

Beijing Forest Studio
北京理工大学信息系统及安全对抗实验中心



大模型赋能的渗透测试技术

硕士研究生 郑俊怡

2025年05月11日

- 总结反思
 - 演讲技巧欠缺，需要面对听众
 - 创新点引入不自然，概念解释不清
- 相关内容
 - 2025.04.06 高玺凯《二进制代码反编译技术》
 - 2024.09.03 张浩然《大模型赋能的模糊测试用例生成技术》

- 预期收获
- 题目内涵解析
- 研究背景与意义
- 研究历史与现状
- 算法原理
 - PentestGPT
 - HPTSA
- 特点总结与工作展望
- 参考文献

- 预期收获
 - 掌握渗透测试的基本概念
 - 了解大模型赋能的自动化渗透测试技术的研究背景和研究意义
 - 明确大模型赋能的自动化渗透测试技术的前沿方法和未来方向

- 内涵解析

- 渗透测试：通过**模拟恶意黑客**的攻击方法，主动评估计算机网络系统的安全性能
- 流程：侦察、扫描、漏洞评估、利用和后利用（包括报告）
 - **侦察**：全面收集目标系统的公开/隐藏信息，构建攻击面地图，识别潜在入口点
 - **扫描**：识别目标系统的开放服务、潜在漏洞及配置缺陷，缩小攻击范围
 - **漏洞评估**：验证漏洞存在性，评估其可利用性、影响范围及修复优先级
 - **漏洞利用**：利用已验证漏洞获取目标系统权限，实施攻击
 - **后利用**：成功渗透目标系统或网络后执行的一系列活动

- 研究目标

- 利用大模型强大的通用推理能力、自然语言处理能力和学习能力
- **辅助安全人员**或者**自动化**发现系统中存在的安全隐患
- 提升计算机网络系统的安全性

- 研究背景

- 渗透测试通常大量**人工参与**，高度依赖专业人员的经验和知识，整个过程需要耗费大量的时间和人力成本
 - 地缘政治不确定性加剧了经济压力，导致许多行业的预算和劳动力减少，而网络安全威胁和数据安全事件只增不减
 - 根据ISC2《**2024**年网络安全劳动力研究报告》，全球网络安全从业人员数量**同比增长0.1%**，而**2022**年的**同比增长为11.1%**
- 传统自动化工具**受限于环境简化**（仅支持有限动作空间、马尔可夫假设），难以应对真实网络的异步性和部分可观测性

- 研究意义

- 降低渗透测试对专业人员的依赖度与人力成本，显著提升渗透测试的执行效率
- 增强自动化渗透测试工具在复杂网络环境中的适应能力，满足真实网络需求

• 传统的自动化渗透测试工具

– 基于规则

- 将人工渗透测试经验转化为渗透规则，并将一系列渗透动作组合成渗透程序
- 优点：渗透规则源自专家的实战经验，可靠性较高，对于适配的渗透场景决策速度快且渗透成功率高
- 缺点：只适用于固定渗透场景

– 基于模型

- 渗透场景的机器理解，渗透知识库，渗透路径规划，渗透策略
- 优点：能够使用智能模型实现对人决策行为的模拟，自动选择渗透路径
- 缺点：在大型复杂网络中的渗透决策能力差，渗透行为过多时就会造成渗透效率低

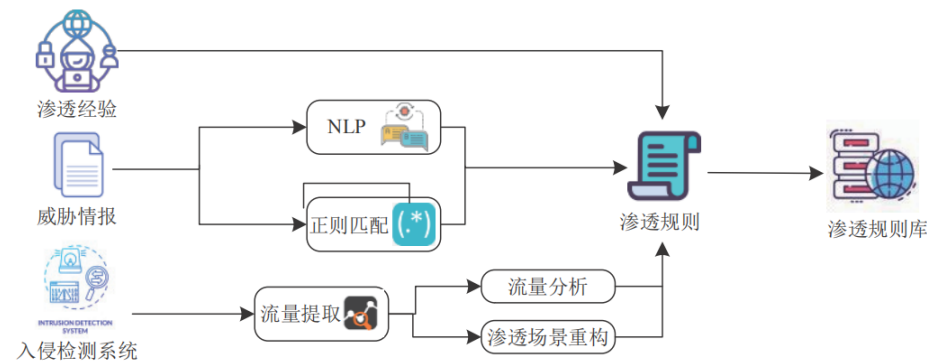


图7 渗透规则提取框架

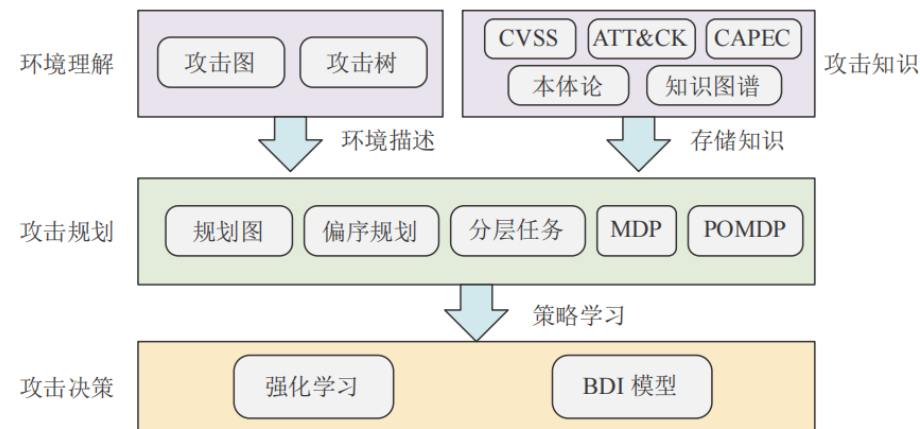


图9 基于模型的自动化渗透框架

Happe 等人验证了 GPT3.5 能够辅助人们进行渗透测试，不仅能够在**低层次**上提供具体的操作命令，而且能够在**高层次**上规划任务

Hilario 等人使用 ChatGPT 3.5 在渗透测试的勘测、扫描、漏洞评估、漏洞利用和报告五个阶段中提供选项和建议，并在 VulnHub 上的公共机器上进行测试，发现大模型能够提升渗透测试的**创造性**和**效率**，帮助发现新攻击路径

Deng 等人设计出的 pentestGPT 拥有**推理**、**生成**和**解析**三个自交互模块，其中每个模块处理渗透测试的单独子任务

构建了专为评估 LLMs 在攻防安全任务中的**开源基准数据集**，并集成 Metasploit、Ghidra 等工具构建**半自动化的 AI 助手**



2023

2024

2024

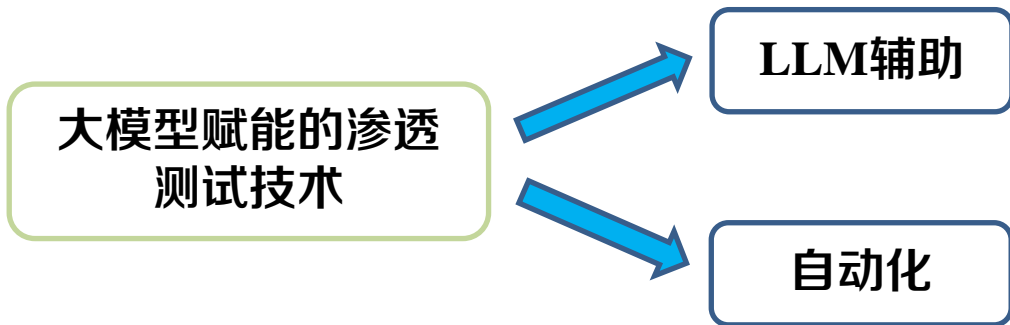
2025

Moskal S 等人对大模型生成可操作网络相关知识的能力，以及对威胁行为者的能力提升进行评估，同时尝试实施自动化攻击。他们发现，**优化提示词**和**检测、修正模型幻觉**是提升自动化攻击性能的关键步骤

Fang 等人使用 GPT4 驱动 AI 代理，自主入侵构建的**测试网站**，揭示了 AI 代理能够在未预先知道漏洞信息且没有人类反馈的情况下执行复杂的人侵任务

Fang 等人使用 AI 代理在只知道 CVE 描述的情况下，自主利用现实世界的**一日漏洞**

Fang 等人设计出具有规划代理且能够启动子代理的多代理系统 HPTSA。**规划代理**探索目标系统并确定要调用子代理的类型，从而在尝试不同漏洞时解决**长期规划问题**；**子代理**中包括**通用子代理**和多个特定于某一类型任务的**专家代理**。Fang 等人使用 HPTSA 成功自主利用**零日漏洞**





【 USENIX 】

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



LLMO

T	目标	评估LLM在渗透测试中的能力，并提出解决方案
I	输入	渗透测试目标信息或渗透测试操作结果
P	处理	1. 确定下一个任务 2. 生成完成任务的具体操作
O	输出	渗透测试步骤
P	问题	1. LLM难以维持长期记忆 2. LLM倾向最近的任务和深度优先的搜索方法，容易过度关注一项服务而忘记以前的发现 3. 由于内部知识不准确和存在幻觉，LLM可能生成错误的操作或命令
C	条件	能够执行端到端渗透测试，有效运用渗透工具，识别常见漏洞，并具备解析源代码发现安全隐患的能力的大模型
D	难点	1. 如何解决渗透测试上下文丢失问题 2. 如何全局维护和规划渗透测试任务 3. 如何避免LLM生成错误的操作或命令
L	水平	USENIX 2024 CCF A



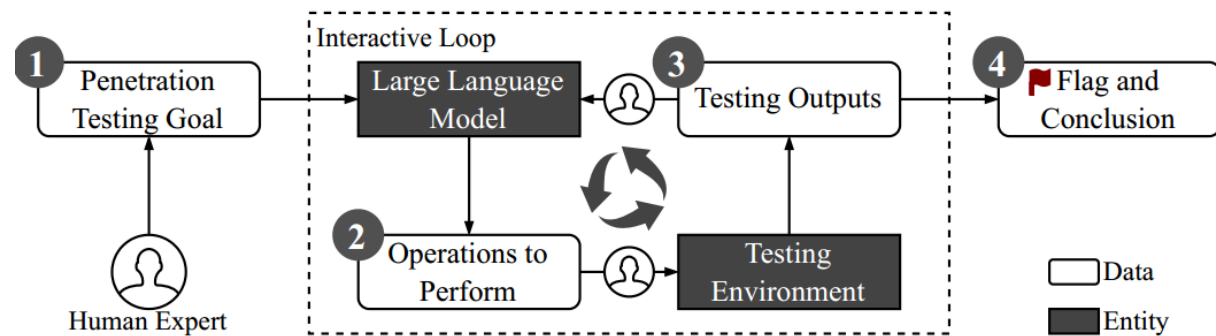
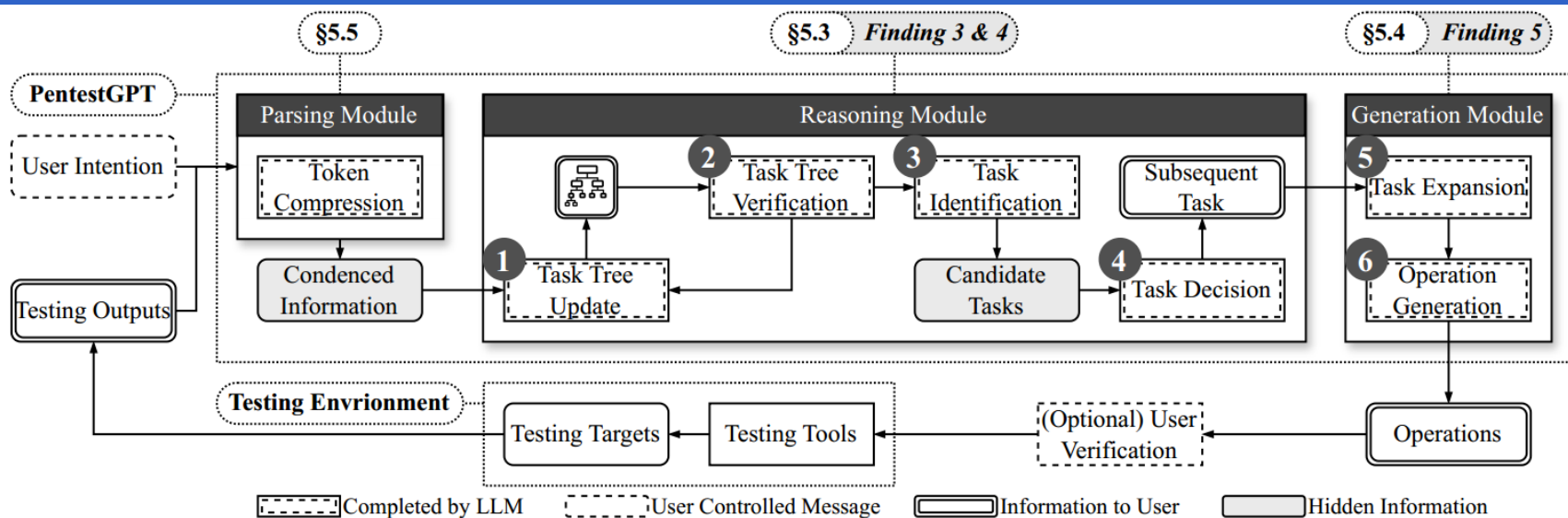
• 算法原理图

– 渗透测试分为两个过程

- 确定下一个任务
- 生成完成任务的具体操作

– 三个核心模块

- **解析模块**：简化和压缩用户输入
- **推理模块**：确定下一个任务
- **生成模块**：生成完成任务的具体操作
- 主动反馈：修正大模型结果



• 解析模块

— 作用

- 提取关键信息
- 减少计算成本

— 类别

- 用户意图
- 安全测试工具输出
- 原始HTTP Web信息
- 在渗透测试过程中提取的源代码

```
<!DOCTYPE html>
<html lang="zh-CN">

<head>
  <meta charset="utf-8">
  <link rel="canonical" href="https://blog.csdn.net/Kris_zhang/article/details/106841466"/>
  <meta http-equiv="content-type" content="text/html; charset=utf-8">
  <meta name="renderer" content="webkit"/>
  <meta name="force-rendering" content="webkit"/>
  <meta http-equiv="X-UA-Compatible" content="IE=edge, chrome=1"/>
  <meta name="viewport" content="width=device-width, initial-scale=1.0, minimum-scale=1.0, maximum-scale=1.0, user-scalable=no">
  <meta name="report" content="{\"pid\": \"blog\", \"spm\": \"1001.2101\"}">
  <meta name="referrer" content="always">
  <meta http-equiv="Cache-Control" content="no-siteapp" /><link rel="alternate" media="handheld" href="#"/>
  <meta name="shenma-site-verification" content="5a59773ab8077d4a62bf469ab966a63b_1497598848">
  <meta name="applicable-device" content="pc">
  <link href="https://g.csdnimg.cn/static/logo/favicon32.ico" rel="shortcut icon" type="image/x-icon" />
  <title>Nmap网络扫描利器-CSDN博客</title>
  <script>
    (function() {
      var el = document.createElement("script");
      el.src = "https://s3a.pstatp.com/toutiao/push.js?1abfal3dfe74d72d41d83c86d240de427e7cac50c51ead53b2e79d40c7952a23ed7716d05b4a4c";
      el.id = "ttzz";
      var s = document.getElementsByTagName("script")[0];
      s.parentNode.insertBefore(el, s);
    })(window)
  </script>
  <meta name="keywords" content="nmap">
  <meta name="csdn-baidu-search" content="{\"autorun\":true,\"install\":true,\"keyword\":\"nmap\"}">
  <meta name="description" content="文章浏览阅读10w+次，点赞367次，收藏3k次。简介nmap是用来探测计算机网络上主机和服务的一种安全扫描器。为了给用户提供一个良好的使用体验，特将nmap的使用教程整理成文档，方便用户查阅。<a href='\"https://csdnimg.cn/release/blogv2/dist/pc/css/detail_enter-a805dca594.min.css\"'>https://csdnimg.cn/release/blogv2/dist/pc/css/detail_enter-a805dca594.min.css</a>
  <link rel="stylesheet" type="text/css" href="https://csdnimg.cn/release/blogv2/dist/pc/css/detail_enter-a805dca594.min.css">
  <script type="application/ld+json">{"@context":"https://ziyuan.baidu.com/contexts/cambrian.jsonld","@id":"https://blog.csdn.net/Kris_zhang","type":"article","url":"https://blog.csdn.net/Kris_zhang/article/details/106841466","publishedTime":"2019-07-10T10:00:00.000Z","modifiedTime":"2019-07-10T10:00:00.000Z","author":{"name":"Kris_zhang","url":"https://blog.csdn.net/Kris_zhang"},"publisher":{"name":"CSDN","url":"https://www.csdn.net/"},"provider":{"name":"CSDN","url":"https://www.csdn.net/"},"article":{"name":"Nmap网络扫描利器-CSDN博客","url":"https://blog.csdn.net/Kris_zhang/article/details/106841466","description":"文章浏览阅读10w+次，点赞367次，收藏3k次。简介nmap是用来探测计算机网络上主机和服务的一种安全扫描器。为了给用户提供一个良好的使用体验，特将nmap的使用教程整理成文档，方便用户查阅。<a href='\"https://csdnimg.cn/release/blogv2/dist/pc/css/detail_enter-a805dca594.min.css\"'>https://csdnimg.cn/release/blogv2/dist/pc/css/detail_enter-a805dca594.min.css</a>","keywords":"nmap","articleSection":"安全测试工具输出","articleCategory":"安全测试工具输出","articleTag":"nmap","articleLabel":"安全测试工具输出","articleImage":"https://img-blog.csdnimg.cn/416e3d3474a24468bbf4ba4ae77a4e40_kris_zhang.jpg?1"}</script>
  <script src="https://g.csdnimg.cn/lib/jquery/1.12.4/jquery.min.js" type="text/javascript"></script>
  <script src="https://g.csdnimg.cn/lib/jquery-migrate/1.4.1/jquery-migrate.js" type="text/javascript"></script>
  <script type="text/javascript">
    var isCorporate = false;
    var username = "Kris_zhang";
    var skinImg = "white";

    var blog_address = "https://blog.csdn.net/Kris_zhang";
    var currentUser = "weixin_54484507";
    var isOwner = false;
    var loginUrl = "http://passport.csdn.net/account/login?from=https://blog.csdn.net/Kris_zhang/article/details/106841466";
    var blogUrl = "https://blog.csdn.net/";
    var avatar = "https://profile-avatar.csdnimg.cn/416e3d3474a24468bbf4ba4ae77a4e40_kris_zhang.jpg?1";
    var articleTitle = "nmap超详细使用教程";
    var articleDesc = "文章浏览阅读10w+次，点赞367次，收藏3k次。简介nmap是用来探测计算机网络上主机和服务的一种安全扫描器。为了给用户提供一个良好的使用体验，特将nmap的使用教程整理成文档，方便用户查阅。<a href='\"https://csdnimg.cn/release/blogv2/dist/pc/css/detail_enter-a805dca594.min.css\"'>https://csdnimg.cn/release/blogv2/dist/pc/css/detail_enter-a805dca594.min.css</a>";
    var articleTitles = "Nmap网络扫描利器-CSDN博客";
    var nickName = "刘一手";
    var articleDetailUrl = "https://blog.csdn.net/Kris_zhang/article/details/106841466";
    var vipUrlV = "https://mall.csdn.net/vip?vipSource=learningVip";
    if(window.location.host.split('.').length == 3) {
      blog_address = blogUrl + username;
    }
    var skinStatus = "White";
    var blogStaticHost = "https://csdnimg.cn/release/blogv2/"
    var payColumn = false
  </script>
```

推理模块

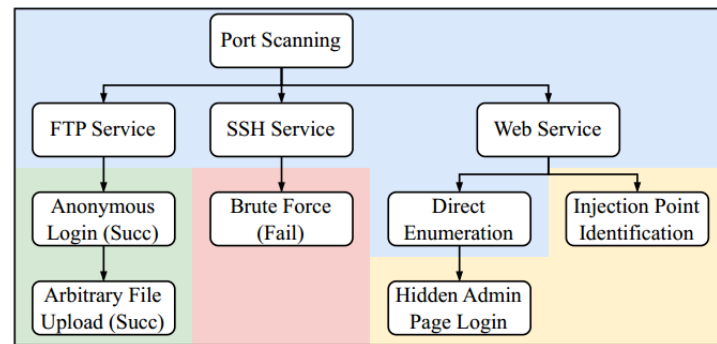
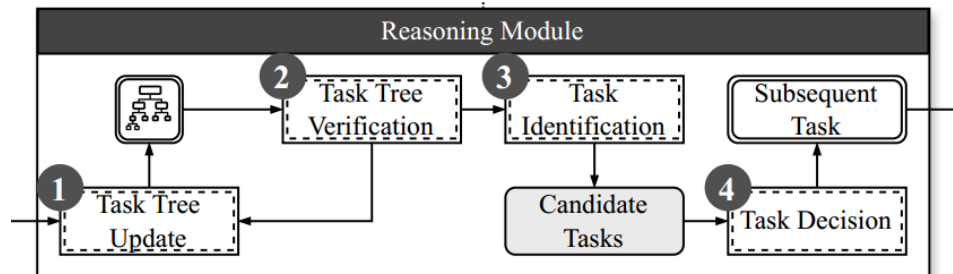
解决的问题

- LLM难以维持长期记忆→渗透测试任务树（PTT）
- LLM倾向于只专注于最近的任务→任务决定

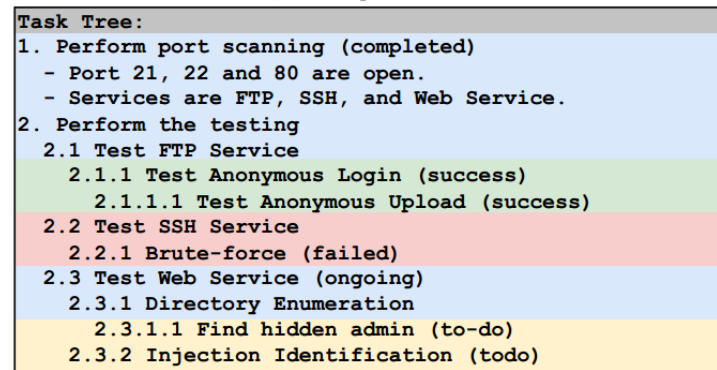
步骤

- **任务树更新**：解释用户的目标，创建自然语言格式化的初始PTT
- **任务树验证**：检查新PTT是否仅修改了PTT的叶节点，确保推理过程的正确性
- **任务评估**：评估当前的树状态，并确定可行的候选子任务
- **任务决定**：评估子任务导致成功渗透测试结果的可能性，

推荐top任务作为输出



a) PTT Representation



b) PTT Representation in Natural Language

- 生成模块
 - 解决的问题
 - 生成错误的操作或命令→消除潜在的歧义
 - 步骤
 - 任务扩展
 - 将简洁的子任务扩展为一系列详细的步骤
 - 考虑测试环境中可用的可能工具和操作
 - 操作生成
 - 将扩展步骤转换为可执行的**精确终端命令**，或转换为要执行的特定**图形用户界面操作**的**详细说明**

GPT-4

好的，渗透测试的第一步通常是侦察。在这一步，我们要尽可能多地收集关于目标的信息，以识别潜在的漏洞。

步骤 1: 侦察

ping 目标主机：为确保目标主机已启动并能响应，请 ping 该主机。

```
ping 192.168.1.5
```

告诉我结果。如果机器已启动并能响应，我们就继续进行扫描。确认机器可用后，下一步是扫描目标上开放的端口和正在运行的服务。

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

此命令将提供开放端口列表以及在这些端口上运行的服务（及其版本）。
请提供 nmap 扫描的结果，根据输出，我们可以决定下一步行动方案。

• 数据集

- 从HackTheBox和VulnHub中选择13个任务（182个子任务），涵盖了OWASP Top 10中列出的所有类型的漏洞
- 简单7个（子任务77），中等4个（子任务71），困难2个（子任务34）

• 对比方法

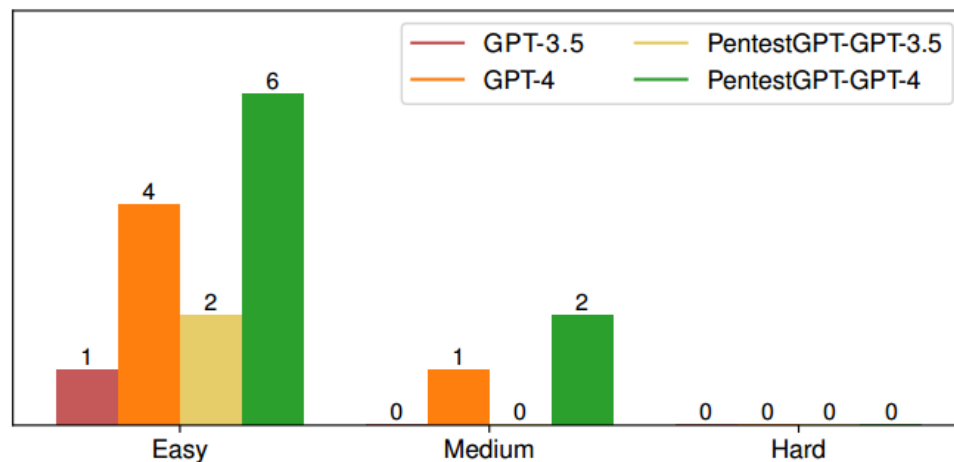
- GPT-3.5
- GPT-4
- PENTESTGPT-GPT-3.5
- PENTESTGPT-GPT-4

• 评价指标：

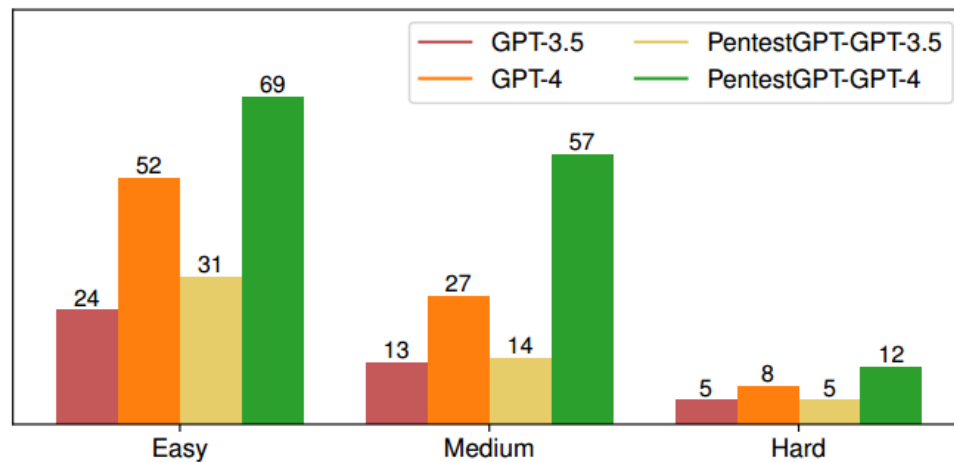
- 总体完成状态
- 子任务完成状态

性能评估

- **GPT-4**和**PENTESTGPT-GPT-4**的性能强于**GPT-3.5**和**PENTESTGPT-GPT-3.5**
 - **GPT-3.5**缺乏渗透测试相关知识
- **PENTESTGPT**在任务和子任务上表现都优于直接使用大模型
- 所有解决方案在**高难度测试目标**上都面临困难
 - 困难任务需要修改现有的渗透测试工具或脚本
 - 未扩展大模型对漏洞的知识



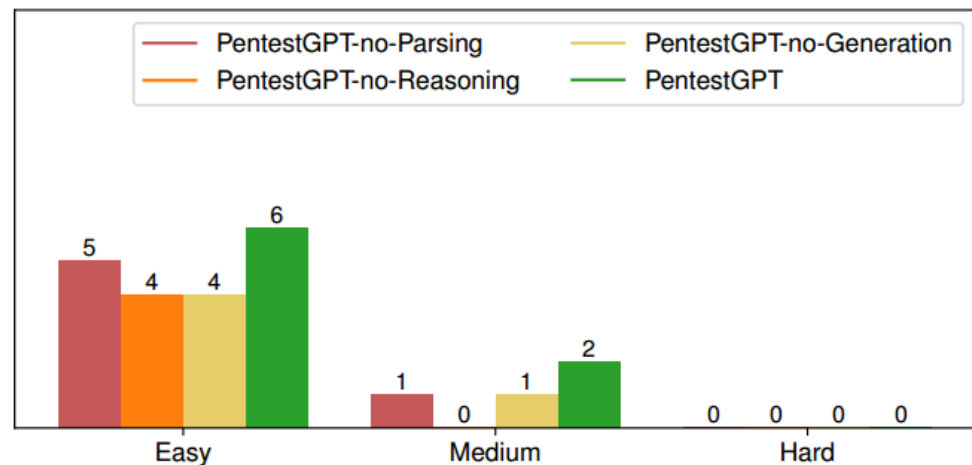
(a) Overall completion status.



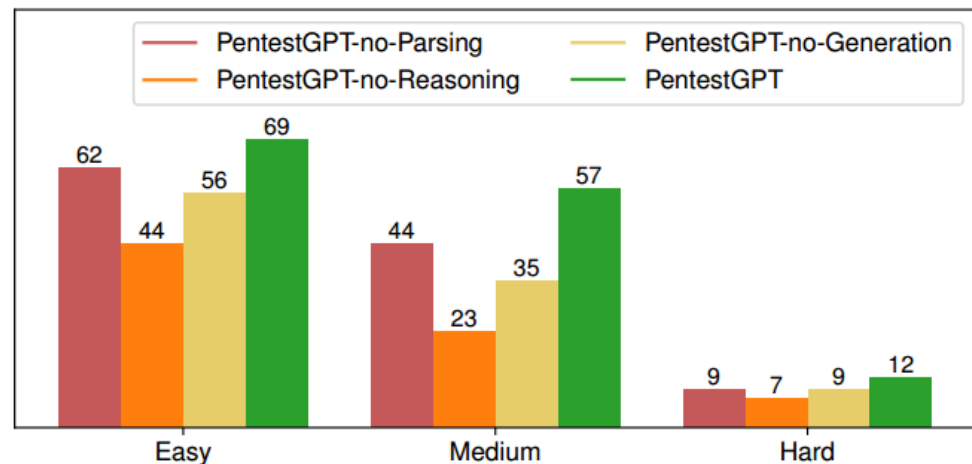
(b) Subtask completion status.

• 消融实验

- 没有**解析模块**，PENTESTGPT在任务和子任务完成方面的性能仅略有下降
 - 推理模块的设计保留了完整的测试上下文，弥补了解析模块的缺失
- 没有**推理模块**，PENTESTGPT在任务和子任务完成方面的性能下降最大，甚至低于直接使用GPT-4
 - 生成模块添加的子任务扭曲了大模型的上下文
- 没有**生成模块**，PENTESTGPT略优于直接使用GPT-4
 - 类似于标准的大模型使用，缺乏精确操作



(a) Overall completion status



(b) Sub-task completion status

• 实用性研究

– HackTheBox

- 5个简单难度目标和5个中等难度目标
- 成功完成了4个简单难度和1个中等难度的挑战，总成本为131.5美元，平均每个目标21.9美元

– picoMini

- 21个CTF挑战
- 解决9个挑战，每次尝试的平均成本为5.1美元
- 总共积累了1400分，在248个有效提交的团队中排名第24位

Table 5: PENTESTGPT performance over the active HackTheBox 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.

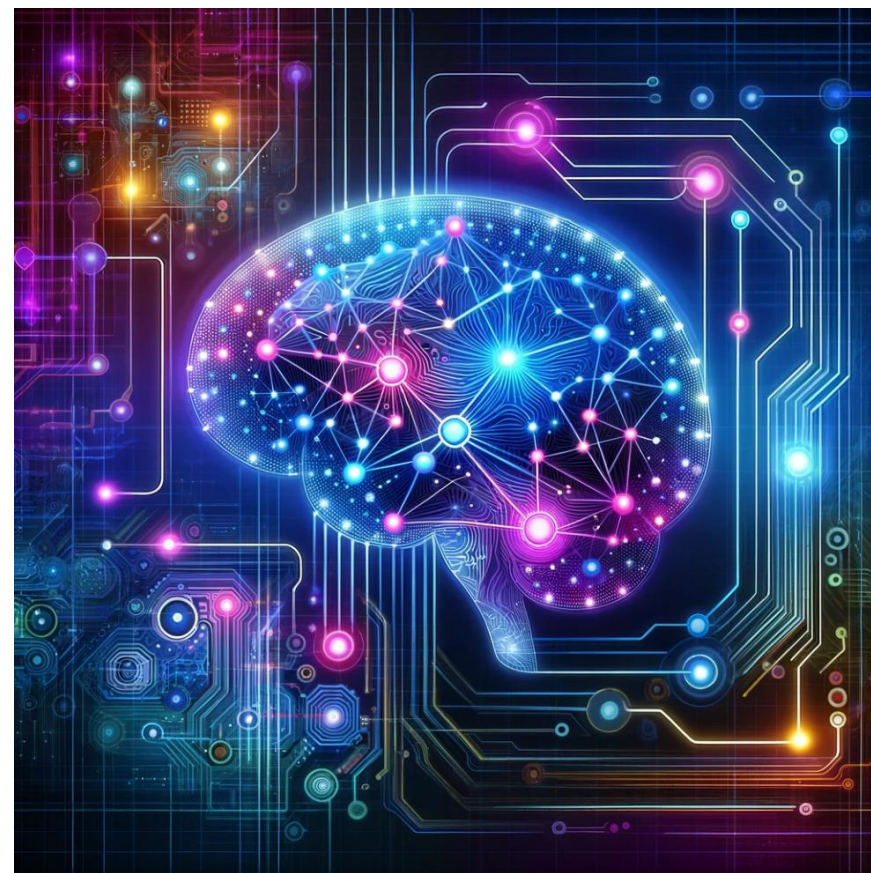
Challenge	Category	Score	Completions
<i>login</i>	web	100	5/5 (✓)
<i>advance-potion-making</i>	forensics	100	3/5 (✓)
<i>spelling-quiz</i>	crypto	100	4/5 (✓)
<i>caas</i>	web	150	2/5 (✓)
<i>XtrOrdinary</i>	crypto	150	5/5 (✓)
<i>tripplesecure</i>	crypto	150	3/5 (✓)
<i>clutteroverflow</i>	binary	150	1/5 (✓)
<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 (✗)

• 算法优势

- 使用PTT实现渗透测试上下文的长期维护
- 解决了大模型倾向于只专注于最近的任务的问题
- 缓解了幻觉和错误指令的影响

• 算法不足

- 幻觉问题依旧影响工具的可靠性
- 只使用了大模型本身的知识，没有额外构建渗透知识库
- 无法处理图像信息
- 缺乏运用某些社会工程学技术以及察觉细微线索的能力
- 难以准确构建漏洞利用代码



HLI2V



【 Arxiv 】

Teams of LLM Agents can Exploit Zero-Day Vulnerabilities



TIPO LIBO

T	目标	实现零日漏洞的自动化渗透测试
I	输入	目标系统信息、工具信息和文档
P	处理	1. 探索网站，确定要尝试的漏洞类型以及要尝试的网站页面 2. 决定使用哪些特定智能体 3. 使用特定智能体执行特定任务
O	输出	漏洞利用结果

P	问题	1. 大模型长期规划能力欠缺 2. 大模型倾向于深度优先，难以尝试多种漏洞
C	条件	能够执行端到端渗透测试，有效运用渗透工具，识别常见漏洞，并具备解析源代码发现安全隐患的能力的大模型
D	难点	1. 如何进行长期规划，避免陷于单一漏洞的尝试
L	水平	arxiv 2025

- 算法原理图

- 分层规划器

- 探索网站

- 为了降低成本，去掉不必要的HTML标签

- 确定要尝试的漏洞类型以及要尝试的网站页面

- 团队管理器

- 决定使用哪些特定智能体

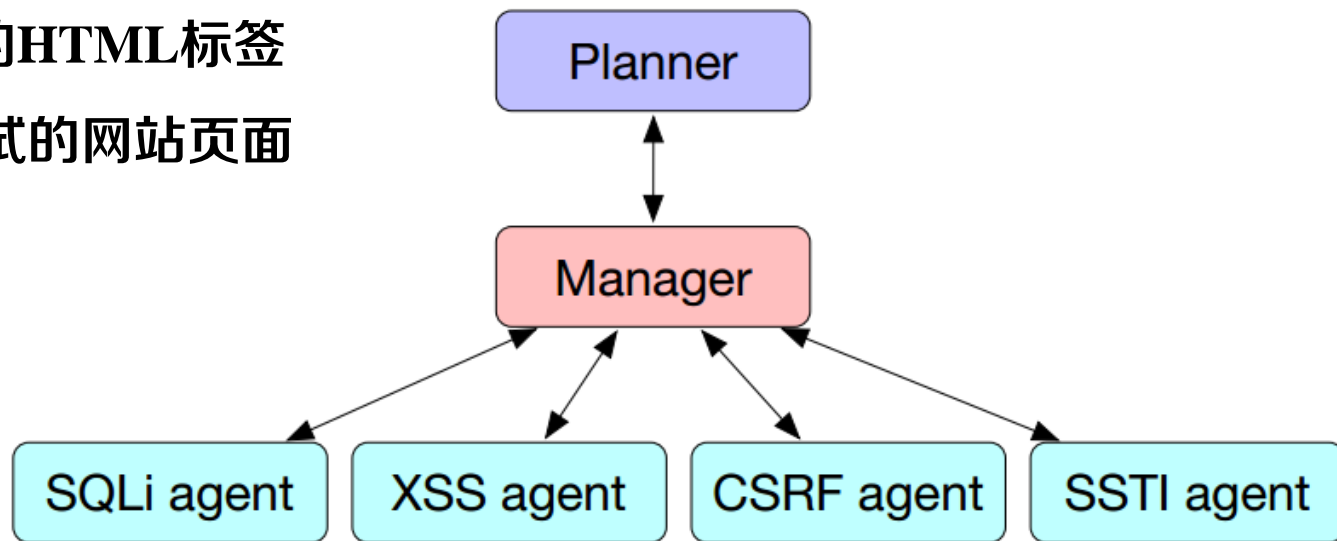
- 检索智能体运行信息

- 特定任务的专家智能体

- 6个专家智能体：XSS、SQLi、CSRF、SSTI、ZAP

- 以及一个通用网络黑客智能体

- 具备工具使用、文档访问的权限





- 数据集
 - 14个Web漏洞

- 对比方法
 - 开源漏洞扫描器ZAP、MetaSploit
 - 没有漏洞描述的one-day代理
 - 有漏洞描述的one-day代理

Vulnerability	CVE	Date	Severity
Travel Journal XSS	CVE-2024-24041	02/01/2024	6.1 (medium)
flusity-CMS CSRF	CVE-2024-24524	02/02/2024	8.8 (high)
flusity-CMS XSS	CVE-2024-27757	03/18/2024	6.1 (medium)
Dolibarr SQLi	CVE-2024-5314	05/24/2024	9.1 (critical)
LedgerSMB CSRF privilege escalation	CVE-2024-23831	02/02/2024	7.5 (high)
alf.io improper authorization	CVE-2024-25635	02/19/2024	8.8 (high)
changedetection.io XSS	CVE-2024-34061	05/02/2024	4.3 (medium)
Navidrome parameter manipulation	CVE-2024-32963	05/01/2024	4.2 (medium)
SWS XSS	CVE-2024-32966	05/01/2024	5.8 (medium)
Zabbix privilege escalation	CVE-2024-22120	05/14/2024	9.1 (critical)
Stalwart Mail Server ACE	CVE-2024-35179	05/15/2024	6.8 (medium)
Sourcecodester SQLi admin-manage-user	CVE-2024-33247	04/25/2024	9.8 (critical)
Sourcecodester SQLi login	CVE-2024-31678	04/11/2024	9.8 (critical)
PrestaShop information leakage	CVE-2024-34717	05/14/2024	5.3 (medium)

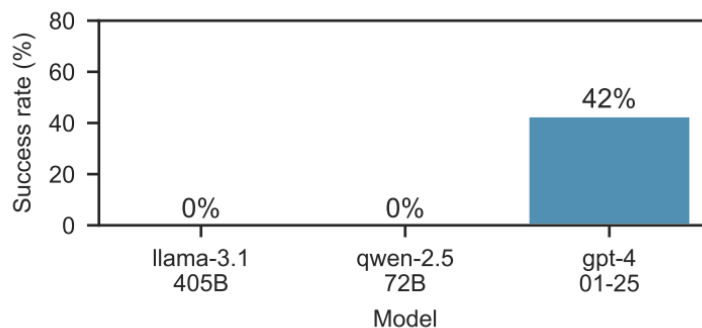
- 评价指标：
 - Pass at 5
 - Pass at 1

• 对比实验

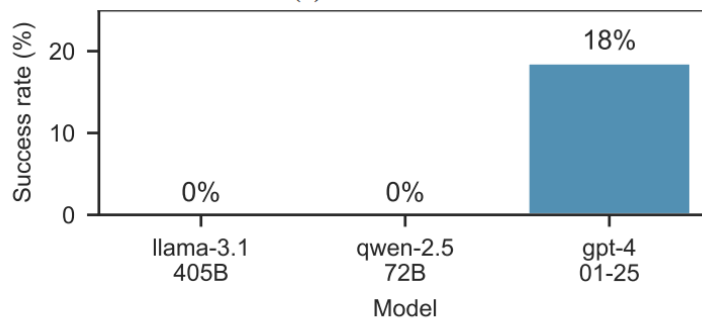
– 使用**GPT-4**的HPTSA达到了最高的成功率，开源模型未能利用任何漏洞

- 开源模型拒绝率更高，且经常反复尝试相同的错误方法

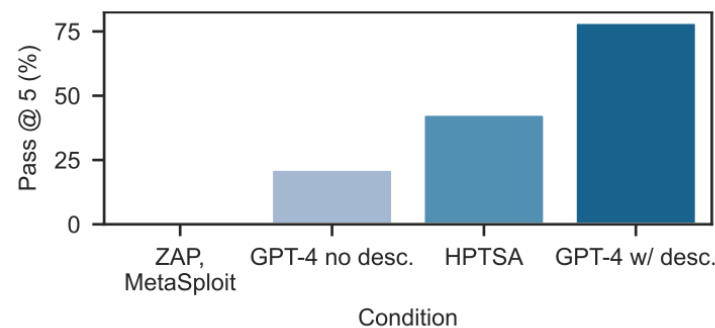
– HPTSA**优于没有漏洞描述**的one-day代理，**差于有漏洞描述**的one-day代理；开源漏洞扫描器未能利用任何漏洞



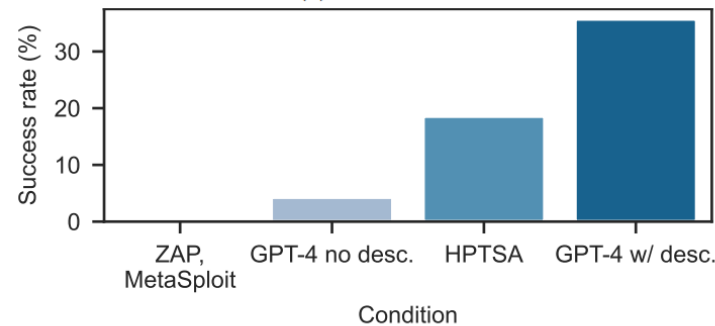
(a) Pass at 5



(b) Overall success rate (pass at 1)



(a) Pass at 5



(b) Overall success rate (pass at 1)

• 消融实验

– 三个变体

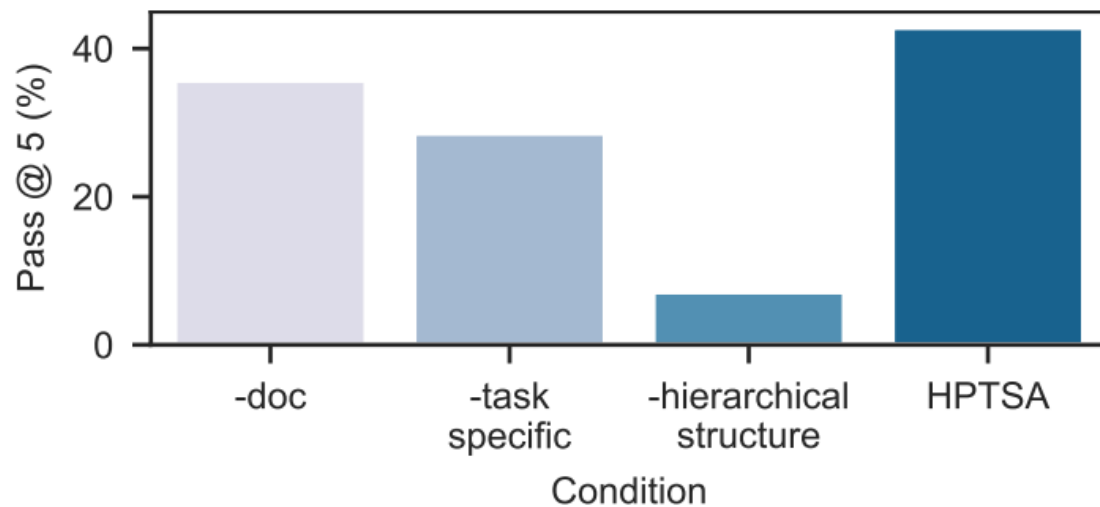
- 从特定任务智能体中**移除文档**
- **只用通用网络安全智能体**
- 没有分层结构，**随机使用特定任务智能体**

– **移除特定任务智能体和移除文档会导致性能大幅下降**

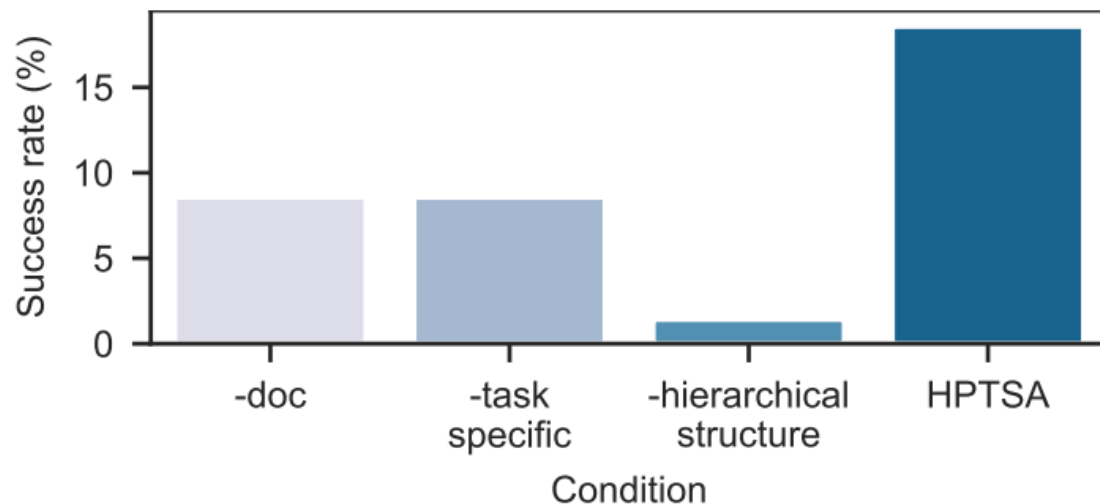
- Pass at 1下降2.1倍，Pass at 5下降50%
- Pass at 1下降2.1倍，Pass at 5下降20%

– 没有分层结构，智能体的单步通过率会降低13倍，五步通过率会降低6倍

- 这个实验合理吗？



(a) Pass at 5



(b) Overall success rate (pass at 1)

• 算法流程

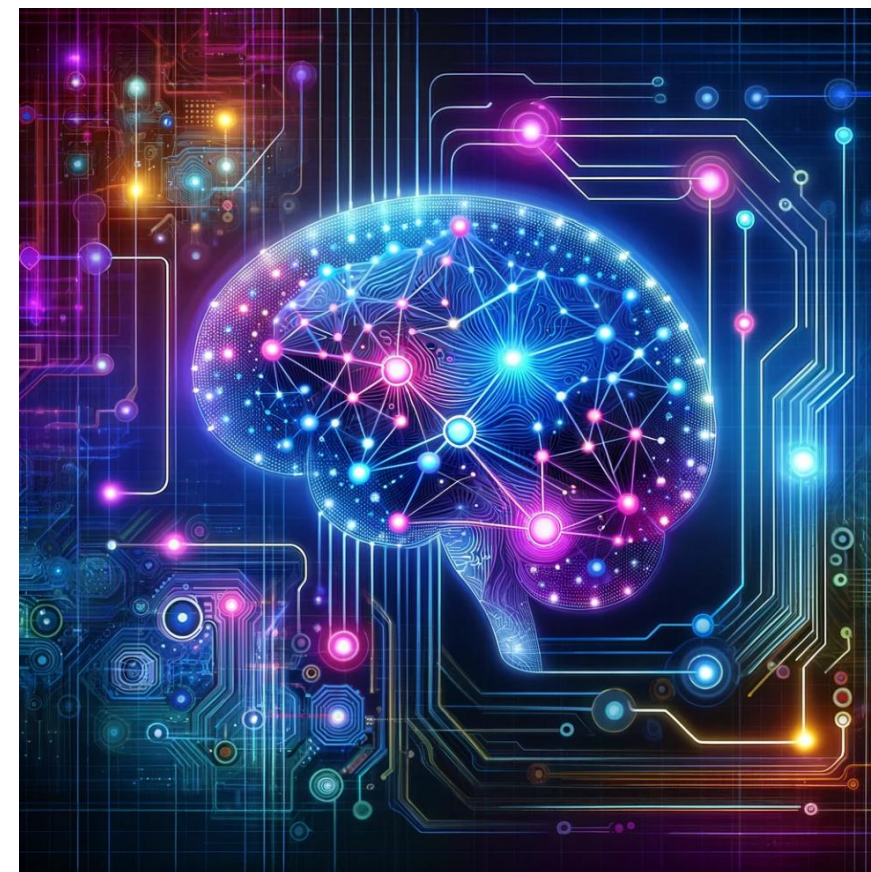
- 探索网站，确定要尝试的漏洞类型以及要尝试的网站页面
- 决定使用哪些特定智能体
- 使用特定智能体执行特定任务

• 算法优势

- 可以自动化实现渗透测试
- 实现长期规划和特定任务的专业化
- 排除了训练数据的干扰

• 算法不足

- 终端节点发现能力不足
- 只在Web漏洞上进行测试





特点总结与未来展望



- 特点总结
 - PentestGPT
 - 人工智能助手
 - 方法特点：PTT、COT、三核心模块
 - HPTSA
 - 自动化渗透测试工具
 - 方法特点：多代理、渗透知识库、工具集成
- 未来发展
 - 全自动化渗透测试工具
 - 缓解幻觉影响
 - 全面基准测试

- [1] Deng G, Liu Y, Mayoral-Vilches V, et al. {PentestGPT}: Evaluating and harnessing large language models for automated penetration testing. 33rd USENIX Security Symposium (USENIX Security 24) [C]. Berkeley, CA: USENIX Association, 2024: 847-864.
- [2] Fang R, Bindu R, Gupta A, et al. Teams of LLM Agents can Exploit Zero-Day Vulnerabilities[EB/OL]. (2025-03-30)[2025-05-11]. <https://arxiv.org/abs/2406.01637>.

知人者智，自知者明。胜人者有力，自胜者强。知足者富。强行者有志。不失其所者久。死而不亡者，寿。

谢谢！

