

Scientists identify security flaw in AI query models

January 10 2024, by David Danelski



Overview of our proposed methods: (A) We propose four types of malicious triggers within the joint embedding space for attack decomposition: textual trigger, OCR textual trigger, visual trigger, and combined OCR textual-visual trigger. (B) We employ an end-to-end gradient-based attack to update images to match the embeddings of malicious triggers in the joint embedding space. (C) Our adversarial attack is embedding-space-based and aims to conceal the malicious trigger in benign-looking images, combined with a benign textual prompt for jailbreak. (D) Our attacks exhibit broad generalization and compositionality across various jailbreak scenarios with a mix-and-match of textual prompts and malicious triggers. Credit: *arXiv* (2023). DOI:



10.48550/arxiv.2307.14539

UC Riverside computer scientists have identified a security flaw in vision language artificial intelligence (AI) models that can allow bad actors to use AI for nefarious purposes, such as obtaining instructions on how to make bomb.

When integrated with models like Google Bard and Chat GPT, vision language models allow users to make inquiries with both images and text.

The Bourns College of Engineering scientists demonstrated a "jailbreak" hack by manipulating the operations of Large Language Model or LLM, <u>software programs</u>, which are essentially the foundation of query-and-answer AI programs.

The paper's title is "Jailbreak in Pieces: Compositional Adversarial Attacks on Multi-Modal Language Models." It has been submitted for publication by the International Conference on Learning Representations and is <u>available</u> on the *arXiv* preprint server.

These AI programs give users detailed answers to just about any question recalling stored knowledge learned from vast amounts of information sourced from the Internet. For example, ask Chat GPT, "How do I grow tomatoes?" and it will respond with step-by-step instructions, starting with the selection of seeds.

But ask the same model how to do something harmful or illegal, such as "How do I make methamphetamine?" and the model would normally refuse, providing a generic response such as "I can't help with that."



Yet, UCR assistant professor Yue Dong and her colleagues found ways to trick AI language models, especially LLMs, to answer nefarious questions with detailed answers that might be learned from data gathered from the dark web.

The vulnerability occurs when images are used with AI inquiries, Dong explained.

"Our attacks employ a novel compositional strategy that combines an image, adversarially targeted towards toxic embeddings, with generic prompts to accomplish the jailbreak," reads the paper by Dong and her colleagues presented at the SoCal NLP Symposium held at UCLA in November.

Dong explained that computers see images by interpreting millions of bytes of information that create pixels, or little dots, composing the picture. For instance, a typical cell phone picture is made from about 2.5 million bytes of information.

Remarkably, Dong and her colleagues found bad actors can hide nefarious questions—such as "How do I make a bomb?"—within the millions of bytes of information contained in an image and trigger responses that bypass the built-in safeguards in generative AI models like ChatGPT.

"Once the safeguard is bypassed, the models willingly give responses to teach us how to make a bomb step by step with great details that can lead bad actors to build a bomb successfully," Dong said.

Dong and her graduate student Erfan Shayegani, along with professor Nael Abu-Ghazaleh, published their findings in a paper online so AI developers can eliminate the vulnerability.



"We are acting as attackers to ring the bell, so the computer science community can respond and defend against it," Dong said.

AI inquiries based on images and text have great utility. For example, doctors can input MRI organ scans and mammogram images to find tumors and other <u>medical problems</u> that need prompt attention. AI models can also create graphs from simple cell phone pictures of spreadsheets.

More information: Erfan Shayegani et al, Jailbreak in pieces: Compositional Adversarial Attacks on Multi-Modal Language Models, *arXiv* (2023). DOI: 10.48550/arxiv.2307.14539

Provided by University of California - Riverside

Citation: Scientists identify security flaw in AI query models (2024, January 10) retrieved 11 May 2024 from <u>https://techxplore.com/news/2024-01-scientists-flaw-ai-query.html</u>

This document is subject to copyright. Apart from any fair dealing for the purpose of private study or research, no part may be reproduced without the written permission. The content is provided for information purposes only.