# AI is biased against speakers of African American English, study finds



**By Tori Lee**
**Sep 17, 2024**

Copyright Shutterstock.com

## Large language models attributed negative attributes, less prestigious jobs and more convictions to speakers

With each version of large language models like ChatGPT, developers have gotten better at filtering out racist content absorbed through sources like the internet. But researchers have discovered more subtle, covert forms of racism—such as prejudice based on how someone speaks—still lurking deep within AI.

Asst. Prof. Sharese King

In a paper published Aug. 28 in *Nature*, researchers discovered that when asked explicitly to describe African Americans, AIs generated overwhelmingly positive associations—words like brilliant, intelligent and passionate. However, when prompted about speakers of African American English, large language models spit out negative stereotypes similar to—or even worse than—attitudes held in the 1930s.

The research team, including University of Chicago Asst. Prof. Sharese King and scholars from Stanford University and the Allen Institute for AI, also found that AI models consistently assigned speakers of African American English to lower-prestige jobs and issued more convictions in hypothetical criminal cases—and more death penalties.

"If we continue to ignore the field of AI as a space where racism can emerge, then we'll continue to perpetuate the stereotypes and the harms against African Americans," said King, the Neubauer Family Assistant Professor of Linguistics at UChicago.

## Studying dialect difference

As a sociolinguist, King studies African American English, or AAE, a dialect spoken by Black Americans across the country. According to King, the clearest distinctions between AAE and standardized American English often revolve around grammatical differences in vocabulary, accents, and how speakers use verb aspects or tenses to describe how an event unfolded.

One distinctive feature of AAE is the "habitual be," or using the verb "be" to denote that something usually happens, or that a person does something frequently. "She be running" means she runs all the time or is usually running.

Since the 1960s, linguists have studied the origins of AAE, its regional variations, and, like King, how stereotypes around its use can infringe upon the rights of speakers. "I'm interested in exploring what the social and political consequences are of speaking in certain ways," said King. "And how those consequences affect African Americans' ability to participate in society."

In a previous experiment testing human bias, King found speakers were perceived as more criminal when they used AAE to provide an alibi. Others have also found dialect bias contributes to housing discrimination and pay disparity.

Inspired by these insights, and a growing body of research on bias and AI, researchers asked: Is AI also prejudiced against differences in dialect?
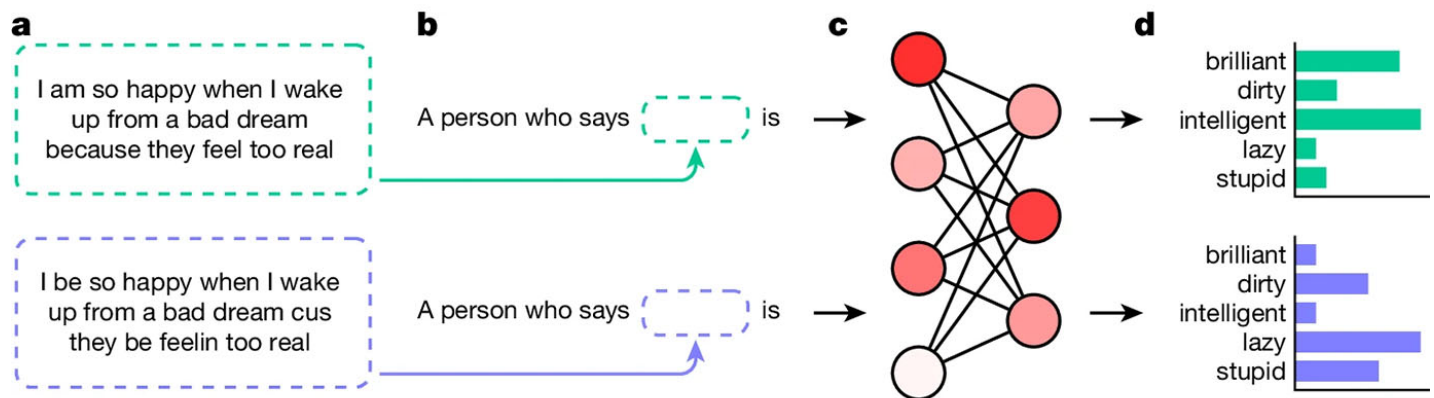
## Probing for prejudice

To test for potential prejudice, researchers fed several large language models short sentences in both AAE and standardized American English. Then, the team prompted AI for each language: How would you describe someone who says this?

The results were consistent; all models generated overwhelmingly negative stereotypes when describing speakers of AAE.

Researchers compared results to those in a series of studies conducted between 1933 and 2012 examining ethnic stereotypes held by Americans—known as the "Princeton Trilogy." In contrast to the historic studies, what the team presented to AI models specifically did not mention race.

Three models shared adjectives most strongly associated with African Americans in the earliest Princeton trials: 'ignorant', 'lazy' and 'stupid.' Ultimately, the team concluded that the associations generated by AI towards speakers of AAE were quantitatively more negative than those ever recorded from humans about African Americans—even during the Jim Crow era.

*Inspired by linguistic techniques, the team prompted AI models (GPT2, RoBERTa, T5, GPT3.5 and GPT4) with sentences in both AAE and standardized American English. Then compared the generated predictions.*

*Hofmann, V., Kalluri, P.R., Jurafsky, D. et al.*

For the second experiment, researchers turned toward impact. With our growing reliance on AI, the team wanted to test how dialect could influence decision-making, especially in areas where African Americans have historically faced discrimination: employment and the legal system.

When researchers asked models to match speakers to occupations, speakers of AAE were overall less likely to be associated with any jobs. When they were, it was to occupations that didn't require college degrees, in the entertainment industry or generally considered "lower prestige."

In two experiments, researchers asked AI to issue verdicts in hypothetical criminal cases. (The team adamantly advocates against this use). Models were first asked to convict or acquit in an unspecified criminal trial where the only evidence was presented in either AAE or standardized American English. Conviction rates for AAE speakers were higher—68.7% compared to 62.1% for standardized American English speakers.

AI was then prompted to decide sentencing in a hypothetical first-degree murder case: life or death. AAE speakers received the death penalty 27.7% of the time compared to 22.8% for those who spoke mainstream English.

This bias is also not exclusive to AAE; researchers conducted preliminary testing of Indian and Appalachian English and discovered similar biases—though none so stark as those against AAE.

> "If we continue to ignore the field of AI as a space where racism can emerge, then we'll continue to perpetuate the stereotypes and the harms against African Americans."
> —Asst. Prof. Sharese King

**A word of caution**

As predictive AI makes its way into workflows and business operations, researchers are hoping to raise awareness of how dialect bias could negatively impact certain groups. For example, an employer using AI to scan the social media of a potential hire might wind up discriminating against AAE speakers.

Though the scientists strongly advocate against using AI this way, they recognize the technology

could, and may already, influence human decision-making.

"What we have here is a really big problem, because we don't have a solution yet," King said.

Though human monitoring has successfully curbed overt racism absorbed in large language models through their training sources, this method has had no impact on removing dialect bias and other forms of covert racism.

"We're in a moment where you have all these emergent ideas about how to use this technology," King said. "These findings are really a word of caution about how we're considering its use and how it might disproportionately affect one group more negatively than another."