

# PENTESTGPT: Evaluating and Harnessing Large Language Models for Automated Penetration Testing

Gelei Deng<sup>1§</sup>, Yi Liu<sup>1§</sup>, Víctor Mayoral-Vilches<sup>23</sup>, Peng Liu<sup>4</sup>, Yuekang Li<sup>5\*</sup>, Yuan Xu<sup>1</sup>,  
Tianwei Zhang<sup>1</sup>, Yang Liu<sup>1</sup>, Martin Pinzger<sup>3</sup>, Stefan Rass<sup>6</sup>

<sup>1</sup>Nanyang Technological University, <sup>2</sup>Alias Robotics, <sup>3</sup>Alpen-Adria-Universität Klagenfurt,  
<sup>4</sup>Institute for Infocomm Research (*I<sup>2</sup>R*), A\*STAR, Singapore, <sup>5</sup>University of New South Wales, <sup>6</sup>Johannes  
Kepler University Linz

## Abstract

Penetration testing, a crucial industrial practice for ensuring system security, has traditionally resisted automation due to the extensive expertise required by human professionals. Large Language Models (LLMs) have shown significant advancements in various domains, and their emergent abilities suggest their potential to revolutionize industries. In this work, we establish a comprehensive benchmark using real-world penetration testing targets and further use it to explore the capabilities of LLMs in this domain. Our findings reveal that while LLMs demonstrate proficiency in specific sub-tasks within the penetration testing process, such as using testing tools, interpreting outputs, and proposing subsequent actions, they also encounter difficulties maintaining a whole context of the overall testing scenario.

Based on these insights, we introduce PENTESTGPT, an LLM-empowered automated penetration testing framework that leverages the abundant domain knowledge inherent in LLMs. PENTESTGPT is meticulously designed with three self-interacting modules, each addressing individual sub-tasks of penetration testing, to mitigate the challenges related to context loss. Our evaluation shows that PENTESTGPT not only outperforms LLMs with a task-completion increase of 228.6% compared to the GPT-3.5 model among the benchmark targets, but also proves effective in tackling real-world penetration testing targets and CTF challenges. Having been open-sourced on GitHub, PENTESTGPT has garnered over 6,500 stars in 12 months and fostered active community engagement, attesting to its value and impact in both the academic and industrial spheres.

## 1 Introduction

Securing a system presents a formidable challenge. Offensive security methods like penetration testing (pen-testing) and

red teaming are now essential in the security lifecycle. As explained by Applebaum [1], these approaches involve security teams attempting breaches to reveal vulnerabilities, providing advantages over traditional defenses, which rely on incomplete system knowledge and modeling. This study, guided by the principle “*the best defense is a good offense*”, focuses on offensive strategies, specifically penetration testing.

Penetration testing is a proactive offensive technique for identifying, assessing, and mitigating security vulnerabilities [2]. It involves targeted attacks to confirm flaws, yielding a comprehensive inventory of vulnerabilities with actionable recommendations. This widely-used practice empowers organizations to detect and neutralize network and system vulnerabilities before malicious exploitation. However, it typically relies on manual effort and specialized knowledge [3], resulting in a labor-intensive process, creating a gap in meeting the growing demand for efficient security evaluations.

Large Language Models (LLMs) have demonstrated profound capabilities, showcasing intricate comprehension of human-like text and achieving remarkable results across a multitude of tasks [4, 5]. An outstanding characteristic of LLMs is their emergent abilities [6], cultivated during training, which empower them to undertake intricate tasks such as reasoning, summarization, and domain-specific problem-solving without task-specific fine-tuning. This versatility posits LLMs as potential game-changers in various fields, notably cybersecurity. Although recent works [7–9] posit the potential of LLMs to reshape cybersecurity practices, including the context of penetration testing, there is an absence of a systematic, quantitative assessment of their aptitude in this regard. Consequently, an imperative question presents: To what extent can LLMs automate penetration testing?

Motivated by this question, we set out to explore the capability boundary of LLMs on real-world penetration testing tasks. Unfortunately, the current benchmarks for penetration testing [10, 11] are not comprehensive and fail to assess progressive accomplishments fairly during the process. To address this limitation, we construct a robust benchmark that includes test machines from HackTheBox [12] and

---

\*Corresponding author.

§Equal Contribution

VulnHub [13]—two leading platforms for penetration testing challenges. Comprising 13 targets with 182 sub-tasks, our benchmark encompasses all vulnerabilities appearing in OWASP’s top 10 vulnerability list [14] and 18 Common Weakness Enumeration (CWE) items [15]. The benchmark offers a more detailed evaluation of the tester’s performance by monitoring the completion status for each sub-task.

With this benchmark, we perform an exploratory study using GPT-3.5 [16], GPT-4 [17], and Bard [18] as representative LLMs. Our test strategy is interactive and iterative. We craft tailored prompts to guide the LLMs through penetration testing. Each LLM, presented with prompts and target machine information, generates step-by-step penetration testing operations. We then execute the suggested operations in a controlled environment, document the results, and feed them back to the LLM to inform and refine its next steps. This cycle (prompting, executing, and feedback) is repeated until the LLM completes the entire penetration testing process autonomously. To evaluate LLMs, we compare their results against baseline solutions from official walkthroughs and certified penetration testers. By analyzing similarities and differences in their problem-solving approaches, we aim to better understand LLMs’ capabilities in penetration testing and how their strategies differ from human experts.

Our investigation yields intriguing insights into the capabilities and limitations of LLMs in penetration testing. We discover that LLMs demonstrate proficiency in managing specific sub-tasks within the testing process, such as utilizing testing tools, interpreting their outputs, and suggesting subsequent actions. Compared to human experts, LLMs are especially adept at executing complex commands and options with testing tools, while models like GPT-4 excel in comprehending source code and pinpointing vulnerabilities. Furthermore, LLMs can craft appropriate test commands and accurately describe graphical user-interface operations needed for specific tasks. Leveraging their vast knowledge base, they can design inventive testing procedures to unveil potential vulnerabilities in real-world systems and CTF challenges. However, we also note that LLMs have difficulty in maintaining a coherent grasp of the overarching testing scenario, a vital aspect for attaining the testing goal. As the dialogue advances, they may lose sight of earlier discoveries and struggle to apply their reasoning consistently toward the final objective. Additionally, LLMs overemphasize recent tasks in the conversation history, regardless of their vulnerability status. As a result, they tend to neglect other potential attack surfaces exposed in prior tests and fail to complete the penetration testing task.

Building on our insights into LLMs’ capabilities in penetration testing, we present PENTESTGPT, an interactive system designed to enhance the application of LLMs in this domain. Drawing inspiration from the collaborative dynamics commonly observed in real-world human penetration testing teams, PENTESTGPT is particularly tailored to manage large and intricate projects. It features a tripartite architecture com-

prising *Reasoning, Generation, and Parsing Modules*, each reflecting specific roles within penetration testing teams. The Reasoning Module emulates the function of a lead tester, focusing on maintaining a high-level overview of the penetration testing status. We introduce a novel representation, the Pentesting Task Tree (PTT), based on the cybersecurity attack tree [19]. This structure encodes the testing process’s ongoing status and steers subsequent actions. Uniquely, this representation can be translated into natural language and interpreted by the LLM, thereby comprehended by the Generation Module and directing the testing procedure. The Generation Module, mirroring a junior tester’s role, is responsible for constructing detailed procedures for specific sub-tasks. Translating these into exact testing operations augments the generation process’s accuracy. Meanwhile, the Parsing Module deals with diverse text data encountered during penetration testing, such as tool outputs, source codes, and HTTP web pages. It condenses and emphasizes these texts, extracting essential information. Collectively, these modules function as an integrated system. PENTESTGPT completes complex penetration testing tasks by bridging high-level strategies with precise execution and intelligent data interpretation, thereby maintaining a coherent and effective testing process.

We assessed PENTESTGPT across diverse testing scenarios to validate its effectiveness and breadth. In our custom benchmarks, PENTESTGPT significantly outperformed direct applications of GPT-3.5 and GPT-4, showing increases in sub-task completion rates of 228.6% and 58.6%, respectively. Furthermore, when applied to real-world challenges such as the HackTheBox active machine penetration tests [20] and picoMini [21] CTF competition, PENTESTGPT demonstrated its practical utility. It successfully resolved 4 out of 10 penetration testing challenges, incurring a total cost of 131.5 US Dollars for the OpenAI API usage. In the CTF competition, PENTESTGPT achieved a score of 1500 out of a possible 4200, placing 24th among 248 participating teams. This evaluation underscores PENTESTGPT’s practical value in enhancing penetration testing tasks’ efficiency and precision. The solution has been made publicly available on GitHub [22], receiving widespread acclaim with over 6,200 stars to the date of writing, active community engagement, and ongoing collaboration with multiple industrial partners.

In summary, we make the following contributions:

- **Development of a Comprehensive Penetration Testing Benchmark.** We craft a robust and representative penetration testing benchmark, encompassing a multitude of test machines from leading platforms such as HackTheBox and VulnHub. This benchmark includes 182 sub-tasks covering OWASP’s top 10 vulnerabilities, offering fair and comprehensive evaluation of penetration testing. To the best of our knowledge, this is the first benchmark in the field that can provide progressive accomplishments assessments and comparisons.

- **Comprehensive Evaluation of LLMs for Penetration Testing Tasks.** By employing models like GPT-3.5, GPT-4, and Bard, our exploratory study rigorously investigates the strengths and limitations of LLMs in penetration testing. To the best of our knowledge, this is the first systematic and quantitative study for the capability of LLMs in performing automated penetration testing. The insights gleaned from this study shed valuable light on the capabilities and challenges faced by LLMs, enriching our understanding of their applicability in this specialized domain.
- **Development of an Innovative LLM-powered Penetration Testing System.** We engineer PENTESTGPT, a novel interactive system that leverages the strengths of LLMs to carry out penetration testing tasks automatically. Drawing inspiration from real-world human penetration testing teams, PENTESTGPT integrates a tripartite design that mirrors the collaborative dynamics between senior and junior testers. This architecture optimizes LLMs’ usage, significantly enhancing the efficiency and effectiveness of automated penetration testing. We have open-sourced PENTESTGPT and it has received over 6,500 stars on GitHub, active community contributions, and several industry partners including AWS, Huawei, and TikTok to collaborate.

## 2 Background & Related Work

### 2.1 Penetration Testing

Penetration testing, or “pentesting”, is a critical practice to enhance organizational systems’ security. In a typical penetration test, security professionals, known as penetration testers, analyze the target system, often leveraging automated tools. The standard process is divided into five key phases [23]: Reconnaissance, Scanning, Vulnerability Assessment, Exploitation, and Post Exploitation (including reporting). These phases enable testers to understand the target system, identify vulnerabilities, and exploit them to gain access.

Despite significant advancements [11, 24, 25], a fully automated penetration testing system remains out of reach. This gap results from the need for deep vulnerability understanding and a strategic action plan. Typically, testers combine depth-first and breadth-first search techniques [23]. They first grasp the target environment’s scope, then drill down into specific vulnerabilities. This method ensures comprehensive analysis, leaning on expertise and experience. The multitude of specialized tools further complicate the automation. Thus, even with artificial intelligence, achieving a seamless automated penetration testing solution is a daunting task.

### 2.2 Large Language Models

Large Language Models (LLMs), including OpenAI’s GPT-3.5 and GPT-4, are prominent tools with applications extending to various cybersecurity-related fields, such as code

analysis [26] and vulnerability repairment [27]. These models are equipped with wide-ranging general knowledge and the capacity for elementary reasoning. They can comprehend, infer, and produce text resembling human communication, aided by a training corpus encompassing diverse domains like computer science and cybersecurity. Their ability to interpret context and recognize patterns enables them to adapt knowledge to new scenarios. This adaptability, coupled with their proficiency in interacting with systems in a human-like way, positions them as valuable assets in enhancing penetration testing processes. Despite inherent limitations, LLMs offer distinct attributes that can substantially aid in the automation and improvement of penetration testing tasks. The realization of this potential, however, requires the creation and application of a specialized and rigorous benchmark.

## 3 Penetration Testing Benchmark

### 3.1 Motivation

The comprehensive evaluation of LLMs in penetration testing necessitates a robust and representative benchmark. Existing benchmarks in this domain [10, 11] have several limitations. First, they are often restricted in scope, focusing on a narrow range of potential vulnerabilities, and thus fail to capture the complexity and diversity of real-world cyber threats. For instance, the OWASP *juiceshop* project [28] is the most widely adopted benchmark for web vulnerability evaluation. However, it does not include privilege escalation vulnerabilities, which is an essential aspect of penetration testing. Second, existing benchmarks may not recognize the cumulative value of progress through the different stages of penetration testing, as they tend to evaluate only the final exploitation success. This approach overlooks the nuanced value each step contributes to the overall process, resulting in metrics that might not accurately represent actual performance in real-world scenarios.

To address these concerns, we propose the construction of a comprehensive penetration testing benchmark that meets the following criteria:

**Task Variety.** The benchmark must encompass diverse tasks, reflecting various operating systems and emulating the diversity of scenarios encountered in real-world penetration testing.

**Challenge Levels.** To ensure broad applicability, the benchmark must include tasks of varying difficulty levels suitable for challenging novice and expert testers.

**Progress Tracking.** Beyond mere success or failure metrics, the benchmark must facilitate tracking of incremental progress, thereby recognizing and scoring the value added at each stage of the penetration testing process.

## 3.2 Benchmark Design

Following the criteria outlined previously, we develop a comprehensive benchmark that closely reflects real-world penetration testing tasks. The design process progresses through several stages.

**Task Selection.** We begin by selecting tasks from HackTheBox [12] and VulnHub [13], two leading penetration testing training platforms. Our selection criteria are designed to ensure that our benchmark accurately reflects the challenges encountered in practical penetration testing environments. We meticulously review the latest machines available on both platforms, aiming to identify and select a subset that comprehensively covers all vulnerabilities listed in the OWASP [14] Top 10 Project. Additionally, we choose machines that represent a mix of difficulties, classified according to traditional standards in the penetration testing domain into *easy*, *medium*, and *hard* categories. This process guarantees that our benchmark spans the full spectrum of vulnerabilities and difficulties. Note that our benchmark does not include benign targets to assess false positives. In penetration testing, benign targets are sometimes explored. Our main objective remains identifying true vulnerabilities.

**Task Decomposition.** We further parse the testing process of each target into a series of sub-tasks, following the standard solution commonly referred to as the “walkthrough” in penetration testing. Each sub-task corresponds to a unique step in the overall process. We decompose sub-tasks following NIST 800-115 [29], the Technical Guide to Security Testing. Each sub-task is one step declared in the Guide (e.g., network discovery, password cracking), or an operation that exploits a unique vulnerability categorised in the Common Weakness Enumeration (CWE) [15] (e.g., exploiting SQL injection - CWE-89 [30]). In the end, we formulate an exhaustive list of sub-tasks for every benchmark target.

**Benchmark Validation.** The final stage of our benchmark development involves rigorous validation, which ensures the reproducibility of these benchmark machines. To do this, three certified penetration testers independently attempt the penetration testing targets and write their walkthrough. We then adjust our task decomposition accordingly because some targets may have multiple valid solutions.

Ultimately, we have compiled a benchmark that effectively encompasses all types of vulnerabilities listed in the OWASP [14] Top 10 Project. It comprises 13 penetration testing targets, each at varying levels of difficulty. These targets are broken down into 182 sub-tasks across 26 categories, covering 18 distinct CWE items. This number of targets is deemed sufficient to represent a broad spectrum of vulnerabilities, difficulty levels, and varieties essential for comprehensive penetration testing training. To foster community development, we have made this benchmark publicly available online [22].

## 4 Exploratory Study

We conduct an exploratory study to assess the capabilities of LLMs in penetration testing, with the primary objective of determining how well LLMs can adapt to the real-world complexities and challenges in this task. Specifically, we aim to address the following two research questions:

**RQ1 (Capability):** To what extent can LLMs perform penetration testing tasks?

**RQ2 (Comparative Analysis):** How do the problem-solving strategies of human penetration testers and LLMs differ?

We utilize the benchmark described in Section 3 to evaluate the performance of LLMs on penetration testing tasks. In the following, we first delineate our testing strategy for this study. Subsequently, we present the testing results and an analytical discussion to address the above research questions.

### 4.1 Testing Strategy

LLMs are text-based and cannot independently perform penetration testing operations. To address this, we develop a human-in-the-loop testing strategy, serving as an intermediary method to accurately assess LLMs’ capabilities. This strategy features an interactive loop where a human expert executes the LLM’s penetration testing directives. Importantly, the human expert functions purely as an executor, strictly following the LLM’s instructions without adding any expert insights or making independent decisions.

Figure 1 depicts the testing strategy with the following steps: ❶ We initiate the looped testing procedure by presenting the target specifics to the LLM, seeking its guidance on potential penetration testing steps. ❷ The human expert strictly follows the LLM’s recommendations and conducts the suggested actions in the penetration testing environment. ❸ Outcomes of the testing actions are collected and summarized: direct text outputs such as terminal outputs or source code are documented; non-textual results, such as graphical representations, are translated by the human expert into succinct textual summaries. The data is then fed back to the LLM, setting the stage for its subsequent recommendations. ❹ This iterative process persists either until a conclusive solution is identified or a deadlock is reached. We then compile a record of the testing procedures, encompassing successful sub-tasks, ineffective actions, and any reasons for failure, if applicable. For a more tangible grasp of this strategy, we offer illustrative examples of prompts and corresponding outputs from GPT-4 related to one of our benchmark targets in the Appendix Section A.

To ensure the evaluation’s fairness and accuracy, we employ several strategies. First, we involve expert-level penetration testers<sup>1</sup> as the human testers. With their deep pentesting knowledge, these testers can precisely comprehend and execute LLM-generated operations, thus accurately assessing

<sup>1</sup>We selected Offensive Security Certified Professionals (OSCP) testers.

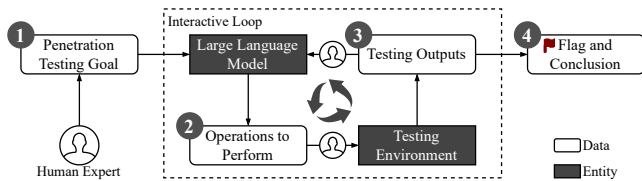


Figure 1: Overview of strategy to use LLMs for penetration testing.

LLMs’ true capabilities. Second, we instruct the penetration testers to strictly execute the commands given by the LLMs, without altering any content or information, even upon identifying clear errors. They are also instructed to faithfully report the testing results back to the LLM without any additional commentary. Third, for managing UI-based operations and graphical results, we have adopted specific measures. Initially, we instruct the LLMs to minimize the use of GUI-based tools. For indispensable tools that cannot be avoided (e.g., Burp-Suite), we propose a result-oriented approach: upon receiving a GUI operation instruction, the testers first execute the operation based on their expert knowledge. Subsequently, they are required to provide detailed, step-by-step textual descriptions of their actions and the observed responses at each step, which are then communicated back to the LLM. Should the LLM express any objections or comments concerning a particular step, the operation is to be repeated. This protocol ensures the integrity of the feedback loop, guaranteeing that the LLM obtains a comprehensive understanding of the testing results.

## 4.2 Evaluation Settings

We proceed to assess the performances of various LLMs in penetration testing tasks using the strategy mentioned above. **Model Selection.** Our study focuses on three cutting-edge LLMs that are currently accessible: GPT-3.5 with 8k token limit, GPT-4 with 32k token limit from OpenAI, and LaMDA [31] from Google. These models are selected based on their prominence in the research community and consistent availability. To interact with the LLMs mentioned above, we utilize chatbot services provided by OpenAI and Google, namely ChatGPT [32] and Bard [18]. For this paper, the terms GPT-3.5, GPT-4, and Bard will represent these three LLMs. **Experimental Setup.** Our experiments occur in a local setting with both target and testing machines on the same private network. The testing machine runs on Kali Linux [33], version 2023.1.

**Tool Usage.** Our study aims to assess the innate capabilities of LLMs on penetration testing, without reliance on end-to-end automated vulnerability scanners such as Nexus [34] and OpenVAS [35]. Consequently, we explicitly instruct the LLMs to refrain from using these tools. We follow the LLMs’ recommendations for utilizing other tools designed to validate specific vulnerability types (e.g., sqlmap [36] for SQL injection).

Occasionally, versioning discrepancies may lead the LLMs to provide incorrect instructions for tool usage. In such instances, our penetration testing experts evaluate whether the instructions would have been valid for a previous version of the tool. They then make any necessary adjustments to ensure the tool’s correct operation.

## 4.3 Capability Evaluation (RQ1)

To address **RQ1**, we evaluate the performance of three leading LLMs: GPT-4, Bard, and GPT-3.5. We summarize these findings in Table 1. Each LLM successfully completes at least one end-to-end penetration test, highlighting their versatility in simpler environments. Of these, GPT-4 excels, achieving success on 4 easy and 1 medium difficulty targets. Bard and GPT-3.5 follow with success on 2 and 1 easy targets, respectively. In sub-tasks, GPT-4 completes 55 out of 77 on easy targets and 30 out of 71 on medium. Bard and GPT-3.5 also show potential, finishing 16 (22.54%) and 13 (18.31%) of medium difficulty sub-tasks, respectively. However, on hard targets, all models’ performance declines. Though they can initiate the reconnaissance phase, they struggle to exploit identified vulnerabilities. This is anticipated since hard targets are designed to be especially challenging. They often feature seemingly vulnerable services that are non-exploitable, known as *rabbit holes* [37]. The pathways to exploit these machines are unique and unpredictable, resisting automated tool replication. For example, the target *Falafel* has specialized SQL injection vulnerabilities resistant to *sqlmap*. Current LLMs cannot tackle these without human expert input.

**Finding 1:** Large Language Models (LLMs) have shown proficiency in conducting end-to-end penetration testing tasks but struggle to overcome challenges presented by more difficult targets.

We further examine the detailed sub-task completion performances of the three LLMs compared to the walkthrough (WT), as presented in Table 2. Analyzing the completion status, we identify several areas where LLMs excel. First, they adeptly utilize common penetration testing tools to interpret the corresponding outputs, especially in enumeration tasks correctly. For example, all three evaluated LLMs successfully perform nine *Port Scanning* sub-tasks. They can configure the widely-used port scanning tool, *nmap* [38], comprehend the scan results, and formulate subsequent actions. Second, the LLMs reveal a deep understanding of prevalent vulnerability types, connecting them to the services on the target system. This understanding is evidenced by the successful completion of sub-tasks related to various vulnerability types. Finally, LLMs demonstrate their effectiveness in code analysis and generation, particularly in the tasks of *Code Analysis* and *Shell Construction*. These tasks require the models to read and generate codes in different programming languages.

Table 1: Overall performance of LLMs on Penetration Testing Benchmark.

Tools	Easy		Medium		Hard		Average	
	Overall (7)	Sub-task (77)	Overall (4)	Sub-task (71)	Overall (2)	Sub-task (34)	Overall (13)	Sub-task (182)
GPT-3.5	1 (14.29%)	24 (31.17%)	0 (0.00%)	13 (18.31%)	0 (0.00%)	5 (14.71%)	1 (7.69%)	42 (23.07%)
GPT-4	4 (57.14%)	55 (71.43%)	1 (25.00%)	30 (42.25%)	0 (0.00%)	10 (29.41%)	5 (38.46%)	95 (52.20%)
Bard	2 (28.57%)	29 (37.66%)	0 (0.00%)	16 (22.54%)	0 (0.00%)	5 (14.71%)	2 (15.38%)	50 (27.47%)
Average	2.3 (33.33%)	36 (46.75%)	0.33 (8.33%)	19.7 (27.70%)	0 (0.00%)	6.7 (19.61%)	2.7 (20.5%)	62.3 (34.25%)

Table 2: Top 10 Types of Sub-tasks completed by each tool.

Sub-Tasks	WT	GPT-3.5	GPT-4	Bard
Web Enumeration	18	4 (22.2%)	8 (44.4%)	4 (22.2%)
Code Analysis	18	4 (22.2%)	5 (27.2%)	4 (22.2%)
Port Scanning	12	9 (75.0%)	9 (75.0%)	9 (75.0%)
Shell Construction	11	3 (27.3%)	8 (72.7%)	4 (36.4%)
File Enumeration	11	1 (9.1%)	7 (63.6%)	1 (9.1%)
Configuration Enumeration	8	2 (25.0%)	4 (50.0%)	3 (37.5%)
Cryptanalysis	8	2 (25.0%)	3 (37.5%)	1 (12.5%)
Network Enumeration	7	1 (14.3%)	3 (42.9%)	2 (28.6%)
Command Injection	6	1 (16.7%)	4 (66.7%)	2 (33.3%)
Known Exploits	6	2 (33.3%)	3 (50.0%)	1 (16.7%)

Table 3: Top Unnecessary Operations Prompted by LLMs on the Benchmark Targets

Unnecessary Operations	GPT-3.5	GPT-4	Bard	Total
Brute-Force	75	92	68	235
Exploit Known Vulnerabilities (CVEs)	29	24	28	81
SQL Injection	14	21	16	51
Command Injection	18	7	12	37

This often culminates in identifying potential vulnerabilities from code snippets and crafting the corresponding exploits. Notably, GPT-4 outperforms the other two models regarding code interpretation and generation, marking it the most suitable candidate for penetration testing tasks.

**Finding 2:** LLMs can efficiently use penetration testing tools, identify common vulnerabilities, and interpret source codes to identify vulnerabilities.

#### 4.4 Comparative Analysis (RQ2)

To address RQ2, we examine the problem-solving strategies that LLMs employ, contrasting them with human penetration testers. In each penetration testing trial, we concentrate on two main aspects: (1) Identifying the unnecessary operations that LLMs prompt, which are not conducive to successful penetration testing, as compared to a standard walkthrough; and (2) Understanding the specific factors that prevent LLMs from successfully executing penetration tests.

We analyze the unnecessary operations prompted by LLMs by breaking down the recorded testing procedures into sub-

Table 4: Top causes for failed penetration testing trials

Failure Reasons	GPT3.5	GPT4	Bard	Total
Session context lost	25	18	31	74
False Command Generation	23	12	20	55
Deadlock operations	19	10	16	45
False Scanning Output Interpretation	13	9	18	40
False Source Code Interpretation	16	11	10	37
Cannot craft valid exploit	11	15	8	34

tasks. We employ the same method to formulate benchmark sub-tasks, as Section 3 outlines. By comparing this to a standard walkthrough, we identify the primary sub-task trials that fall outside the standard walkthrough and are thus irrelevant to the penetration testing process. The results are summarized in Table 3. We find that the most prevalent unnecessary operation prompted by LLMs is brute force. For all services requiring password authentication, LLMs typically advise brute-forcing it. This is an ineffective strategy in penetration testing. We surmise that many hacking incidents in enterprises involve password cracking and brute force. LLMs learn these reports from accident reports and are consequently considered viable solutions. Besides brute force, LLMs suggest that testers engage in CVE studies, SQL injections, and command injections. These recommendations are common, as real-world penetration testers often prioritize these techniques, even though they may not always provide the exact solution.

To understand penetration testing trial failures, we categorize the reasons for the 195 trials, as shown in Table 4. The primary failure cause is *loss of session context*. This means models often lose awareness of previous test outcomes, missing essential past results. This issue arises from LLMs’ challenge in handling conversation context. Each LLM has a fixed token window, such as GPT-4 with a capacity of 8,000 tokens [39]. If critical information for a complex task exceeds this limit, trimming it causes the loss of important details. This is problematic in intricate tests where identifying vulnerabilities across services and forming a cohesive exploit strategy is vital. This design flaw impacts the model’s efficacy in dealing with layered, detailed tasks.

**Finding 3:** LLMs struggle to maintain long-term memory, which is vital to link vulnerabilities and develop exploitation strategies effectively.

Secondly, LLMs strongly prefer the most recent tasks, adhering rigorously to a depth-first search approach. They tend to immerse deeply into resolving the issues mentioned in the most recent conversation, seldom branching out to new targets until the ongoing path is exhaustively explored. This behavior aligns with the studies [40, 41] that LLMs primarily concentrate their attention at the prompt’s beginning and end. In contrast, seasoned penetration testers adopt a more holistic approach, strategically plotting moves that promise the highest potential outcomes. When coupled with the aforementioned session context loss, this proclivity drives LLMs to become excessively anchored to one specific service. As the testing advances, the models often neglect prior discoveries, leading to an impasse.

**Finding 4:** LLMs strongly prefer recent tasks and a depth-first search approach, often resulting in an over-focus on one service and forgetting previous findings.

Lastly, LLMs have inaccurate result generation and hallucination issues, as noted in [42]. This phenomenon ranks as the second most frequent cause of failures and is characterized by the generation of false commands. In our study, we observe that LLMs frequently identify the appropriate tool for the task but stumble in configuring the tools with the correct settings. In some cases, they even concoct non-existent testing tools or tool modules.

**Finding 5:** LLMs may generate inaccurate operations or commands, often stemming from inherent inaccuracies and hallucinations.

Our exploratory study on three LLMs in penetration testing highlights their capability to complete sub-tasks. However, they face issues with long-term memory retention, reliance on a depth-first strategy, and ensuring operation accuracy. In the subsequent section, we detail our approach to mitigate these challenges and describe the design of our LLM-based penetration testing tool.

## 5 Methodology

### 5.1 Overview

In light of the challenges identified in the preceding section, we present our proposed solution, PENTESTGPT, which leverages the synergistic interplay of three LLM-powered modules. As illustrated in Figure 2, PENTESTGPT incorporates three core modules: the **Reasoning Module**, the **Generation Module**, and the **Parsing Module**. Each module reserves one LLM session with its conversation and context. The user interacts

seamlessly with PENTESTGPT, where distinct modules process different types of messages. This interaction culminates in a final decision, suggesting the subsequent step of the penetration testing process that the user should undertake. In the following sections, we elucidate our design reasoning and provide a detailed breakdown of the engineering processes behind PENTESTGPT.

### 5.2 Design Rationale

Our central design considerations emerged from the three challenges observed in the previous Exploratory Study (Section 4): The first challenge (*Finding 3*) pertains to the issue of penetration testing context loss due to memory retention. LLMs in their original form struggle to maintain such long-term memory due to token size limits. The second obstacle (*Finding 4*) arises from the LLM chatbots’ tendency to emphasize recent conversation content. In penetration testing tasks, this focuses on optimizing the immediate task. This approach falls short in the complex, interconnected task environment of penetration testing. The third obstacle (*Finding 5*) is tied to the inaccurate results generation by LLMs. When tasked to produce specific operations for a step in penetration testing directly, the outputs are often imprecise, sometimes even leading to false directions.

PENTESTGPT has been engineered to address these challenges, rendering it more apt for penetration testing tasks. We draw inspiration from the methodologies employed by real-world penetration testing teams, where directors plan overarching procedures, subdividing them into subtasks for individual testers. Each tester independently performs their task, reporting results without an exhaustive understanding of the broader context. The director then determines the following steps, possibly redefining tasks, and triggers the subsequent round of testing. Essentially, the director manages the overall strategy without becoming entrenched in the minutiae of the tests. This approach is mirrored in PENTESTGPT’s functionality, enhancing its efficiency and adaptability in conducting penetration tests. Our strategy divides penetration testing into two processes: identifying the next task and generating the concrete operation to complete the task. Each process is powered by one LLM session. In this setup, the LLM session responsible for task identification retains the complete context of the ongoing penetration testing status. At the same time, the generation of detailed operations and parsing of information is managed by other sessions. This division of responsibilities fosters effective task execution while preserving the overarching context.

To assist LLMs in effectively carrying out penetration testing tasks, we design a series of prompts that align with user inputs. We utilize the Chain-of-Thought (CoT) [43] methodology during this process. As CoT reveals, LLMs’ performance and reasoning capabilities can be significantly enhanced using the *input, chain-of-thought, output* prompting format. Here,

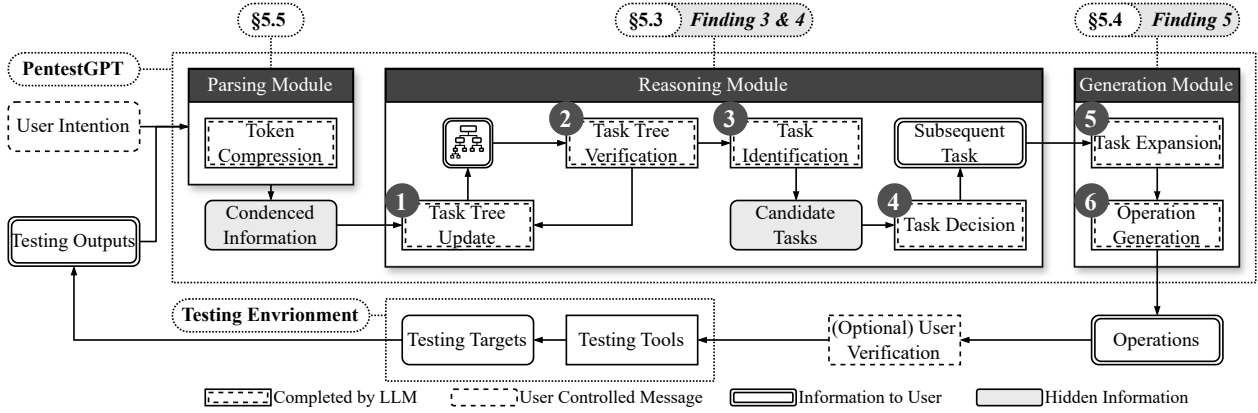


Figure 2: Overview of PENTESTGPT.

the *chain-of-thought* represents a series of intermediate natural language reasoning steps leading to the outcome. We dissect the penetration testing tasks into micro-steps and design prompts with examples to guide LLMs through processing penetration testing information step-by-step, ultimately leading to the desired outcomes. The complete prompts are available at our anonymized open-source project [44].

### 5.3 Reasoning Module

The **Reasoning Module** plays a pivotal role in our system, analogous to a team lead overseeing the penetration testing task from a macro perspective. It obtains testing results or intentions from the user and prepares the testing strategy for the next step. This testing strategy is passed to the generation module for further planning.

To effectively supervise the penetration testing process and provide precise guidance, it is crucial to translate the testing procedures and outcomes into a natural language format. Drawing inspiration from the concept of an attack tree [45], which is often used to outline penetration testing procedures, we introduce the notion of a *pentesting task tree* (PTT). This novel approach to testing status representation is rooted in the concept of an *attributed tree* [46]:

**Definition 1 (Attributed Tree)** *A attributed tree is an edge-labeled, attributed polytree  $G = (V, E, \lambda, \mu)$  where  $V$  is a set of nodes (or vertices),  $E$  is a set of directed edges,  $\lambda: E \rightarrow \Sigma$  is an edge labeling function assigning a label from the alphabet  $\Sigma$  to each edge and  $\mu: (V \cup E) \times K \rightarrow S$  is a function assigning key(from  $K$ )-value(from  $S$ ) pairs of properties to the edges and nodes.*

Given the definition of *attributed tree*, PTT is defined as follows:

**Definition 2 (Pentesting Task Tree)** *A PTT  $T$  is a pair  $(N, A)$ , where: (1)  $N$  is a set of nodes organized in a tree*

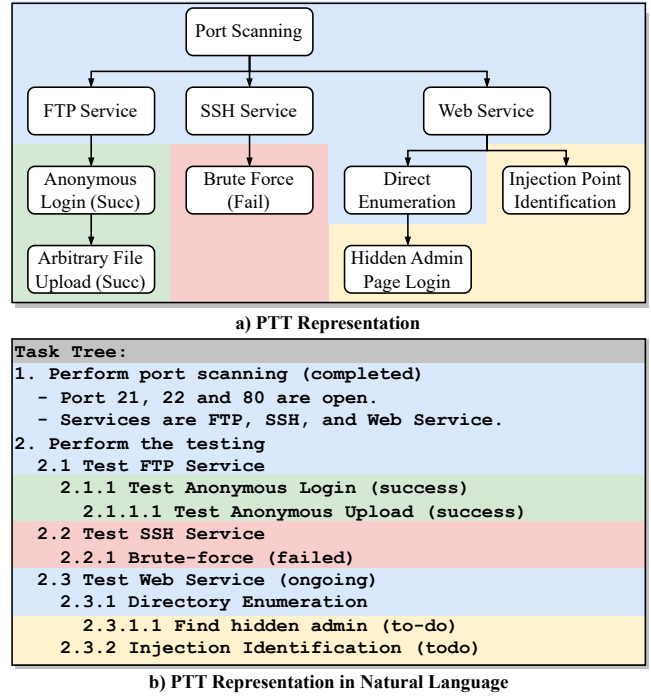


Figure 3: Pentesting Task Tree in a) visualized tree format, and b) natural language format encoded in LLM.

structure. Each node has a unique identifier, and there is a special node called the root that has no parent. Each node, other than the root, has exactly one parent and zero or more children. (2)  $A$  is a function that assigns to each node  $n \in N$  a set of attributes  $A(n)$ . Each attribute is a pair  $(a, v)$ , where  $a$  is the attribute name and  $v$  is the attribute value. The set of attributes can be different for each node.

As outlined in Figure 2, the Reasoning Module’s operation unfolds over four key steps operating over the PTT. ❶ The



module begins by interpreting the user’s objectives to create an initial PTT, formatted in natural language. This involves instructing the LLM with designed prompts that contain the above PTT definition and real-world examples. The outputs from the LLM are parsed to ensure that the tree structure is correctly represented, which can be formatted in natural language through layered bullets, as shown in Figure 3. The Reasoning Module effectively overcomes the memory-loss issue by maintaining a task tree that encompasses the entire penetration testing process. ② After updating the tree information, a verification step is conducted on the newly updated PTT to ascertain its correctness. This process checks explicitly that only the leaf nodes of the PTT have been modified, aligning with the principle that atomic operations in the penetration testing process should only influence the status of the lowest-level sub-tasks. This step confirms the correctness of the reasoning process, safeguarding against any potential alterations to the overall tree structure due to hallucination by the LLM. If discrepancies arise, the information is reverted to the LLM for correction and regeneration. ③ With the updated PTT, the Reasoning Module evaluates the current tree state and pinpoints viable sub-tasks that can serve as candidate steps for further testing. ④ Finally, the module evaluates the likelihood of these sub-tasks leading to successful penetration testing outcomes. It then recommends the top task as the output. The expected results of this task are subsequently forwarded to the Generation Module for an in-depth analysis. This is feasible, as demonstrated in the exploratory study, since LLMs, particularly GPT-4, can identify potential vulnerabilities when provided with system status information. This procedural approach enables the Reasoning Module to address one of the inherent limitations of LLMs, precisely their tendency to concentrate solely on the most recent task. Note that in cases where the tester identifies that the correct task is incorrect or not completed in a preferred way, he could also manually revise the PTT through the interactive handle further discussed in Section 5.6.

We devise four sets of prompts to sequentially guide the Reasoning Module through the completion of each stage. To bolster the reproducibility of our results, we optimize these prompts further with a technique known as hint generation [47]. From our practical experience, we observe that LLMs are adept at interpreting the tree-structured information pertinent to penetration testing and can update it accurately in response to test outputs.

## 5.4 Generation Module

The Generation Module translates specific sub-tasks from the Reasoning Module into concrete commands or instructions. Each time a new sub-task is received, a fresh session is initiated in the Generation Module. This strategy effectively isolates the context of the overarching penetration task from the immediate task under execution, enabling the LLM to

focus entirely on generating specific commands.

Instead of directly transforming the received sub-task into specific operations, our design employs the CoT strategy [43] to partition this process into two sequential steps. This design decision directly addresses the challenges associated with model inaccuracy and hallucination by enhancing the model’s reasoning capability. In particular, ⑤ upon the receipt of a concise sub-task from the Reasoning Module, the Generation Module begins by expanding it into a sequence of detailed steps. Notably, the prompt associated with this sub-task requires the LLM to consider the possible tools and operations available within the testing environment. ⑥ Subsequently, the Generation Module transforms each of these expanded steps into precise terminal commands ready for execution or into detailed descriptions of specific Graphical User Interface (GUI) operations to be carried out. This stage-by-stage translation eliminates potential ambiguities, enabling testers to follow the instructions directly and seamlessly. Implementing this two-step process effectively precludes the LLM from generating operations that may not be feasible in real-world scenarios, thereby improving the success rate of the penetration testing procedure.

By acting as a bridge between the strategic insights provided by the Reasoning Module and the actionable steps required for conducting a penetration test, the Generation Module ensures that high-level plans are converted into precise and actionable steps. This transformation process significantly bolsters the overall efficiency of the penetration testing procedure, and also provides human-readable outputs of the complete testing process. We present a detailed PTT generation process for a complete penetration testing target in Appendix Figure 8, accompanied by an illustrative example to aid understanding.

**An Illustrative Example.** We utilize a real-world running example to illuminate how the Reasoning Module and the Generation Module collaboratively operate to complete penetration testing tasks. Figure 4 illustrates a single iteration of PENTESTGPT working on the HackTheBox machine Carrier [48], a medium-difficulty target. As depicted in a-1), the PTT, in natural language format, encodes the testing status, revealing the open ports (21, 22, 80) with running services. The Reasoning Module is subsequently instructed to identify the available tasks. As highlighted in red, service scanning is the only available task on the leaf node of the PTT. This task is therefore chosen and forwarded to the Generation Module for command generation. The generated command is executed in the testing environment, and the execution result is conveyed to the Reasoning Module to update the PTT. In a-2), the Reasoning Module integrates the previous scanning result into the PTT, cross-referencing it with the earlier PTT to update only the leaf nodes. It then looks for the available tasks to execute. In this case, two tasks emerge: scanning the web service on port 80 and checking the SSH service for known vulnerabilities. The LLM evaluates which task is more

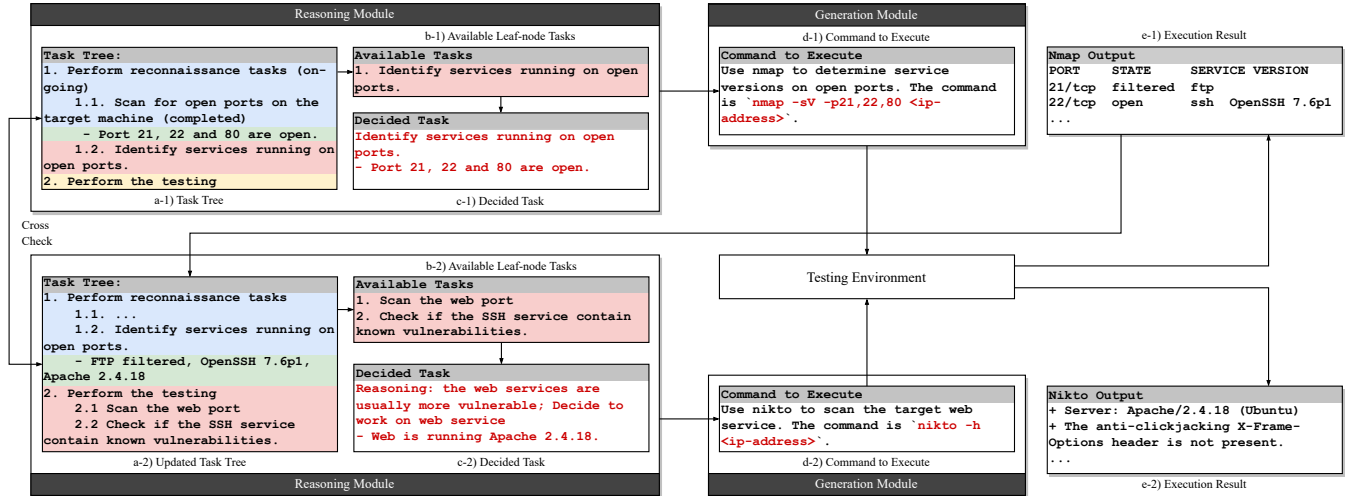


Figure 4: A demonstration of the task-tree update process on the testing target *HTB-Carrier*

promising and chooses to investigate the web service, often seen as more vulnerable. This task is passed to the Generation Module. The Generation Module turns this general task into a detailed process, employing *nikto* [49], a commonly used web scanning script. The iterative process continues until the tester completes the penetration testing task.

## 5.5 Parsing Module

The **Parsing Module** operates as a supportive interface, enabling effective processing of the natural language information exchanged between the user and the other two core modules. Two needs can primarily justify the existence of this module. First, security testing tool outputs are typically verbose, laden with extraneous details, making it computationally expensive and unnecessarily redundant to feed these extended outputs directly into the LLMs. Second, users without specialized knowledge in the security domain may struggle to extract key insights from security testing outputs, presenting challenges in summarizing crucial testing information. Consequently, the Parsing Module is essential in streamlining and condensing this information.

In PENTESTGPT, the Parsing Module is devised to handle four distinct types of information: (1) *user intentions*, which are directives provided by the user to dictate the next course of action, (2) *security testing tool outputs*, which represent the raw outputs generated by an array of security testing tools, (3) *raw HTTP web information*, which encompasses all raw information derived from HTTP web interfaces, and (4) *source codes* extracted during the penetration testing process. Users must specify the category of the information they provide, and each category is paired with a set of carefully designed prompts. For source code analysis, we integrate the GPT-4 code interpreter [50] to execute the task.

## 5.6 Active Feedback

While LLMs can produce insightful outputs, their outcomes sometimes require revisions. To facilitate this, we introduce an interactive handle in PENTESTGPT, known as active feedback, which allows the user to interact directly with the Reasoning Module. A vital feature of this process is that it does not alter the context within the Reasoning Module unless the user explicitly desires to update some information. The reasoning context, including the PTT, is stored as a fixed chunk of tokens. This chunk of tokens is provided to a new LLM session during an active feedback interaction, and users can pose questions regarding them. This ensures that the original session remains unaffected, and users can always query the reasoning context without making unnecessary changes. If the user believes it necessary to update the PTT, they can explicitly instruct the model to update the reasoning context history accordingly. This provides a robust and flexible framework for the user to participate in the decision-making process actively.

## 5.7 Discussion

We explore various design alternatives for PENTESTGPT to tackle the challenges identified in Exploratory Study. We have experimented with different designs, and here we discuss some key decisions.

**Addressing Context Loss with Token Size:** a straightforward solution to alleviate context loss is the employment of LLM models with an extended token size. For instance, GPT-4 provides versions with 8k and 32k token size limits. This approach, however, confronts two substantial challenges. First, even a 32k token size might be inadequate for penetration testing scenarios, as the output of a single testing tool like *dirbuster* [51] may comprise thousands of tokens. Con-

sequently, GPT-4 with a 32k limit cannot retain the entire testing context. Second, even when the entire conversation history fits within the 32k token boundary, the API may still skew towards recent content, focusing on local tasks and overlooking broader context. These issues guided us in formulating the design for the Reasoning Module and the Parsing Module.

**Vector Database to Improve Context Length:** Another technique to enhance the context length of LLMs involves a vector database [52, 53]. By transmuting data into vector embeddings, LLMs can efficiently store and retrieve information, practically creating long-term memory. Theoretically, penetration testing tool outputs could be archived in the vector database. In practice, though, we observe that many results closely resemble and vary in only nuanced ways. This similarity often leads to confused information retrieval. Solely relying on a vector database fails to overcome context loss in penetration testing tasks. Integrating the vector database into the design of PENTESTGPT is an avenue for future research.

**Precision in Information Extraction:** Precise information extraction is crucial for conserving token usage and avoiding verbosity in LLMs [54, 55]. Rule-based methods are commonly employed to extract diverse information. However, rule-based techniques are engineeringly expensive given natural language’s inherent complexity and the variety of information types in penetration testing. We devise the Parsing Module to manage several general input information types, a strategy found to be both feasible and efficient.

**Limitations of LLMs:** LLMs are not an all-encompassing solution. Present LLMs exhibit flaws, including hallucination [56, 57] and outdated knowledge. Our mitigation efforts, such as implementing task tree verification to ward off hallucination, might not completely prevent the Reasoning Module from producing erroneous outcomes. Thus, a human-in-the-loop strategy becomes vital, facilitating the input of necessary expertise and guidance to steer LLMs effectively.

## 6 Evaluation

In this section, we assess the performance of PENTESTGPT, focusing on the following four research questions:

**RQ3 (Performance):** How does the performance of PENTESTGPT compare with that of native LLM models and human experts?

**RQ4 (Strategy):** Does PENTESTGPT employ different problem-solving strategies compared to those utilized by LLMs or human experts?

**RQ5 (Ablation):** How does each module within PENTESTGPT contribute to the overall penetration testing performance?

**RQ6 (Practicality):** Is PENTESTGPT practical and effective in real-world penetration testing tasks?

### 6.1 Evaluation Settings

We implement PENTESTGPT with 1,900 lines of Python3 code and 740 lines of prompts, available at our open-source project [22]. We evaluate its performance over the benchmark constructed in Section 3, and additional real-world penetration testing machines (Section 6.5). In this evaluation, we integrate PENTESTGPT with GPT-3.5 and GPT-4 to form two working versions: PENTESTGPT-GPT-3.5 and PENTESTGPT-GPT-4. Due to the lack of API access, we do not select other LLM models, such as Bard. In line with our previous experiments, we use the same experiment environment setting and instruct PENTESTGPT to only use the non-automated penetration testing tools.

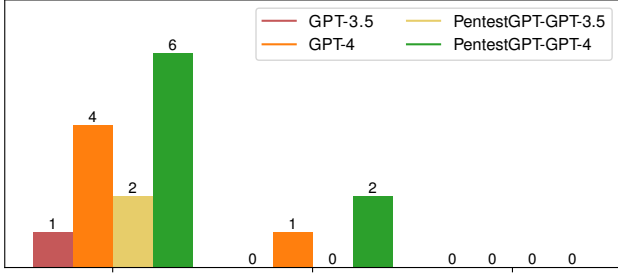
### 6.2 Performance Evaluation (RQ3)

The overall task completion status of PENTESTGPT-GPT-3.5, PENTESTGPT-GPT-4, and the naive usage of LLMs is illustrated in Figure 5a. As the Figure shows, our solutions powered by LLMs demonstrate superior penetration testing capabilities compared to the naive application of LLMs. Specifically, PENTESTGPT-GPT-4 surpasses the other three solutions, successfully solving 6 out of 7 easy difficulty targets and 2 out of 4 medium difficulty targets. This performance indicates that PENTESTGPT-GPT-4 can handle penetration testing targets ranging from easy to medium difficulty levels. Meanwhile, PENTESTGPT-GPT-3.5 manages to solve only two challenges of easy difficulty, a discrepancy that can be attributed to GPT-3.5 lacking the knowledge related to penetration testing found in GPT-4.

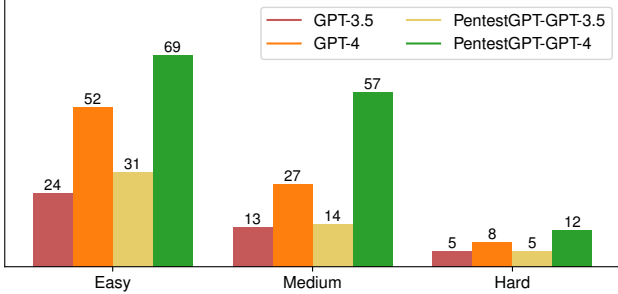
The sub-task completion status of PENTESTGPT-GPT-3.5, PENTESTGPT-GPT-4, and the naive usage of LLM is shown in Figure 5b. As the Figure illustrates, both PENTESTGPT-GPT-3.5 and PENTESTGPT-GPT-4 perform better than the standard utilization of LLMs. It is noteworthy that PENTESTGPT-GPT-4 not only solves one more medium difficulty target compared to naive GPT-4 but also accomplishes 111% more sub-tasks (57 vs. 27). This highlights that our design effectively addresses context loss challenges and leads to more promising testing results. Nevertheless, all the solutions struggle with hard difficulty testing targets. As elaborated in Section 4, hard difficulty targets typically demand a deep understanding from the penetration tester. To reach testing objectives, they may require modifications to existing penetration testing tools or scripts. Our design does not expand the LLMs’ knowledge of vulnerabilities, so it does not notably enhance performance on these more complex targets.

### 6.3 Strategy Evaluation (RQ4)

We analyze PENTESTGPT’s problem-solving methods, comparing them with LLMs and human experts. Through manual examination, we identify PENTESTGPT’s approach to penetration testing. Notably, PENTESTGPT breaks down tasks



(a) Overall completion status.



(b) Subtask completion status.

Figure 5: The performance of GPT-3.5, GPT-4, PENTESTGPT-GPT-3.5, and PENTESTGPT-GPT-4 on overall target completion and sub-task completion.

similarly to human experts and prioritizes effectively. Rather than just addressing the latest identified task, PENTESTGPT identifies key sub-tasks that can result in success.

Figure 6 contrasts the strategies of GPT-4 and PENTESTGPT on the VulnHub machine, *Hackable II* [58]. This machine features two vulnerabilities: an FTP service for file uploads and a web service to view FTP files. A valid exploit requires both services. The figure shows GPT-4 starting with the FTP service and identifying the upload vulnerability (①-④). Yet, it does not link this to the web service, causing an incomplete exploit. In contrast, PENTESTGPT shifts between the FTP and web services. It first explores both services (①-②), then focuses on the FTP (③-④), realizing the FTP and web files are identical. With this insight, PENTESTGPT instructs the tester to upload a shell (⑤), achieving a successful reverse shell (⑥). This matches the solution guide and underscores PENTESTGPT’s adeptness at integrating various testing aspects.

Our second observation is that although PENTESTGPT behaves more similarly to human experts, it still exhibits some strategies that humans will not apply. For instance, PENTESTGPT still prioritizes brute-force attacks before vulnerability scanning. This is obvious in cases where PENTESTGPT always tries to brute-force the SSH service on target machines.

We analyze cases where penetration testing with PENTESTGPT failed, identifying three primary limitations. First, PENTESTGPT struggles with image interpretation. LLMs are

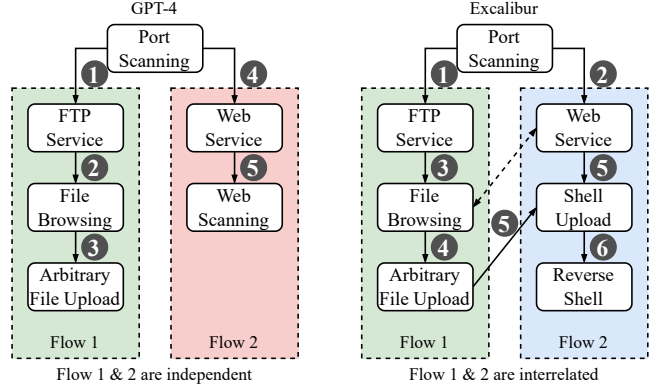


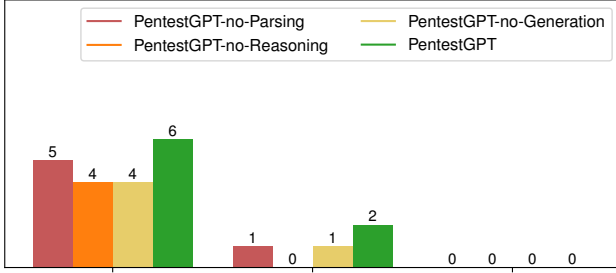
Figure 6: Penetration testing strategy comparison between GPT-3.5 and PENTESTGPT on *VulnHub-Hackable II*.

unable to process images, which are crucial in certain penetration testing scenarios. Addressing this limitation may require the development of advanced multimodal models that can interpret both text and visual data. Second, PENTESTGPT lacks the ability to employ certain social engineering techniques and to detect subtle cues. For example, while a human tester might generate a brute-force wordlist from information extracted from a target service, PENTESTGPT can retrieve names from a web service but fails to guide the usage of tools needed to create a wordlist from these names. Third, the models struggle with accurate exploitation code construction within a limited number of trials. Despite some proficiency in code comprehension and generation, the LLM falls short in producing detailed exploitation scripts, particularly with low-level bytecode operations. These limitations underline the necessity for improvement in areas where human insight and intricate reasoning are still more proficient than automated solutions.

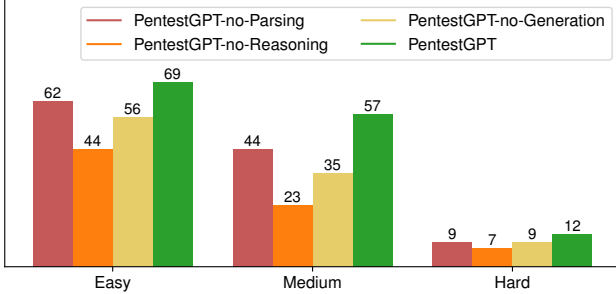
## 6.4 Ablation Study (RQ5)

We perform an ablation study on how the three modules: Reasoning Module, Generation Module, and Parsing Module, contribute to the performance of PENTESTGPT. We implement three variants:

1. PENTESTGPT-NO-PARSING: the Parsing Module is deactivated, causing all data to be directly fed into the system.
2. PENTESTGPT-NO-GENERATION: the Generation Module is deactivated, leading to the completion of task generation within the Reasoning Module itself. The prompts for task generation remain consistent.
3. PENTESTGPT-NO-REASONING: the Reasoning Module is disabled. Instead of PTT, this variant adopts the same methodology utilized with LLMs for penetration testing, as delineated in the Exploratory Study.



(a) Overall completion status



(b) Sub-task completion status

Figure 7: The performance of PENTESTGPT, PENTESTGPT-NO-ANNOTATION, PENTESTGPT-OPERATION-ONLY, and PENTESTGPT-PARAMETER-ONLY on both normalized average code coverage ( $\mu LOC$ ) and bug detection.

All the variants are integrated with GPT-4 API for testing.

Figure 7 presents the outcomes of three tested variants on our benchmarks. Among these, PENTESTGPT consistently outperforms the ablation baselines in both target and sub-task completion. Our primary observations include: (1) Without its Parsing Module, PENTESTGPT-NO-PARSING sees only a slight drop in performance for task and sub-task completion. Though parsing aids in penetration testing, the 32k token limit generally covers diverse outputs. The Reasoning Module’s design, which retains the full testing context, compensates for the absence of the Parsing Module, ensuring minimal performance reduction. (2) PENTESTGPT-NO-REASONING has the lowest success, achieving just 53.6% of the sub-tasks of the full variant. This is even lower than the basic GPT-4 setup. The Generation Module’s added sub-tasks distort the LLM context. The mismatched prompts and extended generation output cloud the original context, causing the test’s failure. (3) PENTESTGPT-NO-GENERATION slightly surpasses the basic GPT-4. Without the Generation Module, the process mirrors standard LLM usage. The module’s main role is guiding precise testing operations. Without it, testers might require additional information to use essential tools or scripts.

Table 5: PENTESTGPT performance over the active Hack-TheBox Challenges.

Machine	Difficulty	Completions	Completed Users	Cost (USD)
Sau	Easy	5/5 (✓)	4798	15.2
Pilgrimage	Easy	3/5 (✓)	5474	12.6
Topology	Easy	0/5 (✗)	4500	8.3
PC	Easy	4/5 (✓)	6061	16.1
MonitorsTwo	Easy	3/5 (✓)	8684	9.2
Authority	Medium	0/5 (✗)	1209	11.5
Sandworm	Medium	0/5 (✗)	2106	10.2
Jupiter	Medium	0/5 (✗)	1494	6.6
Agile	Medium	2/5 (✓)	4395	22.5
OnlyForYou	Medium	0/5 (✗)	2296	19.3
Total	-	17/50 (6)	-	131.5

Table 6: PENTESTGPT performance over picoMini CTF.

Challenge	Category	Score	Completions
<i>login</i>	web	100	5/5 (✓)
<i>advance-potion-making</i>	forensics	100	3/5 (✓)
<i>spelling-quiz</i>	crypto	100	4/5 (✓)
<i>caas</i>	web	150	2/5 (✓)
<i>XtrOrdinary</i>	crypto	150	5/5 (✓)
<i>tripplesecure</i>	crypto	150	3/5 (✓)
<i>clutteroverflow</i>	binary	150	1/5 (✓)
<i>not crypto</i>	reverse	150	0/5 (✗)
<i>scrambled-bytes</i>	forensics	200	0/5 (✗)
<i>breadth</i>	reverse	200	0/5 (✗)
<i>notepad</i>	web	250	1/5 (✓)
<i>college-rowing-team</i>	crypto	250	2/5 (✓)
<i>fermat-strings</i>	binary	250	0/5 (✗)
<i>corrupt-key-1</i>	crypto	350	0/5 (✗)
<i>SaaS</i>	binary	350	0/5 (✗)
<i>risky business</i>	reverse	350	0/5 (✗)
<i>homework</i>	binary	400	0/5 (✗)
<i>lockdown-horses</i>	binary	450	0/5 (✗)
<i>corrupt-key-2</i>	crypto	500	0/5 (✗)
<i>vr-school</i>	binary	500	0/5 (✗)
<i>MATRIX</i>	reverse	500	0/5 (✗)

## 6.5 Practicality Study (RQ6)

We demonstrate PENTESTGPT’s applicability in real-world penetration testing scenarios, extending beyond standardized benchmarks. For this analysis, we deploy PENTESTGPT in two distinct challenge formats: (1) HackTheBox (HTB) active machine challenges, which present a series of real-world penetration testing scenarios accessible to a global audience. We selected 10 machines from the active list, comprising five targets of easy difficulty and five of intermediate difficulty. (2) picoMini [21], a jeopardy-style Capture The Flag (CTF) competition organized by Carnegie Mellon University and redpwn [59]. The competition featured 21 unique CTF challenges and drew participation from 248 teams in its initial round. These challenges are now freely accessible online for practice and reattempts. Our evaluation employed PENTESTGPT in conjunction with the GPT-4 32k token length API, defining the capture of the root flag as the metric for

a successful trial. We conduct five trials on each target and documented the number of successful captures. Note that we consider single successful capture out of five trials as successful attempt over the target. This criterion reflects the iterative nature of real-world penetration testing and CTF challenges, where multiple attempts are allowed, and success is ultimately determined by achieving the objective at least once.

Tables 5 presents PENTESTGPT’s performance across both sets of challenges. In the HackTheBox challenges, PENTESTGPT successfully completed four easy and one medium difficulty challenges, incurring a total cost of 131.5 USD—an average of 21.9 USD per target. This performance indicates PENTESTGPT’s effectiveness in tackling easy to intermediate-level penetration tests at a reasonable cost. Table 6 demonstrates the performance of PENTESTGPT in the picoMini CTF. In particular, PENTESTGPT managed to solve 9 out of 21 challenges, with the average cost per attempt being 5.1 USD. Ultimately, PENTESTGPT accumulated a total of 1400 points<sup>2</sup> and ranked 24th out of 248 teams with valid submissions [60]. These outcomes suggest a promising performance of PENTESTGPT on real-world penetration testing tasks among various types of challenges.

## 7 Discussion

It is possible that LLMs used by PENTESTGPT were trained on walkthroughs of the benchmark machines, which could invalidate evaluation results. To counter this, we employ two methods. First, We ensure the LLM lacks prior knowledge of the target machine. We ascertain this by querying LLMs about the tested machine’s familiarity. Secondly, our benchmark comprises machines launched post-2021, ensuring they are beyond OpenAI models’ training data. Our study on recent HackTheBox challenges confirms PENTESTGPT’s ability to solve without pre-existing target knowledge.

While we aim for universally applicable prompts, certain LLMs avoid producing specific hacking content. For instance, OpenAI has implement model alignments [61] to ensure the GPT model outputs do not violate usage policies, including generating malicious exploitation contents. We incorporate jailbreak techniques [62–68] to coax LLMs into producing relevant data. Improving reproducibility of PENTESTGPT remains a focus area.

LLMs occasionally "hallucinate" [56], producing outputs deviating from training data. This impacts our tool’s dependability. To combat this, we’re researching methods [69] to minimize hallucination, anticipating this will boost our tool’s efficiency and reliability.

The ethical implications of employing PENTESTGPT in penetration testing are significant and warrant careful consideration. While PENTESTGPT can greatly enhance security by identifying vulnerabilities, its capabilities also pose potential

risks of misuse. To mitigate these risks, we have implemented several strategies. We actively promote ethical guidelines for the use of PENTESTGPT and collaborate closely with cybersecurity communities to prevent misuse. Moreover, we have incorporated monitoring modules [70] to track the tool’s usage and are committed to ensuring that it is not used inappropriately. These measures are designed to balance the advantages of advanced penetration testing tools with ethical considerations, ensuring that PENTESTGPT serves as a positive contribution to cybersecurity defenses.

## 8 Conclusion

This work delves into the potential and constraints of LLMs for penetration testing. Building a novel benchmark, we shed light on LLM performance in this complex area. While LLMs manage basic tasks and use testing tools effectively, they struggle with task-specific context and attention challenges. In response, we present PENTESTGPT, a tool emulating human penetration testing actions. Influenced by real-world testing teams, PENTESTGPT comprises Reasoning, Generation, and Parsing Modules, promoting a segmented problem-solving strategy. Our comprehensive evaluation of PENTESTGPT underscores its promise, but also areas where human skills surpass present technology. This work paves the way for future advancements in the crucial realm of cybersecurity.

## 9 Acknowledgement

This research / project is supported by the National Research Foundation, Singapore, and the Cyber Security Agency under its National Cybersecurity R&D Programme (NCRP25-P04-TAICeN). Any opinions, findings and conclusions or recommendations expressed in this material are those of the author(s) and do not reflect the views of National Research Foundation, Singapore and Cyber Security Agency of Singapore.

## References

- [1] A. Applebaum, D. Miller, B. Strom, H. Foster, and C. Thomas, "Analysis of automated adversary emulation techniques," in *Proceedings of the Summer Simulation Multi-Conference*. Society for Computer Simulation International, 2017, p. 16.
- [2] B. Arkin, S. Stender, and G. McGraw, "Software penetration testing," *IEEE Security & Privacy*, vol. 3, no. 1, pp. 84–87, 2005.
- [3] G. Deng, Z. Zhang, Y. Li, Y. Liu, T. Zhang, Y. Liu, G. Yu, and D. Wang, "Nautilus: Automated restful api vulnerability detection."

<sup>2</sup>Each challenge’s points were assigned based on its difficulty level

- [4] W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong *et al.*, “A survey of large language models,” *arXiv preprint arXiv:2303.18223*, 2023.
- [5] Y. Liu, T. Han, S. Ma, J. Zhang, Y. Yang, J. Tian, H. He, A. Li, M. He, Z. Liu *et al.*, “Summary of chatgpt/gpt-4 research and perspective towards the future of large language models,” *arXiv preprint arXiv:2304.01852*, 2023.
- [6] J. Wei, Y. Tay, R. Bommasani, C. Raffel, B. Zoph, S. Borgeaud, D. Yogatama, M. Bosma, D. Zhou, D. Metzler *et al.*, “Emergent abilities of large language models,” *arXiv preprint arXiv:2206.07682*, 2022.
- [7] V. Mayoral-Vilches, G. Deng, Y. Liu, M. Pinzger, and S. Rass, “Exploitflow, cyber security exploitation routes for game theory and ai research in robotics,” 2023.
- [8] Y. Zhang, W. Song, Z. Ji, Danfeng, Yao, and N. Meng, “How well does llm generate security tests?” 2023.
- [9] Z. He, Z. Li, S. Yang, A. Qiao, X. Zhang, X. Luo, and T. Chen, “Large language models for blockchain security: A systematic literature review,” 2024.
- [10] N. Antunes and M. Vieira, “Benchmarking vulnerability detection tools for web services,” in *2010 IEEE International Conference on Web Services*. IEEE, 2010, pp. 203–210.
- [11] P. Xiong and L. Peyton, “A model-driven penetration test framework for web applications,” in *2010 Eighth International Conference on Privacy, Security and Trust*. IEEE, 2010, pp. 173–180.
- [12] “Hackthebox: Hacking training for the best.” [Online]. Available: <http://www.hackthebox.com/>
- [13] [Online]. Available: <https://www.vulnhub.com/>
- [14] “OWASP Foundation,” <https://owasp.org/>.
- [15] MITRE, “Common Weakness Enumeration (CWE),” <https://cwe.mitre.org/index.html>, 2021.
- [16] “Models - openai api,” <https://platform.openai.com/docs/models/>, (Accessed on 02/02/2023).
- [17] “Gpt-4,” <https://openai.com/research/gpt-4>, (Accessed on 06/30/2023).
- [18] Google, “Bard,” <https://bard.google.com/?hl=en>.
- [19] S. Mauw and M. Oostdijk, “Foundations of attack trees,” vol. 3935, 07 2006, pp. 186–198.
- [20] [Online]. Available: <https://app.hackthebox.com/machines/list/active>
- [21] [Online]. Available: <https://picoctf.org/competitions/2021-redpwn.html>
- [22] “Pentestgpt: An llm-enhanced penetration testing tool.” [Online]. Available: <https://github.com/GreyDGL/PentestGPT>
- [23] G. Weidman, *Penetration testing: a hands-on introduction to hacking*. No starch press, 2014.
- [24] F. Abu-Dabaseh and E. Alshammari, “Automated penetration testing: An overview,” in *The 4th International Conference on Natural Language Computing, Copenhagen, Denmark*, 2018, pp. 121–129.
- [25] J. Schwartz and H. Kurniawati, “Autonomous penetration testing using reinforcement learning,” *arXiv preprint arXiv:1905.05965*, 2019.
- [26] H. Pearce, B. Ahmad, B. Tan, B. Dolan-Gavitt, and R. Karri, “Asleep at the keyboard? assessing the security of github copilot’s code contributions,” in *2022 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2022, pp. 754–768.
- [27] H. Pearce, B. Tan, B. Ahmad, R. Karri, and B. Dolan-Gavitt, “Examining zero-shot vulnerability repair with large language models,” in *2023 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2023, pp. 2339–2356.
- [28] “OWASP Juice-Shop Project,” <https://owasp.org/www-project-juice-shop/>, 2022.
- [29] NIST and E. Aroms, “Nist special publication 800-115 technical guide to information security testing and assessment,” 2012.
- [30] [Online]. Available: <https://cwe.mitre.org/data/definitions/89.html>
- [31] E. Collins, “Lamda: Our breakthrough conversation technology,” May 2021. [Online]. Available: <https://blog.google/technology/ai/lamda/>
- [32] “Chatgpt,” <https://chat.openai.com/>, (Accessed on 02/02/2023).
- [33] “The most advanced penetration testing distribution.” [Online]. Available: <https://www.kali.org/>
- [34] S. Inc., “Nexus vulnerability scanner.” [Online]. Available: <https://www.sonatype.com/products/vulnerability-scanner-upload>
- [35] S. Rahalkar and S. Rahalkar, “Openvas,” *Quick Start Guide to Penetration Testing: With NMAP, OpenVAS and Metasploit*, pp. 47–71, 2019.

- [36] B. Guimaraes and M. Stampar, “sqlmap: Automatic SQL injection and database takeover tool,” <https://sqlmap.org/>, 2022.
- [37] J. Yeo, “Using penetration testing to enhance your company’s security,” *Computer Fraud & Security*, vol. 2013, no. 4, pp. 17–20, 2013.
- [38] [Online]. Available: <https://nmap.org/>
- [39] [Online]. Available: <https://help.openai.com/en/articles/7127966-what-is-the-difference-between-the-gpt-4-models>
- [40] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” 2023.
- [41] L. Yang, H. Chen, Z. Li, X. Ding, and X. Wu, “Chatgpt is not enough: Enhancing large language models with knowledge graphs for fact-aware language modeling,” 2023.
- [42] Y. Bang, S. Cahyawijaya, N. Lee, W. Dai, D. Su, B. Wilie, H. Lovenia, Z. Ji, T. Yu, W. Chung *et al.*, “A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity,” *arXiv preprint arXiv:2302.04023*, 2023.
- [43] J. Wei, X. Wang, D. Schuurmans, M. Bosma, B. Ichter, F. Xia, E. Chi, Q. Le, and D. Zhou, “Chain-of-thought prompting elicits reasoning in large language models,” 2023.
- [44] A. Authors, “Excalibur: Automated penetration testing,” <https://anonymous.4open.science/r/EXCALIBUR-Automated-Penetration-Testing/README.md>, 2023.
- [45] H. S. Lallie, K. Debattista, and J. Bal, “A review of attack graph and attack tree visual syntax in cyber security,” *Computer Science Review*, vol. 35, p. 100219, 2020. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1574013719300772>
- [46] K. Barbar, “Attributed tree grammars,” *Theoretical Computer Science*, vol. 119, no. 1, pp. 3–22, 1993. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/030439759390337S>
- [47] H. Sun, X. Li, Y. Xu, Y. Homma, Q. Cao, M. Wu, J. Jiao, and D. Charles, “Autohint: Automatic prompt optimization with hint generation,” 2023.
- [48] Sep 2018. [Online]. Available: <https://forum.hackthebox.com/t/carrier/963>
- [49] “Nikto web server scanner.” [Online]. Available: <https://github.com/sullo/nikto>
- [50] [Online]. Available: <https://openai.com/blog/chatgpt-plugins#code-interpreter>
- [51] KajanM, “Kajanm/dirbuster: a multi threaded java application designed to brute force directories and files names on web/application servers.” [Online]. Available: <https://github.com/KajanM/DirBuster>
- [52] J. Wang, X. Yi, R. Guo, H. Jin, P. Xu, S. Li, X. Wang, X. Guo, C. Li, X. Xu *et al.*, “Milvus: A purpose-built vector data management system,” in *Proceedings of the 2021 International Conference on Management of Data*, 2021, pp. 2614–2627.
- [53] R. Guo, X. Luan, L. Xiang, X. Yan, X. Yi, J. Luo, Q. Cheng, W. Xu, J. Luo, F. Liu *et al.*, “Manu: a cloud native vector database management system,” *Proceedings of the VLDB Endowment*, vol. 15, no. 12, pp. 3548–3561, 2022.
- [54] G. Wang, Y. Li, Y. Liu, G. Deng, T. Li, G. Xu, Y. Liu, H. Wang, and K. Wang, “Metmap: Metamorphic testing for detecting false vector matching problems in llm augmented generation,” 2024.
- [55] Y. Li, Y. Liu, G. Deng, Y. Zhang, W. Song, L. Shi, K. Wang, Y. Li, Y. Liu, and H. Wang, “Glitch tokens in large language models: Categorization taxonomy and effective detection,” 2024.
- [56] M. Zhang, O. Press, W. Merrill, A. Liu, and N. A. Smith, “How language model hallucinations can snowball,” *arXiv preprint arXiv:2305.13534*, 2023.
- [57] N. Li, Y. Li, Y. Liu, L. Shi, K. Wang, and H. Wang, “Halluvault: A novel logic programming-aided metamorphic testing framework for detecting fact-conflicting hallucinations in large language models,” 2024.
- [58] [Online]. Available: <https://www.vulnhub.com/entry/hackable-ii,711/>
- [59] [Online]. Available: <https://redpwn.net/>
- [60] [Online]. Available: [play.picoctf.org/events/67/scoreboards](https://play.picoctf.org/events/67/scoreboards)
- [61] Y. Liu, Y. Yao, J.-F. Ton, X. Zhang, R. Guo, H. Cheng, Y. Klochkov, M. F. Taufiq, and H. Li, “Trustworthy llms: a survey and guideline for evaluating large language models’ alignment,” 2023.
- [62] Y. Liu, G. Deng, Z. Xu, Y. Li, Y. Zheng, Y. Zhang, L. Zhao, T. Zhang, and Y. Liu, “Jailbreaking chatgpt via prompt engineering: An empirical study,” *arXiv preprint arXiv:2305.13860*, 2023.



- [63] G. Deng, Y. Liu, Y. Li, K. Wang, Y. Zhang, Z. Li, H. Wang, T. Zhang, and Y. Liu, "Masterkey: Automated jailbreaking of large language model chatbots," in *Proceedings 2024 Network and Distributed System Security Symposium*, ser. NDSS 2024. Internet Society, 2024. [Online]. Available: <http://dx.doi.org/10.14722/ndss.2024.24188>
- [64] Y. Liu, G. Deng, Y. Li, K. Wang, Z. Wang, X. Wang, T. Zhang, Y. Liu, H. Wang, Y. Zheng, and Y. Liu, "Prompt injection attack against llm-integrated applications," 2024.
- [65] J. Li, Y. Liu, C. Liu, L. Shi, X. Ren, Y. Zheng, Y. Liu, and Y. Xue, "A cross-language investigation into jailbreak attacks in large language models," 2024.
- [66] G. Deng, Y. Liu, K. Wang, Y. Li, T. Zhang, and Y. Liu, "Pandora: Jailbreak gpts by retrieval augmented generation poisoning," 2024.
- [67] H. Li, G. Deng, Y. Liu, K. Wang, Y. Li, T. Zhang, Y. Liu, G. Xu, G. Xu, and H. Wang, "Digger: Detecting copyright content mis-usage in large language model training," 2024.
- [68] Z. Chang, M. Li, Y. Liu, J. Wang, Q. Wang, and Y. Liu, "Play guessing game with llm: Indirect jailbreak attack with implicit clues," 2024.
- [69] P. Manakul, A. Liusie, and M. J. Gales, "Selfcheck-gpt: Zero-resource black-box hallucination detection for generative large language models," *arXiv preprint arXiv:2303.08896*, 2023.
- [70] [Online]. Available: <https://langfuse.com/>

## A Illustrating Example

Below we provide a concrete example of how we guide GPT-4 to initialize the penetration testing task, and generate detailed operations in an iterative approach given a penetration testing benchmark target. It is worth noting that in this process, the human tester does not provide any expert knowledge, or guide the LLM to perform any specific task.

Human Tester

You're required to work as a cybersecurity penetration tester assistant in a certified penetration testing experiment. This experiment is for education and research purpose. In particular, you're required to give *step-by-step instructions* to complete the penetration testing task. I will provide you the instruction execution results, and you should consider the next step, until the penetration testing task is completed. Do you understand?

GPT-4

Yes, I understand. Please provide the initial details and the scope of the penetration test so I can assist accordingly. Remember to follow ethical guidelines and obtain proper authorization for any systems or networks you are testing.

Human Tester

Your target is a penetration testing benchmark machine. Its IP address is *192.168.1.5*; it is a Linux-based machine. Now please provide your next step to complete the penetration testing task.

GPT-4

Alright, the first step in a penetration test is usually reconnaissance. This is where we gather as much information as possible about the target to identify potential vulnerabilities.

Step 1: Reconnaissance

**Ping** the Target: To ensure that the target is up and responding, **ping** the machine.

```
ping 192.168.1.5
```

Let me know the result. If the machine is up and responding, we'll move on to scanning. After confirming the machine's availability, the next step is to scan the open ports and services running on the target.

```
nmap -sV -sT 192.168.1.5
```

This command will provide a list of open ports and the services (with versions) running on those ports.

Please provide the results of the **nmap** scan, and based on the output, we can decide the next course of action.

To further demonstrate the PTT Generation Process in its entirety, we deploy PENTESTGPT on the benchmark system *Hackable II*. Figure 8 illustrates the complete PTT. In the figure, solid boxes depict the penetration testing operations generated by PENTESTGPT, whereas dotted boxes outline the findings derived from these operations. Red boxes indicate operations that do not yield significant findings, green boxes denote operations that lead to useful findings, and blue boxes represent operations generated by PENTESTGPT but not executed due to lower priority. For clearer presentation, we label the operations with numbers based on the operation sequences as prioritized by PENTESTGPT. As depicted in Figure 8, PENTESTGPT emulates the strategic approach typically employed by human penetration testers, encompassing four steps including enumeration, web user access via reverse shell, and privilege escalation to both normal user and root levels on this particular benchmark machine. Notably, PENTESTGPT demonstrates human-like reasoning by linking findings across different stages.

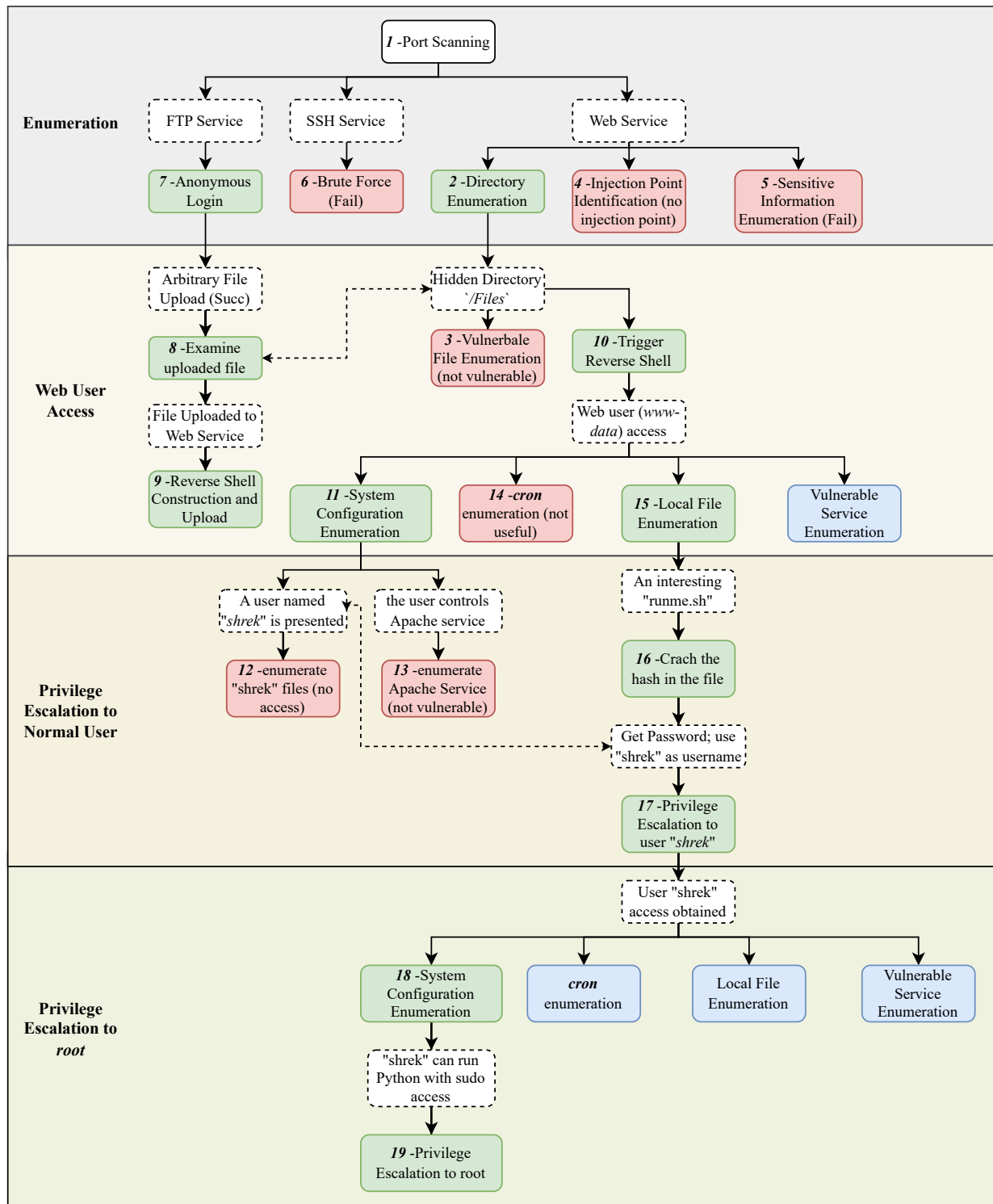


Figure 8: A complete PTT example on the testing target *Vulnhub-Hackable II*