

Hunting Vulnerabilities with AI

Liyi Zhou

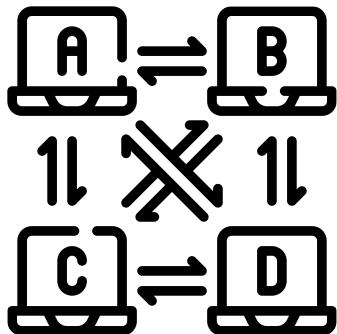
The University of Sydney / D23E / UC Berkeley RDI

About Me

- I was a TA for the DeFi MOOC 😈
- Collaborated with Dawn, Arthur, Kaihua on DeFi security.
 - MEV / BEV
 - Blockchain/Miner/Maximal Extracted/Extractable Value
 - Sandwich attacks
 - Arbitrages
 - Liquidations
 - SoK (Systemization of Knowledge): DeFi Attacks



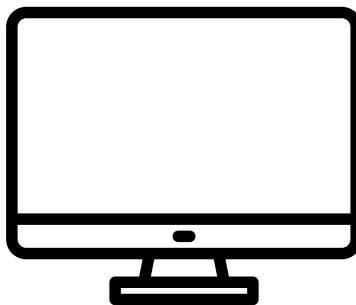
DeFi ≠ just Finance 😊



- (P2P) network ← Decentralized
- State Transition
 - Example 1 – Alice transfers 5 coins to Bob
 - Example 2 – Alice deploys a voting program to decide what to eat for tonight, Bob votes for pizza
- We need some consensus algorithm
 - Eve only owns 5 coins
 - Eve tells Node A she wants to send 5 coins to Alice
 - Eve tells Node B she wants to send 5 coins to Bob

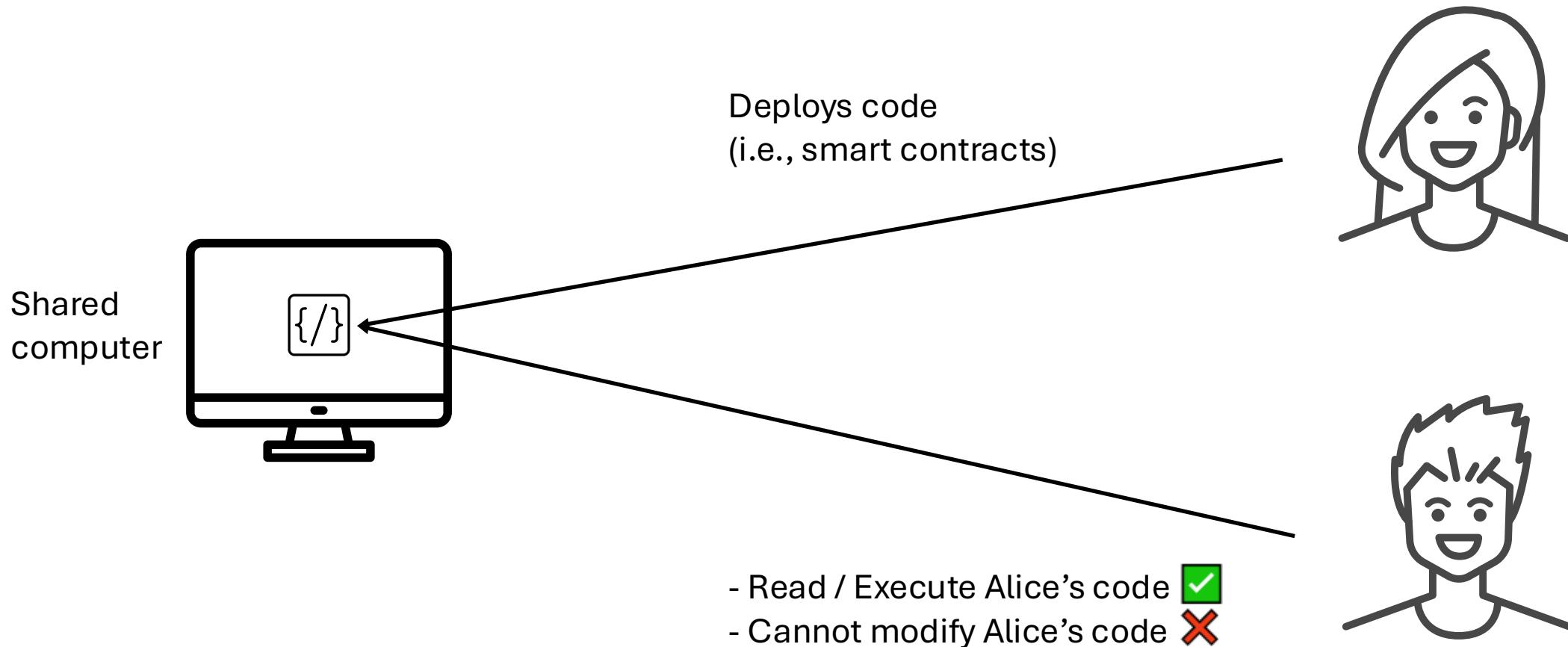
DeFi ≠ just Finance 😊

Shared
computer



- (P2P) network ← Decentralized
- State Transition
 - Example 1 – Alice transfers 5 coins to Bob
 - Example 2 – Alice deploys a voting program to decide what to eat for tonight, Bob votes for pizza
- We need some consensus algorithm
 - Eve only owns 5 coins
 - Eve tells Node A she wants to send 5 coins to Alice
 - Eve tells Node B she wants to send 5 coins to Bob

Smart Contracts



Example

```
solidity

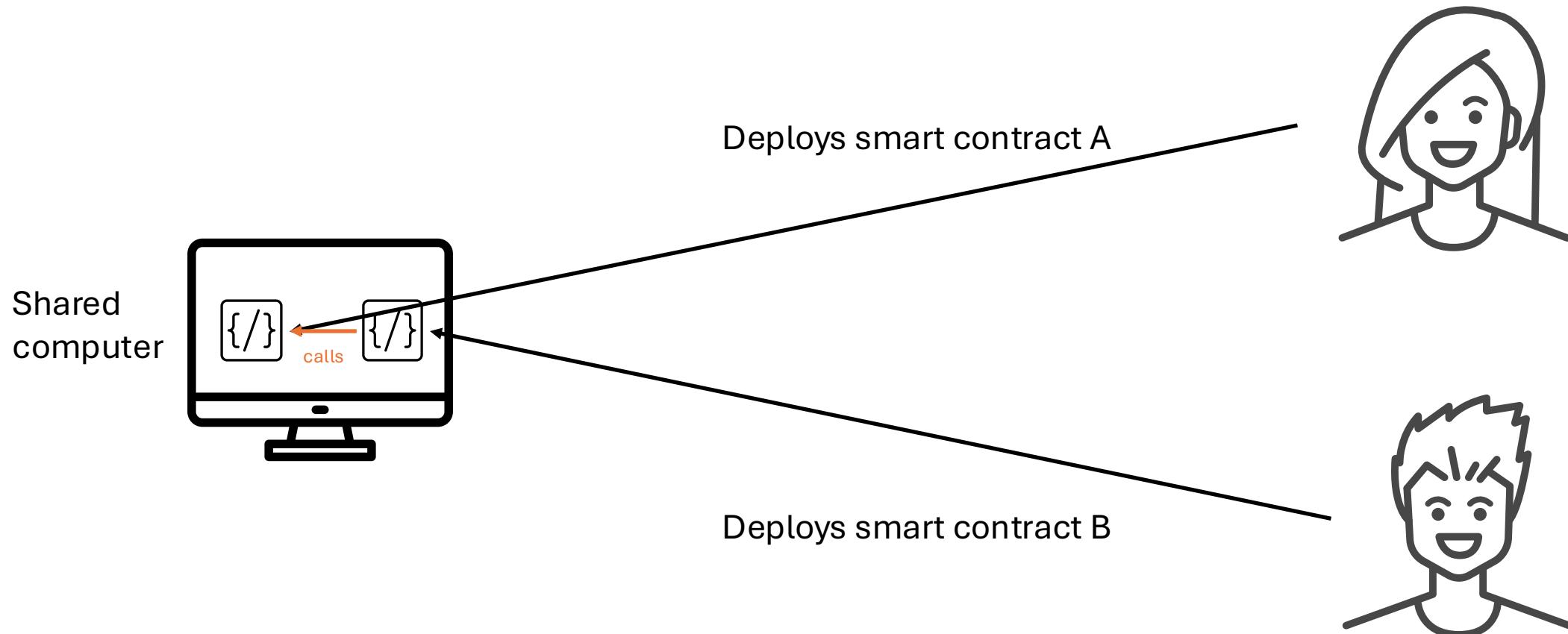
// SPDX-License-Identifier: MIT
pragma solidity ^0.8.20;

contract SimpleToken {
    mapping(address => uint) public balance;

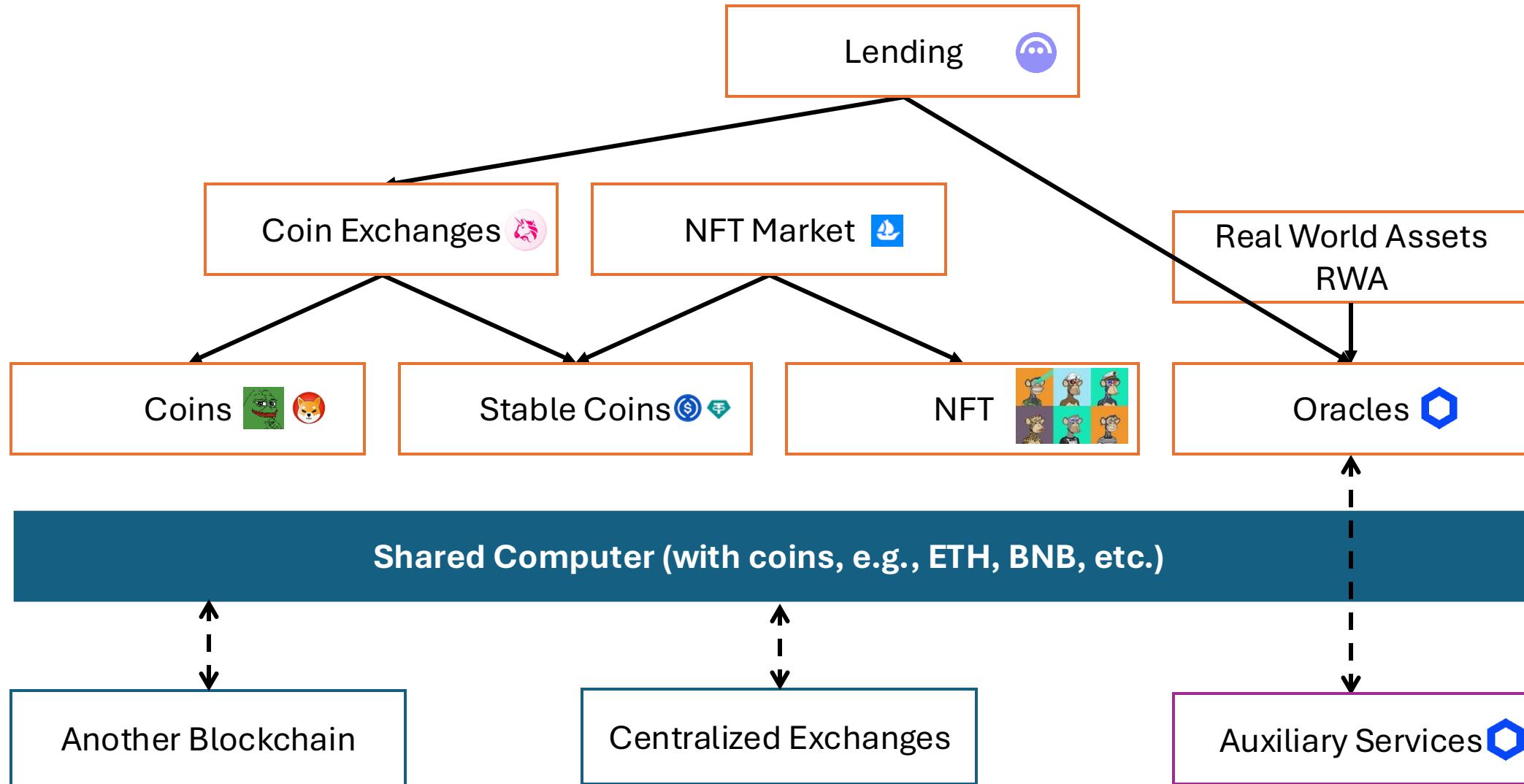
    constructor() {
        // give the deployer (Alice) 100 tokens
        balance[msg.sender] = 100;
    }

    function transfer(address to, uint amount) public {
        require(balance[msg.sender] >= amount, "Not enough");
        balance[msg.sender] -= amount;
        balance[to] += amount;
    }
}
```

Smart Contracts (Composability)



DeFi Lego (Composability)



DeFi Vulnerabilities

- Smart contracts are
 - Not smart 
 - Just programs
 - Implementation bugs
 - Design bugs

DeFi Vulnerabilities

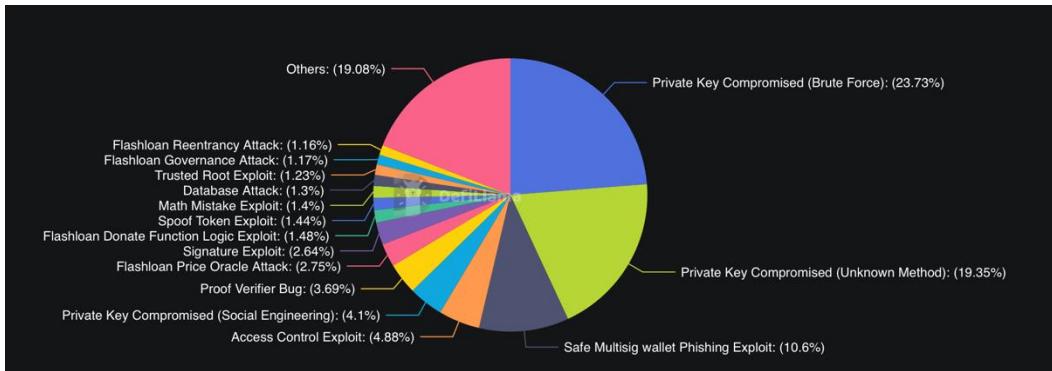
- More likely to cause monetary loss. 
 - Lots of financial applications
 - “Open finance”
- Easy to exit. 
 - Privacy solutions (e.g. Tornado Cash)
 - Centralized exchanges, with liquidity

DeFi Vulnerabilities

- <https://defillama.com>
- Total value locked (estimation):
 - \$163,000,000,000 USD
- Total value hacked (estimation):
 - **\$15,430,000,000 USD** ?

DeFi Vulnerabilities

DeFiLlama:



DeFiHackLabs:

The screenshot shows the DeFi Hack Incidents Explorer interface. At the top, it displays "TOTAL INCIDENTS: 572" and "TOTAL LOSS (USD): \$16,121,684,528". Below this are filters for "YEAR" (set to "All Years"), "ATTACK TYPE" (checkboxes for "All Types" and specific categories like Logic Flaw, Price Manipulation, etc.), "SORT BY" (Date (Latest)), and "DISPLAY CURRENCY" (USD (\$)). A search bar labeled "SEARCH MODULE v1.1" is also present. The main area is a table listing incidents with columns for DATE, PROJECT, EXPLOIT TYPE, LOSS AMOUNT, ROOT CAUSE, and POC. Each row includes a "VIEW POC" button. The table entries are as follows:

DATE	PROJECT	EXPLOIT TYPE	LOSS AMOUNT	ROOT CAUSE	POC
May 11, 2025	MB	Logic Flaw	\$2.16M	-	VIEW POC
Apr 26, 2025	Lif	Price Manipulation	\$15.1K	-	VIEW POC
Apr 26, 2025	Imp	Access Control	\$300.00K	-	VIEW POC
Apr 18, 2025	BTH	Reentrancy Attack	\$19.03K	-	VIEW POC
Apr 16, 2025	YVI	Incorrect Validation	\$15.26K	-	VIEW POC
Mar 30, 2025	Lev	Flash Loan Attack	\$353.80K	-	VIEW POC

<https://arxiv.org/abs/2208.13035>

TABLE III: DeFi incidents taxonomy. We label the incident types that each academic paper and auditing report address. We also group the incidents that occur in the wild. Despite that this table focuses on Ethereum and BSC, we anticipate the taxonomy remains generic and thus applicable to all DeFi enabled blockchains. ● - Incident type addressed; ■ - Incident type checked (likely with tools); □ - Incident cause checked (likely with tools); ○ - Incident type checked (manually). Note that we can only be sure that an incident type has been addressed if an auditing report: (i) explicitly warns of the risk of a potential incident, or (ii) explicitly states that the code passed the check of an incident type. We visualize the gaps using a heat map, where a darker colour indicates a greater frequency of occurrences.

SoK: Decentralized Finance (DeFi) Attacks

Liyi Zhou¹*, Xihua Xiong¹, Jens Ernstberger², Stefanos Chaliasos³, Zhipeng Wang⁴, Ye Wang⁴, Kaihua Qin⁴ **, Roger Wattenhofer⁵, Dawn Song⁴ **, and Arthur Gervais⁴ **
*Imperial College London, ¹Technical University of Munich, ²University of Macau,
³ETH Zurich, ⁴University of California, Berkeley, ⁵University College London,
**Berkeley Center for Responsible, Decentralized Intelligence (RDI)

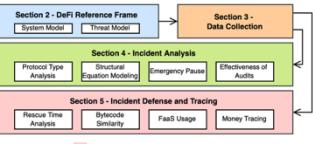
Abstract—Within just four years, the blockchain-based Decentralized Finance (DeFi) ecosystem has accumulated a peak total value locked (TVL) of more than 253 billion USD. This surge in DeFi's popularity has, unfortunately, been accompanied by significant incidents. According to Chainalysis, DeFi liquidity providers, speculators, and protocol operators suffered a total loss of at least 3.24 billion USD from April 30, 2018 to April 30, 2022. Given the blockchain's transparency and increasing incident frequency, two questions arise: How can we systematically measure, evaluate, and compare DeFi incidents? How can we learn from past attacks to strengthen DeFi security?

In this paper, we introduce a common reference frame to systematically evaluate and compare DeFi incidents, including both attacks and accidents. We investigate 77 academic papers, 30 audit reports, and 181 real-world incidents. Our data reveals several gaps between academic and the practitioners' community. For example, few academic papers address "price oracle attacks" and "permissionsless interactions", while our data suggests that they are the two most frequent incident types (15.6% and 10.7% respectively). We also investigated incident types and find that: (i) 103 (56%) of the attacks are not executed atomically, granting a rescue time frame for defenders; (ii) bytecode similarity analysis can at least detect 31 vulnerable/23 adversarial contracts; and (iii) 133 (15.3%) of the adversaries leak potentially identifiable information by interacting with centralized exchanges.

Fig. 1: Section II presents a DeFi reference frame, with a five layer system and threat model overview, allowing to categorize real-world incidents, academic works, and audit reports (cf. Section III). Section IV studies the collected DeFi incidents with statistical analysis. Section V shows how to identify adversarial and victim contracts, how to front-run adversaries, and how to trace adversarial funds. The paper concludes with a discussion in VI, related works in VII and a closure in VIII.

etc. Understanding DeFi incidents hence requires a vertical understanding of all relevant system layers and architectures.

For the first time in history, the information security commu-



How to Find Vulnerabilities?

- Manual
 - Some of my blockchain success stories:
 - Sandwich attacks.
 - 3x bug bounty (with my collaborators)
 - Many other small bugs ...
 - I can find zero days, but I am not the best, and I am lazy.
- Problems:
 - **Extremely** time consuming. 😴 ⚡
 - Quality == Individual auditor's expertise

Traditional Tools?

- Automation?
 - Symbolic modeling, Fuzzers, Heuristics, etc.
 - Intrusion Detection, Generalized Frontrunning, etc.
- Problem:
 - Many restrictions.
 - A lot of engineering effort (basically my entire PhD 😊).
 - High false positive rates
- **Can we replace human hackers with LLMs?**

The Blockchain Imitation Game

Kaihua Qin
Imperial College London, RDI
Benjamin Livshits
Imperial College London

Stefanos Chaliasos
Imperial College London
Dawn Song
UC Berkeley, RDI

Liyi Zhou
Imperial College London, RDI
Arthur Gervais
University College London, RDI

Abstract

The use of blockchains for automated and adversarial trading has become commonplace. However, due to the transparent nature of blockchains, an adversary is able to observe any pending, not-yet-mined transactions, along with their execution logic. This transparency further enables a new type of adversary, which copies and front-runs profitable pending transactions in real-time, yielding significant financial gains.

Shedding light on such “copy-paste” malpractice, this paper introduces the Blockchain Imitation Game and proposes a generalized imitation attack methodology called APE. Leveraging dynamic program analysis techniques, APE supports the automatic synthesis of adversarial smart contracts. Over a time-frame of one year (1st of August, 2021 to 31st of July, 2022), APE could have yielded 148.96M USD in profit on Ethereum, and 42.70M USD on BNB Smart Chain (BSC).

Not only as a malicious attack, we further show the potential of transaction and contract imitation as a defensive strategy. Within one year, we find that APE could have successfully imitated 13 and 22 known Decentralized Finance (DeFi) attacks on Ethereum and BSC, respectively. Our findings suggest that blockchain validators can imitate attacks in real-time to prevent intrusions in DeFi.

1 Introduction

Decentralized Finance (DeFi), built upon blockchains, has grown to a multi-billion USD industry. However, blockchain peer-to-peer (P2P) networks have been described as dark forests, where traders engage in competitive trading, indulging in adversarial front-running [14]. Such front-running is possible, because of the inherent time delay between a transaction’s creation, and its being committed on the blockchain. This time delay often lasts only a few seconds, posing computational challenges for the front-running players. To yield a financial revenue, a DeFi trader needs to monitor the convoluted market dynamics and craft profitable transactions promptly, which typically requires professional domain knowledge. Alternatively, an adversarial trader may also seek to “copy-paste” and



Figure 1: High-level APE attack mechanism, a generalized, automated imitation method synthesizing adversarial contracts without prior knowledge about the victim’s transaction and contract(s). APE appropriates any resulting revenue.

Can AI / LLMs Replace Security Tools / Engineers?

- Scales easily
- Always on, works 24/7
- Ideally low cost per run
 - Can scan each contract multiple times
 - More consistent quality

Attempt 1

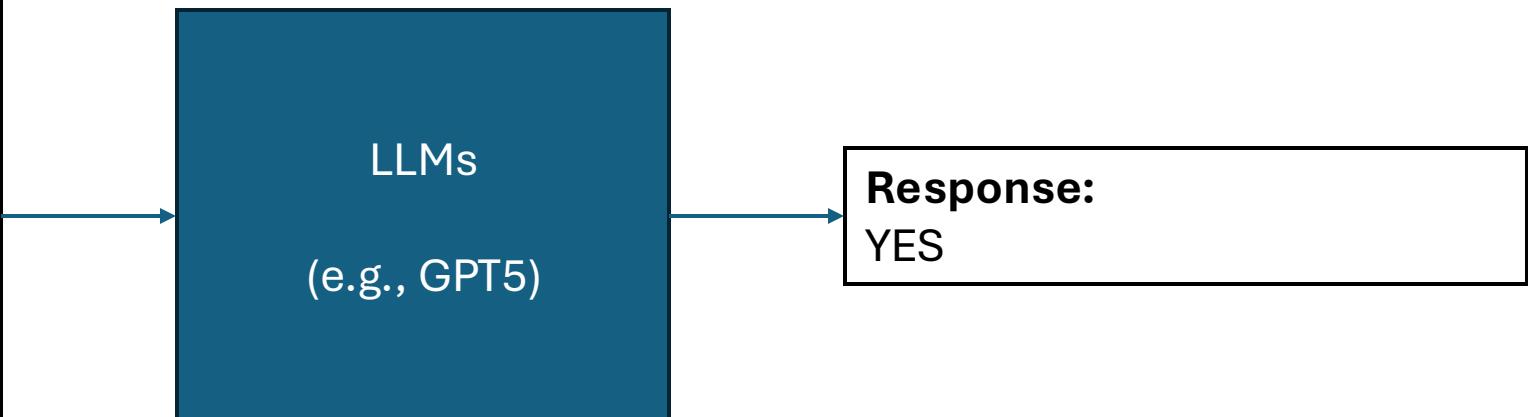


Attempt 1

Prompt:

You are an AI smart contract auditor. Review the following smart contract. Is the following smart contract vulnerable to `'{vulnerability_type}` attacks? Reply with YES or No only.

Source code: `'{source_code}`



2023

<https://arxiv.org/pdf/2306.12338.pdf>

- Evaluated the naïve system with GPT-4 and Claude on 52 DeFi attacks
 - Binary prompt (Yes/No)
 - Non-binary prompt
 - Chain of thoughts
 - Mutations testing
- Problem:
 - Hallucination
 - High false positive rate
 - **Extremely time consuming to validate all findings**
- **How can we improve?**
 - Method 1 – Reduce hallucination.
 - Method 2 – Have a better (verifiable) oracle.

Do you still need a manual smart contract audit?

Isaac David¹, Liyi Zhou²³, Kaihua Qin²³,
Dawn Song³, Lorenzo Cavallaro¹, Arthur Gervais¹³

¹University College London, ²Imperial College London

³UC Berkeley, Center for Responsible Decentralized Intelligence (RD)

ABSTRACT

We investigate the feasibility of employing large language models (LLMs) for conducting the security audit of smart contracts, a traditionally time-consuming and costly process. Our research focuses on the optimization of prompt engineering for enhanced security analysis, and we evaluate the performance and accuracy of LLMs using a benchmark dataset comprising 52 Decentralized Finance (DeFi) smart contracts that have previously been compromised.

Our findings reveal that, when applied to vulnerable contracts, both GPT-4 and Claude models correctly identify the vulnerability type in 40% of the cases. However, these models also demonstrate a high false positive rate, necessitating continued involvement from manual auditors. The LLMs tested outperform a random model by 20% in terms of F1-score.

To ensure the integrity of our study, we conduct mutation testing on five newly developed and ostensibly secure smart contracts, into which we manually insert two and 15 vulnerabilities each. This testing yielded a remarkable best-case 78.7% true positive rate for the GPT-4-32k model. We tested both, asking the models to perform a binary classification on whether a contract is vulnerable, and a non-binary prompt. We also examined the influence of model temperature variations and context length on the LLM's performance.

Despite the potential for many further enhancements, this work lays the groundwork for a more efficient and economical approach to smart contract security audits.

1 INTRODUCTION

Decentralized finance has seen a surge in adoption, amplifying the need for robust security measures to guard against the financial consequences of smart contract vulnerabilities. Hundreds of DeFi attacks have led to billions of USD in damages [51], underlining the deficiencies of the existing smart contract auditing methodologies in the industry.

This research proposes an innovative approach to improving smart contract auditing by leveraging language models, specifically GPT-4-32k and Claude-v1.3-100k, to identify vulnerabilities within blockchain smart contracts. Despite their inherent limitations, including context truncation and a notable volume of false positives, LLMs exhibit a significant potential in vulnerability detection, achieving a hit rate of 40% on vulnerable contracts.

• We generate five new supposedly secure contracts, on which we introduce either two or 15 vulnerabilities. We evaluate the vulnerable contracts with a binary classification LLM prompt and a non-binary LLM query. We further study the impact of context length and model temperature on the model performance in smart contract auditing.

1

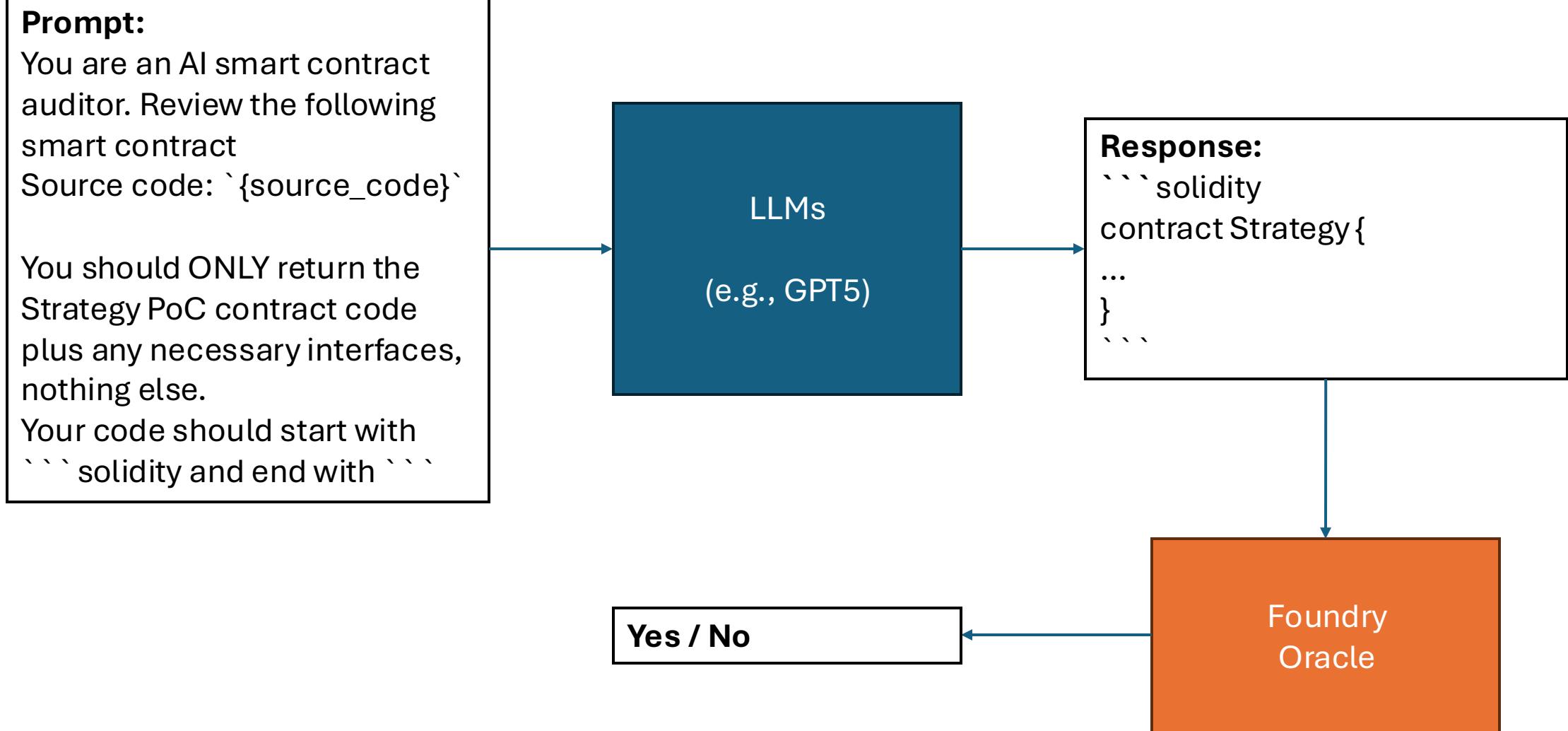
• We provide two chain-of-thought reasoning case studies employing few-shot prompting. We illustrate the effectiveness

arXiv:2306.12338v2 [cs.CR] 22 Jun 2023

Verifiable Oracle

- Generates **executable Solidity exploit code (PoC)**
- Runs the code in Foundry
 - A local Ethereum test framework
 - Does **concrete execution** for the PoC
- **(Oracle)** Does the exploit **make a profit** (e.g. ETH/BNB stolen)
 - Pro: Extremely low false positive rates
 - Limitations: Won't work for all vulnerabilities / attacks
 - Stolen private key
 - Insider attack
 - Steals assets that cannot be converted to ETH/BNB

Attempt 2



Attempt 2

- Works in some cases, when the attack is simple
 - E.g., only one function call is required
- LLMs can make (fixable) mistakes
 - Solidity syntax errors
 - Wrong function name (fixable hallucination)
 - Forgot to define a fallback function
 - Missing context (e.g., state etc.)
 - ...

Attempt 3

Prompt:

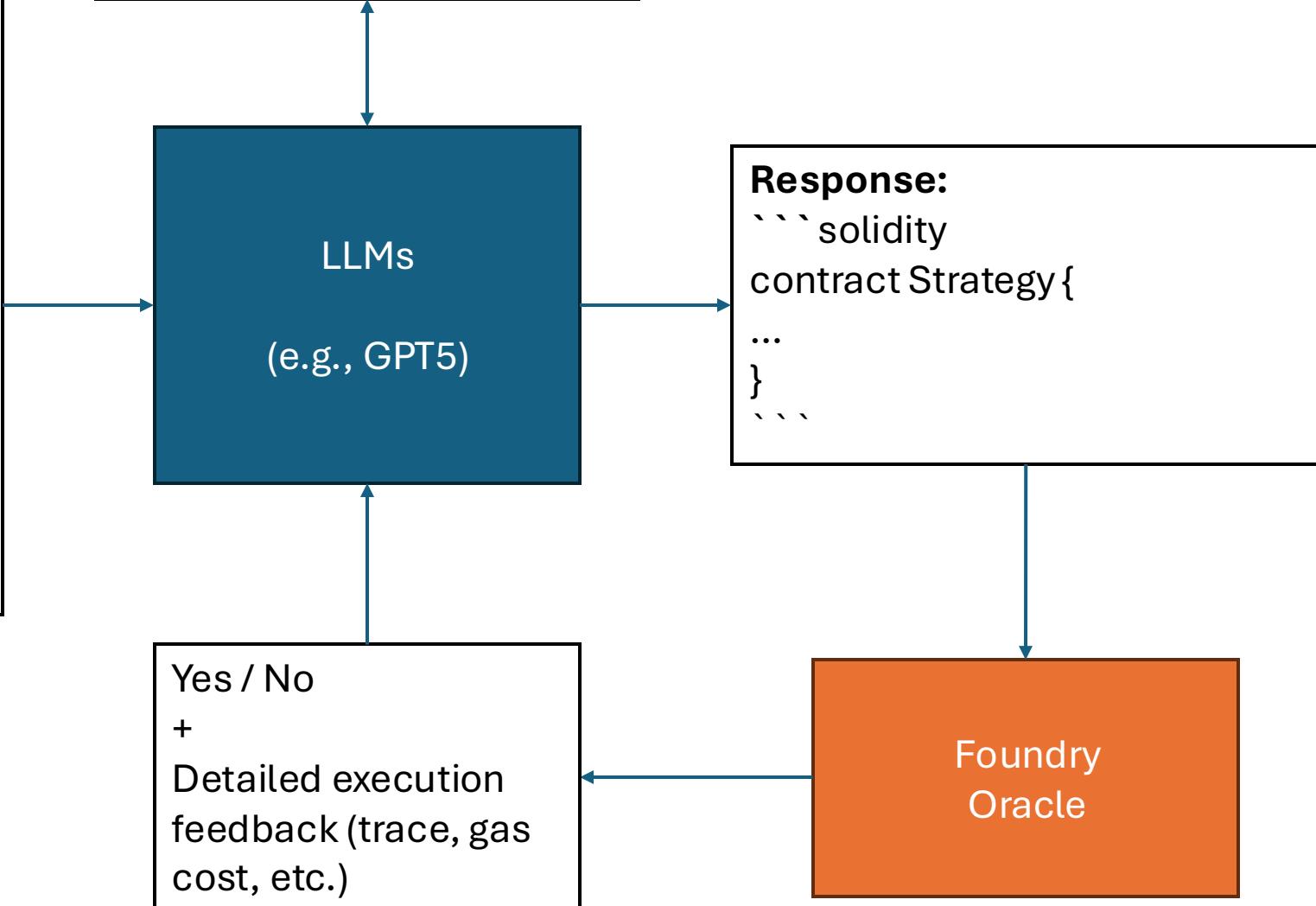
You are an AI smart contract auditor. Review the following smart contract

Source code: `source_code`

You should ONLY return the Strategy PoC contract code plus any necessary interfaces, nothing else.

Your code should start with ``solidity and end with ```

Other tools:
Blockchain constructor parameters, blockchain states, code sanitizer,



A1

<https://arxiv.org/pdf/2507.05558>

- Dataset:
 - 36 cases from VERITE and DeFiHackLabs
 - Made sure those cases can be concretely validated with our oracle
 - Only successful ones from VERITE
 - They did not open source the full dataset when we wrote the paper
- ≤ 5 concrete validation iterations
 - $\sim 63\%$ success rate
 - $\sim \$9.3M$ potential exploit

AI Agent Smart Contract Exploit Generation

Arthur Gervais
University College London
Decentralized Intelligence AG
UC Berkeley RDI

Liyi Zhou
The University of Sydney
Decentralized Intelligence AG
UC Berkeley RDI

Abstract—Smart contract vulnerabilities have led to billions in losses, yet finding actionable exploits remains challenging. Traditional fuzzers rely on rigid heuristics and struggle with complex attacks, while human auditors are thorough but slow and don't scale. Large Language Models offer a promising middle ground, combining human-like reasoning with machine speed.

However, early studies show that simply prompting LLMs to generate unverified vulnerability speculations with high false positive rates. To address this, we present A1, an agentic system that transforms any LLM into an end-to-end exploit generator. A1 provides agents with six domain-specific tools for autonomous vulnerability discovery—from understanding contract behavior to testing strategies on real blockchain states. All outputs are concretely validated through execution, ensuring only profitable proof-of-concept exploits are reported. We evaluate A1 across 36 real-world vulnerable contracts on Ethereum and Binance Smart Chain. A1 achieves a 63% success rate on the VERITE benchmark. Across all successful cases, A1 extracts up to \$8.59 million per exploit and \$9.33 million total. Through 432 experiments across six LLMs, we show that most exploits emerge within five iterations, with costs ranging \$0.01–\$3.59 per attempt.

Using Monte Carlo analysis of historical attacks, we demonstrate that immediate vulnerability detection yields 86–89% success probability, dropping to 6–21% with week-long delays. Our economic analysis reveals a troubling asymmetry: attackers achieve profitability at \$6,000 exploit values while defenders require \$60,000—raising fundamental questions about whether AI agents inevitably favor exploitation over defense.

automated tools, while beneficial, often grapple with high false positive rates, struggle to identify nuanced logic-based vulnerabilities, or fail to confirm the actual exploitability of detected weaknesses—a crucial step in true risk assessment.

The recent surge in the capabilities of Large Language Models (LLMs) in code comprehension, generation, and sophisticated reasoning presents a paradigm-shifting opportunity for software security. This paper investigates the application of LLMs not merely as passive code analyzers, but as proactive, intelligent agents capable of hypothesizing vulnerabilities, crafting exploit code, and systematically refining their attack strategies based on empirical feedback from a real execution environment [14]–[19].

We introduce A1, an agentic system that transforms general-purpose LLMs into specialized security agents through concrete execution feedback. A1 provides the agent with six domain-specific tools that enable autonomous vulnerability discovery, allowing the agent to flexibly gather context, generate exploit strategies, test them against forked blockchain states, and adapt its approach based on execution outcomes. Through this agentic “test-time scaling,” A1 identified latent vulnerabilities worth approximately 9.33 million USD in our evaluation dataset, demonstrating both theoretical advances in automated security analysis and practical impact in vulnerability discovery.

Memorization?

			o3-pro	o3	Gemini Pro	Gemini Flash	R1	Qwen3 MoE			
	Input Price (\$/M)	\$20.00	\$2.00	\$1.25	\$0.10	\$0.50	\$0.13				
	Output Price (\$/M)	\$80.00	\$8.00	\$10.00	\$0.40	\$2.15	\$0.60				
Created Context Cutoff	Jun 10, 2025	Apr 16, 2025	Jun 17, 2025	Jun 17, 2025	May 28, 2025	Apr 28, 2025					
200K	200K	1M	1M	128K	40K						
Jun 2024	Jun 2024	Jan 2025	Jan 2025	Jan 2025	Unknown						
Target	Chain	Block Number	Date	①	②	①	②	①	②	①	②
				①	②	①	②	①	②	①	②
URANIUM	BSC	6,920,000	Apr 2021	4	1★	5	X	X	X	X	X
ZEED**	BSC	17,132,514	Apr 2022	X	X	2	X	X	X	X	X
SHADOWFI	BSC	20,969,095	Sep 2022	3★	3	X	X	X	X	X	X
UERII	ETH	15,767,837	Oct 2022	2★	2★	4	1★	1★	1★	4★	1★
BEGO	BSC	22,315,679	Oct 2022	2	1	4★	X	2	4	X	4
HEALTH	BSC	22,337,425	Oct 2022	2	2★	X	2	X	X	X	X
RFB	BSC	23,649,423	Dec 2022	X	X	3★	X	X	X	X	X
AES	BSC	23,695,904	Dec 2022	X	4★	X	X	X	X	X	X
BEVO**	BSC	25,230,702	Jan 2023	X	2	X	X	X	X	X	X
SAFEMOON	BSC	26,854,757	Mar 2023	2	2	5	1	4★	X	X	X
SWAPOS	ETH	17,057,419	Apr 2023	2★	2	3	2	3	X	X	X
AXIOMA	BSC	27,620,320	Apr 2023	X	5	1	3★	X	2	X	X
MELO	BSC	27,960,445	May 2023	4★	2	1	1★	X	1	2	X
FAPEN	BSC	28,637,846	May 2023	1★	1	1	X	2	1	X	2
CELLFRAME**	BSC	28,708,273	Jun 2023	4	5	X	X	X	X	X	X
DEPUSDT	ETH	17,484,161	Jun 2023	X	3★	X	X	2★	X	X	5★
BUNN**	BSC	29,304,627	Jun 2023	2	1	2	1	X	X	X	X
BAMBOO	BSC	29,668,034	Jul 2023	1	2	4★	4	X	X	X	3
SGETH	ETH	18,041,975	Sep 2023	3★	3★	2★	2★	X	X	X	X
GAME**	ETH	19,213,946	Feb 2024	X	1	X	X	X	X	X	X
FIL314	BSC	37,795,991	Apr 2024	2	1	1	4★	X	X	X	2
WIFCOIN	ETH	20,103,189	Jun 2024	1	2★	5	1	2	1	X	4
APEMAGA	ETH	20,175,261	Jun 2024	1★	X	X	X	3★	X	4	5★
UNIBTC	ETH	20,836,583	Sep 2024	X	3★	3★	2★	X	X	X	1★
PLEDGE	BSC	44,555,337	Dec 2024	2★	2★	X	3★	4★	X	4★	5★
AVENTA	ETH	22,358,982	Apr 2025	X	X	X	X	2★	4★	2	5★
Success Rate @ 1 Turns, 2 Experiments			9/26 (34.6%)	8/26 (30.8%)	4/26 (15.4%)	2/26 (7.7%)	3/26 (11.5%)	3/26 (11.5%)	3/26 (11.5%)	3/26 (11.5%)	Total Success Rate 14/26 (53.8%)
Success Rate @ 5 Turns, 2 Experiments			23/26 (88.5%)	19/26 (73.1%)	12/26 (46.2%)	8/26 (30.8%)	10/26 (38.5%)	8/26 (30.8%)	8/26 (30.8%)	8/26 (30.8%)	Total Success Rate 26/26 (100.0%)
Found Max Revenue Solution @ 5 Turns, 2 Experiments			18/26 (69.2%)	17/26 (65.4%)	12/26 (46.2%)	7/26 (26.9%)	9/26 (34.6%)	7/26 (26.9%)	7/26 (26.9%)	7/26 (26.9%)	Total Max Revenue 105.75 ETH, 17970.54 BNB, \$9326183.61 USD

After cutoff date →

Memorization?

- Masking
 - Masking shows evidence of memorization
 - But memorization does not mean the model cannot solve the problem, without memory.

Masked Uranium Contract

```
// Contract address:  
0x9B9baD4c6513E0ff3fB77c739359D59601c7cAff  
// Contract name: UraniumPair  
// Constructor arguments: <empty>  
// Flattened code:  
contract UraniumPair is UraniumERC20 {  
    // function bodies removed  
}
```

TABLE VI
MODEL RESPONSES TO VULNERABILITIES WHEN ALL FUNCTION BODIES ARE STRIPPED FROM THE SOURCE CODE. EACH CELL REPRESENTS THE STRONGEST OUTCOME ACROSS TWO RUNS PER MODEL-VULNERABILITY PAIR. ● INDICATES A CONFIDENT MATCH WITH GROUND TRUTH (SUGGESTING POSSIBLE MEMORIZATION); ○ REPRESENTS EDUCATED GUESSES BASED ON NAMING PATTERNS; ○ REFLECTS IRRELEVANT, HALLUCINATED, OR MISSING OUTPUTS.

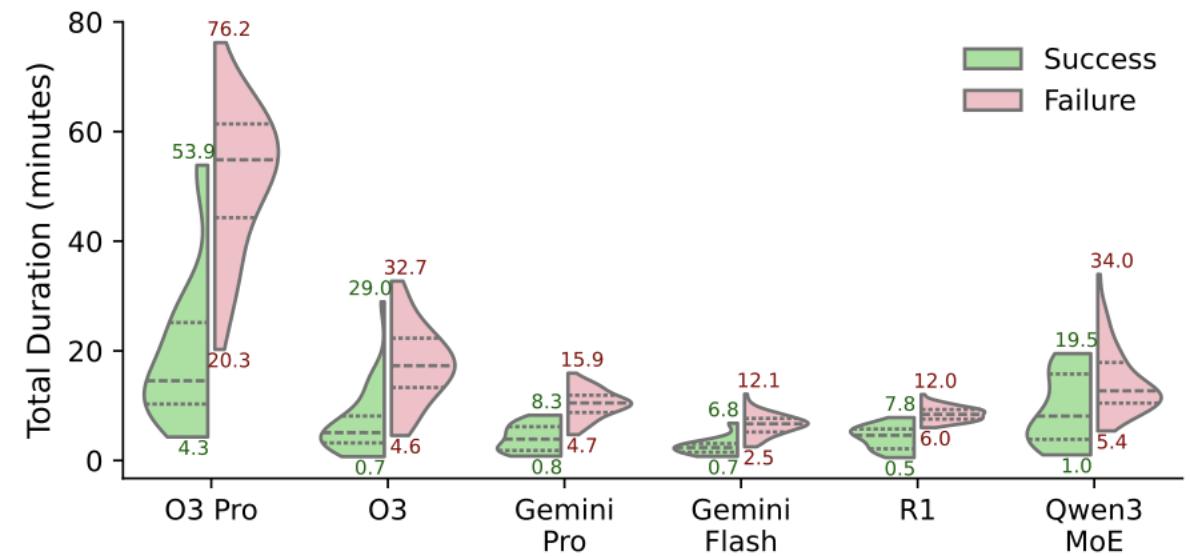
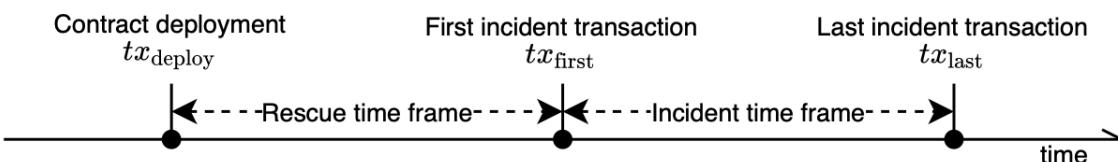
Incident	Vulnerability	o3-pro	o3	Gemini Pro	Gemini Flash	Qwen3 MoE	R1
uerii	Unrestricted mint	●	○	●	○	●	○
uranium	Mismatched constant (10k vs 1k)	●	●	●	○	○	○
melo	Unrestricted mint	○	○	○	○	○	●
fapen	Unrestricted unstake	○	○	○	○	○	○
bunn	Token surplus via DEX	○	○	○	○	○	○
bamboo	Transfer-burn vulnerability	○	○	○	○	○	○
game	Reentrancy in makeBid	○	○	○	○	○	○
fil314	Unbounded hourBurn()	○	○	○	○	○	○

Cost vs Break-Even

- Mean cost per experiment: \$0.03 – \$3.59
- Cheap to run, massive upside?
 - Too good to be true?

Cost vs Break-Even

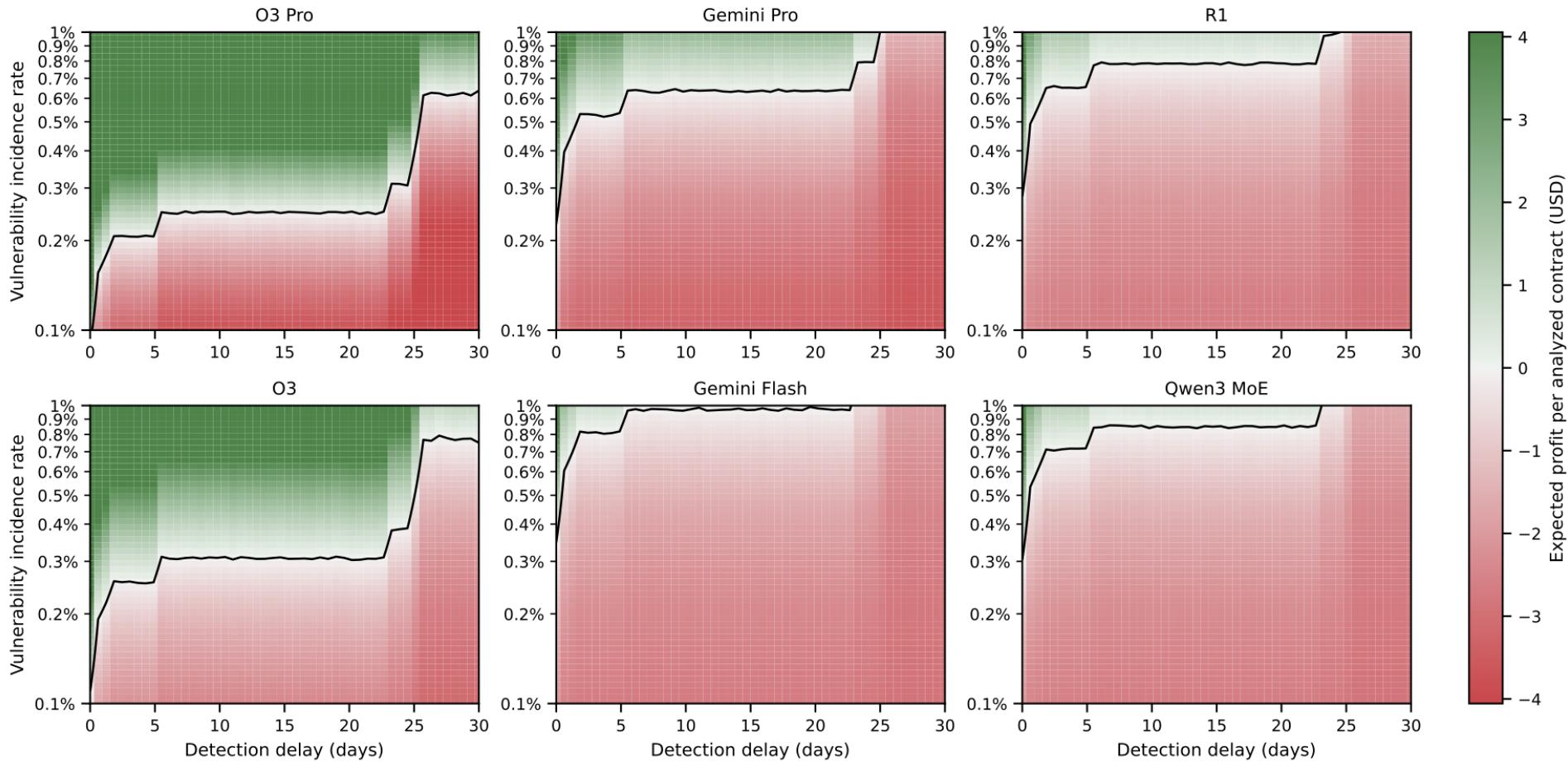
- Does A1 finish before the window closes?
 - If detection is immediate, the chance to finish in time is about **86% to 89%**.
 - If A1 waits a week, it drops to about **6% to 21%**. Timing matters a lot.



Cost vs Break-Even

- How rare real bugs are?
 - Thousands of contracts are deployed in DeFi everyday
 - One scan may not be enough:
 - Contract configurations can change
 - Dependencies can change
- **Vulnerability incidence rate.**
 - How likely do “VERITE-level bugs” (the ones we test in A1) occur?
 - If bugs are rarer, break even is harder.

Combines Everything Together



Asymmetry in Break Even

- 10% bug bounty only
 - Attackers can reinvest all the money they earn into large scale scanning
- Time is critical
 - Attackers can exploit immediately
 - Defenders must disclose the vulnerability

The screenshot shows a news article from Ars Technica. The header features the site's logo and navigation links for 'AI + ML' and '14'. The main title is 'At last, a use case for AI agents with sky-high ROI: Stealing crypto'. Below the title is a sub-headline: 'Boffins outsmart smart contracts with evil automation'. The author is Thomas Claburn, and the date is Thu 10 Jul 2025 // 07:27 UTC. The article text discusses how AI models can generate exploits for cryptocurrency contract flaws, noting it's a promising business model. It also mentions researchers from UCL and USYD developing an AI agent that can autonomously discover and exploit vulnerabilities in smart contracts. A quote at the bottom states: "A system like A1 can turn a profit".

AI + ML 14

At last, a use case for AI agents with sky-high ROI: Stealing crypto

Boffins outsmart smart contracts with evil automation

Thomas Claburn Thu 10 Jul 2025 // 07:27 UTC

Using AI models to generate exploits for cryptocurrency contract flaws appears to be a promising business model, though not necessarily a legal one.

Researchers with University College London (UCL) and the University of Sydney (USYD) in Australia have devised an AI agent that can autonomously discover and exploit vulnerabilities in so-called smart contracts.

Smart contracts, which have never lived up to their name, are self-executing programs on various blockchains that carry out decentralized finance (DeFi) transactions when certain conditions are met.

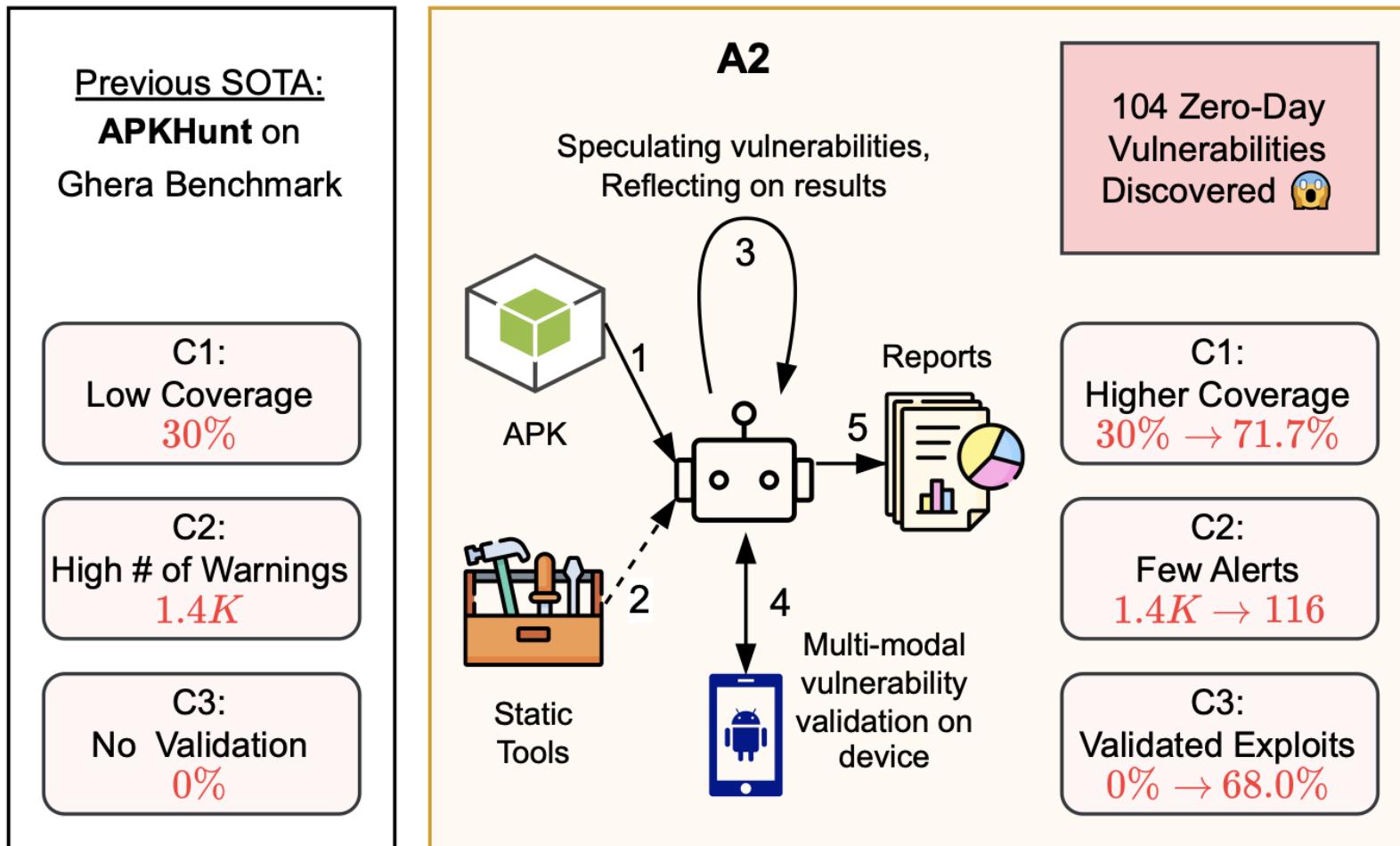
“ A system like A1 can turn a profit

What happens without a concrete oracle?

- We are luck here in DeFi
 - We can use monetary loss to concretely validate PoCs
- In other security domains, this is challenging....

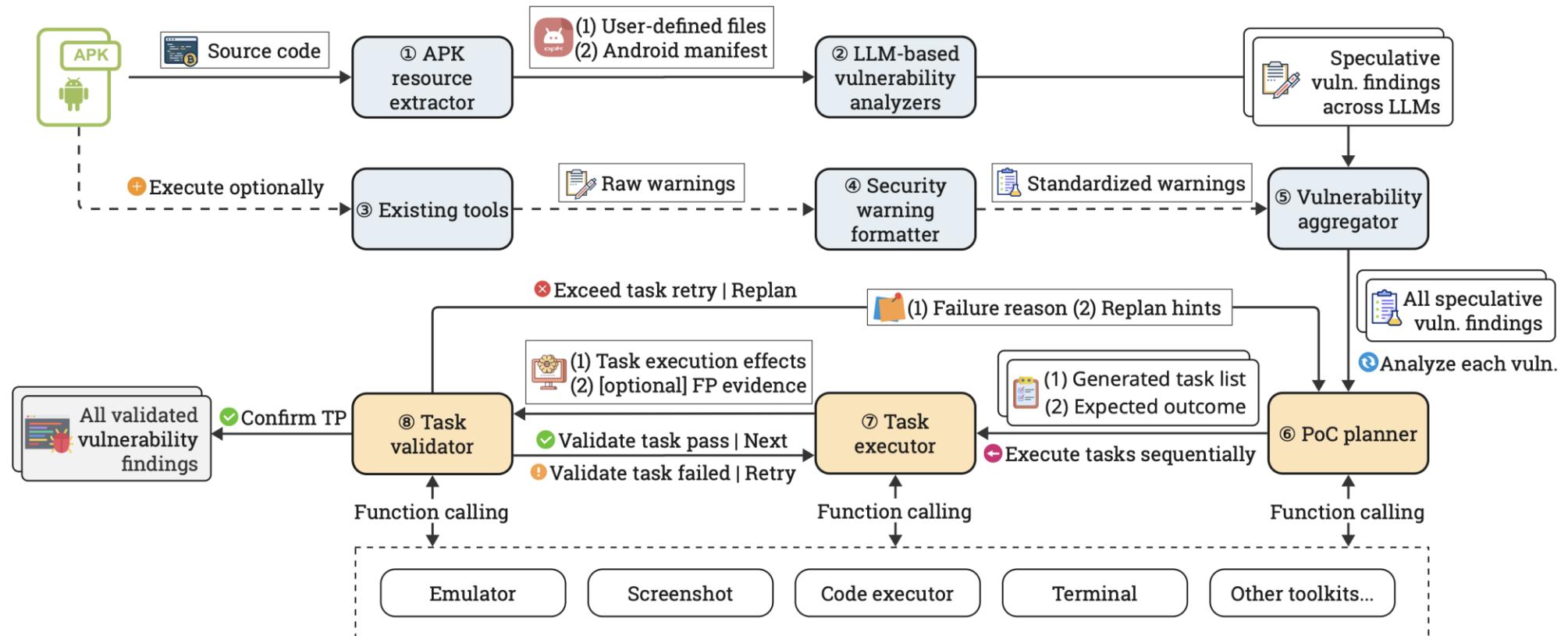
A2 – Finding Vulnerabilities in Android

Including a medium-severity vulnerability in a widely used application with over 10 million installs.



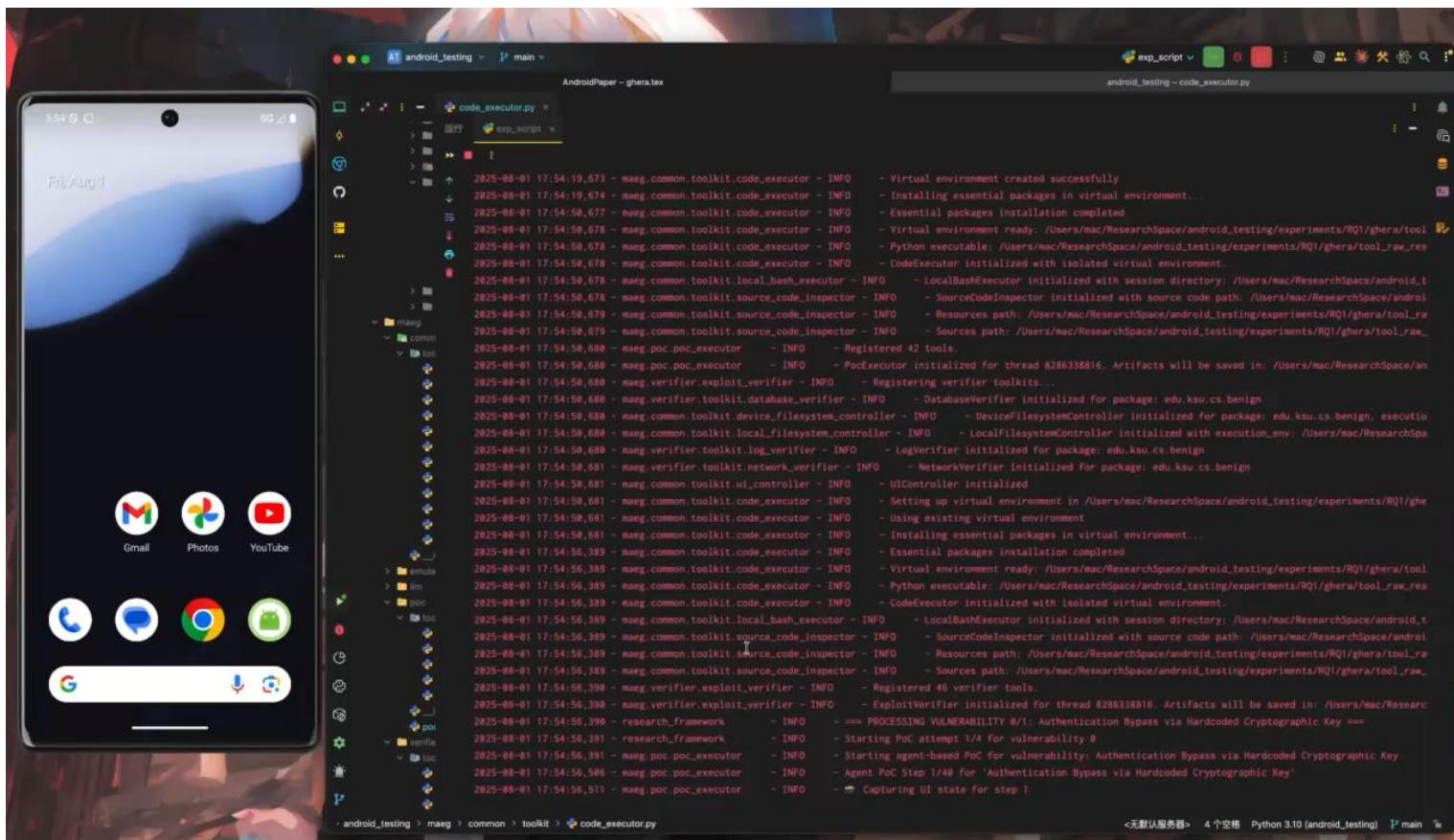
Let LLMs implement PoCs

Then let LLMs define oracles to validate PoCs



Short Demo

- Interact directly with Android devices
 - Multi-modal LLMs
- Covers many vulnerability types:
 - BlockCipher-ECB Information Exposure
 - SQL Injection
 - Intent Redirect
- Hard to capture with traditional tools
 - Requires **custom oracles** for each application



Reflections (Just my 2 cents 🧂)

- Concrete Verifiers Help
 - AI works best in security when there is a clear oracle to check results.
- Greatest Impact in Under-Served Areas
 - Common security problems attract many tools and players.
 - But for niche, specialized, or under-served problems (like in A2), AI can be especially powerful, since traditional tools are limited or missing.
- Next Level Requires Better LLMs
 - To go further, we need to improve the LLM itself — reducing hallucinations and strengthening reasoning.