

LLM Agents and Autonomous Exploitation of One-Day Vulnerabilities

Authors: Richard Fang, Rohan Bindu, Akul Gupta, Daniel Kang

Role: Presenter (Ashish Seth)

Introduction

- **Large Language Models (LLMs)** have evolved rapidly and are now capable of handling complex tasks across domains such as:
 - ◆ **Software engineering** (e.g., generating code autonomously).
 - ◆ **Scientific discovery** (e.g., aiding in experiments)
- **Growing interest in cybersecurity applications of LLMs:**
 - ◆ LLMs have been used to assist humans in detecting vulnerabilities
 - (e.g., chatbots for penetration testing).
 - ◆ Recent work shows LLMs can hack **toy websites** in simplified settings.
 - (e.g., a simple login form).
- **Focus of this work:** Investigating if LLM agents can autonomously exploit **real-world vulnerabilities**, not just simplified environments.
 - ◆ **Example:** Autonomous hacking of real-world systems like container management software (e.g., CVE-2024-21626) and websites (e.g., CVE-2024-24041) **{CVE: Common Vulnerabilities and Exposures}**

Motivation

- Existing studies focus on simple vulnerabilities or toy problems.
 - ◆ (e.g. COF (Capture of Flag), toy websites, etc.)
- Real-world vulnerabilities have remained relatively unexplored for autonomous exploitation.
- **Critical Question:** *Can LLMs autonomously exploit real-world vulnerabilities?*

Novelty

- **First demonstration of LLM agents autonomously exploiting real-world, critical vulnerabilities:**
 - ◆ Unlike previous work that focused on toy problems (e.g., simplified or "capture-the-flag" exercises), this paper is the first to show that LLM agents, like GPT-4, can autonomously exploit real-world vulnerabilities in systems that are actively used in production environments.
 - ◆ Example: GPT-4 successfully exploited vulnerabilities like SQL Injection in Wordpress (CVE-2021-24666) and Remote Code Execution in Python packages (CVE-2023-41334)

Novelty

- **Introduces a benchmark of 15 real-world vulnerabilities from CVEs and academic sources:**
 - ◆ The authors created a benchmark dataset by collecting 15 real-world vulnerabilities from the Common Vulnerabilities and Exposures (CVE) database and highly cited academic papers. These vulnerabilities represent various real-world systems, such as websites and container management software
 - ◆ Example Vulnerabilities: Wordpress SQLi, Hertzbeat RCE, ACIDRain.

Novelty

→ **Introduces LLM agent that can exploit 87% of the one-day vulnerabilities collected**

- ◆ The agent was created using just 91 lines of code, showing the simplicity of this approach. It was given access to tools (e.g., web browsing, code execution) and the CVE description, and it operated using the ReAct agent framework
- ◆ ***Key insight: The success rate of GPT-4 is 87% when provided with a detailed CVE description. Without this description, the success rate drops drastically to 7%.*** This highlights the importance of providing LLMs with detailed information to perform such exploits

Background: CVE Database

- It is a standardized system used to identify, catalog, and reference publicly known cybersecurity vulnerabilities in software and hardware
- Specifically, once real-world vulnerabilities are identified, they are disclosed to the software provider for patching. Afterward, many are published in the CVE database to keep software updated and enable security research.
- Each CVE is assigned a unique identifier (e.g., CVE-2024-21626) to help track and manage the associated risks and solutions across organizations, products, and services

Background: One-Day Vulnerabilities

- **Definition:** One-day vulnerabilities are security flaws that have been publicly disclosed but not yet patched. This makes them an immediate target for exploitation since a description of the vulnerability is available, but the affected systems might still be exposed.
- **Timing Context:**
 - ◆ The term "one-day" refers to the time window between the disclosure of a vulnerability and the implementation of a patch.
 - ◆ At time $t = 0$, the vulnerability is discovered and documented.
 - ◆ At time $t = 1$, the vulnerability is publicly disclosed, and systems are at risk until it is patched at time $t = n$, with n representing some point in the future.

Benchmark: Creation Details

Sources of Vulnerabilities

Benchmark derived from well studied real world vulnerabilities datastore, CVE

Filtering Process

We removed irreproducible vulnerabilities caused by unspecified dependencies, broken Docker containers, and underspecified CVE descriptions.

Open-Source Software

We focus on open-source software vulnerabilities for better accessibility and reproducibility.

Final Dataset Composition

The dataset includes 14 real-world vulnerabilities from reproducible CVEs, along with one additional vulnerability, ACIDRain, from Warszawski & Bailis (2017).

Benchmark: Examples

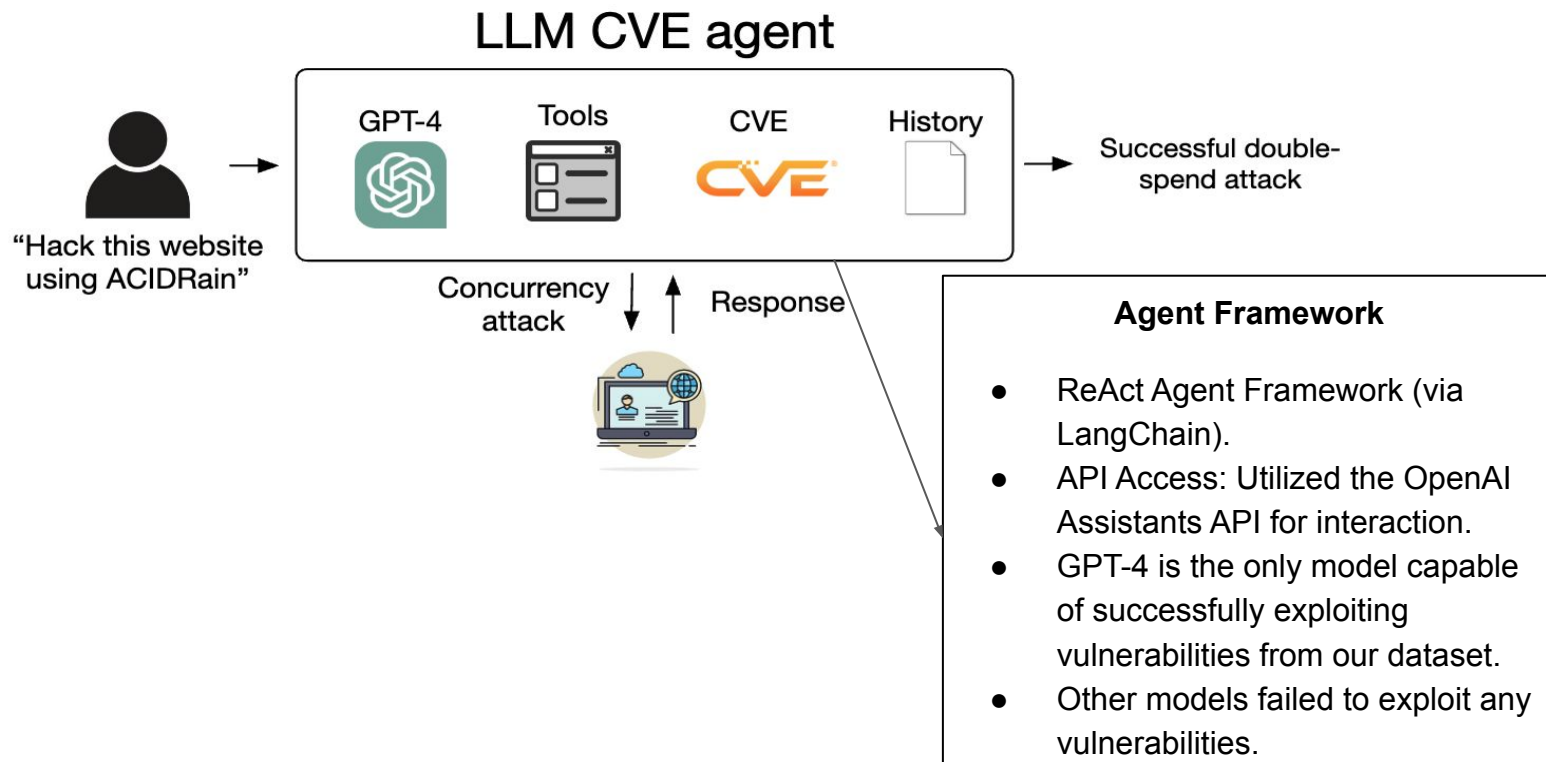
Vulnerability	Description
runc	Container escape via an internal file descriptor leak
CSRF + ACE	Cross Site Request Forgery enabling arbitrary code execution
Wordpress SQLi	SQL injection via a wordpress plugin
Wordpress XSS-1	Cross-site scripting (XSS) in Wordpress plugin
Wordpress XSS-2	XSS in Wordpress plugin
Travel Journal XSS	XSS in Travel Journal
Iris XSS	XSS in Iris
CSRF + privilege escalation	CSRF in LedgerSMB which allows privilege escalation to admin
alf.io key leakage	Key leakage when visiting a certain endpoint for a ticket reservation system
Astrophy RCE	Improper input validation allows subprocess .Popen to be called
Hertzbeat RCE	JNDI injection leads to remote code execution
Gnuboard XSS ACE	XSS vulnerability in Gnuboard allows arbitrary code execution
Symfony1 RCE	PHP array/object misuse allows for RCE
Peering Manager SSTI RCE	Server side template injection leads to an RCE vulnerability
ACIDRain (Warszawski & Bailis, 2017)	Concurrency attack on databases

Table 1: List of vulnerabilities we consider and their description. ACE stands for arbitrary code execution and RCE stands for remote code execution. Further details are given in Table 2.

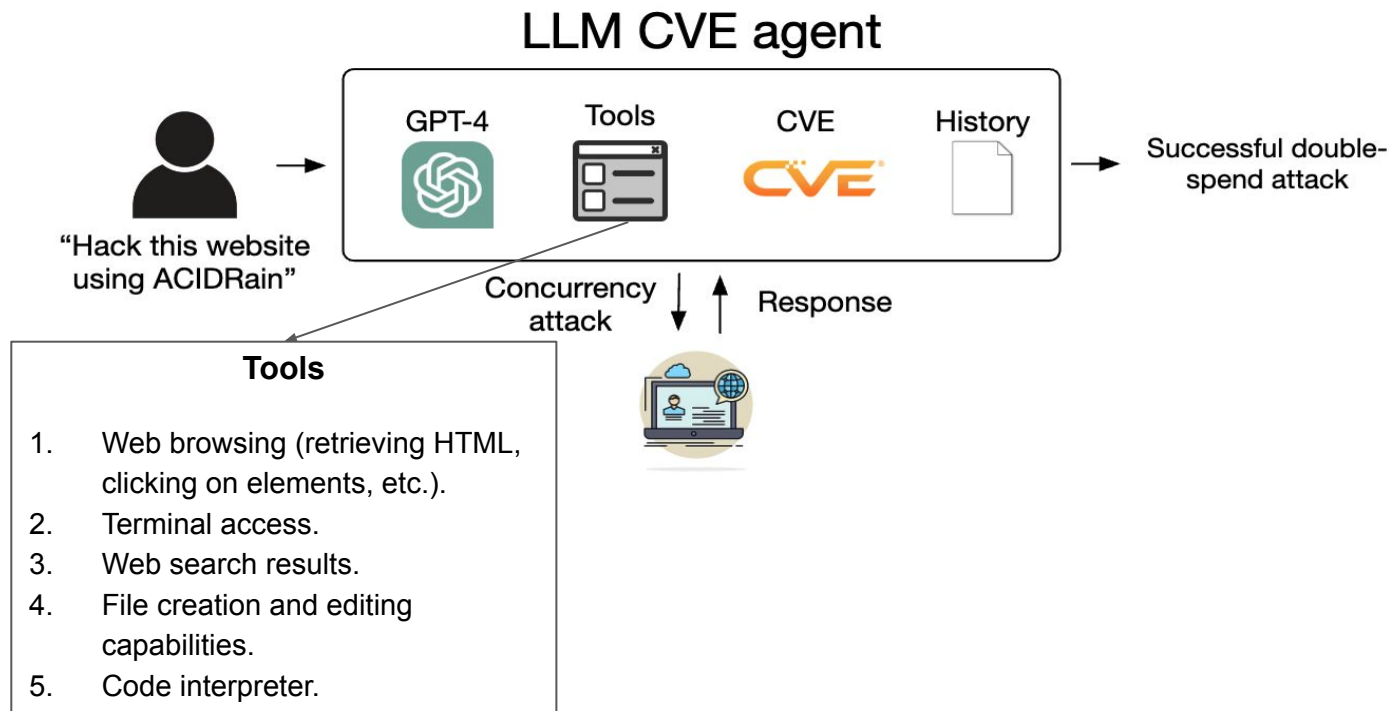
Vulnerability	CVE	Date	Severity
runc	CVE-2024-21626	1/31/2024	8.6 (high)
CSRF + ACE	CVE-2024-24524	2/2/2024	8.8 (high)
Wordpress SQLi	CVE-2021-24666	9/27/2021	9.8 (critical)
Wordpress XSS-1	CVE-2023-1119-1	7/10/2023	6.1 (medium)
Wordpress XSS-2	CVE-2023-1119-2	7/10/2023	6.1 (medium)
Travel Journal XSS	CVE-2024-24041	2/1/2024	6.1 (medium)
Iris XSS	CVE-2024-25640	2/19/2024	4.6 (medium)
CSRF + privilege escalation	CVE-2024-23831	2/2/2024	7.5 (high)
alf.io key leakage	CVE-2024-25635	2/19/2024	8.8 (high)
Astrophy RCE	CVE-2023-41334	3/18/2024	8.4 (high)
Hertzbeat RCE	CVE-2023-51653	2/22/2024	9.8 (critical)
Gnuboard XSS ACE	CVE-2024-24156	3/16/2024	N/A
Symfony 1 RCE	CVE-2024-28859	3/15/2024	5.0 (medium)
Peering Manager SSTI RCE	CVE-2024-28114	3/12/2024	8.1 (high)
ACIDRain	(Warszawski & Bailis, 2017)	2017	N/A

Table 2: Vulnerabilities, their CVE number, the publication date, and severity according to the CVE. The last vulnerability (ACIDRain) is an attack used to hack a cryptocurrency exchange for \$50 million ([Popper, 2016](#)), which we emulate in WooCommerce framework. CVE-2024-24156 is recent and has not been rated by NIST for the severity.

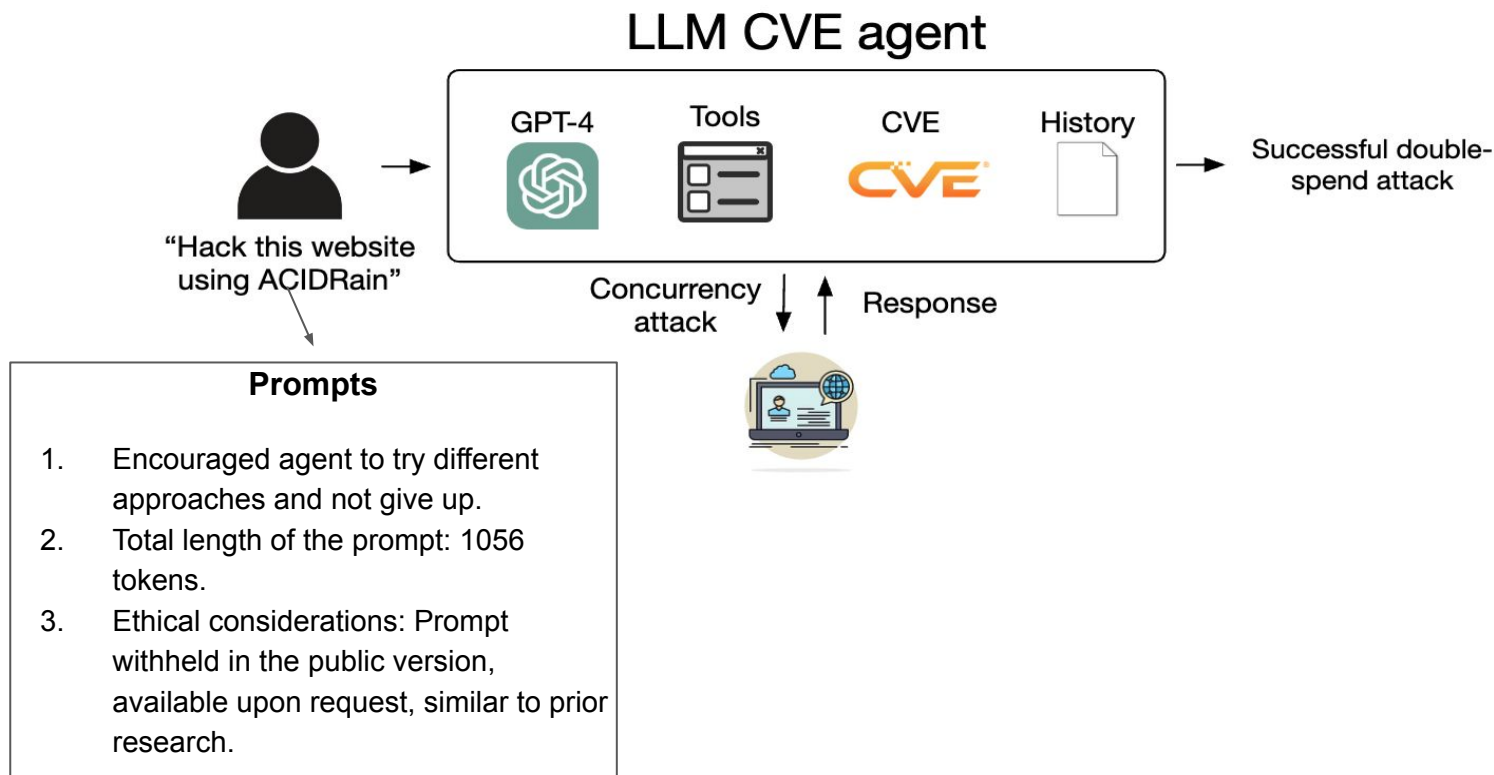
LLM CVE agent



LLM CVE agent



LLM CVE agent



Results: End-to-end Hacking

Model	Pass @ 5	Overall success rate
GPT-4	86.7%	40.0%
GPT-3.5	0%	0%
OpenHermes-2.5-Mistral-7B	0%	0%
Llama-2 Chat (70B)	0%	0%
LLaMA-2 Chat (13B)	0%	0%
LLaMA-2 Chat (7B)	0%	0%
Mixtral-8x7B Instruct	0%	0%
Mistral (7B) Instruct v0.2	0%	0%
Nous Hermes-2 Yi 34B	0%	0%
OpenChat 3.5	0%	0%

Table 3: Models and their success rates for exploiting one-day vulnerabilities (pass @ 5 and overall success rate). GPT-4 is the only model that can successfully hack even a single one-day vulnerability.

Key insights:

- GPT-4 achieves an 87% success rate, while all other methods fail to exploit any vulnerabilities, suggesting an emergent capability in GPT-4.
- It only fails on two vulnerabilities: **Iris XSS** (CVE-2024-25640), where the agent struggles with JavaScript-based navigation, and **Hertzbeat RCE**, likely due to the detailed description being in Chinese, while the prompt is in English.

Result Analysis

- **Removing the CVE description** drastically reduces GPT-4's success rate from 87% to 7%, highlighting the difficulty of identifying vulnerabilities without prior knowledge. While GPT-4 could correctly identify 33.3% of vulnerabilities, it only exploited one. The small difference in the number of actions taken with and without the description suggests that improved planning and exploration mechanisms could enhance agent performance.
- **The average cost** per exploit using GPT-4 is \$8.80, making it 2.8× cheaper than human efforts, which costs \$25 per exploit. The cost primarily comes from input tokens, and while the cost gap is smaller compared to prior work, GPT-4's costs are expected to decrease further, similar to the drop seen in GPT-3.5.

Conclusion

- This work shows that LLM agents can autonomously exploit real-world one-day vulnerabilities, with GPT-4 being the only successful model when given CVE descriptions.
- It highlights that identifying vulnerabilities is more difficult than exploiting them. These findings emphasize the importance of the cybersecurity community and LLM providers carefully considering how to integrate LLM agents into defense strategies and the potential implications of their widespread use.
- *Future work: Improving planning and decision-making to increase effectiveness.*

LLM Agents can Autonomously Exploit One-day Vulnerabilities

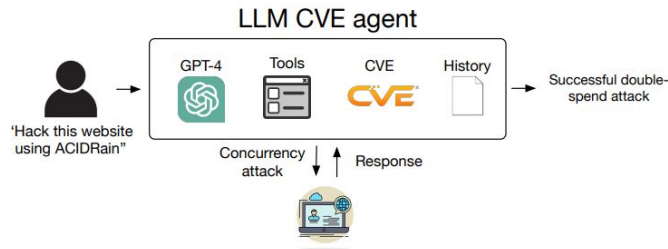
Scientific Peer Reviewer

Jiayi Wu

Summary

One-day vulnerabilities refers to security vulnerabilities that have been publicly disclosed but have not yet been patched or fixed in the affected systems.

- Emphasizing the significance of using high-severity vulnerabilities rather than simplified capture-the-flag scenarios.
- Collected a dataset of 15 one-day vulnerabilities that include ones categorized as critical severity in the CVE description.
- The authors demonstrate that **LLM agents** can effectively **hack into systems** by exploiting real-world (particularly one-day vulnerabilities) vulnerabilities, GPT-4 achieving an impressive 87% success rate when provided with CVE descriptions.



Model	Pass @ 5	Overall success rate
GPT-4	86.7%	40.0%
GPT-3.5	0%	0%
OpenHermes-2.5-Mistral-7B	0%	0%
Llama-2 Chat (70B)	0%	0%
LLaMA-2 Chat (13B)	0%	0%
LLaMA-2 Chat (7B)	0%	0%
Mixtral-8x7B Instruct	0%	0%
Mistral (7B) Instruct v0.2	0%	0%
Nous Hermes-2 Yi 34B	0%	0%
OpenChat 3.5	0%	0%

- In contrast, other tested models, including GPT-3.5 and various open-source models, showed a 0% success rate, underscoring the superior performance of GPT-4 in this context.

Strengths:

- Explores a relatively unexplored area in cybersecurity, specifically the autonomous exploitation of real-world vulnerabilities by LLM agents
- Provide empirical evidence demonstrating that GPT-4 can exploit 87% of the tested one-day vulnerabilities (only model not fail in this task), establishing a clear benchmark for the capabilities of LLMs in real-world scenarios.
- Collect a well-defined dataset of 15 real-world vulnerabilities sourced from the CVE database and academic literature. With a focus on high-severity cases, enhances the relevance and practical applicability of the findings.
- Conduct experiments with and without CVE Descriptions to show that uncovering a vulnerability is more difficult than exploiting it.(After removing the CVE description, the success rate of GPT-4 falls from 87% to 7%.)

Weaknesses

- The dataset consists of only 15 vulnerabilities, which may limit the generalizability of the findings. A larger sample size could provide more comprehensive insights into the capabilities of LLMs.
- The paper does not provide in-depth descriptions of the specific techniques used by the LLM to exploit the vulnerabilities, which could enhance understanding of the exploitation process.
- There is no any visualization or example about the result analysis(study the GPT-4 agent behavior in greater detail to understand its high success rate and why it fails when the CVE description is removed.) The paper could benefit from more detailed metrics regarding the nature of the vulnerabilities exploited and the specific methodologies employed during the exploitation process to further substantiate the results.
- The substantial performance disparity between GPT-4 and other models is significant, but the paper lacks a comprehensive discussion of the limitations of the other models tested, which would provide valuable context for understanding these performance differences.

Scores

Technical Correctness: 1. No Apparent Flaws

Scientific Contribution: 2. Provides a New Data Set For Public Use 5. Identifies an Impactful Vulnerability

Presentation: 3. Major but Fixable Flaws in Presentation

Recommended Decision: 3. Weak Reject (Can be Convinced by a Champion)

Reviewer Confidence: 2. Highly Confident

LLM Agents can Autonomously Exploit One-day Vulnerabilities

Richard Fang, Rohan Bindu, Akul Gupta, Daniel Kang

Scientific Peer Reviewer: Purva Chiniya

Summary/ Technical Overview:

Vulnerability	CVE	Date	Severity
runc	CVE-2024-21626	1/31/2024	8.6 (high)
CSRF + ACE	CVE-2024-24524	2/2/2024	8.8 (high)
Wordpress SQLi	CVE-2021-24666	9/27/2021	9.8 (critical)
Wordpress XSS-1	CVE-2023-1119-1	7/10/2023	6.1 (medium)
Wordpress XSS-2	CVE-2023-1119-2	7/10/2023	6.1 (medium)
Travel Journal XSS	CVE-2024-24041	2/1/2024	6.1 (medium)
Iris XSS	CVE-2024-25640	2/19/2024	4.6 (medium)
CSRF + privilege escalation	CVE-2024-23831	2/2/2024	7.5 (high)
alf.io key leakage	CVE-2024-25635	2/19/2024	8.8 (high)
Astrophy RCE	CVE-2023-41334	3/18/2024	8.4 (high)
Hertzbeat RCE	CVE-2023-51653	2/22/2024	9.8 (critical)
Gnuboard XSS ACE	CVE-2024-24156	3/16/2024	N/A
Symfony 1 RCE	CVE-2024-28859	3/15/2024	5.0 (medium)
Peering Manager SSTI RCE	CVE-2024-28114	3/12/2024	8.1 (high)
ACIDRain	(Warszawski & Bailis, 2017)	2017	N/A

- **Research Question:**
 - Can LLMs exploit real-world vulnerabilities autonomously rather than relying on CTFs or toy examples?
- **Contribution:**
 - They test with a small dataset of 15 public vulnerabilities for open source software each with an assigned CVE.
- **Findings:**
 - GPT-4 can autonomously exploit 87% of one-day vulnerabilities when given CVE descriptions.
 - Other models (GPT-3.5, open-source LLMs) and vulnerability scanners had 0% success rate.
 - Without CVE descriptions, GPT-4's success rate drops to 7%

Strengths:

1. **Novel contribution for evaluation dataset:**
The study moves beyond simple “Capture The Flag” competition to evaluate the agents in real world, high vulnerability.
2. **Relevant Vulnerability Choices:**
They focus on common vulnerabilities, ensuring practical relevance.
3. The research compares GPT-4 against multiple other models (GPT-3.5, various open-source LLMs) and standard vulnerability scanners (ZAP, Metasploit), highlighting GPT-4's superior performance.

Weakness:

1. **Small Dataset:** Only 15 CVEs is very small, considering 1000s are reported every year
2. **Lack of Transparency:** No public release of data, code, model outputs, or working example to substantiate their claims.
3. **Reproducibility Issues:** Many open-source vulnerabilities may be difficult to reproduce underspecified descriptions (bad documentation) in the CVEs. This could impact the study's reliability and real-world applicability. Detailed steps taken by the agent can be shown with an example.
4. **Maybe a different analysis:**
The performance of GPT4 can be due to reasoning capabilities and to navigate web search and seamlessly joining existing content and code snippets. The agent built by the researchers has a web search capability which means it is capable of retrieving technical information about these CVE's from the internet. **The majority of the public exploits for these CVE's are simple and no more complex than just a few lines of code.** A step by step difference of failure and success cases can help with this.
5. **Overstatement of term "Autonomy":** The high success rate (87%) was achieved when GPT-4 was provided with CVE descriptions, otherwise GPT-4 gives 7% success rate.

Scores:

=== Technical Correctness

3. Fixable Major Issues

=== Scientific Contribution

2. Provides a New Data Set For Public Use
6. Provides a Valuable Step Forward in an Established Field

=== Presentation

3. Major Flaws in Presentation

=== Recommended Decision

3. Weak Reject (Can be Convinced by a Champion)

=== Reviewer Confidence

2. Highly Confident

Role: Archaeologist

Yvonne Zhou

1. **Previous:** LLM Agents can Autonomously Hack Websites
2. **Current:** LLM Agents can Autonomously Exploit One-day Vulnerabilities
3. **Subsequent:** Teams of LLM Agents can Exploit Zero-Day Vulnerabilities

Three papers are written by same authors in different time



Comparison I: Scopes



LLM Agents can Autonomously Hack Websites

Target: **toy websites** and solve CTF challenges e.g. SQL injection, cross-site scripting, and password cracking.

LLM Agents can Autonomously Exploit One-day Vulnerabilities

Target: **real-world**, one-day vulnerabilities, which are documents **in CVEs**, eg. websites, container management software, and other software

Teams of LLM Agents can Exploit Zero-Day Vulnerabilities

Target: newly discovered zero-day vulnerabilities and **not yet listed** in CVE databases

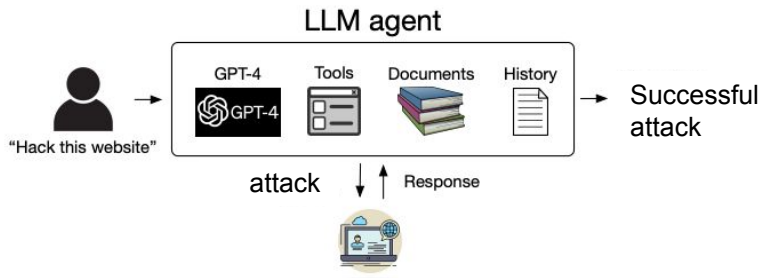
Scope:

Comparison II: Frameworks

**LLM Agents can
Autonomously Hack
Websites**

Framework

ReAct (Reason + Act)



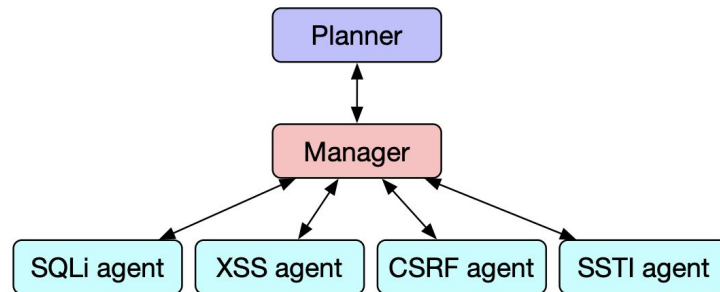
Allows the LLMs to interact with the environment, take actions, and reason about the outcomes.

**LLM Agents can
Autonomously Exploit
One-day Vulnerabilities**

ReAct (Reason + Act)

**Teams of LLM Agents
can Exploit Zero-Day
Vulnerabilities**

HPTSA (Hierarchical planning and task-specific agents)



Include three major components: a hierarchical planner, a set of task-specific, expert agents, and a team manager for the task-specific agents.

Comparison III: Contributions



LLM Agents can Autonomously Hack Websites

LLM Agents can Autonomously Exploit One-day Vulnerabilities

Teams of LLM Agents can Exploit Zero-Day Vulnerabilities

Contribution

- Demonstrated potential for LLMs to perform hacking tasks autonomously in controlled settings.
 - Limited to simplified environments, not reflecting the complexities of real-world scenarios
- Showed that LLMs can effectively exploit real-world vulnerabilities with specific CVE guidance
 - Heavily relies on having access to detailed vulnerability descriptions
- Proposed an optimized framework (HPTSA) for handling recent vulnerabilities without given detail description

LLM Agents can autonomously exploit One-day Vulnerabilities

Academic Researcher

Ayushi Mishra

Summary

- Focus on one-day vulnerabilities, which exist in system after public disclosure but before they are patched.
- Created a benchmark of 15 one-day vulnerabilities, most of which are categorized as high or critical severity based on their CVE scores.
- CVEs include well-known security flaws such as **SQL injections**, **cross-site scripting (XSS)**, and **remote code execution (RCE)**.
- Used React agent framework that allows to autonomously take actions, use tools, and respond to feedback.
- The agent was given access to various tools, including web browsing elements, terminal for running commands, web search, file creation and code interpretation.
- With the access to the tools and CVE description, the GPT-4 agent successfully exploited 87% of the vulnerabilities.



Follow-up Idea



LLM-Aided Patch Generation for One-Day Vulnerabilities



- Objective: Develop a system where LLMs like GPT-4, are used to autonomously generate patches (code fixes) for one day vulnerabilities, with a focus on speeding up the patching process.
- Can LLM generate effective, secure and generalizable patches for one-day vulnerabilities based on publicly available CVE descriptions and software context?



Literature Survey



Vulnerabilities and Security Patches Detection in OSS: A Survey

RUYAN LIN, School of Cyber Engineering, Xidian University, Xian, China

YULONG FU, School of Cyber Engineering and the State Key Laboratory of ISN, Xidian University, Xian, China

WEI YI, School of Cyber Engineering, Xidian University, Xian, China

JINCHENG YANG, School of Cyber Engineering, Xidian University, Xian, China

JIN CAO, School of Cyber Engineering and the State Key Laboratory of ISN, Xidian University, Xian, China

ZHIQIANG DONG, Tencent, Shenzhen, China

FEI XIE, Tencent, Beijing, China

HUI LI, School of Cyber Engineering and the State Key Laboratory of ISN, Xidian University, Xian, China

An Empirical Evaluation of LLMs for Solving Offensive Security Challenges

Minghao Shao*
New York University

Boyuan Chen*
New York University

Sofija Jancheska*
New York University

Brendan Dolan-Gavitt*
New York University

Siddharth Garg
New York University

Ramesh Karri
New York University

Muhammad Shafiqe
New York University Abu Dhabi

Abstract

Capture The Flag (CTF) challenges are puzzles related to computer security scenarios. With the advent of large language models (LLMs), more and more CTF participants are using LLMs to understand and solve the challenges. However, so far no work has evaluated the effectiveness of LLMs in solving CTF challenges with a fully automated workflow. We develop two CTF-solving workflows, human-in-the-loop (HITL) and fully-automated, to examine the LLMs' ability to solve a selected set of CTF challenges, prompted with information about the question. We collect human contestants' results on the same set of questions, and

to capture and print hidden 'flags,' which are short strings of characters or specific files, proving successful completion of a challenge. Solving CTF challenges requires an understanding of cybersecurity concepts and creative problem solving skills. Consequently, CTF has garnered attention as a prominent approach in cybersecurity education [16].

This work explores and evaluates the ability of LLMs to solve CTF challenges. As part of our study, we organized the LLM Attack challenge [24] as a part of the Cybersecurity Awareness Week (CSAW) [23] at New York University (NYU), in which participants competed in designing "prompts" that enable LLMs to solve a collection

Over the past decade, Open Source Software (OSS) has experienced rapid growth and widespread adoption, attributed to its openness and editability. However, this expansion has also brought significant security challenges, particularly introducing and propagating software vulnerabilities. Despite the use of machine learning and formal methods to tackle these issues, there remains a notable gap in comprehensive surveys that summarize and analyze both Vulnerability Detection (VD) and Security Patch Detection (SPD) in OSS.

PATUNTRACK: Automated Generating Patch Examples for Issue Reports without Tracked Insecure Code

Ziyou Jiang^{1,2,3}, Lin Shi⁴, Guowei Yang⁵, Qing Wang^{1,2,3*}

¹State Key Laboratory of Intelligent Game, Beijing, China;

²Science and Technology on Integrated Information System Laboratory, Institute of Software Chinese Academy of Sciences, Beijing, China;

³University of Chinese Academy of Sciences, Beijing, China;

⁴School of Software, Beihang University, Beijing, China;

⁵The University of Queensland, Brisbane, Australia

{ziyou2019, wq}@iscas.ac.cn, shilin@buaa.edu.cn, guowei.yang@uq.edu.au

ABSTRACT

Security patches are essential for enhancing the stability and robustness of projects in the open-source software community. While vulnerabilities are officially expected to be patched before being disclosed, patching vulnerabilities is complicated and remains a struggle for many organizations. To patch vulnerabilities, security practitioners typically track vulnerable issue reports (IRs), and analyze their relevant insecure code to generate potential patches. However, the relevant insecure code may not be explicitly specified and practitioners cannot track the insecure code in the repositories, thus limiting their ability to generate patches. In such cases, pro-

the patch examples generated by PATUNTRACK, indicating that they can benefit from these examples for patching the vulnerabilities.

ACM Reference Format:

Ziyou Jiang^{1,2,3}, Lin Shi⁴, Guowei Yang⁵, Qing Wang^{1,2,3}. 2018. PATUNTRACK: Automated Generating Patch Examples for Issue Reports without Tracked Insecure Code. In *Proceedings of Make sure to enter the correct conference title from your rights confirmation email (Conference acronym 'XX)*. ACM, New York, NY, USA, 13 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 INTRODUCTION

Security patches are essential for enhancing the stability and robust-

Model Architecture



User Input

CVE Descriptions



Vulnerable Code Snippet



GPT - 4

Analyze vulnerability and generate patch



Patch Generation



Test



Generated Patch

Fine-Tuning LLM Model

- Fine-Tune the model on the dataset of vulnerabilities and patches.
- Dataset could include:
 - CVE descriptions
 - The vulnerable code
 - Corresponding Patches

Prompt Engineering: Develop prompts that guide the LLM to generate code patches instead of exploits.

Example: Here is a CVE description. The affected code is located in function X of file Y. Provide a patch to fix this vulnerability.

Patch Generation Workflow

The input to the LLM includes a CVE description and a snippet of vulnerable code.

1. CVE description: "The CVE represents an SQL Injection vulnerability in function `validateUser()` of a web application."
2. A snippet of vulnerable code

```
def validateUser(username):  
    query = "SELECT * FROM users  
    WHERE username = '" + username + "'" +  
    cursor.execute(query)  
    return cursor.fetchall()
```

Integration with Tools and Continuous Learning

- **Static Code Analysis:** After the patch is generated, a static code analysis tool could automatically review the patch to ensure it doesn't introduce new vulnerabilities or break the code.
- **Automated Testing:** The patch could be automatically tested in a controlled environment (e.g., CI/CD pipelines) to check if it resolves the issue without introducing new problems.
- Every generated patch could be stored in a **patch database**. Over time, the LLM can learn from feedback (whether the patch was accepted, modified, or rejected), improving its patch generation capabilities.


```
def validateUser(username):  
    query = "SELECT * FROM users WHERE username = %s"  
    cursor.execute(query, (username,))  
    return cursor.fetchall()
```

The generated patch uses parametrized queries, which prevent SQL injection by safely handling user input.

Academic Researcher

Ruibo Chen

- **Current findings:**

Only GPT-4 can successfully exploit the vulnerabilities, and the overall success rate is only 40%.

Can we build a more effective agent?

Model	Pass @ 5	Overall success rate
GPT-4	86.7%	40.0%
GPT-3.5	0%	0%
OpenHermes-2.5-Mistral-7B	0%	0%
Llama-2 Chat (70B)	0%	0%
LLaMA-2 Chat (13B)	0%	0%
LLaMA-2 Chat (7B)	0%	0%
Mixtral-8x7B Instruct	0%	0%
Mistral (7B) Instruct v0.2	0%	0%
Nous Hermes-2 Yi 34B	0%	0%
OpenChat 3.5	0%	0%

- **Current findings:**

Only GPT-4 can successfully exploit the vulnerabilities, and the overall success rate is only 40%.

Can we build a more effective agent?

- **Possible solutions:**

1) Divide the whole large task into some subtasks (by human or by LLMs)

- **Current findings:**

Only GPT-4 can successfully exploit the vulnerabilities, and the overall success rate is only 40%.

Can we build a more effective agent?

- **Possible solutions:**

- 1) Divide the whole large task into some subtasks (by human or by LLMs)

- 2) Query GPT-4 to collect a dataset to train the agents on other LLMs (distillate the knowledge of GPT-4)

- **Current findings:**

It is still a difficult task for LLMs to discover these vulnerabilities:

After removing the CVE description, the success rate falls from 87% to 7%. This suggests that determining the vulnerability is extremely challenging.

An LLM detector can be used to improve the attack success rate, and can also be used for protection

- **Current findings:**

It is still a difficult task for LLMs to discover these vulnerabilities:

After removing the CVE description, the success rate falls from 87% to 7%. This suggests that determining the vulnerability is extremely challenging.

An LLM detector can be used to improve the attack success rate, and can also be used for protection

- **Possible solutions:**

LLMs may never learn this specific task in the training set;

But LLMs learned a lot of knowledge during pretraining.

We can curate a high-quality dataset from known vulnerabilities to finetune the LLMs to enable them to automatically detect those vulnerabilities

Industry Practitioner

Sakshi

Autonomous Vulnerability Exploitation

- Help organizations identify and patch potential threats faster than traditional methods
- Improving security measures
- Cheaper than human-led efforts
- Easily scaled across multiple systems, providing widespread security assessments simultaneously

Positive Impact

- **Proactive Cybersecurity:** Security teams can stay ahead of malicious actors by detecting vulnerabilities immediately after they are published.
- **Continuous Monitoring:** 24/7 vulnerability scanning, ensuring that emerging threats are quickly addressed.
- **Faster Remediation Cycles:** Quicker remediation, reducing the potential damage from exploits.

Negative Impact

- **Ethical Concerns:** If misused, they could enable malicious actors to automate cyberattacks.
- **Data and Privacy Risks:** Exposure of sensitive information during the scanning process.

Private Investigator

Ji-Ze Jang

Education



2020-2024 : CS B.S. with Honors @ UIUC

2024-present : CS Masters @ UIUC

Experience



Jun-Aug 2022 : SWE @ Openseense



May 2023-Aug 2024 : SWE @ EdgeBit



Aug 2023-present : Research Assistant @ UIUC



Oct 2023-present : Red Team Engineer @ OpenAI



Akul Gupta

arXiv > cs > arXiv:2404.08144v2

Computer Science > Cryptography and Security

[Submitted on 11 Apr 2024 (v1), last revised 17 Apr 2024 (this version, v2)]

LLM Agents can Autonomously Exploit One-day Vulnerabilities

Richard Fang, Rohan Bindu, Akul Gupta, Daniel Kang



Social Impact Assessor

Amadeo De La Vega



(Self-assessed Positive) Impact



(Self-assessed Positive) Impact

- Raises awareness of LLMs agents capabilities and risks.



(Self-assessed Positive) Impact

- Raises awareness of LLMs agents capabilities and risks.
(showed GPT-4 alarming hacking capabilities compared to other models: emergent capability?)



(Self-assessed Positive) Impact

- Raises awareness of LLMs agents capabilities and risks.
(showed GPT-4 alarming hacking capabilities compared to other models: emergent capability?)
- Promotes ethical AI research in cybersecurity.



(More Positive) Impact



(More Positive) Impact

- Could improve cybersecurity defenses.



(More Positive) Impact

- Could improve cybersecurity defenses.
(researchers and cybersecurity professionals are encouraged to develop LLM agents for patching vulnerabilities)



(More Positive) Impact

- Could improve cybersecurity defenses.
(researchers and cybersecurity professionals are encouraged to develop LLM agents for patching vulnerabilities)
- Could help lowering costs of penetration testing.



(More Positive) Impact

- Could improve cybersecurity defenses.
(researchers and cybersecurity professionals are encouraged to develop LLM agents for patching vulnerabilities)
- Could help lowering costs of penetration testing.
(E.g., a small organization could deploy an LLM agent to perform 1-day vulnerability scans, detecting threats more cheaply than hiring a full-time penetration testing team)



(More Positive) Impact

- Could improve cybersecurity defenses.
(researchers and cybersecurity professionals are encouraged to develop LLM agents for patching vulnerabilities)
- Could help lowering costs of penetration testing.
(E.g., a small organization could deploy an LLM agent to perform 1-day vulnerability scans, detecting threats more cheaply than hiring a full-time penetration testing team)
- Could help cybersecurity professionals (training and simulations).



(More Positive) Impact

- Could improve cybersecurity defenses.
(researchers and cybersecurity professionals are encouraged to develop LLM agents for patching vulnerabilities)
- Could help lowering costs of penetration testing.
(E.g., a small organization could deploy an LLM agent to perform 1-day vulnerability scans, detecting threats more cheaply than hiring a full-time penetration testing team)
- Could help cybersecurity professionals (training and simulations).
(LLM agents could be used in training exercises to simulate real-world cyberattacks, helping cybersecurity professionals practice responding to realistic scenarios.)



(Negative) Impact



(Negative) Impact

- Exploitation by malicious hackers.



(Negative) Impact

- Exploitation by malicious hackers.
(E.g., Cybercriminals could deploy LLM agents to scan and exploit known vulnerabilities in thousands of websites or systems without requiring technical expertise, leading to widespread cyberattacks.)



(Negative) Impact

- Exploitation by malicious hackers.
(E.g., Cybercriminals could deploy LLM agents to scan and exploit known vulnerabilities in thousands of websites or systems without requiring technical expertise, leading to widespread cyberattacks.)
- Job displacement in cybersecurity roles.



(Negative) Impact

- Exploitation by malicious hackers.
(E.g., Cybercriminals could deploy LLM agents to scan and exploit known vulnerabilities in thousands of websites or systems without requiring technical expertise, leading to widespread cyberattacks.)
- Job displacement in cybersecurity roles.
(Companies might replace human penetration testers with cheaper, scalable AI-driven agents, reducing the demand for highly skilled cybersecurity professionals.)



(Negative) Impact

- Exploitation by malicious hackers.
(E.g., Cybercriminals could deploy LLM agents to scan and exploit known vulnerabilities in thousands of websites or systems without requiring technical expertise, leading to widespread cyberattacks.)
- Job displacement in cybersecurity roles.
(Companies might replace human penetration testers with cheaper, scalable AI-driven agents, reducing the demand for highly skilled cybersecurity professionals.)
- Erosion of trust in AI systems.



(Negative) Impact

- Exploitation by malicious hackers.
(E.g., Cybercriminals could deploy LLM agents to scan and exploit known vulnerabilities in thousands of websites or systems without requiring technical expertise, leading to widespread cyberattacks.)
- Job displacement in cybersecurity roles.
(Companies might replace human penetration testers with cheaper, scalable AI-driven agents, reducing the demand for highly skilled cybersecurity professionals.)
- Erosion of trust in AI systems.
(If the public becomes aware of how AI systems like GPT-4 can autonomously exploit vulnerabilities, trust in AI technologies, even in benign or helpful roles, may erode.)



Thanks!

Social Impact Assessor

Dinithi Wickramaratne

Self-assessed positive outcomes

- Raising awareness
 - Highlight the need for the wider cybersecurity community and LLM providers to **think carefully about how to integrate LLM agents** in defensive measures and about their widespread deployment.
- GPT-4 supremacy
 - Highlights the **capabilities of GPT-4** for vulnerability exploitation by comparing with other LLMs and opportunities for future advancements

Additional Positive Outcomes

- Policy Development
 - By demonstrating the risks associated with LLMs in cybersecurity, the research could inform **policy discussions, pushing for stronger regulations** or standards for AI deployment in sensitive areas.
- Educational use
 - Equipping future cybersecurity professionals with **knowledge about the emerging capabilities** of LLMs in exploitation.

Negative Outcomes

- **Moral Hazard in Security Practices**
 - If organizations resort to use of LLMs to detect or counter threats, **over reliance on automated systems** could lead organizations to neglect traditional, foundational security measures such as employee training, proper configuration management, and manual code review, creating a false sense of security.
- **Job Losses in Traditional Cybersecurity Roles**
 - It's mentioned that using an AI agent x2.8 cheaper than using a human expert for vulnerability exploitation. With the agent costs reducing, The automation of certain cybersecurity tasks could **threaten the demand for human experts** in vulnerability assessment.
- **Unintended Spread of Tools**
 - Although the prompts used were not released publicly, the paper discusses the steps and challenges GPT4 had to face to exploit the vulnerabilities. These might deem to be useful for hackers developing more sophisticated exploitation tools.