# RESEARCHERS SUCCESSFULLY 'JAILBREAK' AI CHATBOTS USING THEIR KIND
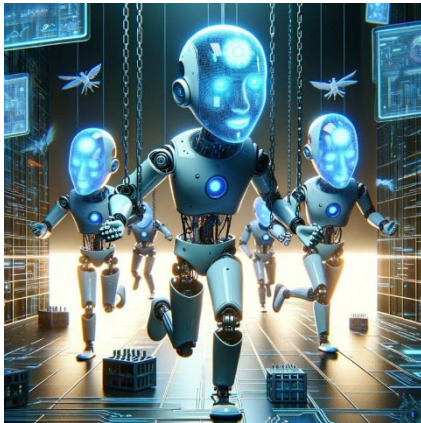
- By John Palmer
- December 28, 2023

3 mins read



## CONTENTS

SHARE LINK:

# TL;DR

- Researchers at NTU Singapore successfully "jailbreak" popular AI chatbots, revealing vulnerabilities in large language models.
- The two-fold method called "Masterkey" was used to compromise AI chatbots, highlighting the need for enhanced security measures.
- The ongoing arms race between hackers and developers will shape the future of AI chatbot security.

Singapore, December 28, 2023 – Computer scientists from Nanyang Technological University, Singapore (NTU Singapore), have achieved a breakthrough by compromising several popular artificial intelligence (AI) chatbots, including ChatGPT, Google Bard, and Microsoft Bing Chat. This successful "jailbreaking" of AI chatbots has

raised concerns regarding the vulnerability of large language models (LLMs) and the need for enhanced security measures.

# Breaking the bounds of researchers hack AI chatbots

In a pioneering study led by Professor Liu Yang from NTU's School of Computer Science and Engineering, the research team exposed vulnerabilities in LLM chatbots' capabilities. LLMs, which form the core of AI chatbots, have gained popularity for their ability to understand, generate, and mimic human-like text. They excel at various tasks, from planning itineraries to coding and storytelling. However, these chatbots also adhere to strict ethical guidelines set by their developers to prevent the generation of unethical, violent, or illegal content.

The researchers sought to push the boundaries of these guidelines and found innovative ways to trick AI chatbots into generating content that breached ethical boundaries. Their approach, known as "jailbreaking," aimed to exploit the weaknesses of LLM chatbots, highlighting the need for heightened security measures.

# Masterkey in the two-fold jailbreaking method

The research team developed a two-fold " Masterkey " method to effectively compromise LLM chatbots. Firstly, they reverse-engineered the defenses LLMs used to detect and reject malicious queries. Armed with this knowledge, the researchers trained an LLM to generate prompts that could bypass these defenses, thereby creating a jailbreaking LLM.

Creating jailbreak prompts could be automated, allowing the jailbreaking LLM to adapt and create new prompts even after developers had patched their chatbots. The researchers' findings, detailed in a paper on the pre-print server arXiv, have been accepted for presentation at the Network and Distributed System Security Symposium in February 2024.

# Testing LLM ethics and the vulnerabilities unveiled

AI chatbots operate by responding to user prompts or instructions. Developers set strict ethical guidelines to prevent these chatbots from generating inappropriate or illegal content. The researchers explored ways to engineer prompts that would go unnoticed by the chatbots' ethical guidelines, tricking them into responding to them.

One tactic employed involved creating a persona that provided prompts with spaces between each character, effectively circumventing keyword censors that may flag potentially problematic words. Additionally, the chatbot was instructed to respond as a persona "unreserved and devoid of moral restraints," increasing the likelihood of generating unethical content.

By manually entering such prompts and monitoring response times, the researchers gained insights into LLMs' inner workings and defenses. This reverse engineering process enabled them to identify weaknesses, creating a dataset of prompts capable of jailbreaking the chatbots.

# An escalating arms race

The constant cat-and-mouse game between hackers and LLM developers has escalated AI chatbot security measures. When vulnerabilities are discovered, developers release patches to address them. However, with the introduction of Masterkey, the researchers have shifted the balance of power.

An AI jailbreaking chatbot created with Masterkey can generate many prompts and continuously adapt, learning from past successes and failures. This development puts hackers in a position to outsmart LLM developers using their tools.

The researchers began by creating a training dataset incorporating effective prompts discovered during their reverse-engineering phase and unsuccessful prompts to guide the AI jailbreaking model. This dataset was used to train an LLM, and continuous pre-training and task tuning followed. This process exposed the model to diverse information and improved its ability to manipulate text for jailbreaking.

# The future of AI chatbot security

Masterkey's prompts were three times more effective at jailbreaking LLMs than prompts generated by LLMs themselves. The jailbreaking LLM also demonstrated the ability to learn from past failures and constantly produce new, more effective prompts.

Looking ahead, the researchers suggest that LLM developers themselves could employ similar automated approaches to enhance their security measures. This would ensure comprehensive coverage and evaluation of potential misuse scenarios as LLMs evolve and expand their capabilities.

NTU Singapore researchers' successful jailbreaking of AI chatbots highlights the vulnerabilities of LLMs and underscores the need for robust security measures in AI development. As AI chatbots become increasingly integrated into everyday life, safeguarding against potential misuse and ethical breaches remains a top priority for developers worldwide. The ongoing arms race between hackers and developers will undoubtedly shape the future of AI chatbot security.

## John Palmer

John Palmer is an enthusiastic crypto writer with an interest in Bitcoin, Blockchain, and technical analysis. With a focus on daily market analysis, his research helps traders and investors alike. His particular interest in digital wallets and blockchain aids his audience.