Software Engineering, ASE*, 2018.
- [38] G. Grieco, W. Song, A. Cygan, J. Feist, and A. Groce, "Echidna: effective, usable, and fast fuzzing for smart contracts," in *ACM SIGSOFT International Symposium on Software Testing and Analysis, ISSTA*, 2020.
- [39] X. Lin, Q. Xie, B. Zhao, Y. Tian, S. Zonouz, N. Ruan, J. Li, R. Beyah, and S. Ji, "Promfuzz: Leveraging llm-driven and bug-oriented composite analysis for detecting functional bugs in smart contracts," in *IEEE/ACM International Conference on Automated Software Engineering, ASE*, 2025.
- [40] 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 *IEEE/ACM International Conference on Software Engineering, ICSE*, 2024.
- [41] S. Xu, J. Dong, W. Cai, J. Li, A. Shaghaghi, N. Sun, and S. Ma, "Enhancing security in third-party library reuse - comprehensive detection of 1-day vulnerability through code patch analysis," in *Network and Distributed System Security Symposium, NDSS*, 2025.
- [42] F. Zhang, L. Fan, S. Chen, M. Cai, S. Xu, and L. Zhao, "Does the vulnerability threaten our projects? automated vulnerable API detection for third-party libraries," *IEEE Transactions on Software Engineering*, vol. 50, no. 11, pp. 2906–2920, 2024.
- [43] X. Zhan, L. Fan, S. Chen, F. Wu, T. Liu, X. Luo, and Y. Liu, "ATVHUNTER: reliable version detection of third-party libraries for vulnerability identification in android applications," in *IEEE/ACM International Conference on Software Engineering, ICSE*, 2021.
- [44] J. Zhang, A. R. Beresford, and S. A. Kollmann, "Libid: reliable identification of obfuscated third-party android libraries," in *ACM SIGSOFT International Symposium on Software Testing and Analysis, ISSTA*, 2019.
- [45] M. Backes, S. Bugiel, and E. Derr, "Reliable third-party library detection in android and its security applications," in *ACM SIGSAC Conference on Computer and Communications Security, CCS*, 2016.
- [46] J. Jiang, J. Wu, X. Ling, T. Luo, S. Qu, and Y. Wu, "App-miner: Detecting API misuses via automatically mining API path patterns," in *IEEE Symposium on Security and Privacy, SP*, 2024.
- [47] J. Liu, Y. Yang, K. Chen, and M. Lin, "Generating API parameter security rules with LLM for API misuse detection," in *Network and Distributed System Security Symposium, NDSS*, 2025.
- [48] Z. Li, A. Machiry, B. Chen, M. Naik, K. Wang, and L. Song, "AR-BITRAR: user-guided API misuse detection," in *IEEE Symposium on Security and Privacy, SP*, 2021.
- [49] T. Lv, R. Li, Y. Yang, K. Chen, X. Liao, X. Wang, P. Hu, and L. Xing, "Rtfm! automatic assumption discovery and verification derivation from library document for API misuse detection," in *ACM SIGSAC Conference on Computer and Communications Security, CCS*, 2020.
- [50] M. Wen, Y. Liu, R. Wu, X. Xie, S. Cheung, and Z. Su, "Exposing library API misuses via mutation analysis," in *International Conference on Software Engineering, ICSE*, 2019.
- [51] I. Yun, C. Min, X. Si, Y. Jang, T. Kim, and M. Naik, "Apsan: Sanitizing API usages through semantic cross-checking," in *USENIX Security Symposium*, 2016.
- [52] Y. J. Kang, B. Ray, and S. Jana, "Apex: automated inference of error specifications for C apis," in *IEEE/ACM International Conference on Automated Software Engineering, ASE*, 2016.

## Appendix A.

### Distribution-based Weight Inference for Composite Signature Retrieval

Each SCR composite signature is represented as a four-dimensional tuple  $Sign = (feature_1, feature_2, feature_3, feature_4)$ , where  $feature_1$  and  $feature_2$  correspond to the *contract name* and *function name*, and  $feature_3$  and  $feature_4$  correspond to the *number of parameters* and *return value*, respectively. The first two features are continuous in lexical semantics, while the latter two are discrete. Accordingly, we assign equal group priors,  $W_C = 0.5$  and  $W_D = 0.5$ , where  $C = \{1, 2\}$  and  $D = \{3, 4\}$  denote the continuous and discrete feature groups.

To capture the inherent stability of each feature, we infer its weight from the empirical distribution of string lengths collected from the knowledge base. For each feature  $i$ , the observed length samples are denoted as  $\ell_i = \{\ell_i^{(1)}, \ell_i^{(2)}, \dots, \ell_i^{(N)}\}$ . Each sample is normalized by its total length, as defined in Equation 2.

$$\hat{\ell}_i^{(n)} = \frac{\ell_i^{(n)}}{\sum_{m=1}^N \ell_i^{(m)}}, \quad n = 1, \dots, N. \quad (2)$$

The sample variance of the normalized lengths is then computed using Equation 3, which measures the dispersion of the length distribution:

$$s_i^2 = \frac{1}{N-1} \sum_{n=1}^N (\hat{\ell}_i^{(n)} - \bar{\ell}_i)^2. \quad (3)$$

To prevent numerical instability when a feature exhibits extremely low variance, we introduce a small stability constant  $\epsilon$ , as shown in Equation 4.

$$\epsilon = 0.01 \cdot \text{Median}(\{s_j^2\}_{j \in \text{group}(i)}). \quad (4)$$

The intra-group weight  $\tilde{w}_i$  of each feature is then determined through inverse-variance normalization, as defined in Equation 5. This operation assigns larger weights to features with more stable length distributions.

$$\tilde{w}_i = \frac{(s_i^2 + \epsilon)^{-1}}{\sum_{j \in \text{group}(i)} (s_j^2 + \epsilon)^{-1}}. \quad (5)$$

Finally, the global weight  $w_i$  is obtained by combining the group priors with the intra-group coefficients, as described in Equation 6.

$$w_i^s = \begin{cases} W_C \cdot \tilde{w}_i, & i \in C, \\ W_D \cdot \tilde{w}_i, & i \in D. \end{cases} \quad (6)$$

In our empirical evaluation, the resulting weights are  $\omega^s = (0.079, 0.421, 0.313, 0.187)$ .

## Appendix B.

### Sensitivity Analysis of Weights and Thresholds in the Similarity-Based Checker

To evaluate the soundness of both weights and the thresholds in the similarity-based checker, we conducted a sensitivity analysis.

For each group of weights  $\omega = (\omega_n, \omega_c)$  satisfying  $\omega_n + \omega_c = 1$ , we varied  $\omega_n$  within the interval  $[\omega_n - 0.05, \omega_n + 0.05]$  with a step size of 0.01, while  $\omega_c$  is adjusted accordingly to maintain the normalization constraint. This setting provides a sufficiently fine granularity to capture parameter sensitivity without overfitting to small fluctuations, which is appropriate for discrete-task evaluation where performance metrics change non-continuously. For each threshold  $\tau$ , we varied its value within  $\pm 5\%$  of its original setting with a step of 1%. This range ensures a reasonable local perturbation that reflects potential deviations caused by empirical estimation errors or environmental variance, while keeping the perturbation narrow enough to preserve the detection semantics of the model. Table 8 provides the detailed configuration of each parameter.

TABLE 8: The configuration of weights and thresholds in sensitivity analysis.

Parameter		Value	Interval	Step
$\omega^t$	$\omega_n^t$	0.42	$[\omega_n^t - 0.05, \omega_n^t + 0.05]$	0.01
	$\omega_c^t = 1 - \omega_n^t$	0.58	$[\omega_c^t + 0.05, \omega_c^t - 0.05]$	0.01
$\omega^o$	$\omega_n^o$	0.27	$[\omega_n^o - 0.05, \omega_n^o + 0.05]$	0.01
	$\omega_c^o = 1 - \omega_n^o$	0.73	$[\omega_c^o + 0.05, \omega_c^o - 0.05]$	0.01
$\tau_t$		0.68	$[\tau_t - 5\% \tau_t, \tau_t + 5\% \tau_t]$	$1\% \tau_t$
$\tau_o$		0.92	$[\tau_o - 5\% \tau_o, \tau_o + 5\% \tau_o]$	$1\% \tau_o$
$\tau_l$		0.90	$[\tau_l - 5\% \tau_l, \tau_l + 5\% \tau_l]$	$1\% \tau_l$

In the evaluation, we measure the Precision, Recall, and F1-score, and compute the average variation of these metrics within the corresponding perturbation intervals. As shown in Table 9, the average change across the interval remained minor, indicating that the value of parameters is stable and appropriate. Therefore, the sensitivity analysis verifies that the selected weights and thresholds provide robust performance under the perturbations.

TABLE 9: Sensitivity analysis for parameters.

Parameters	$\Delta Precision$	$\Delta Recall$	$\Delta F1 - Score$
$\omega^t$	1.19%	0	0.68%
$\omega^o$	0.70%	1.05%	0.90%
$\tau_t$	2.20%	0	1.35%
$\tau_o$	3.35%	3.16%	2.96%
$\tau_l$	0.83%	0.31%	0.57%