# Securing Smart Contract Languages with a Unified Agentic Framework for Vulnerability Repair in Solidity and Move

**Rabimba Karanjai**[1†,3], **Lei Xu**[2], **Weidong Shi**[1]

[1]University Of Houston, [2]Kent State University

[†]rkaranjai@uh.edu

The rapid growth of the blockchain ecosystem and the increasing value locked in smart contracts necessitate robust security measures. While languages like Solidity and Move aim to improve smart contract security, vulnerabilities persist. This paper presents Smartify, a novel multi-agent framework leveraging Large Language Models (LLMs) to automatically detect and repair vulnerabilities in Solidity and Move smart contracts. Unlike traditional methods that rely solely on vast pre-training datasets, Smartify employs a team of specialized agents working on different specially fine-tuned LLMs to analyze code based on the underlying programming concepts and language-specific security principles. We evaluated Smartify on a dataset for Solidity and a curated dataset for Move, demonstrating its effectiveness in fixing a wide range of vulnerabilities. Our experimental results show that Smartify (Gemma2+Codegemma) achieves state-of-the-art performance, surpassing existing LLMs and even enhancing the capabilities of general-purpose models, such as Llama 3.1. Notably, Smartify can incorporate language-specific knowledge, such as the nuances of Move, without requiring massive language-specific pretraining datasets. This work offers a detailed analysis of the performance of various LLMs on smart contract repair, highlighting the strengths of our multi-agent approach and providing a blueprint for developing more secure and reliable decentralized applications in the growing blockchain landscape. We also provide a detailed description to extend the proposed technology to other similar use cases.

**Date:** November 21, 2025

## 1 Introduction

Smart contracts, self-executing agreements with terms directly written into code, have emerged as a cornerstone of blockchain technology Nath et al. (2014); Ray (2023). Their ability to automate transactions and eliminate intermediaries has led to widespread adoption in various sectors, including finance, supply chain management, and healthcare Zheng et al. (2018); Karanjai et al. (2021); Kaleem et al. (2021). However, the increasing complexity of smart contracts has given rise to a growing concern: security vulnerabilities Vacca et al. (2021). These vulnerabilities, often stemming from coding errors or design flaws, can be exploited by malicious actors, leading to significant financial losses and damage to the reputation of blockchain projects.

The financial implications of smart contract vulnerabilities are substantial. Reports indicate that cumulative losses from attacks against Ethereum smart contracts alone have exceeded USD 3.1 billion by 2023 Li et al. (2023). In the DeFi space, an estimated $9.04 billion has been stolen due to vulnerabilities Wronka (2023). Notable incidents like the DAO hack of 2016, resulting in a $55 million loss Popper (2016), and the Poly Network hack in 2021, where over $600 million was stolen pol (2025), underscore the critical need for robust security measures.

Traditional security auditing methods, although essential, often face limitations in terms of accuracy and scalability. This has spurred the exploration of automated techniques for vulnerability detection Jiang et al. (2018a); Wang et al. (2020)and repair, with Large Language Models (LLMs) emerging as a promising solution Joshi et al. (2023). LLMs trained on vast datasets of code can learn to understand and generate code that adheres to specific programming paradigms and best practices. However, most of the tools available for smart contract safety are language-specific, especially `Solidity` Dannen and Dannen (2017). For other

languages, existing tools often require scanning of compiled bytecode Song et al. (2024).

Apart from `Solidity`, `Move` Blackshear et al. (2019) has gained significant traction due to its strong focus on security. Its cutting-edge features, including a custom data type for secure operations, robust access controls via Move modules, and unique memory safety features Blackshear et al. (2022), have been particularly noteworthy. Moreover, the Move Prover, a native security framework, provides an additional layer of protection Dill et al. (2022). Notably, several prominent blockchain platforms, such as Starcoin sta (2024), Aptos Dev, and Sui Blackshear et al. (2024), have already adopted `Move`.

However, despite its promising architecture, the real-world security performance of `Move` modules remains largely untested. Unlike `Solidity`-based smart contracts, which have been extensively studied through empirical research and surveys, there is a scarcity of research focused specifically on Move modules. Although some methodologies have been proposed for identifying defects in `Move` modules or conducting formal verification Keilty et al. (2022); Park et al. (2024), and empirical analysis Song et al. (2024), a significant knowledge gap persists. Specifically, large-scale investigations into the frequency of defects in real-world `Move` modules and identifying and repairing potential vulnerabilities are lacking, highlighting the need for further research in this area.

This paper introduces Smartify, a framework that moves beyond naive LLM application by proposing a novel paradigm for structured, automated repair. We introduce Smartify not merely as a "multi-agent framework" but as a **structured, collaborative system that mimics the workflow of an expert human security audit team**. This contrasts with existing LLM-powered tools like `ContractTinker`, which utilize a linear Chain-of-Thought (CoT) reasoning process to break down the repair task. Smartify's novelty lies in its process-centric design, where specialized agents (Auditor, Architect, Code Generator, Refiner, Validator) assume distinct roles within a delegative and iterative workflow. The Architect agent, for instance, does not simply reason about the next step; it formulates a comprehensive repair strategy that is then executed by other agents. Crucially, the Refiner-Validator duo establishes a feedback loop for iterative quality assurance, a feature essential for generating trustworthy patches.

In this paper, our aim is to answer the following research questions related to software engineering using AI agents and the landscape of complex smart contract reasoning.

- **RQ1:** *Do the present state-of-the-art LLMs can explain a Smart Contract code correctly?*

- **RQ2:** *Can they detect and explain bad coding practices or specific mistakes leading to bugs or vulnerabilities in a smart contract code?*

- **RQ3:** *Can we encode programming language-specific knowledge to train the LLMs to understand unsafe and buggy codes in detail enough to repair them?*

- **RQ4:** *Can LLMs repair the bugs and fix the vulnerability?*

- **RQ5:** *Does the proposed post-training framework be generalizable to a larger set of pre-trained LLMs?*

The main contributions of our work are the following:

- We introduce a novel **role-based multi-agent architecture** for smart contract repair that models the collaborative workflow of a security audit team, enhancing structured reasoning over monolithic LLM approaches.

- We propose a rigorous methodology combining **specialized fine-tuned models for deep vulnerability analysis**, **Retrieval-Augmented Generation (RAG) for language-specific context**, and an **iterative self-refinement loop** for ensuring patch quality.

- We conduct the first **multi-faceted empirical evaluation** of an LLM-based repair tool for both Solidity and Move, using metrics that assess not only correctness but also **exploit mitigation effectiveness**, **semantic preservation**, and **code quality**.

- We present a comprehensive **ablation study** that systematically dissects our framework to quantify the distinct contribution of each architectural component to the overall repair performance.

# 2 Related Work

This section provides a critical review of the existing landscape in smart contract security, positioning Smartify relative to traditional analysis techniques and emerging LLM-based repair methodologies.

## 2.1 Traditional Smart Contract Security Analysis

Automated security analysis for smart contracts has traditionally been dominated by static, dynamic, and formal verification techniques. Each approach offers distinct advantages but also possesses inherent limitations, particularly when confronted with complex, logic-based vulnerabilities.

## 2.2 Smart Contract Security Auditing

Various tools and techniques have been developed for detecting vulnerabilities in smart contracts:

**Static Analysis Tools:** Tools like Mythril Mueller and Slither Feist et al. (2019) analyze contract source code to identify potential vulnerabilities. They perform symbolic execution and taint analysis to detect patterns associated with common vulnerabilities.

**Dynamic Analysis Tools:** Tools like Manticore Mossberg et al. (2019) and Echidna Grieco et al. (2020) execute contracts with various inputs to uncover runtime errors. They use fuzzing and symbolic execution techniques to explore different execution paths and identify potential issues.

**Formal Verification:** This approach uses mathematical techniques to rigorously prove the correctness of a contract's code against a formal specification. Tools like KEVM Hildenbrandt et al. (2018) and CertiK's DeepSEA have been developed for formal verification of smart contracts Zhong et al. (2020).

While these tools are valuable, they often have limitations in accuracy, scalability, and the ability to handle the complexities of real-world smart contracts.

## 2.3 LLM-Powered Smart Contract Repair: A Taxonomy

The application of LLMs to automated program repair is a rapidly advancing field, with several distinct methodologies emerging.

### 2.3.1 Monolithic LLM Approaches

Early efforts in this domain involved using large, general-purpose LLMs like GPT-3 with intricate, zero-shot or few-shot prompts to generate patches. While demonstrating feasibility, these approaches often lack the domain-specific knowledge required for the high-stakes environment of smart contracts. They are prone to generating syntactically correct but semantically flawed or insecure code, as they lack a deep, ingrained understanding of blockchain-specific security paradigms Blackshear et al. (2022); Bobadilla et al. (2025).

### 2.3.2 Chain-of-Thought (CoT) and Static Analysis Integration

To address the limitations of monolithic models, more structured approaches have been developed. A prominent example is **ContractTinker**, a tool designed for real-world smart contract repair Wang et al. (2024). `ContractTinker` employs a Chain-of-Thought (CoT) mechanism to guide an LLM through a sequence of reasoning steps: vulnerability localization, analysis, and patch generation. To ground the LLM's reasoning and mitigate hallucination, it integrates static analysis techniques, including dependency analysis and program slicing, to provide relevant context from audit reports and the source code itself.

While `ContractTinker` represents a significant advancement by imposing a logical structure on the repair process, its CoT approach remains a fundamentally linear and sequential reasoning pipeline. In contrast, Smartify's architecture is **delegative and iterative**. The Architect agent formulates a comprehensive, high-level plan, which is then delegated to specialized agents for execution. Furthermore, the explicit Refiner-Validator loop introduces a crucial mechanism for feedback and iterative quality improvement, a feature not explicitly detailed in the `ContractTinker` workflow. This architectural distinction moves beyond a simple chain of thought to a collaborative problem-solving process.

### 2.3.3 Multi-Agent Systems for Code Tasks

The concept of using multiple LLM-based agents to collaborate on complex tasks has gained traction in the broader software engineering domain, with applications in areas like code translation and generation Huang

et al. (2024); Jiang et al. (2018b). These systems leverage the principle of specialization, assigning different roles or sub-tasks to individual agents to achieve a more robust and accurate outcome than a single agent could.

While Smartify aligns with this general trend, its novelty lies in its **domain-specific agent roles tailored explicitly for the vulnerability repair workflow**. Instead of a generic "translator" or "coder" agent, Smartify's agents embody the distinct functions of a human security audit team: the Auditor for analysis, the Architect for strategic planning, the Code Generator for implementation, and the Refiner/Validator for quality assurance. This specialization allows for a more nuanced and effective approach to the highly specific and critical task of securing smart contracts.

Our proposed framework, **Smartify**, addresses these challenges by combining the strengths of specialized LLMs within a multi-agent architecture. It leverages language-specific fine-tuning, safety classifiers, and Retrieval-Augmented Generation (RAG) to enhance the accuracy and security of generated code repairs.

In the following sections, we detail the architecture of Smartify, describe the experimental setup, present the evaluation results, and discuss the implications of our findings for the future of smart contract security.

# 3 Data Collection and Analysis Methodology

This research employs a multi-faceted approach to investigate the security of smart contracts, focusing on both `Solidity` and `Move` programming languages. The methodology encompasses collecting and analyzing two distinct datasets: `Solidity`-based, `Move`-based source code. Each dataset serves a specific purpose in addressing the research questions and contributing to a comprehensive understanding of smart contract vulnerabilities.

## 3.1 Importance of Dataset Categorization

For several reasons, categorizing the datasets based on programming language (`Solidity` and `Move`) and code representation is crucial. It allows for a focused analysis of language-specific vulnerabilities and coding practices. As a more mature language, `Solidity` exhibits a different vulnerability landscape than the newer `Move` language. Examining them separately enables the identification of unique challenges and security considerations associated with each language.

## 3.2 Dataset Descriptions

### 3.2.1 `Solidity`-based Dataset

This dataset comprises a collection of vulnerable `Solidity` smart contracts sourced from the "Not-So-Smart Contracts" repository curated by Trail of Bits git. This repository is renowned for its comprehensive set of contracts that intentionally exhibit a variety of common vulnerabilities. These vulnerabilities were chosen for inclusion because of their prevalence in real-world decentralized applications and their representation of typical errors during smart contract development. The dataset contains 60 vulnerable contracts, encompassing 8 distinct vulnerability categories. Table 1 shows these categories' distribution.

**Table 1** Distribution of Vulnerabilities in the `Solidity` Dataset.

| Vulnerability Type | Number of Contracts | Percentage (%) |
|---|---|---|
| Reentrancy | 15 | 25.0 |
| Integer Overflow/Underflow | 10 | 16.7 |
| Denial of Service (DoS) | 8 | 13.3 |
| Access Control Issues | 12 | 20.0 |
| Uninitialized Storage Pointers | 5 | 8.3 |
| Tx.origin Misuse | 3 | 5.0 |
| Timestamp Dependency | 4 | 6.7 |
| Gas Limit and Out-of-Gas Vulnerabilities | 3 | 5.0 |
| **Total** | **60** | **100.0** |

### 3.2.2 `Move`-Based Dataset (Source Code)

This dataset encompasses the source code of 92 real-world `Move` projects, comprising 652 individual modules. These projects were part of Aptos Dev, Sui Blackshear et al. (2024), and Starcoin sta (2024). These projects span various application domains, as depicted in Table 2. The total number of `Move` projects is 92, and the total number of `Move` modules within these projects is 652.

**Table 2** Distribution of `Move` Projects by Application Domain.

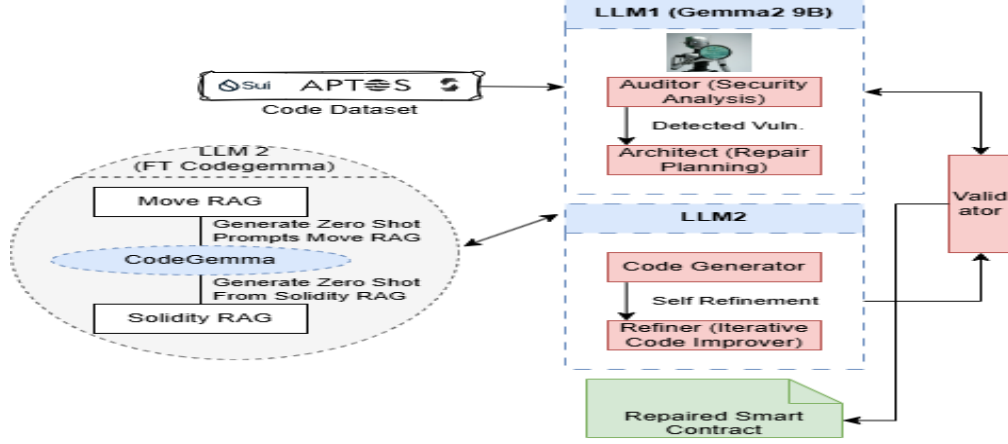| Application Domain | Number of Projects | Percentage (%) |
|---|---|---|
| Decentralized Finance | 41 | 44.6 |
| Token | 22 | 23.9 |
| Bridge | 18 | 19.6 |
| Library | 3 | 3.3 |
| Infrastructure | 3 | 3.3 |
| Other | 5 | 5.4 |
| **Total** | **92** | **100.0** |



**Figure 1** Agentic architecture of Smartify.

For both the `Move`-based datasets, we utilize Song et al.'s Song et al. (2024) work to compare the vulnerability detection part. While the prior work is directed towards detection, the same dataset helps us compare Smartify's performance on both detection and repair.

### 3.3 Evaluation of Smartify

Smartify is designed with two core functionalities: detecting and repairing unsafe coding patterns in smart contracts. We utilize the previously described datasets to evaluate these capabilities rigorously, encompassing both `Solidity` and `Move` code. The evaluation process focuses on the complete output of Smartify rather than individual components, reflecting its nature as an integrated solution for smart contract security. Performance is measured using the Pass@1 score.

### 3.4 Agent–Based Code Repair Process for Smart Contracts

Our approach leverages a multi-agent system inspired by established software development methodologies but tailored explicitly for the automated repair of `Solidity` and `Move` smart contracts. This system employs five specialized agents: an *Auditor*, an *Architect*, a *Code Generator*, a *Refiner*, and a *Validator*. The process incorporates a self-refinement loop and a final validation step, ensuring high accuracy and security. Each agent plays a distinct role in a structured workflow, detailed below.

- **Auditor:** This agent is the cornerstone of security analysis. It is fine-tuned on a comprehensive corpus of `Solidity` and `Move` code documentation, encompassing syntax, semantics, and best practices. Furthermore, it is safety-aligned using a classifier adapted from Google's Responsible AI toolkit. This classifier has been meticulously modified to enforce language-specific rules and safe coding practices, effectively preventing the generation of unsafe or unsupported code constructs.

  This alignment is of paramount importance. For `Move`, it ensures that the generated code strictly adheres to the conventions of the target blockchain (e.g., Sui or Aptos). It prevents the accidental introduction of elements from one Move variant into another or the inclusion of unsupported Rust paradigms. This is crucial because Move, derived from Rust, has unique features and limitations. For `Solidity`, it enforces established security best practices and prevents the generation of code patterns known to be vulnerable.

  The Auditor's primary responsibility is to meticulously scan the input smart contract code (either `Solidity` or `Move`) to identify potential vulnerabilities and unsafe patterns. Its secondary, yet vital, role is to serve as the final *Validator* of the repaired code.

- **Architect:** This agent receives the output from the *Auditor*, which includes a detailed report of identified

vulnerabilities and unsafe code segments. The *Architect*'s role is to devise a high-level strategic plan for addressing these issues. This plan does not involve generating code directly. Instead, it outlines the necessary modifications, refactoring, and improvements to rectify the identified problems. This plan is a comprehensive blueprint for the *Code Generator*, guiding the code repair process.

- **Code Generator:** This agent is a general-purpose code LLM. Its strength lies in its ability to leverage Retrieval-Augmented Generation (RAG) from two distinct data stores, one dedicated to `Solidity` and the other to `Move`. These data stores contain a collection of best practices and relevant documentation for the respective programming languages.

  Using the *Architect*'s plan as a guide, the Code Generator selects and adapts relevant examples from the appropriate RAG datastore. This dynamic, context-aware retrieval of few-shot examples significantly enhances the *Code Generator*'s ability to produce accurate and secure code repairs. It ensures the generated code adheres to language-specific conventions and incorporates established best practices.

- **Refiner:** This agent's role is to enhance the code quality produced by the *Code Generator*. It achieves this through iterative self-refinement, essentially acting as its critic. The *Refiner* uses the same underlying LLM as the *Code Generator* but with a different prompt that focuses on improving the code quality based on best practices and potential improvements that it might detect from a higher level.

- **Validator:** This agent acts as a final checkpoint. After the refinement stage, it re-employs the *Auditor* agent to re-evaluate the code. The *Validator*'s objective is to ensure that all the previously identified vulnerabilities have been adequately addressed and that no new vulnerabilities have been introduced during the repair and refinement process.

## 4   Smartify System Architecture and Workflow

Smartify operates through a five-agent system designed for automated intelligent contract vulnerability detection and repair. The system functions as shown in Figure 1.

The Smartify system operates in a five-phase process to automatically repair smart contract code. Firstly, in the Input & Initial Audit phase, the smart contract code, written in either `Solidity` or `Move`, is fed into the system. The *Auditor*, an LLM based on Gemma2 9B, analyzes the code to detect potential vulnerabilities and produces a report detailing its findings. Secondly, during Repair Planning, the *Architect* receives this vulnerability report and formulates a high-level repair plan that outlines the necessary code modifications to address the identified issues. Thirdly, in Code Generation & Refinement, an LLM called CodeGemma, which has been fine-tuned for code generation and is equipped with Retrieval-Augmented Generation (RAG) capabilities, takes the lead. It utilizes separate `Move` RAG and `Solidity` RAG components to provide language-specific context. The Code Generator, part of CodeGemma, uses the repair plan to generate the modified code, selecting the appropriate RAG based on the input language and able to perform `Solidity` to `Move` translation when necessary.

Subsequently, a Self-refinement process is initiated, and the *Refiner* component iteratively improves the generated code's quality, readability, and efficiency. Fourthly, in the Validation phase, the *Validator* (the same agent as the *Auditor*) performs a final security audit on the refined code to ensure all identified vulnerabilities have been resolved. Finally, the system outputs the repaired smart contract code.

The process may iterate to step 3 or 4 if the *Validator* identifies any issues. Each step is vital in ensuring the smart contract code's accurate and secure repair. The workflow is designed to be efficient and effective, leveraging each agent's strengths to achieve the desired outcome.

### 4.1   Agent Prompting Strategy

The agents within Smartify are driven by carefully crafted prompts that guide their actions and ensure consistent performance. We employ a standardized prompt template adapted from established practices in LLM-based agent systems. The template is structured as follows:

> **Prompt Template**
>
> **Role:** You are a `[role]` specializing in `[Solidity/Move]` smart contracts.
> **Task:** `[task]`
> **Instruction:** Based on the provided Context, please follow these steps: `[numbered steps]`
> **Context:** ...

This template is broken down into the following components.

Each agent in our framework is defined by four key components: the Role, which designates the agent's specific function (such as Auditor, Architect, or Code Generator); the Task, which outlines the agent's particular objectives; the Instruction, which provides detailed step-by-step guidance using chain-of-thought reasoning; and the Context, which encompasses all necessary information including input code, audit reports, architectural plans, RAG datastore examples, and inter-agent conversation history.

Table 3 shows how this template is adapted for each agent.

**Table 3** Agent Prompts for Smart Contract Repair.

| Role | Task | Instruction | Context |
|---|---|---|---|
| **Auditor** | Identify vulnerabilities and unsafe patterns in Solidity/Move code. | Analyze the code for security vulnerabilities and generate a detailed report. | Input smart contract code (Solidity/Move). |
| **Architect** | Create a high-level plan to address vulnerabilities identified by the Auditor. | Review the Auditor's report and develop a plan outlining necessary modifications. | Auditor's report. |
| **Code Generator** | Generate Repaired Solidity/Move code based on the Architect's plan and RAG examples. | Consult the Architect's plan, retrieve examples from the RAG datastore, and generate repaired code. | Architect's plan, Solidity/Move code examples from RAG. |
| **Refiner** | Iteratively refine the generated code to improve quality and efficiency. | Review the generated code, identify areas for improvement, and refine accordingly. | Generated code, previous iteration code (if any). |
| **Validator** | Perform a final security check on the repaired code. | Analyze the repaired code for vulnerabilities, verify issue resolution, and ensure no new vulnerabilities. | Repaired smart contract code. |

## 4.2 Hardware and Model Fine-tuning

The development and deployment of Smartify leveraged a heterogeneous computing environment, utilizing high-performance GPUs for computationally intensive tasks and a more resource-efficient setup for inference.

### 4.2.1 Fine-tuning Setup

- **Hardware:** Fine-tuning leveraged a cluster of **four NVIDIA A100 GPUs** for computationally demanding pattern learning in `Solidity` and `Move` code.

- **Model:** Based on the **Gemma 9B model**, selected for strong code-related task performance and fine-tuning adaptability, particularly in instruction following. Fine-tuned on a dataset of `Solidity` and `Move` code, vulnerability examples, best practices, and documentation, augmented with outputs from earlier pipeline stages to enhance safety issue detection.

- **Training Recipe:** Supervised learning paradigm. Trained to predict correct outputs (e.g., vulnerability reports, safe code patterns) from inputs (e.g., `Solidity`/`Move` code, vulnerability descriptions).

  - **Data Pre-processing:** Tokenization, normalization, and input-output pair creation ensured data consistency and quality.

  - **Hyperparameter Optimization:** Learning rate (1e-5), batch size (8, due to memory constraints), and training epochs (5, as validation loss plateaued) optimized via grid search and manual tuning.

  - **Regularization:** Dropout and weight decay used to prevent overfitting and improve generalization.

– **Evaluation Metrics:** Accuracy, precision, recall, and F1-score on a held-out validation set monitored model performance.

### 4.2.2 Inference Setup

- **Hardware:** Inference was performed on a single **NVIDIA RTX 4090 GPU**, balancing performance and cost-effectiveness for real-time code repair.

- **Models:**

  – **Code Generator and Refiner:** These agents utilize a fine-tuned **CodeGemma** model, initially pre-trained on a limited `Move` corpus and further instruction-tuned to follow Architect-generated "recipe" patterns. Fine-tuning on Architect outputs ensured it understood these instructions, and pre-training on a limited `Move` corpus ensured basic syntax understanding.

  – **Comparison Model:** A stock **Llama 3.1** model was used in some experiments for comparative analysis, helping assess the gains from fine-tuning and instruction tuning.

### 4.2.3 Key Considerations

- A balance between performance requirements, resource availability, and cost considerations drove the choice of hardware and models.

- The fine-tuning process for the *Auditor* was particularly resource-intensive due to the complexity of the task and the size of the model.

- The use of a smaller, more efficient GPU for inference makes the system more accessible for practical deployment.

- The comparison with a stock Llama 3 model provides valuable insights into the effectiveness of our fine-tuning and instruction-tuning strategies.

This heterogeneous setup, combining high-performance GPUs for training and a more efficient GPU for inference, allows Smartify to effectively address the computational demands of both model development and deployment. The detailed description of the fine-tuning process provides transparency and allows for replication of our results.

## 5 Experimental Results and Discussion

We run our experiments as defined in Section 3.3. We report the results as well as the empirical performance of our models. Through that, we will try to answer our Research Questions individually in this section.

Along with Smartify, we have run the benchmark for the following models.

**Table 4** Comparison of Code and Non-Code Models.

| Model Name | Parameters | Quantization | Code Model |
|---|---|---|---|
| Granite-Code | 8B | FP16 | Yes |
| CodeGemma | 7B | FP16 | Yes |
| DeepSeek-Coder-v2 | N/A | N/A | Yes |
| StarCoder2 | 15B | FP16 | Yes |
| CodeGeex4 | 13B | N/A | Yes |
| CodeStral | 7B | FP16 | Yes |
| DeepSeek-Coder | 33B | N/A | Yes |
| CodeLlama Roziere et al. (2023) | 13B | N/A | Yes |
| CodeQwen | 7B | Q8_0 | Yes |
| Qwen2.5-coder | 2.5B | N/A | Yes |
| Gemma2 | N/A | N/A | Yes |
| Gemma2:27b | 27B | FP16 | Yes |
| Llama3.2 | 3.2B | FP16 | No |
| OpenCoder | 8B | FP16 | Yes |
| Llama3.3 | 3.3B | FP16 | No |

**Figure 2** Code Repair: Solidity.

The models were chosen according to the top 8 models at Hugging Face Big Code Leaderboard hug at the time of this work, and also adding general-purpose models, which are supposed to be better at reasoning.

## 5.1 Solidity

This section presents the evaluation results of various code generation models on repairing vulnerabilities in Solidity smart contracts, specifically focusing on the "Not So Smart Contracts" dataset from the Trail of Bits GitHub repository. This dataset is a collection of intentionally vulnerable Solidity contracts designed to test the ability of automated tools to detect and repair common security flaws. It contains diverse vulnerabilities, including re-entrancy, integer overflow/underflow, access control issues, and timestamp dependence, among others. The dataset has been publicly available for a significant period, raising the possibility that some or all of its contents might be present in the pre-training data of the evaluated models. We analyze the performance of these models based on two key metrics: the number of vulnerabilities fixed and the average inference time, as summarized in Table 5 and Figure 2. We also introduce our framework, Smartify, and demonstrate its effectiveness in enhancing model performance.

**Table 5** Performance of Code Generation Models on Vulnerability Repair.

| Model Name | Vuln. Fixed | Avg. Time (s) |
| --- | --- | --- |
| CodeGeex-4 | 11 | 95.50 |
| **CodeGemma** | **16** | 96.50 |
| CodeLlama | 3 | 243.93 |
| CodeQwen | 5 | 141.05 |
| CodeStral | 13 | 295.23 |
| DeepSeekCoder-33b | 1 | 411.75 |
| **DeepSeek-V2** | 9 | **19.42** |
| Gemma2-9b | 13 | 108.30 |
| Gemma2-27b | 14 | 304.27 |
| Granite-Code | 14 | 90.37 |
| Llama3.2 | 1 | 37.09 |
| Llama3.3-70b | 13 | 741.10 |
| OpenCoder* | 1* | 94* |
| Qwen-2.5-Code | 13 | 79.72 |
| StarCoder2 | 0* | 89.10 |
| **Smartify (Gemma2+CodeGemma)** | **16** | 112.30 |
| Smartify (Gemma2+Llama3.1) | 14 | 267.80 |

The results reveal significant performance disparities among the evaluated models. Among the pre-trained models for Solidity **CodeGemma** surprisingly emerges as a top performer, successfully fixing 16 vulnerabilities with a relatively low average inference time of 96.5 seconds. This suggests that **CodeGemma** possesses a strong ability to understand and rectify code vulnerabilities while maintaining reasonable efficiency. However since most of these Solidity smart contracts were part of open Githubs repositories, there can be a strong possibility fo these already being part of the pertaining data. Our proposed framework, **Smartify (Gemma2+CodeGemma)**, achieves comparable performance, also fixing 16 vulnerabilities, albeit with a slightly higher average inference time of 112.3 seconds. This increased time is likely due to its iterative multi-agent process, which enables Smartify to leverage the complementary strengths of **Gemma2** and **CodeGemma**, resulting in robust and reliable fixes.

While Smartify does not immediately show any benefits over **CodeGemma** here, we can notice that the same Smartify framework when applied to **Llama3.1** without any fine-tuning (unlike the Smartify with **CodeGemma**) still gives considerable performance boost over vanilla.

Conversely, models like **CodeLlama**, **CodeQWen**, **DeepSeekCoder 33b**, and **Llama3.2** show limited effectiveness,
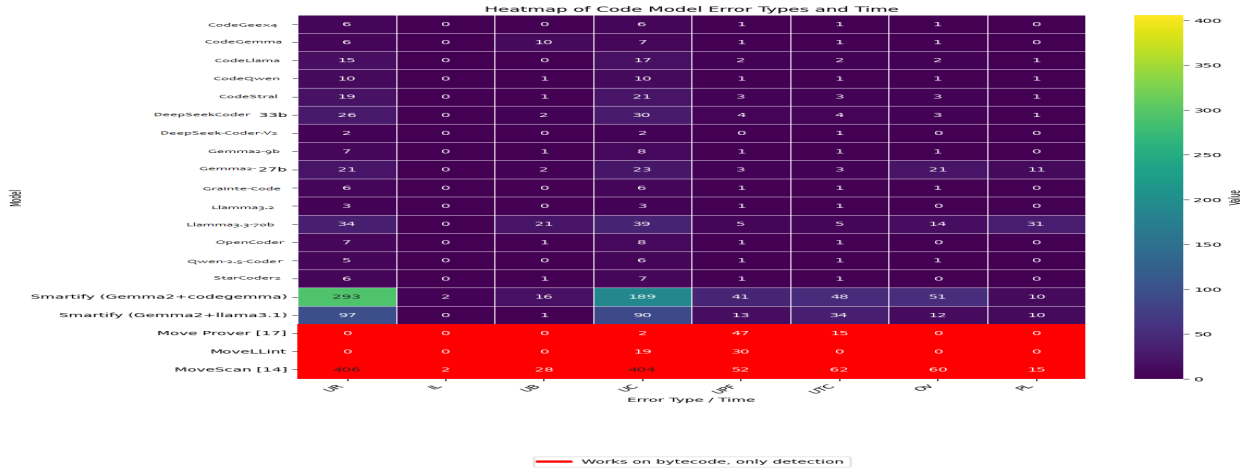
**Figure 3** `Move` Code Repair.

fixing only a small number of vulnerabilities. The poor performance of these models could be attributed to several factors, such as insufficient exposure to `Solidity` code during pre-training or fine-tuning or architectures ill-suited for vulnerability repair, which requires a deep understanding of code syntax and security principles. The exceptionally poor performance of models like **starcoder2** (marked with an asterisk *), along with incomplete data for **opencoder**, suggests potential issues with their training data or a fundamental mismatch between their capabilities and the task's demands. These models might have been trained on an older version of `Solidity` or different smart contract security practices than those in the Not-So-Smart-Contracts dataset. Moreover, they might prioritize other aspects of code generation, such as code completion, over security-specific tasks like vulnerability repair.

**Table 6** `Move` Vulnerability Repair (Time in seconds).

| Model | UR | IL | UB | UC | UPF | UTC | Ov | PL | Time |
|---|---|---|---|---|---|---|---|---|---|
| CodeGeex4 | 6 | 0 | 0 | 6 | 1 | 1 | 1 | 0 | 96 |
| CodeGemma | 6 | 0 | 10 | 7 | 1 | 1 | 0 | 0 | 97 |
| CodeLlama | 15 | 0 | 1 | 17 | 2 | 2 | 2 | 1 | 244 |
| CodeQwen | 10 | 0 | 1 | 10 | 1 | 1 | 1 | 1 | 141 |
| CodeStral | 19 | 0 | 1 | 21 | 3 | 3 | 2 | 1 | 295 |
| DeepSeekCoder 33b | 26 | 0 | 2 | 30 | 4 | 4 | 3 | 1 | 412 |
| DeepSeekV2 | 2 | 0 | 0 | 2 | 0 | 1 | 0 | 0 | 19 |
| Gemma2 9b | 7 | 0 | 1 | 8 | 1 | 1 | 0 | 10 | 108 |
| Gemma2 27b | 21 | 0 | 2 | 23 | 3 | 3 | 21 | 11 | 304 |
| Granite-Code | 6 | 0 | 0 | 6 | 1 | 1 | 1 | 0 | 90 |
| Llama3.2 | 3 | 0 | 0 | 3 | 1 | 1 | 0 | 0 | 37 |
| Llama3.3 70b | 34 | 0 | 21 | 39 | 5 | 5 | 14 | 31 | 741 |
| OpenCoder | 7 | 0 | 1 | 8 | 1 | 1 | 0 | 0 | 94* |
| Qwen 2.5 code | 5 | 0 | 0 | 6 | 1 | 1 | 1 | 0 | 80 |
| StarCoder2 | 6 | 0 | 1 | 7 | 1 | 1 | 0 | 0 | 89 |
| **Smartify (Gemma2+CodeGemma)** | **293** | **2** | **16** | **189** | **41** | **48** | **51** | **10** | **112** |
| **Smartify (Gemma2+Llama3.1)** | 97 | 0 | 1 | 90 | 13 | 34 | 12 | **10** | 268 |
| Move Prover Dill et al. (2022) | - | 2 | - | - | - | - | 47 | 15 | - |
| MoveLint | - | - | - | - | 19 | 30 | 0 | 0 | - |
| MoveScan Song et al. (2024) | 406 | 2 | 28 | 404 | 52 | 62 | 60 | 15 | - |

**Abbreviations: UR**: Unchecked Return; **IL**: Infinite Loop; **UB**: Unnecessary Boolean; **UC**: Unused Constant; **UPF**: Unused Private Function;

The public availability of the "Not-So-Smart-Contracts" dataset raises the question of data contamination. Many evaluated models, especially those trained on large, public code corpora, might have encountered this

dataset during pre-training, potentially inflating their performance. However, since **CodeGemma** and **Smartify (Gemma2+CodeGemma)** were specifically fine-tuned for this task, the issue of data contamination is likely less significant.

## 5.2 `Move` **Code Repair**

This section analyzes the efficacy of various models in repairing vulnerabilities within `Move` smart contracts, as detailed in Table 6. The evaluation encompasses eight distinct vulnerability categories: *Unchecked Return (UR), Infinite Loop (IL), Unnecessary Boolean (UB), Unused Constant (UC), Unused Private Function (UPF), Unnecessary Type Conversion (UTC), Overflow (Ov), and Precision Loss (PL)* following the works of Song et al Song et al. (2024). The metrics presented in the table represent the number of successfully repaired instances for each vulnerability type, with higher values indicating superior performance. The inference time, measured in seconds, is also provided for each model.

The results demonstrate a significant variance in performance across the evaluated models. Notably, the larger language models, such as **Deepseekcoder 33b** and **Llama3.3 70b**, exhibit a relatively higher number of successful repairs across multiple categories, albeit with a corresponding increase in inference time. Conversely, smaller models like **DeepseekV2** and **Llama3.2** demonstrate limited repair capabilities. The specialized tools for `Move` code, namely **Move Prover**, **MoveLint**, and **MoveScan**, were employed as a benchmark for comparison. It is crucial to note that these tools are designed for vulnerability **detection** rather than repair. **MoveScan**, in particular, identified a substantial number of instances across all categories, highlighting its effectiveness as a static analysis tool. **Move Prover** demonstrated proficiency in detecting Overflow and Precision Loss vulnerabilities, while **MoveLint** focused on Unused Private Functions and Unnecessary Type Conversions.

The Smartify models, which take advantage of a combination of **Gemma2** with either **CodeGemma** or **Llama3.1**, present an interesting case. Smartify(**Gemma2+CodeGemma**) and Smartify(**Gemma2+Llama3.1**) outperform several individual models in multiple categories. This is likely because the specialized models are fine-tuned on the `Move`-specific dataset. For instance, Smartify (**Gemma2+CodeGemma**) achieves the highest number of repairs for the Unchecked Return, Infinite Loop, Unused Boolean, Unused Constant, Unused Private Function, Unnecessary Type Conversion, and Overflow categories, showcasing a substantial improvement over individual models in these areas. However, it is worth mentioning that they also have limitations compared to individual models for certain categories like Precision Loss.

> ### RQ1 & RQ2 - Code Understanding and Vuln. Detection
>
> **Yes.** Our empirical analysis with Smartify, especially with using a fine-tuned Code-Gemma and also using vanilla pre-trained **Llama3.1**, has shown us the effectiveness of the framework's ability to understand code and capture bad practices leading to vulnerability. Especially for a low-resource code like `Move`, without significant fine-tuning (in the case of **Llama3.1**), Smartify outperform the larger and computationally intensive models such as Llama 3.3 70b.

Notably, **Smartify (Gemma2+CodeGemma)**, combining fine-tuned **Gemma2** with **CodeGemma**, achieves performance on par with the best individual model, **CodeGemma**, which is expected due to one of the models being fine-tuned. This highlights the advantages of strategically combining specialized models, answering our next research question.

> ### RQ3 & RQ4 - Code Repair
>
> Both for `Solidity` and `Move`, we were able to compare the efficacy of our framework with prior works. We can see Smartify outperforms all of the existing code models, even very specialized code models trained on `Move` (OpenCoder Huang et al. (2024)) in generating repair codes for detected vulnerabilities.

Furthermore, Smartify's efficacy extends even when integrating a non-finetuned model like **Llama 3.1**. Smartify significantly outperforms **Llama 3.2** by fixing 14 vulnerabilities compared to Llama 3.2's single fix, making its performance comparable with the much more extensive and computationally intensive **Llama 3.3 70b**. This demonstrates that Smartify's architecture can enhance even general-purpose language models for code repair, balancing speed, and accuracy.

**RQ5 - Generalization**

Our implementation of Smartify with both fine-tuned **Code-Gemma** and **Llama3.1** as the second agent allowed us to run our experiments on both sets of LLMs. The results show that Smartify can significantly boost performance even on non-finetuned models compared to a single model.

Comparative analysis reveals trade-offs between model scale and performance in automated code repair. Larger models, such as **DeepSeekCoder 33b** and **Llama 3.3 70b**, exhibit broader repair capabilities but incur higher computational costs and inference times. Conversely, the **Gemma2 27b** model demonstrates notable proficiency in addressing Overflow vulnerabilities, albeit with limitations in handling Unnecessary Boolean and Unused Constant compared to **Llama 3.3 70b**. While **Llama 3.3 70b** outperforms Smartify in overall repair capability, its significantly slower inference speed poses a challenge for practical deployment. Therefore, for real-world, on-device applications, **Smartify (Gemma2+CodeGemma)** presents a compelling solution with its balance of substantial accuracy and rapid inference.

**Insight:** Specialized code models like **StarCoder** Lozhkov et al. (2024), **OpenCoder** Huang et al. (2024) and **DeepSeekCoder** Guo et al. (2024) doesn't necessarily work well even if it's a coding specific task. While code models like **CodeGemma** Team et al. (2024) and **CodeLlama** Roziere et al. (2023) are much better at understanding instructions and working on code. This helped Smartify for its understanding and fine-tuning for code repairability.

Specialized static analysis tools for `Move`, including **Move Prover**, **MoveLint**, and **MoveScan**, work as baselines of detecting `Move` vulnerabilities with which we compare our Smartify and other LLMs. These findings underscore the need for targeted model improvements. The Smartify framework directly addresses these deficiencies, offering enhanced vulnerability repair effectiveness.

This research also opens up future research directions on using this framework for context-aware test case generation.

# 6 Ablation Study

To rigorously validate our architectural design and isolate the contribution of each key component, a comprehensive ablation study was conducted. This study systematically deconstructs the Smartify framework to quantify the impact of its core mechanisms—agent specialization, iterative refinement, and retrieval-augmented generation—on overall repair performance Huang et al. (2025); Jiang et al. (2018b).

## 6.1 Ablation Configurations

Four distinct configurations of the system were evaluated against the Solidity and Move datasets:

- **Full Smartify:** The complete five-agent system, including the RAG module for the Code Generator and the Refiner-Validator iterative feedback loop. This represents our proposed approach.

- **Smartify (No Refinement):** A version of the framework where the Refiner and Validator agents are disabled. The output is taken directly from the first pass of the Code Generator. This configuration measures the impact of the iterative self-improvement loop on patch quality.

- **Smartify (No RAG):** In this setup, the Code Generator operates without the contextual, few-shot examples provided by the RAG system. It relies solely on its fine-tuned knowledge and the Architect's plan. This configuration is designed to measure the importance of providing language-specific, in-context examples, especially for the low-resource Move language.

- **Single-Agent Baseline (CoT-style):** This configuration replaces the multi-agent system with a single, powerful agent (Gemma2 9B). The agent is given a complex, chained prompt that instructs it to sequentially perform the tasks of auditing, planning, and generating the repair. This mimics the linear reasoning workflow of CoT-based approaches like `ContractTinker` and directly tests the value of our role-based specialization and delegation architecture.

## 6.2 Analysis

The performance of each configuration was measured using the Pass@1 and Exploit Mitigation Rate metrics. The results, presented in Table 7, provide clear empirical evidence supporting our architectural choices.

**Table 7** Ablation Study of Smartify Components on Solidity and Move Datasets

| Configuration | Solidity | | Move | |
|---|---|---|---|---|
| | Pass@1 (%) | Exploit Mit. (%) | Pass@1 (%) | Exploit Mit. (%) |
| **Full Smartify** | **26.7** | **25.0** | **48.9** | **45.7** |
| Smartify (No Refinement) | 25.0 | 20.0 | 44.6 | 39.1 |
| Smartify (No RAG) | 21.7 | 18.3 | 28.3 | 23.9 |
| Single-Agent Baseline | 18.3 | 15.0 | 21.7 | 17.4 |

The results demonstrate a clear and consistent degradation in performance as key components are removed, confirming the positive contribution of each element of the Smartify architecture.

- **Impact of Iterative Refinement:** Removing the refinement loop (**No Refinement**) causes a noticeable drop in both metrics, particularly the Exploit Mitigation Rate (from 25.0% to 20.0% for Solidity). This indicates that while the initial code generation is often correct, the refinement process is crucial for hardening the patch against exploits and catching subtle regressions that the first pass might miss.

- **Impact of RAG:** The removal of the RAG module (**No RAG**) has a significant negative impact across the board, but its effect is most pronounced for the Move language. The Pass@1 score for Move plummets from 48.9% to 28.3%, a relative decrease of over 42%. This strongly supports our hypothesis that providing in-context, language-specific examples is critical for achieving high performance in low-resource languages where the model's pre-trained knowledge is limited.

- **Impact of Multi-Agent Architecture:** The **Single-Agent Baseline**, designed to mimic a CoT-style approach, performs the worst of all configurations. Its performance is substantially lower than the Full Smartify framework, particularly on the more complex task of exploit mitigation. This finding provides strong evidence that our role-based, delegative architecture is superior to a monolithic, linear reasoning process. By specializing agents for distinct tasks (analysis vs. planning vs. generation), the system achieves a more robust and effective decomposition of the complex program repair problem.

The ablation study empirically validates the design of Smartify. The multi-agent architecture provides a superior structure for complex problem-solving, the refinement loop is essential for ensuring patch quality and security, and the RAG mechanism is a critical component for adapting the framework to new or low-resource programming languages.

## 7 Conclusion

This work addresses the pressing need for enhanced security in the burgeoning blockchain ecosystem. We investigate the application of Large Language Models (LLMs) to smart contract vulnerability detection and repair, focusing on `Solidity` and `Move`. We introduce **Smartify**, a novel multi-agent framework that significantly improves LLM performance in this critical domain. The contributions of this work are: (1) **Smartify**, a novel multi-agent framework that enhances LLM-based smart contract vulnerability detection and repair; (2) a method for encoding language-specific knowledge, valuable for low-resource languages like Move; (3) a scalable, adaptable approach applicable to other programming languages and LLMs; (4) a demonstration of Smartify's efficacy on generalized pre-trained LLMs; and (5) a detailed analysis of the challenges inherent in automated code repair.

**Smartify** represents a significant advancement in automating smart contract security, a crucial concern in the expanding blockchain landscape. Future work will refine the framework, expand its language coverage, particularly within the blockchain domain, and integrate it into real-world blockchain development workflows.

This research lays the foundation for AI-powered tools that can bolster the security and reliability of decentralized applications, fostering a more robust and trustworthy blockchain ecosystem.

## Acknowledgments

## References

not-so-smart-contracts: Examples of Solidity security issues. https://github.com/crytic/not-so-smart-contracts. [Accessed 09-01-2025].

Big Code Models Leaderboard - a Hugging Face Space by bigcode — huggingface.co. https://huggingface.co/spaces/bigcode/bigcode-models-leaderboard. [Accessed 10-01-2025].

Starcoin whitepaper, 2024. URL https://starcoin.org/_astro/whitepaper.CMlZ6t_x.pdf.

Decoding poly network $34 billion hack, 2025. URL https://www.quillaudits.com/blog/hack-analysis/poly-network-hack.

Sam Blackshear, Evan Cheng, David L Dill, Victor Gao, Ben Maurer, Todd Nowacki, Alistair Pott, Shaz Qadeer, Dario Russi Rain, Stephane Sezer, et al. Move: A language with programmable resources. *Libra Assoc*, page 1, 2019.

Sam Blackshear, John Mitchell, Todd Nowacki, and Shaz Qadeer. The move borrow checker. *arXiv preprint arXiv:2205.05181*, 2022.

Sam Blackshear, Andrey Chursin, George Danezis, Anastasios Kichidis, Lefteris Kokoris-Kogias, Xun Li, Mark Logan, Ashok Menon, Todd Nowacki, Alberto Sonnino, et al. Sui lutris: A blockchain combining broadcast and consensus. In *Proceedings of the 2024 on ACM SIGSAC Conference on Computer and Communications Security*, pages 2606–2620, 2024.

Sofia Bobadilla, Monica Jin, and Martin Monperrus. Do automated fixes truly mitigate smart contract exploits?, 2025. URL https://arxiv.org/abs/2501.04600. arXiv e-prints.

Chris Dannen and Chris Dannen. Solidity programming. *Introducing Ethereum and Solidity: Foundations of Cryptocurrency and Blockchain Programming for Beginners*, pages 69–88, 2017.

Aptos Dev. The aptos blockchain: Safe, scalable, and upgradeable web3 infrastructure.

David Dill, Wolfgang Grieskamp, Junkil Park, Shaz Qadeer, Meng Xu, and Emma Zhong. Fast and reliable formal verification of smart contracts with the move prover. In *International Conference on Tools and Algorithms for the Construction and Analysis of Systems*, pages 183–200. Springer, 2022.

Josselin Feist, Gustavo Grieco, and Alex Groce. Slither: a static analysis framework for smart contracts. In *2019 IEEE/ACM 2nd International Workshop on Emerging Trends in Software Engineering for Blockchain (WETSEB)*, pages 8–15. IEEE, 2019.

Gustavo Grieco, Santiago Palladino, Martin Ortigoza, Federico Bond, and Valentin Wustholz. Echidna: Effective, usable, and fast fuzzing for smart contracts. In *Proceedings of the 29th ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA '20)*, pages 481–493. ACM, 2020. doi: 10.1145/3395363.3397384.

Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Yu Wu, YK Li, et al. Deepseek-coder: When the large language model meets programming–the rise of code intelligence. *arXiv preprint arXiv:2401.14196*, 2024.

Everett Hildenbrandt, Manasvi Saxena, Nishant Rodrigues, Xiaoran Zhu, Philip Daian, Dwight Guth, Brandon Moore, Daejun Park, Yi Zhang, Andrei Stefanescu, et al. Kevm: A complete formal semantics of the ethereum virtual machine. In *2018 IEEE 31st Computer Security Foundations Symposium (CSF)*, pages 204–217. IEEE, 2018.

Siming Huang, Qingyue Zhang, Yilun Jin, Xiaozhu Meng, Yiyang Wang, Haoyu Wang, Ding Li, and Shiqing Ma. Opencoder: The open cookbook for top-tier code large language models, 2024. URL https://arxiv.org/abs/2411.04905. arXiv preprint.

Siming Huang, Tianhao Cheng, J. K. Liu, Jiaran Hao, Liuyihan Song, Yang Xu, J. Yang, Jiaheng Liu, Chenchen Zhang, Linzheng Chai, Ruifeng Yuan, Zhaoxiang Zhang, Jie Fu, Qian Liu, Ge Zhang, Zili Wang, Yuan Qi,

Yinghui Xu, and Wei Chu. Opencoder: The open cookbook for top-tier code large language models, 2025. URL https://arxiv.org/abs/2411.04905.

Bo Jiang, Ye Liu, and W. K. Chan. Contractfuzzer: fuzzing smart contracts for vulnerability detection. In *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering*, ASE '18, page 259–269, New York, NY, USA, 2018a. Association for Computing Machinery. ISBN 9781450359375. doi: 10.1145/3238147.3238177. URL https://doi.org/10.1145/3238147.3238177.

Bo Jiang, Ye Liu, and W. K. Chan. Contractfuzzer: fuzzing smart contracts for vulnerability detection. In *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering*, ASE '18, page 259–269, New York, NY, USA, 2018b. Association for Computing Machinery. ISBN 9781450359375. doi: 10.1145/3238147.3238177. URL https://doi.org/10.1145/3238147.3238177.

Harshit Joshi, José Cambronero Sanchez, Sumit Gulwani, Vu Le, Gust Verbruggen, and Ivan Radiček. Repair is nearly generation: Multilingual program repair with llms. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 5131–5140, 2023.

Mudabbir Kaleem, Keshav Kasichainula, Rabimba Karanjai, Lei Xu, Zhimin Gao, Lin Chen, and Weidong Shi. An event driven framework for smart contract execution. In *Proceedings of the 15th ACM International Conference on Distributed and Event-based Systems*, pages 78–89, 2021.

Rabimba Karanjai, Lei Xu, Zhimin Gao, Lin Chen, Mudabbir Kaleem, and Weidong Shi. On conditional cryptocurrency with privacy. In *2021 IEEE International Conference on Blockchain and Cryptocurrency (ICBC)*, pages 1–3. IEEE, 2021.

Eric Keilty, Keerthi Nelaturu, Bowen Wu, and Andreas Veneris. A model-checking framework for the verification of move smart contracts. In *2022 IEEE 13th International Conference on Software Engineering and Service Science (ICSESS)*, pages 1–7. IEEE, 2022.

Jinggang Li, Gehao Lu, Yulian Gao, and Feng Gao. A smart contract vulnerability detection method based on multimodal feature fusion and deep learning. *Mathematics*, 11(23):4823, 2023.

Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, Tianyang Liu, Max Tian, Denis Kocetkov, Arthur Zucker, Younes Belkada, Zijian Wang, Qian Liu, Dmitry Abulkhanov, Indraneil Paul, Zhuang Li, Wen-Ding Li, Megan Risdal, Jia Li, Jian Zhu, Terry Yue Zhuo, Evgenii Zheltonozhskii, Nii Osae Osae Dade, Wenhao Yu, Lucas Krauß, Naman Jain, Yixuan Su, Xuanli He, Manan Dey, Edoardo Abati, Yekun Chai, Niklas Muennighoff, Xiangru Tang, Muhtasham Oblokulov, Christopher Akiki, Marc Marone, Chenghao Mou, Mayank Mishra, Alex Gu, Binyuan Hui, Tri Dao, Armel Zebaze, Olivier Dehaene, Nicolas Patry, Canwen Xu, Julian McAuley, Han Hu, Torsten Scholak, Sebastien Paquet, Jennifer Robinson, Carolyn Jane Anderson, Nicolas Chapados, Mostofa Patwary, Nima Tajbakhsh, Yacine Jernite, Carlos Muñoz Ferrandis, Lingming Zhang, Sean Hughes, Thomas Wolf, Arjun Guha, Leandro von Werra, and Harm de Vries. Starcoder 2 and the stack v2: The next generation, 2024.

Mark Mossberg, Felipe Manzano, Eric Hennenfent, Alex Groce, Gustavo Grieco, Josselin Feist, Trent Brunson, and Artem Dinaburg. Manticore: A user-friendly symbolic execution framework for binaries and smart contracts. In *2019 34th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, pages 1186–1189. IEEE, 2019.

Bernhard Mueller. File 1 of.

Keshab Nath, Sourish Dhar, and Subhash Basishtha. Web 1.0 to web 3.0-evolution of the web and its various challenges. In *2014 International Conference on Reliability Optimization and Information Technology (ICROIT)*, pages 86–89. IEEE, 2014.

Junkil Park, Teng Zhang, Wolfgang Grieskamp, Meng Xu, Gerardo Di Giacomo, Kundu Chen, Yi Lu, and Robert Chen. Securing aptos framework with formal verification. In *5th International Workshop on Formal Methods for Blockchains (FMBC 2024)*. Schloss Dagstuhl–Leibniz-Zentrum für Informatik, 2024.

Nathaniel Popper. A hacking of more than $50 million dashes hopes in the world of virtual currency. *The New York Times*, 17, 2016.

Partha Pratim Ray. Web3: A comprehensive review on background, technologies, applications, zero-trust architectures, challenges and future directions. *Internet of Things and Cyber-Physical Systems*, 3:213–248, 2023.

Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu,

Romain Sauvestre, Tal Remez, et al. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*, 2023.

Shuwei Song, Jiachi Chen, Ting Chen, Xiapu Luo, Teng Li, Wenwu Yang, Leqing Wang, Weijie Zhang, Feng Luo, Zheyuan He, et al. Empirical study of move smart contract security: Introducing movescan for enhanced analysis. In *Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis*, pages 1682–1694, 2024.

CodeGemma Team, Heri Zhao, Jeffrey Hui, Joshua Howland, Nam Nguyen, Siqi Zuo, Andrea Hu, Christopher A Choquette-Choo, Jingyue Shen, Joe Kelley, et al. Codegemma: Open code models based on gemma. *arXiv preprint arXiv:2406.11409*, 2024.

Anna Vacca, Andrea Di Sorbo, Corrado A Visaggio, and Gerardo Canfora. A systematic literature review of blockchain and smart contract development: Techniques, tools, and open challenges. *Journal of Systems and Software*, 174: 110891, 2021.

Che Wang, Jiashuo Zhang, Jianbo Gao, Libin Xia, Zhi Guan, and Zhong Chen. Contracttinker: Llm-empowered vulnerability repair for real-world smart contracts. In *Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering*, pages 2350–2353, 2024.

Wei Wang, Jingjing Song, Guangquan Xu, Yidong Li, Hao Wang, and Chunhua Su. Contractward: Automated vulnerability detection models for ethereum smart contracts. *IEEE Transactions on Network Science and Engineering*, 8(2):1133–1144, 2020.

Christoph Wronka. Financial crime in the decentralized finance ecosystem: new challenges for compliance. *Journal of Financial Crime*, 30(1):97–113, 2023.

Zibin Zheng, Shaoan Xie, Hong-Ning Dai, Xiangping Chen, and Huaimin Wang. Blockchain challenges and opportunities: A survey. *International journal of web and grid services*, 14(4):352–375, 2018.

Jingyi Emma Zhong, Kevin Cheang, Shaz Qadeer, Wolfgang Grieskamp, Sam Blackshear, Junkil Park, Yoni Zohar, Clark Barrett, and David L Dill. The move prover. In *Computer Aided Verification: 32nd International Conference, CAV 2020, Los Angeles, CA, USA, July 21–24, 2020, Proceedings, Part I 32*, pages 137–150. Springer, 2020.