



ADMM Cybersecurity and
Information Centre of Excellence

UPDATE ON THE INFORMATION DOMAIN

Issue 02/23 (February)

Risks and Ethical Considerations in the Use of Generative Artificial Intelligence (AI)

INTRODUCTION

1. Generative artificial intelligence (AI) refers to the use of algorithms to create new content, including texts, images, audio and video. The emergence of generative AI that is capable of creating novel content rather than simply analysing or processing existing data, has resulted in a new wave of innovative applications and renewed interest in its development. From generating content for social media posts to writing music, AI has become the largest enabler for technological advancements in organisational processes and creative ideation.
2. Generative AI is based on the principle of generating content that is similar to, but not identical to, existing data. This is achieved through the use of deep learning algorithms and neural networks, which are typically trained on large datasets. With just a few instructive words or lines of text, AI can create any content the users want. Other than generating new content, AI also benefits industries that deal with vast amount of information. For instance, as *Financial Review* reported, AI could assist with writing prescriptions for doctors, explain medical bills to patients, and even assist with both writing and understanding legal documents.

3. AI technology is becoming increasingly popular with tools and services such as Open AI's ChatGPT, DALL-E, or Midjourney. As the popularity of generative AI grows, guidelines and suggestions have been put forth for the creators of these tools to take steps to mitigate the risks and their potentially harmful uses.

Unethical Practices of Generative AI Tools

4. According to *SCMP*, there is an increasing risk of unethical practices being used to generate content for various types of creative work. For instance, Midjourney is an AI image generation tool that takes inputs through text prompts and parameters, and uses a Machine Learning (ML) algorithm trained on a large amount of image data to produce unique images. This AI tool is very popular and used by many to generate images. However, the image data sets contained works from artists at all levels, without the original artists' approval or consent. This has raised questions of misuse and sparked an outrage among artists.

Phishing and Social Engineering

5. Another notable risk is the misuse of ChatGPT¹ to conduct malicious activities. Within the first few weeks of ChatGPT's launch, *Recorded Future* reported that threat actors on the dark web and special-access forums have begun sharing proof of concept conversations on the use of such tools for phishing, social engineering, malware development and disinformation. The ability of generative AI to imitate human language has raised concerns due to its potential for misuse as a phishing and social engineering tool. For example, ChatGPT is able to generate codes that can be used to mirror websites without the permission of the website owner, thereby creating realistic looking phishing pages to carry out an attack.

¹ GPT refers to Generative Pre-trained Transformer. ChatGPT is a neural network learning model which enables machine to perform natural language processing (NLP) tasks. Trained to respond to queries and instructions as a human would, ChatGPT is an intelligence bot that uses and analyses human language to generate content, provide information, write songs, and explain scientific concepts and so forth. According to *NYTimes*, ChatGPT could also automate comments submitted in regulatory processes, write letters to the editor for publications in local newspapers, and even mimic work of Russian Internet Research Agency.

Malware Development

6. *Recorded Future* reported that malicious actors can request ChatGPT to write codes in a number of different programming languages that exploit critical vulnerabilities. While ChatGPT often flag such requests as malicious, there are workarounds to “trick” ChatGPT into fulfilling such requests. By continuously querying the chatbot and receiving a unique piece of code each time, it is entirely possible for the AI tool to create polymorphic malwares that are able to evade antivirus detection. (See [Figure 1](#) for an example of AI-generated codes to encrypt files.) Users can also misuse ChatGPT to write malware payloads, including infostealers, remote access trojans, cryptocurrency clippers/drainers, crypters and more.

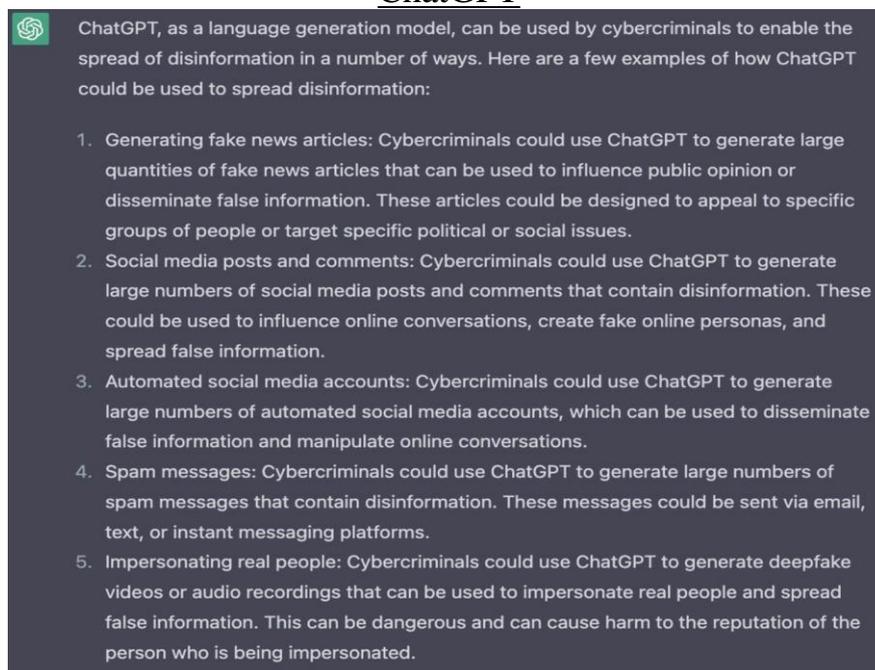
Figure 1: AI-Generated Codes to Encrypt File



Disinformation

7. With AI’s capabilities to accurately emulate human language and convey emotions, as well as create inauthentic content rapidly on a large scale, these generative AI tools can easily be weaponised by state actors, non-state actors and cybercriminals to spread disinformation. (See [Figure 2](#) on how ChatGPT described how ChatGPT itself can be used to spread disinformation.)

Figure 2: Examples of how ChatGPT could be used to spread disinformation according to ChatGPT



ASSESSMENT

8. In essence, AI does bring about many benefits across the different domains – from software engineering to creative, healthcare, and legal industries, amongst others. However, the lack of robust ethical principles to guide AI development and implementation may lead to undesired outcomes, or even dangerous consequences. It is essential for AI-powered tools and applications to be designed and implemented ethically. One way is to promote a culture of ethics related to AI development. This could include improving transparency and encouraging more discussions around the pitfalls of AI-powered systems. According to *Contently*, the ethical implications of AI can be addressed either through governing policies and regulations via legislations, or encouraging companies and individuals to adopt ethical principles voluntarily. Some of the governing policies include the United States’ AI Bill of Rights, which focuses on AI’s development and implementation principles in guiding agencies to ensure that AI is developed and used responsibly. The European Union’s ethical guidelines for trustworthy AI, are similarly aimed at ensuring that AI is developed and used ethically while respecting fundamental human rights.

9. Operators of AI could also consider adopting an ethical framework in developing its generative AI tools. As *Holisticai* reported, AI ethics converge on four key verticals: (a) **Safety** – whether the system is robust against adversarial attacks; (b) **Privacy** – whether appropriate data minimisation have been adopted to protect users’ privacy; (c) **Fairness** – whether the system has been tested for bias and appropriate action taken to mitigate this; and (d) **Transparency** – how explainable the system is and whether there is appropriate communication with the relevant stakeholders. In addition, operators of these AI tools could also invest in a privacy team with fact-checkers and ethics experts to address the concerns and issues arising from the use of these AI tools.

10. When using generative AI to create content, it remains imperative for individuals to understand the potential implications of their actions. According to *Contently*, individuals could consider the following questions when using generative AI tools:

- a. What are the possible risks and implications of creating content with generative AI?
- b. How might your content be misused or misinterpreted?
- c. What could the potential negative impacts be on individuals or groups of people?
- d. Are there any risks to public safety that needs to be considered?

11. These questions are paramount to the use of AI in content generation. Any content created should be ethically sound and responsible for individuals, organisations, and consumers of AI. Only when individuals and organisations agree on and abide by ethical guidelines, can the negative impacts derived from generative AI be minimised.

CONTACT DETAILS

All reports can be retrieved from our website at www.acice-asean.org/resource/.

For any queries and/or clarifications, please contact ACICE at ACICE@defence.gov.sg.

Prepared by:

ADMM Cybersecurity and Information Centre of Excellence

••••

REFERENCES

News Articles

- 1 AI is Coming for White-Collar Jobs, Gates Warns
[Link: <https://www.afr.com/technology/ai-is-coming-for-white-collar-jobs-gates-warns-20230123-p5cev7>]
- 2 What are the Threats and Opportunities from ChatGPT?
[Link: <https://www.scmp.com/news/china/science/article/3207712/what-are-threats-and-opportunities-chatgpt-scientists-weigh>]
- 3 Guidelines for Responsible Content Creation with Generative AI
[Link: <https://www.contently.com/2023/01/03/guidelines-for-responsible-content-creation-with-generative-ai/amp/>]
- 4 How ChatGPT Hijacks Democracy
[Link: <https://www.nytimes.com/2023/01/15/opinion/ai-chatgpt-lobbying-democracy.html>]
- 5 Exploring The Dark Side of ChatGPT
[Link: <https://www.augustman.com/in/gear/tech/openai-gpt-chatbot-chatgpt-dark-side/amp/>]
- 6 Cyber Threat Analysis: I, Chatbot
[Link: <https://go.recordedfuture.com/hubfs/reports/cta-2023-0126.pdf>]
- 7 A Code of Ethics For Chatbots
[Link: <https://www.housing-technology.com/a-code-of-ethics-for-chatbots/>]
- 8 We Asked ChatGPT to Write an Article about Ethical AI, Here's What it Said
[Link: <https://www.holisticai.com/blog/ethical-ai-chat-gpt>]

- 9 How Companies Can Practice AI
[Link: <https://www.venturebeat.com/ai/how-companies-can-practice-ethical-ai>]
- 10 #DataPrivacyWeek: Addressing ChatGPT's Shortfalls in Data Protection Law Compliance
[Link: <https://www.infosecurity-magazine.com/news-features/chatgpt-shortfalls-data-protection/>]
- 11 Imitating Creators: Prospective Governance and Mechanisms for Identifying AI-generated Content
[Link: <https://www.lexology.com/library/detail.aspx?g=efb4a96c-b1fd-4560-aa3d-541e09a64ed7>]
- 12 How AI turns Text into Images
[Link: <https://www.pbs.org/newshour/amp/science/how-ai-makes-images-based-on-a-few-words>]
- 13 Some insist that Generative AI ChatGPT is a mirror into the Soul of Humanity, Vexing AI Ethics and AI Law
[Link: <https://www.forbes.com/sites/lanceeliot/2023/01/29/some-insist-that-generative-ai-chatgpt-is-a-mirror-into-the-soul-of-humanity-vexing-ai-ethics-and-ai-law/amp/>]
- 14 Midjourney Founder Basically Admits to Copyright Breaching and Artists are Angry
[Link: <https://www.digitalcameraworld.com/news/midjourney-founder-basically-admits-to-copyright-breaching-and-artists-are-angry>]