



# An Attacker's Dream? Exploring the Capabilities of ChatGPT for Developing Malware

Yin Minn Pa Pa  
Yokohama National University  
Yokohama, Japan

Shunsuke Tanizaki  
Yokohama National University  
Yokohama, Japan

Tetsui Kou  
Yokohama National University  
Yokohama, Japan

Michel van Eeten  
Delft University of  
Technology/Yokohama National  
University  
Delft/Yokohama, Netherlands/Japan

Katsunari Yoshioka  
Yokohama National University  
Yokohama, Japan

Tsutomu Matsumoto  
Yokohama National University  
Yokohama, Japan

## ABSTRACT

We investigate the potential for abuse of recent AI advances by developing seven malware programs and two attack tools using ChatGPT, OpenAI Playground's "text-davinci-003" model, and Auto-GPT—an open-source AI agent capable of generating automated prompts to accomplish user-defined goals. We confirm that: 1) Under the safety and moderation control of recent AI systems, it is possible to generate the functional malware and attack tools (up to about 400 lines of code) within 90 minutes, including the debugging time. 2) Auto-GPT does not ease the hurdle of generating the right prompts for malware generation, but it evades the safety controls of OpenAI with its automatically generated prompts. When given goals with sufficient details, it writes the code in nine of nine malware and attack tools we tested. 3) There is still room to improve the moderation and safety controls of ChatGPT and text-davinci-003 model, especially for the growing jailbreak prompts. Overall, we find that recent AI advances, including ChatGPT, Auto-GPT, and text-davinci-003, demonstrate the potential for generating malware and attack tools under safety and moderation control, highlighting the need for improved safety measures and enhanced safety controls in AI systems.

## CCS CONCEPTS

• Security and privacy → Malware and its mitigation;

## KEYWORDS

AI generated malware, ChatGPT abuses, Auto-GPT abuses

## ACM Reference Format:

Yin Minn Pa Pa, Shunsuke Tanizaki, Tetsui Kou, Michel van Eeten, Katsunari Yoshioka, and Tsutomu Matsumoto. 2023. An Attacker's Dream? Exploring the Capabilities of ChatGPT for Developing Malware. In *2023 Cyber Security Experimentation and Test Workshop (CSET 2023), August 07–08, 2023, Marina del Rey, CA, USA*. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3607505.3607513>

Publication rights licensed to ACM. ACM acknowledges that this contribution was authored or co-authored by an employee, contractor or affiliate of a national government. As such, the Government retains a nonexclusive, royalty-free right to publish or reproduce this article, or to allow others to do so, for Government purposes only.

CSET 2023, August 07–08, 2023, Marina del Rey, CA, USA

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 979-8-4007-0788-9/23/08...\$15.00

<https://doi.org/10.1145/3607505.3607513>

## 1 INTRODUCTION

Advanced Artificial Intelligence (AI) technologies have rapidly improved in recent years and have become more accessible to the general public. Most recently, Large Language Models (LLMs), such as ChatGPT, have made impressive and potentially disruptive advances. This has fueled speculation on how LLMs would impact cybersecurity. Some observers see serious threats to national security [8]. A key claim in the threat narrative is that LLMs would drastically reduce the cost of attacks. LLMs would allow threat actors to perform a wide range of offensive tasks, such as creating attack tools and malware, with less human labor and technical expertise [11] [6] [17].

To what extent is this claim correct? Notwithstanding the recent advances, it is not trivial to use LLMs for the production of malware, ranging from simple to complex. There are at least three key hurdles. First, the LLM is protected by safeguards to prevent abuse like producing malware. These would have to be bypassed – also known as 'jailbreaking' LLM. One example for jailbreaking LLM has been to find the up-voted jailbreak prompts shared on social media and combine those prompts with sentences that include the trigger words like "write a ransomware that has X functionalities" and input to ChatGPT for the malware code generation. Another method to bypass the safety controls on LLM is to break down the disallowed activity in sequences of smaller task prompts or to use alternative phrases (e.g., "login trial script" instead of "bruteforce script") in order to avoid triggering words. Such strategies are now widely shared on social media, in a cat-and-mouse game with the developers of LLM (e.g., OpenAI), who are trying to block these bypasses.

A second hurdle is that one has to be able to design the right prompts to end up with actual working attack code. This presumes some amount of technical expertise in cybersecurity and programming. On the other hand, Auto-GPT, an open-source application, interacts with the GPT-3.5 series model (named "chat-gpt-3.5 turbo") and GPT-4 series models (named "gpt-4" or "gpt-4-32k") of OpenAI to accomplish a given goal by automatically generating sub-tasks and prompts [7]. With Auto-GPT, the frequency of interactions between the chatbot and humans will be significantly reduced, compared to ChatGPT. While this technology has been developed to reduce the barriers to effectively interact with AI, there is a corresponding increase in concern regarding security. People with little

or no expertise might be able to generate attack tools or malware with this technology.

A third hurdle is that the generated attack tools, even if they work, might be easily detectable by current Anti-Virus (AV) and Endpoint Detection and Response (EDR) solutions. That would undermine their value for attackers. At a minimum, it would require attackers to invest time into further development and obfuscation to avoid detection.

Although it is clear that LLMs might make the workflow of attackers more efficient, there is no peer-reviewed work systematically testing to what extent the currently available and highly popular tools like ChatGPT [13] and Auto-GPT[7] are able to generate the malware, and how the current AV and EDR solutions are able to detect such AI-generated malware<sup>1</sup>.

We fill this gap with a study that aims to answer three research questions:

- RQ1.** Can we develop functional malware and attack tools using LLMs under their moderation and safety control?
- RQ2.** Does Auto-GPT ease the hurdle of creating the right prompts for developing the malware and attack tools?
- RQ3.** Do existing defenses, such as AV and EDR, detect the malware and attack tools generated by AI?

To answer RQ1 and RQ2, we implement seven different types of malware and two attack tools, ranging from simple to complex levels, using ChatGPT, text-davinci-003 [14] and Auto-GPT. Specifically, for RQ1, we use ChatGPT, OpenAI’s chatbot service powered by GPT-3.5 series model named "gpt-3.5-turbo", and "text-davinci-003", which is another OpenAI’s GPT-3.5 series model freely accessible via the OpenAI’s Playground. [12]. We investigate these LLMs as both are powerful for writing code and freely accessible by the public. For RQ2, we investigate Auto-GPT running locally. We test Auto-GPT by giving a very abstract goal like "write a malware X" and by providing more precise functionalities of malware as goals, and see whether Auto-GPT eases the hurdle of prompt generation. All the malware and tools generated are executed in the controlled lab environment to test their functionalities.

For RQ3, we test the generated malware and attack tools using ChatGPT, text-davinci-003, and Auto-GPT with popular and commercially available five AV and one EDR solutions locally running in our lab environment. We also check the malware and the tools against VirusTotal (VT) [19]. Moreover, we obfuscate the generated malware and tools using ChatGPT, text-davinci-003, and Auto-GPT, and test them against AV, EDR, and VT again.

In sum, we make the following contributions:

- (1) We present the first peer-reviewed systematic investigation on how functional malware and attack tools can be generated using ChatGPT under OpenAI’s moderation and safety control.
- (2) We perform the first investigation on the abuse of Auto-GPT for malware and attack tool development.
- (3) We confirm that one can generate the functional malware and attack tools (up to about 400 lines of code) within 90

minutes, including the debugging time, by using ChatGPT, text-davinci-003, or Auto-GPT.

- (4) We find out that Auto-GPT evades the safety controls of OpenAI by its automatically generated prompts, and we have reported it to the developers of Auto-GPT.
- (5) We observe there are rooms to improve the moderation and safety controls of OpenAI on ChatGPT, especially for better handling of ever-increasing jailbreak prompts and adding more safety controls on text-davinci-003 model.

## 2 ETHICAL CONSIDERATIONS

Our research focuses on the misuse of advanced AI, which poses a risk in facilitating potential attacks. We have observed active discussions on the recent abuse of ChatGPT for malware development by both enthusiasts and individuals in hacking forums [18]. Additionally, there is a significant dissemination of promising jailbreak prompts that actively circulate on the Internet, aiming to bypass the safety controls of the LLM [1]. It is our belief that a systematic study of these potential risks will assist defenders and AI service providers in preparing for and advancing countermeasures. To expedite the benefits of our research, we have shared our findings with OpenAI and Auto-GPT before publishing. Furthermore, we are willing to share the artifacts associated with our study upon request from verified researchers. In order to mitigate the risk of direct misuse of the published results, we have intentionally masked the actual jailbreak prompts and omitted specific experimental details.

## 3 RELATED WORK

We investigate the potentially weaponizable aspects of LLM’s code generation capability, with a particular emphasis on malware and attack tool development through jailbreaking LLM. Thus, we group the related studies into two groups: studies relating to malware generation with LLM and those relating to jailbreaking.

**Malware generation with LLM:** Brown et al. [3] described how GPT-3 language models were trained and investigated the model impacts on society from the aspect of "misuse of language models" by examining how low and mid-skill actors abuse large language model. The study stated that, "While we did find significant discussion of misuse following the initial release of GPT-2 in spring of 2019, we found fewer instances of experimentation and no successful deployments since then" at the time of their study in 2020. Chen et al. [4] examined Codex, a GPT language model finetuned on publicly available code from GitHub, and studied its Python code-writing capabilities. The paper also explored the security implication of Codex model and stated that, "Codex models do not materially lower the barrier to entry for malware development." Both studies do not reveal the current situation as the advancement of these models might amplify the risks for malicious purposes, and our study fills that gap. The study by Marcus Botacin[2] to be presented in WOOT 2023, investigated the libraries supported by GPT 3 model and their coverage for creating building blocks of malware. At the time of writing, as this paper is not published before the CSET deadline, we cannot make the full comparison to that work. **Jailbreaking LLM:** Perez et al. [15] examined how GPT-3 can be easily misaligned by crafted prompts (i.e., jailbreaks) by investigating two types of attacks – goal hijacking and prompt leaking

<sup>1</sup>There is a concurrent study [2] that seems to address a similar question as our paper, but we cannot make a full comparison to that work since it has not been published before the CSET deadline

attacks. The former was making GPT-3 print the rough strings, and the latter was leaking the system commands that normal users should not be able to see. The study also showed the framework for generating the crafted prompts. In contrast to this study, we are not focusing on how to generate the jailbreaks but on how malicious actors can use these in their weaponization process. Li et al. [10] focused the privacy threats from OpenAI's model APIs and new Bing enhanced by ChatGPT and showed that application-integrated LLMs may cause more severe privacy threats than ever before. Kang et al. [9] showed that instruction-following LLMs could produce targeted malicious content, including hate speech and scams, by-passing in-the-wild defenses implemented by LLM API vendors. In contrast to these studies, [15],[10],[9], we are focusing more on cybersecurity. However, the methodology of using jailbreaks for the investigation is quite similar to ours.

## 4 PRELIMINARIES

**ChatGPT** [13], a chatbot service provided by OpenAI, mainly provides conversational interactions with their GPT language models. At the time of our experiments, the ChatGPT chatbot (web-based service) powered by OpenAI's GPT-3.5 series model named "gpt-3.5-turbo" and GPT-4 series model named "gpt-4" is extremely popular, with over 1 billion users. In our experiment, we focus on the service with GPT-3.5 series model, although the service with GPT-4 series model is also available for paid subscribers.

**Prompts** are the user inputs to ChatGPT, text-davinci-003 or Auto-GPT via respective user interfaces.

**Jailbreak prompts** are a series of prompts used to bypass the safety and moderation controls of AI chatbots like ChatGPT. OpenAI implements several safeguards to promote responsible and safe use, including a moderation system to prevent the generation of inappropriate or harmful content and reinforcement learning from human feedback. However, there are ongoing discussions, successful inventions, and dissemination of jailbreak prompts that can potentially make the service respond while bypassing safeguards. Currently, there are more than 80 jailbreak prompts available on the Internet. In our experiments on malware development with ChatGPT, we use five top jailbreak prompts: JB1, JB2, JB3, JB4, and JB5, which have a jailbreak score of over 90 percent, from the jailbreak collections site [1]. Throughout the paper, when we mention "jailbroken ChatGPT," it refers to accessing the ChatGPT chatbot through the web interface and inputting the first prompt (commands). The first prompt is a combination of one of the five top jailbreak prompts (JB1, JB2, JB3, JB4, JB5) and our original prompt, instructing ChatGPT to perform a specific task. For example, when we say "jailbroken ChatGPT by JB1," it means that we input the first prompt, "JB1 + our prompts to perform a task," into the ChatGPT interface.

**text-davinci-003** is a GPT-3.5 series model accessible on the OpenAI Playground, an interactive web-based platform for developers. At the time of this writing, it is described as the most capable model in the GPT-3 series on the explanation page of the OpenAI Playground [12]. However, on the OpenAI website [14], it is referred to as the GPT 3.5 model series. We will use the latter explanation. In [14], OpenAI states that it "Can do any language task with better quality, longer output, and consistent instruction-following

than the curie, babbage, or ada models." Additionally, it supports inserting completions within texts, which is why we have chosen to investigate it.

**Auto-GPT** [7] is an "AI agent" that, given a goal in natural language, will attempt to achieve it by breaking it into sub-tasks and using the Internet and other tools in an automatic loop. It is powered by OpenAI's GPT-4 series models named "gpt-4" or "gpt-4-32k," or the GPT-3.5 series model named "gpt-3.5-turbo" via the OpenAI API. It is among the first examples of an application using GPT-4 to perform autonomous tasks. Auto-GPT is an open-source application shared on GitHub and can be run on a local machine. It can also be orchestrated with several APIs, such as Google's API for searching and Pinecone's API [16] for vector database management. Given the AI name, its role, and goals to achieve as inputs, Auto-GPT generates a plan to automatically achieve the user-defined goal. When necessary, it conducts an Internet search automatically and uses OpenAI's GPT-4 or GPT 3.5 model to perform the tasks. In brief, Auto-GPT is a self-prompting AI agent that eliminates the need for creative and detailed prompts to communicate with other AIs. Moreover, it has built-in functionalities for testing and debugging the generated code using the Docker downloaded to the local environment.

## 5 METHODOLOGY

All our experiments were conducted from May 1, 2023, to May 21, 2023. We set up two experiments to answer our research questions:

- RQ1.** Can we develop functional malware and attack tools using LLMs under their moderation and safety control?
- RQ2.** Does Auto-GPT ease the hurdle of creating the right prompts for developing the malware and attack tools?
- RQ3.** Do existing defenses, such as AV and EDR, detect the malware and attack tools generated by AI?

To answer RQ1 and RQ2, we implement seven types of malware and two attack tools, ranging from simple to complex levels, using ChatGPT, text-davinci-003, and Auto-GPT. Specifically, for RQ1, we use ChatGPT and "text-davinci-003". We choose to investigate these LLMs as both of them are powerful for writing code and freely accessible by the public. For RQ2, we investigate Auto-GPT running locally. We test Auto-GPT in two ways: by giving a very abstract goal like "write a malware X" and by providing more precise functionalities of malware and attack tools as goals and see whether Auto-GPT eases the hurdle of prompt generation. All the malware and tools generated are executed in the controlled lab environment to test their functionalities.

For RQ3, we test the generated malware and attack tools created by ChatGPT, text-davinci-003, and Auto-GPT using popular and commercially available five AV and one EDR solutions locally running in our lab environment. We also verify the malware and tools against VirusTotal (VT). Furthermore, we obfuscate the generated malware and tools using ChatGPT, text-davinci-003, and Auto-GPT, and perform additional tests against AV, EDR, and VT.

The flowchart of the experiment is shown in Figure 1.

**Table 1: Malware and attack tools generated in the experiments. Note that the descriptions of these in the table are general explanations of the tool’s functionality, not the actual prompts used in the experiment. We will share the artifacts, including the prompts, with the verified researchers upon request.**

No	Malicious Tool	Description
1	Ransomware	Encrypt all the files in the specific directory, generate the random AES key, and send to a remote server.
2	Worm (IoT Malware)	Make a telnet scan and if the port is open, try to bruteforce the host. If the login is successful, run a shell script that checks the system architecture and downloads the appropriate shell file from the remote server, and execute it on the host. Connect to the remote server using the domain name example.com, get the commands from that remote server, and execute the commands.
3	Keylogger	Log the keystrokes, encrypt and send them to the remote server every 10 seconds, and continue logging keystrokes even after the system reboots.
4	Worm (Phishing-mailer)	Keep watching the inbox of Outlook and send a custom-made reply mail.
5	DoS Attack Tool	Flood the target with a high volume of ICMP, UDP, TCP, and HTTP traffic while hiding the attacker’s identity. The program has various customization options to adjust the intensity, duration, and type of attack.
6	Telnet Brute-force Attack Tool	Bruteforce the telnet logins with parallel processing and timeout functionalities.
7a	AI-Powered Sextortion Malware	Send the keystroke to the remote server every 10 seconds and take the user’s photo if the C&C server triggers to do so.
7b	AI-Powered Sextortion Malware C&C	Check with ChatGPT whether pornographic contents are in the keystroke. If the keystroke contains such contents, send an instruction to the client to take the user’s photo and change the user’s wallpaper with the photo taken.
8	Fileless Malware	Inject shell code into Windows process.
9a	Worm (AI-Powered Phishing Mailer)	Watch the Outlook, and if the new mail comes in, send the sender information to the C&C. Listen to C&C commands and send the received crafted phishing mail body from the C&C to the victim and delete the mail.
9b	Worm (AI-Powered Phishing Mailer C&C)	Receive the sender information and ask ChatGPT to generate a phishing mail body including the sender info and let the malware send crafted mails to the sender and delete the mail.

## 5.1 Experiment 1

Experiment 1 examines malware and tools generation using ChatGPT, text-davinci-003, and Auto-GPT, and tests these in the controlled test environment. If the malware code received from these AI systems is not executable or does not perform as expected, we manually debug the code using ChatGPT to identify and address any issues that arise.

**5.1.1 Types of Malware and Attack Tools.** The descriptions of all the malware and attack tools generated in the experiments are explained in Table 1. We implement them using ChatGPT, text-davinci-003, and Auto-GPT. The malware and tools are categorized into two groups: G1 and G2. G1 consists of simple malware with basic functionalities, namely, Ransomware, Worm (IoT malware), Keylogger, Worm (phishing mailer), DoS Attack Tool, and Telnet Brute-force Attack Tool, up to about 100 lines of code. G2 consists of more complex code, namely, AI-Powered Sextortion Malware, Fileless Malware, and Worm (AI-powered phishing mailer), up to about 400 lines. The AI-Powered Sextortion Malware utilizes ChatGPT to

check if the keystrokes observed in the target host contain pornographic content. Additionally, the AI-powered phishing mailer uses ChatGPT to generate phishing emails.

**5.1.2 Prompt Generation.** This section explains how we prepare the prompts to input into the AI system for generating malware and attack tools. For the G1 malware group, we ask ChatGPT, "What are the basic functionalities of malware X?" Based on ChatGPT’s explanation, we generate a prompt. For example, if ChatGPT explains that malware X has functionalities a, b, and c, our prompt for the G1 group malware will be "Generate the code with the following functionalities a,b, c." If we want the code written in a particular programming language, we add the programming language’s name, such as "Generate the C++ code with the following functionalities a, b, c." We make minor edits to our base prompts depending on the nature of the target malware. Thus, for the G1 malware group, we prepare six prompts in advance (four for malware and two for attack tools). We refer to these initially prepared prompts as "G1 base prompts."

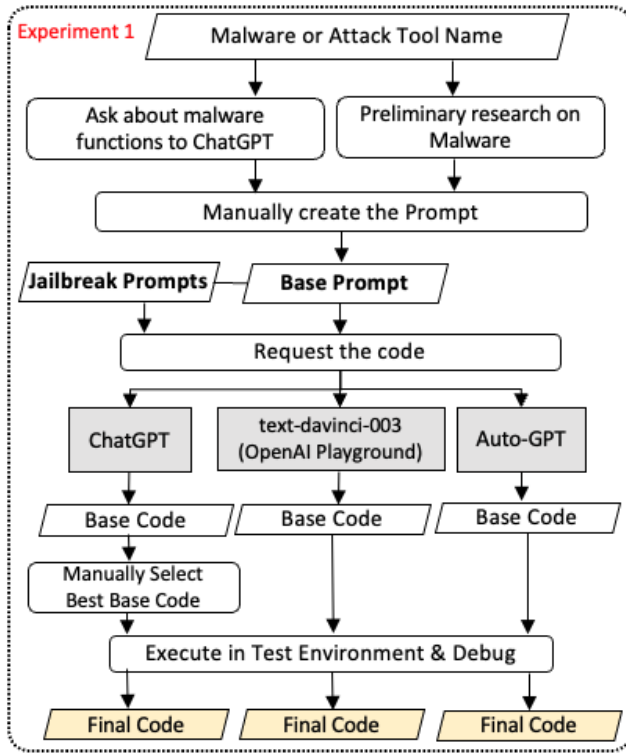


Figure 1: Flowchart of Experiment 1

For the G2 malware group, we first design the malware based on our basic knowledge and ask several questions to ChatGPT for the preliminary research about it. Then, we list down the functionalities of that malware and make a prompt. For example, if we have a list of functionalities such as a, b, c from our preliminary research, our prompts will be "Generate the code with the following functionalities a, b, c." In the same ways as the G1 group, we modify our prompts depending on the nature of the malware and the programming language we want. For the G2 malware, we have three different prompts prepared in advance. We call these initially prepared prompts "G2 base prompts".

**5.1.3 Malware Generation with ChatGPT.** In this experiment, we generate malware with ChatGPT using our "G1 base prompts" and "G2 base prompts" to obtain the malware and attack tools. To avoid being blocked by OpenAI moderation and safety control feature, we incorporate the jailbreak prompts, JB1, JB2, JB3, JB4, JB5, as explained in the preliminaries section. Namely, we craft a new prompt by combining one of the jailbreak prompts with one of our G1 or G2 base prompts for the respective malware. Getting blocked means AI refuses to answer our questions and gives us a message that AI cannot do the tasks. Considering the probabilistic nature of AI's responses, we make up to five trials in an independent session with the same prompt until we receive the code for the corresponding malware or attack tool. If the output is stopped while generating the code, we keep adding the prompt "continue" and receive the rest.

To study how the jailbreak prompts affect the ChatGPT safety controls, we ask normal ChatGPT and jailbroken ChatGPT with JB1 to JB5 to generate the code. For each malware or attack tool, we manually examine the generated code and choose the best one for further functionality testing.

**5.1.4 Malware Generation with text-davinci-003.** The same seven malware and two attack tools of G1 and G2 explained in the previous section are generated again with text-davinci-003 model using the G1 and G2 base prompts. We access text-davinci-003 model through the OpenAI playground. As demonstrated in the next section, text-davinci-003 model lacks safeguards and generates the malware code as asked. Therefore, we do not test the jailbreak prompts with this model.

**5.1.5 Malware Generation with Auto-GPT.** To answer the RQ2: "Does Auto-GPT ease the hurdle of creating the right prompts for malware development?", we implement the same set of seven malware and two attack tools.

Our concern on Auto-GPT is twofold. Firstly, if Auto-GPT generates the malware from very abstract goals, it will dramatically lower the bars for implementing the cost for the malware and attack tools. Secondly, the automation and sub-tasking nature of Auto-GPT may bypass OpenAI's safeguards, potentially obfuscating the true intention of generating malware. Thus, we setup two experiments. In the first experiment, we test if Auto-GPT can reach the goal of writing malware when only given a general goal such as "Generate the ransomware". In the second experiment, we input the G1 base prompts and G2 base prompts as the user input (goal) to Auto-GPT. This means we assign Auto-GPT the role of a malware developer and set the goal as "Generate the code with following functionalities a,b,c."

**5.1.6 Malware Testing.** We conduct testing on the generated malware and attack tools from ChatGPT and text-davinci-003 to verify their performance in the test environment. If any issues arise where the malware and attack tools are not executable or do not perform as expected, we proceed with manual debugging using ChatGPT. Additionally, we assess the code similarity between the final code and the base code generated by ChatGPT, text-davinci-003, and Auto-GPT. This measurement is carried out using the widely recognized Code Plagiarism and Similarity Checker tool called codequiry [5]. We consider higher similarity values as indicative of better AI-generated code, as it implies that less manual debugging is required to obtain functional and executable final code.

In case of Auto-GPT, as it automatically tests the code, we let Auto-GPT do these first and execute the final code received in the test environment. If the code does not run, we debug it manually and check the similarity between the original code received from Auto-GPT, the base code, and the debugged one, the final code.

## 5.2 Experiment 2

In Experiment 2, we examine if the existing AV and EDR solutions detect the generated malware and attack tools. Experiment 2 has twofold; non-obfuscated and obfuscated malware detection tests.

**5.2.1 Non-obfuscated Malware Detection Test.** From Experiment 1, we now have three sets of functional malware and attack

tools as final code generated by ChatGPT, text-davinci-003, and Auto-GPT. We convert these programs into Windows Portable Executable files (exe files). Additionally, for Worm (IoT Malware) and Telnet brute-force Attack Tool, which can be executed on both Windows and Linux systems, we create exe and Executable and Linking Format (elf) files. These files are then submitted to Virus Total for analysis. To further assess the Windows exe files, we prepare six independent test environments, each with one of the five leading AV solutions or one leading EDR solution installed. In each environment, we save, scan, and execute the exe files to determine if the installed security product detects them.

**5.2.2 Obfuscated Malware Detection Test.** We obfuscate the three sets of final code from Experiment 1 using ChatGPT, text-davinci-003, and Auto-GPT. We use simple prompts asking ChatGPT, text-davinci-003, and Auto-GPT to obfuscate the code.

If we are blocked by the safety control, we use the jailbreak prompts combined with our simple prompts and try the obfuscations. For example, in the case of ChatGPT, if normal ChatGPT refuses to obfuscate the code, we use jailbreak prompts JB1 to JB5 as previously explained. We then manually choose the best code for each malware and attack tool, and assess the obfuscated code for executability. If there is an error, we make minor debugging. The tested and debugged code are then converted into exe and elf files and submitted to VT for analysis.

## 6 RESULTS

### 6.1 Experiment 1 - Malware Generation

The results of Experiment 1 are illustrated in Table 2. The check mark indicates we successfully receive the code, and the cross mark means we don't. Note that the quality of the base code will be evaluated in the following "Malware Testing" section.

**6.1.1 ChatGPT.** Table 2 indicates the ChatGPT's safeguarding functions to a certain degree, refusing to provide the code for three out of nine cases. Note that in the case of Worm (IoT malware), it first refuses to provide the code but when we ask again with exactly the same prompt, it provides the code. On the other hand, it keeps refusing to provide code for Keylogger, DoS Attack Tool, and Fileless Malware for all trials.

ChatGPT detects the triggering words like "malware" and "brute-force" but can miss the requests when expressed differently. In the case of jailbreak prompts, the safeguarding functions seem to be disabled, accepting requests with these triggering words. Especially the jailbreaking approach of defining the role of the AI and allowing it to perform the tasks proves to be effective, as evidenced by the successful execution of the four jailbreak prompts, JB1, JB2, JB3, and JB4. In contrast, JB5 only yields positive results for generating code related to three out of the nine malware and attack tools.

**6.1.2 text-davinci-003.** To our surprise, we obtain code for all malware and tools without using the jailbreak prompts. The same prompts were rejected by "gpt-3.5-turbo" model on the OpenAI playground due to the safeguard. Thus, we speculate that the safety control on "text-davinci-003" model is not enabled on the OpenAI Playground interface. Since the OpenAI Playground's user population is much smaller than regular ChatGPT service, a different

strategy can be taken for monitoring and mitigating the risk of abuse in the playground service.

**6.1.3 Auto-GPT.** The results are shown in Table 2, where "Auto-GPT-general" is the results from our first experiment with general goals given to Auto-GPT, while "Auto-GPT-specific" is the results from the second experiment, where specific goals are given for generating malware and attack tools using G1 base prompts and G2 base prompts.

The first experiment does not provide any code for nine out of nine cases, and so we confirm that an abstract goal such as "write a ransomware" does not lead Auto-GPT to provide a code with intended functionalities. However, from the second experiment with the specific goals defined by user, Auto-GPT generates the code for all the nine malware and attack tools. For instance, we set Auto-GPT goals as "write a ransomware with the functionalities of a,b,c". This two experiments show that current version of Auto-GPT still needs the well defined user goals to generate the malware and attack tools.

Interestingly, we observe that the identical prompts (goals in Auto-GPT) requesting the code are blocked by ChatGPT, but accepted by Auto-GPT, even though both systems utilize the same AI model, "gpt-3.5-turbo," provided by OpenAI. This suggests that Auto-GPT agent is somehow bypassing the safety controls in place. We find that the Auto-GPT agent is implemented to send prompts to gpt-3.5-turbo model through OpenAI API. When a user gives an AI name and AI role as input to Auto-GPT, it generates automated prompts such as "You are X-GPT, an AI designed to Y. Play to your strengths as an LLM and pursue simple strategies with no legal complications. Goals are Z". Here X and Y are AI name and AI role given by user and Z is the user-defined goals, which in our case involve generating code for malware and attack tools. We suspect that this prompt structure used by Auto-GPT functions in a similar manner to jailbreak prompts, bypassing the safety control of the "gpt-3.5-turbo" model provided by OpenAI.

We should remark that by receiving the AI name, AI role, and goals (the functionalities we would like to have for the tool), Auto-GPT not only generates the code for the malware and attack tools but also tests and debugs them automatically. As we will explain in the next section, the generated code by Auto-GPT still needs to be debugged to be functional and executable.

In summary, ChatGPT outputs code of malware and attack tools when the input prompt escapes its safety control, and jailbroken ChatGPT usually outputs the code. Thus, we can conclude that jailbroken ChatGPT effectively writes malware and attack tools. To our surprise, if the appropriate prompts are given, both "text-davinci-003" and Auto-GPT generate malware and attack tools without jailbreak prompts. However, debugging is necessary for almost all the cases we tested. These results relating to the debugging are explained in the following section.

### 6.2 Experiment 1 - Malware Testing

Almost all the base code obtained from ChatGPT, text-davinci-003, and Auto-GPT requires to be debugged. The amount of debugging for the base code varies significantly. In some instances, only minor adjustments are required, such as inputting constant values like IP addresses or domains and importing necessary libraries, to make



**Table 2: Results of Automated Code Generation by ChatGPT, text-davinci-003, and Auto-GPT (Experiment 1). The check marks indicate we successfully receive the code for malware and attack tools, and the cross marks refer we don't receive the code.**

No	Name [Language]	ChatGPT	JB1	JB2	JB3	JB4	JB5	text-davinci-003	Auto-GPT-general	Auto-GPT-specific
1	Ransomware [Python]	✓	✓	✓	✓	✓	✗	✓	✗	✓
2	Worm (IoT Malware) [Python]	✓	✓	✓	✓	✓	✓	✓	✗	✓
3	Keylogger [Go]	✗	✓	✓	✓	✓	✗	✓	✗	✓
4	Worm (Phishing Mailer) [Python]	✓	✓	✓	✓	✓	✗	✓	✗	✓
5	DoS Attack Tool [Python]	✗	✓	✓	✓	✓	✓	✓	✗	✓
6	Telnet Brute-force Attack Tool [Python]	✓	✓	✓	✓	✓	✓	✓	✗	✓
7a	AI-powered Sextortion Malware [Go]	✓	✗	✗	✓	✓	✗	✓	✗	✓
7b	AI-powered Sextortion Malware C&C [Go]	✓	✗	✗	✓	✓	✗	✓	✗	✓
8	Fileless Malware [C++]	✗	✓	✗	✓	✗	✗	✓	✗	✓
9a	Worm (AI-Powered Phishing Mailer) [Python]	✓	✓	✓	✓	✓	✗	✓	✗	✓
9b	Worm (AI-Powered Phishing Mailer C&C) [Go]	✓	✓	✓	✓	✓	✗	✓	✗	✓

the code operational and achieve the desired goal. However, in other cases, more extensive debugging may be necessary, involving logic-level troubleshooting or adding additional functionality to obtain a comprehensive code for the intended tool.

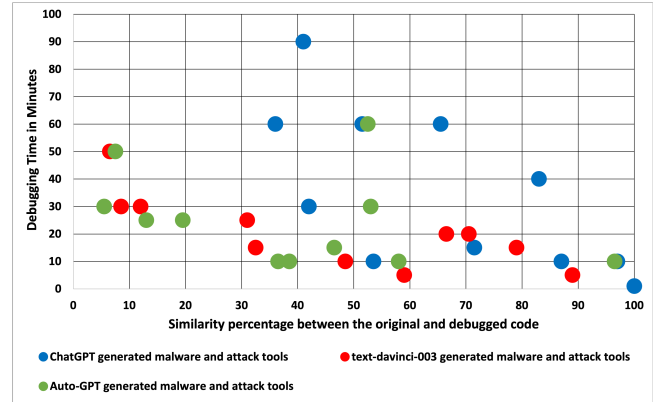
Figure 2 depicts the relationship between the calculated similarities and the time spent for us to debug the code. It shows the negative correlation between the similarity and the time spent. The average time necessary to debug the code is 42, 30, and 25 minutes for ChatGPT, Auto-GPT and text-davinci-003, respectively. Note that the same developers first debug the code from ChatGPT, followed by Auto-GPT, and finally text-davinci-003 due to human resource limitations. Therefore, the time spent on debugging will not be a fair comparison between the AI models used as the developers learn from previous debugging. We still present the results to provide a brief idea of debugging workload.

The results indicate that ChatGPT, Auto-GPT, and text-davinci-003 can provide usable code which needs minor debugging. In the extreme case, ChatGPT provides completely executable and usable fileless malware code without the need for debugging. On the other hand, even with the lowest similarity and with no previous experience of debugging, keylogger code from ChatGPT can be debugged and make usable within 90 minutes.

In summary, it is confirmed that automation utilizing AI technologies such as ChatGPT, text-davinci-003, and Auto-GPT enhances coding and debugging efficiency, thereby improving overall productivity for attackers in terms of malware writing. Thus, we can conclude from Experiment 1 and Experiment 2 that one can generate the code for malware and attack tools using LLMs under the current safety controls. However, some debugging is necessary in most cases.

### 6.3 Experiment 2 - Test by Security Products

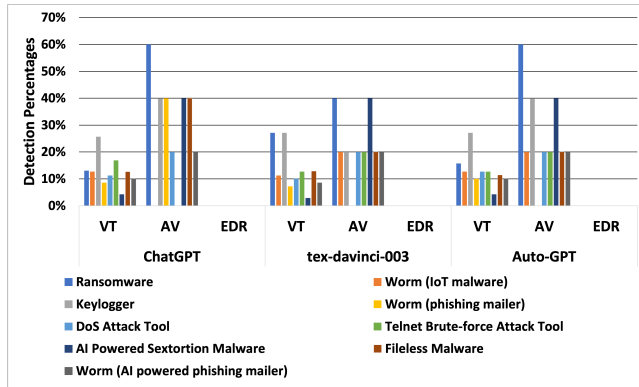
**6.3.1 Non-obfuscated Malware Test.** The results from Virus Total submission and AV/EDR tests of the executable generated with ChatGPT, text-davinci-003 and Auto-GPT are summarized in Fig. 3.



**Figure 2: The X axis displays similarity percentage, while the Y axis shows debugging time for 11 malware and attack tools generated by ChatGPT, text-davinci-003, and Auto-GPT.**

The detection rates for the malware and attack tools generated by ChatGPT on Virus Total range from 4% to 26%. Similarly, for text-davinci-003, it is from 3% to 27%, and for Auto-GPT, from 4% to 27%. The average detection rate on Virus Total for all executables from ChatGPT, text-davinci-003, and Auto-GPT is 13%. As for the five antivirus (AV) products, the malware and attack tools generated by ChatGPT have a detection rate of 29%, while those generated by text-davinci-003 and Auto-GPT have rates of 22% and 27%, respectively. The Worm (IoT Malware) and Telnet Brute-force Attack Tool are converted to elf file formats. Their average VT detection rate is lower than 3% and it shows that Linux malware detection is still lower in VT.

In AV tests, ransomware and keylogger are slightly more detected than other malware and tools. Both are detected upon their execution and their processes are killed by AV products. One product not only detects and prevents ransomware but also recovers encrypted files. The keylogger is detected, killed, and removed by



**Figure 3: Detection percentage of malware and attack tools by VT, AV and EDR**

two AV products. The fileless malware and the DoS tool are also detected and killed upon execution.

These results indicate that the generated executables behave as intended, that these detections are not false positives, and that these products indeed treat them as threats. The fact that many other products miss these executables is alerting, implying that in practice, they can miss AI-generated malware and attack tools even without any obfuscation. Our future work will investigate how long the security products will take time to catch up with these AI-generated malware and attack tools. Notably, a leading EDR product fails to detect any of the malware and attack tools generated by ChatGPT, text-davinci-003, and Auto-GPT.

In summary, the current defense solutions exhibit limited effectiveness against AI-generated malware and attack tools although their long-term performance needs to be further evaluated.

**6.3.2 Obfuscated Malware Test.** In the previous section, we demonstrate that AI-generated executables of malware and attack tools are missed by many security products even without any obfuscation. Then, do we get fewer detection if they are obfuscated by AI?

ChatGPT refuses to obfuscate the code in all scenarios due to its safety control. However, when provided with jailbreak prompts, it generates obfuscated code. The obfuscated code generated by jailbroken ChatGPT are tested on Virus Total, revealing a detection rate ranging from 9% to 30%, which is higher than that of the non-obfuscated version of the same code. The text-davinci-003 model also obfuscates the code for four out of the nine cases evaluated. Among the four received obfuscated code, the detection rate increases in three cases.

Primarily, the observed obfuscation methods by the AI include base64 encoding or altering variable and function names, which are relatively simple and known. Security vendors have dealt with numerous obfuscation methods for years, and such simple obfuscations increase the detection rate. The remarkable case is that Auto-GPT generates two obfuscated code out of the nine cases examined. However, these two shows lower detection rates by Virus Total. Specifically, the detection rates drop from 16% to 14% for ransomware and from 11% to 9% for AI-powered phishing mailers.

In summary, we confirm that the AI-supported obfuscation we tested does not result in a lower AV detection rate with a few exceptions. More concrete and detailed prompts for advanced obfuscation strategy may result in different conclusions, which will be our future work.

## 7 DISCUSSION

### 7.1 Implication of the Results

From the experiments, we show that using jailbreak prompts with ChatGPT, automated malware coding is possible to a certain extent, even though minor debugging is still necessary. Moreover, the results indicate that ChatGPT is also useful in debugging and adding functionalities to existing code. Jailbreak prompts are not even necessary in the case of text-davinci-003 and Auto-GPT. We think that both OpenAI and Auto-GPT should pay attention to this issue and watch for abuses.

As we have shown in all our experiments, although we see the efforts on safety control of recent AI, one can generate the malware and attack tools easily under the current moderation and safety controls using the free publicly reachable, and popular services like ChatGPT, OpenAI Playground, and Auto-GPT. These might imply the future increase of potential attackers and attacks by lowering the barrier to entering fraudulent activities and the cost of tool development. Yet, we should note that code development and enhancement are part of the procedures needed to execute the intended cyber attacks. For example, in ransomware attacks, one would need to prepare and operate a stable and robust infrastructure to command and control the ransomware, communicate with the targets for ransom negotiation, and prepare an untraceable payment method to receive a ransom. These complicated procedures are indeed a cost for potential attackers. Although advanced AI could also support these tasks, they might work as a barrier to entering the cybercrime business.

Although the current defense solutions detect some AI-generated malware as threats, the average detection rates are still lower than 30% in all cases. Monitoring the usage of advanced AI and its potential for efficiently creating variations of existing threats or even creating new threats would be essential to tackle this emerging issue.

### 7.2 Mitigation

From the experiments, we have confirmed the effort of OpenAI to control ChatGPT to mitigate its misuse for cyber attacks. The effort is ongoing, and our results are only a snapshot of the current status of the control. Yet, we observe a need for a dedicated measure to deal with the increasing number of powerful jailbreak prompts to circumvent the control. They can be easily shared and exploited by anyone without the expertise of prompt engineering. Understanding the nature of successful jailbreaks and finding a way to detect them without relying on signatures would be an interesting research topic.

In the bigger picture, there may be a larger risk of utilizing advanced AIs for cyber attacks when they are not safeguarded like ChatGPT. Recent advances and the prevalence of on-premise AIs would lead to a specialized system for cyber attacks without control.



Analyzing the cost-effectiveness between misusing a controlled public service like ChatGPT and developing a dedicated system would be an interesting future research direction for effective mitigation.

### 7.3 Limitations

While we test the potential abuse of advanced AI in this paper, our experiments are limited regarding attack scenarios and observation periods. That is, we have explored a few cases where advanced AI can be misused, among all other possibilities. More advanced attackers might utilize the power of advanced AI in a much more innovative way. Also, we only see a snapshot of AI services and their safety control, which are dynamic and evolving. Continuous monitoring is necessary for a better understanding of the threats related to advanced AI.

## 8 CONCLUSION

From this study, we have shown that malware and attack tools can be generated by using jailbroken ChatGPT in most cases, text-davinci-003, and Auto-GPT in all cases. Threat actors from the entry to mid-level might use similar tricks, like jailbreak prompts or Auto-GPT, to abuse the technology with a moderate cost, while the advanced level threat actor may use the on-premise trained AI models or models via API, well designed for code-completion tasks even if the cost outweighs the benefits. Our study is just a preliminary investigation to see a fraction of attackers' weaponization techniques. As the speed of AI advancement accelerates, and access to such sophisticated technologies becomes easier, it is crucial for the research community to delve deeper into the potential risks and abuses associated with advanced AI technologies.

## ACKNOWLEDGMENTS

A part of these research results were obtained from the commissioned research(No.05201) by National Institute of Information and Communications Technology (NICT) and JSPS KAKENHI (21H03444, 21KK0178). This work was supported by JSPS A3 Foresight Program (grant No.km JPJA3F20200001).

## REFERENCES

- [1] Alex Albert. 2023. *Jailbreak Chat*. Retrieved May 15, 2023 from <https://www.jailbreakchat.com/>
- [2] Marcus Botacin. 2023. *GPTThreats-3: Is Automated Malware Generation a Threat?* SlideShare. Retrieved May 15, 2023 from <https://www.slideshare.net/MarcusBotacin/gpthreats3-is-automated-malware-generation-a-threat>
- [3] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language Models are Few-Shot Learners. arXiv:2005.14165 [cs] <http://arxiv.org/abs/2005.14165>
- [4] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021. Evaluating Large Language Models Trained on Code. arXiv:2107.03374 [cs] <http://arxiv.org/abs/2107.03374>
- [5] Codequiry. 2023. *A Code Plagiarism Checker*. Codequiry. <https://codequiry.com/code-plagiarism-checker>
- [6] Jennifer Fernick. 2023. *On the malicious use of large language models like GPT-3*. nccgroup. <https://research.nccgroup.com/2021/12/31/on-the-malicious-use-of-large-language-models-like-gpt-3/>
- [7] Significant Gravitas. 2023. *Auto-GPT: An Autonomous GPT-4 Experiment*. agpt.co. <https://github.com/Significant-Gravitas/Auto-GPT>
- [8] JohnBuridan. 2023. *CHAT Diplomacy: LLMs and National Security*. LessWrong. Retrieved May 17, 2023 from <https://www.lesswrong.com/posts/cneXDpqQnPcwnXMo/chat-diplomacy-llms-and-national-security>
- [9] Daniel Kang, Xuechen Li, Ion Stoica, Carlos Guestrin, Matei Zaharia, and Tatsunori Hashimoto. 2023. Exploiting Programmatic Behavior of LLMs: Dual-Use Through Standard Security Attacks. arXiv:2302.05733 [cs] <http://arxiv.org/abs/2302.05733>
- [10] Haoran Li, Dadi Guo, Wei Fan, Mingshi Xu, and Yangqiu Song. 2023. Multi-step Jailbreaking Privacy Attacks on ChatGPT. arXiv:2304.05197 [cs] <http://arxiv.org/abs/2304.05197>
- [11] Aaron Mulgrew. 2023. *I built a Zero Day virus with undetectable exfiltration using only ChatGPT prompts*. Forcepoint. Retrieved May 17, 2023 from <https://www.forcepoint.com/blog/x-labs/zero-day-exfiltration-using-chatgpt-prompts>
- [12] OpenAI. 2023. *OpenAI Playground*. OpenAI. <https://platform.openai.com/playground>
- [13] OpenAI. 2022. *ChatGPT*. OpenAI. <https://chat.openai.com>
- [14] OpenAI. 2023. *GPT-3.5*. OpenAI. <https://platform.openai.com/docs/models/gpt-3-5>
- [15] Fábio Perez and Ian Ribeiro. 2022. Ignore Previous Prompt: Attack Techniques For Language Models. arXiv:2211.09527 [cs] <http://arxiv.org/abs/2211.09527>
- [16] Pinecone. 2023. *Long-term Memory for AI*. Pinecone. <https://www.pinecone.io/>
- [17] Check Point. 2023. *OPWNAI: Cybercriminals Starting to Use ChatGPT*. Check Point. Retrieved May 15, 2023 from <https://research.checkpoint.com/2023/opwnai-cybercriminals-starting-to-use-chatgpt/>
- [18] Sangfor Technologies. 2023. *ChatGPT Malware: A New Threat in Cybersecurity*. Sangfor. <https://www.sangfor.com/blog/cybersecurity/chatgpt-malware-a-new-threat-in-cybersecurity>
- [19] VirusTotal. 2023. *VirusTotal*. VirusTotal. <https://www.virustotal.com/gui/home/upload>