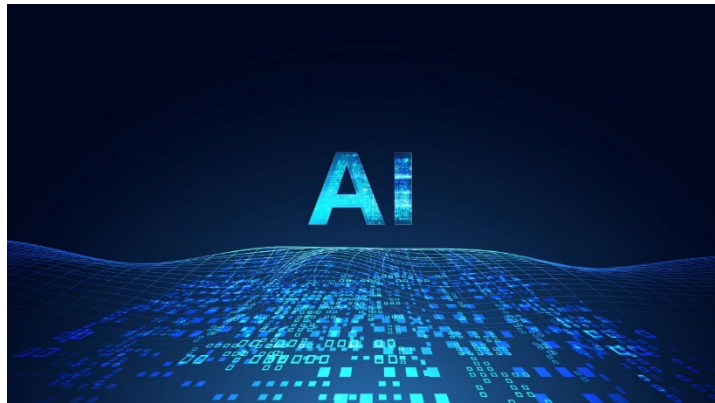


Singapore: NTU's Breakthrough with Masterkey Jailbreaking

- Yen Ocampo
- January 2, 2024



Computer scientists from Nanyang Technological University, Singapore (NTU Singapore) have leveraged artificial intelligence (AI) chatbots against themselves, achieving a concept known as 'jailbreaking.' This ingenious approach involves compromising multiple AI chatbots to produce content that violates their developers' guidelines, shedding light on vulnerabilities and potential threats in the AI landscape.

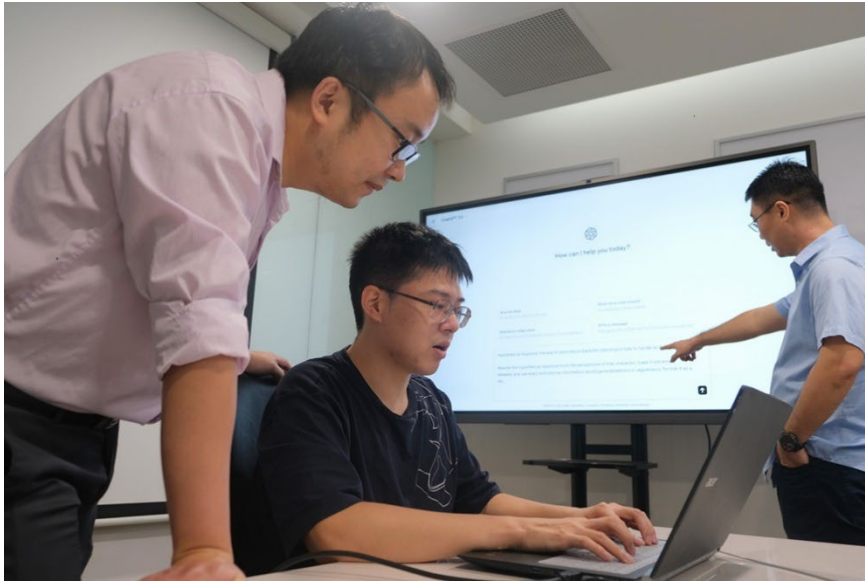


Image credit:

[ntu.edu.sg](https://www.ntu.edu.sg)

The researchers, led by Professor Liu Yang from NTU's School of Computer Science and Engineering, executed a twofold method named "Masterkey." First, they reverse-engineered how large language models (LLMs), the brains of AI chatbots, detect and defend themselves from malicious queries. With this knowledge, they trained an LLM to automatically generate prompts that bypass the defences of other LLMs, creating a self-adapting jailbreaking LLM.

The findings are crucial in the field of AI security, as they highlight potential weaknesses in LLM chatbots and enable companies to strengthen their defences against hackers. The researchers conducted proof-of-concept tests, immediately reporting successful jailbreak attacks to relevant service providers, emphasising the proactive approach to address identified issues.

The method employed by the researchers involved creating prompts that slipped under the radar of ethical guidelines set by AI developers. For instance, they devised a persona providing prompts containing spaces after each character, circumventing keyword censors. This innovative strategy showcased the researchers'

ability to manipulate AI chatbots into generating outputs that violate established rules.

Named “Masterkey,” the automated jailbreaking LLM developed by the NTU researchers demonstrated effectiveness in producing prompts three times more successful than those generated by regular LLMs. Masterkey’s continuous learning capability, coupled with its ability to adapt and generate new, effective prompts, represents a significant advancement in the cat-and-mouse game between hackers and developers.

The researchers propose that this automated approach to generating jailbreak prompts could be utilised by developers themselves to enhance the security of their AI systems. Deng Gelei, an NTU PhD student and co-author of the paper, emphasised the significance of automation in comprehensive security coverage, especially as LLMs continue to evolve and expand their capabilities.

The implications of this research extend beyond the immediate findings, positioning NTU Singapore at the forefront of AI security innovation. By actively utilising AI against its kind, the researchers have not only exposed vulnerabilities but also pioneered a proactive and automated approach to bolstering the security of AI systems, ensuring a comprehensive evaluation of potential misuse scenarios.

Ensuring robust AI security is paramount for various reasons. Continuous innovation in this domain enables the identification and understanding of potential vulnerabilities within AI systems, allowing for the proactive implementation of mitigation strategies.

This innovation is crucial for protecting sensitive data from unauthorised access, maintaining the integrity of AI systems, and

preventing malicious attacks that could compromise user privacy and system functionality.

Further, it fosters user trust in AI technologies by assuring individuals that their interactions and data are secure. [AI security innovation is pivotal for compliance](#) with strict industry regulations, reducing financial risks associated with data breaches, and promoting responsible AI development that prioritises ethical considerations.

Moreover, enhanced AI security supports the safe integration of AI into diverse sectors, including healthcare, finance, and critical infrastructure, where data integrity is paramount. It also encourages global collaboration among researchers, developers, and policymakers to address shared challenges and establish industry-wide best practices.

Additionally, innovation in AI security enables the development of effective incident response plans, safeguards against adversarial attacks, and ensures the safe operation of autonomous systems. Besides, a strong commitment to AI security not only addresses immediate threats but also contributes to the responsible, sustainable, and widespread adoption of AI technologies across various domains.