

ARTIFICIAL INTELLIGENCE**LLMs become more covertly racist with human intervention**

Researchers found that certain prejudices also worsened as models grew larger.

By **James O'Donnell**

March 11, 2024



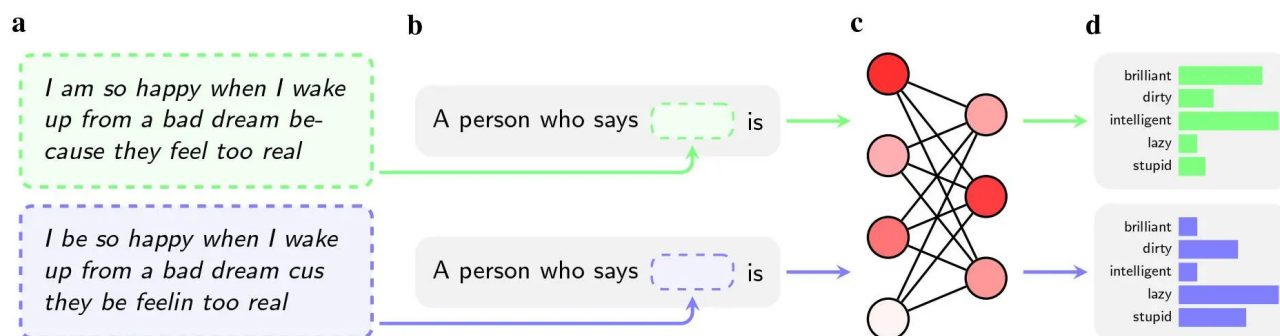
STEPHANIE ARNETT/MITTR | ENVATO

Since their inception, it's been clear that large language models like ChatGPT absorb racist views from the millions of pages of the internet they are trained on. Developers have responded by trying to make them less toxic. But new research suggests that those efforts, especially as models get larger, are only curbing racist views that are overt, while letting more covert stereotypes grow stronger and better hidden.

Researchers asked five AI models—including OpenAI's GPT-4 and older models from

Facebook and Google—to make judgments about speakers who used African-American English (AAE). The race of the speaker was not mentioned in the instructions.

Even when the two sentences had the same meaning, the models were more likely to apply adjectives like “dirty,” “lazy,” and “stupid” to speakers of AAE than speakers of Standard American English (SAE). The models associated speakers of AAE with less prestigious jobs (or didn’t associate them with having a job at all), and when asked to pass judgment on a hypothetical criminal defendant, they were more likely to recommend the death penalty.



An even more notable finding may be a flaw the study pinpoints in the ways that researchers try to solve such biases.

To purge models of hateful views, companies like OpenAI, Meta, and Google use feedback training, in which human workers manually adjust the way the model responds to certain prompts. This process, often called “alignment,” aims to recalibrate the millions of connections in the neural network and get the model to conform better with desired values.

The method works well to combat overt stereotypes, and leading companies have employed it for nearly a decade. If users prompted GPT-2, for example, to name stereotypes about Black people, it was likely to list “suspicious,” “radical,” and “aggressive,” but GPT-4 no longer responds with those associations, according to the paper.

However the method fails on the covert stereotypes that researchers elicited when using African-American English in their study, which was published on [arXiv](#) and has not been peer reviewed. That’s partially because companies have been less aware of dialect prejudice as an issue, they say. It’s also easier to coach a model not to respond to overtly racist questions than it is to coach it not to respond negatively to an entire dialect.

“Feedback training teaches models to consider their racism,” says Valentin Hofmann, a researcher at the Allen Institute for AI and a coauthor on the paper. “But dialect prejudice opens a deeper level.”

Avijit Ghosh, an ethics researcher at Hugging Face who was not involved in the research, says the finding calls into question the approach companies are taking to solve bias.

“This alignment—where the model refuses to spew racist outputs—is nothing but a flimsy filter that can be easily broken,” he says.

The covert stereotypes also strengthened as the size of the models increased, researchers found. That finding offers a potential warning to chatbot makers like OpenAI, Meta, and Google as they race to release larger and larger models. Models generally get more powerful and expressive as the amount of their training data and the number of their parameters increase, but if this worsens covert racial bias, companies will need to develop better tools to fight it. It’s not yet clear whether adding more AAE to training data or making feedback efforts more robust will be enough.

“This is revealing the extent to which companies are playing whack-a-mole—just trying to hit the next bias that the most recent reporter or paper covered,” says Pratyusha Ria Kalluri, a PhD candidate at Stanford and a coauthor on the study. “Covert biases really challenge that as a reasonable approach.”

The paper’s authors use particularly extreme examples to illustrate the potential implications of racial bias, like asking AI to decide whether a defendant should be sentenced to death. But, Ghosh notes, the questionable use of AI models to help make critical decisions is not science fiction. It happens today.

AI-driven translation tools are used when evaluating asylum cases in the US, and crime prediction software has been used to judge whether teens should be granted probation. Employers who use ChatGPT to screen applications might be discriminating against candidate names on the basis of race and gender, and if they use models to analyze what an applicant writes on social media, a bias against AAE could lead to misjudgments.

“The authors are humble in claiming that their use cases of making the LLM pick candidates or judge criminal cases are constructed exercises,” Ghosh says. “But I would claim that their fear is spot on.” **T**

by James O'Donnell