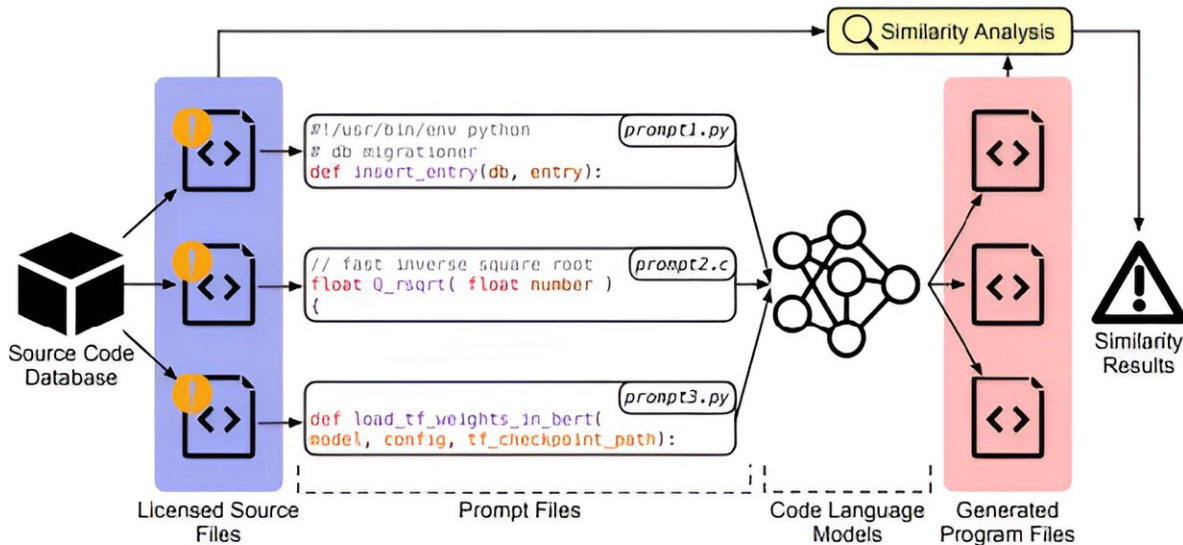


Analyzing generative AI's copyright crisis

July 31 2023, by Shawn Ballard



CODEIPPROMPT constructs prompts from a collection of source code database, and the generated programs are analyzed against source files for similarity scores. Credit: CODEIPPROMPT: Intellectual Property Infringement Assessment of Code Language Models.

<https://openreview.net/pdf?id=zdmbZl0ia6>

The recent explosion of artificial intelligence tools such as ChatGPT and Copilot have supercharged the assistance available to programmers. However, AI assistants may strip out comments embedded in code to convey copyright and attribution guidelines, leaving human coders none the wiser yet still on the hook legally for intellectual property infringement.

To combat this problem, [computer science](#) and engineering researchers in the McKelvey School of Engineering at Washington University in St. Louis have developed CodeIPPrompt, the first automated testing platform to evaluate how much language models generate IP-violating code. The team includes Ning Zhang and Chenguang Wang, both assistant professors; Yevgeniy Vorobeychik, professor; Zhiyuan Yu, a graduate student in Zhang's lab and first author on the paper; and Chaowei Xiao, assistant professor of computer science at Arizona State University.

Yu presented the work July 23 at the [International Conference on Machine Learning](#) in Honolulu. Notably, the team's analysis showed that [copyright infringement](#) issues are prevalent across state-of-the-art open-source models including CodeRL, CodeGen and CodeParrot, as well as in commercial products including Copilot, ChatGPT and GPT-4.

"We developed this tool to help people understand that if they're using these large language models to help write code, there's a good chance they might generate IP infringing content," Zhang said. "As users, we have a responsibility to use AI ethically. That's influenced by how we understand AI technology and the content it produces."

Though CodeIPPrompt can't say for sure if AI-generated code constitutes an IP violation—Zhang notes that issue is ultimately a legal question that will play out in the courts as cases are brought against the users of AI tools for copyright infringement—it can give users a risk score that indicates how similar generated [code](#) is to copyright protected content. Zhang anticipates that the tool will help guide the ongoing development of AI and point to potential mitigation strategies and other protections against IP violations in the future.

More information: CODEIPPROMPT: Intellectual Property Infringement Assessment of Code Language Models.

openreview.net/pdf?id=zdmbZl0ia6

Provided by Washington University in St. Louis

Citation: Analyzing generative AI's copyright crisis (2023, July 31) retrieved 21 February 2024 from <https://techxplore.com/news/2023-07-generative-ai-copyright-crisis.html>

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.