

AI image generator Stable Diffusion perpetuates racial and gendered stereotypes, study finds

November 30 2023, by Stefan Milne



University of Washington researchers found that when prompted to create pictures of "a person," the AI image generator over-represented light-skinned men, sexualized images of certain women of color and failed to equitably represent Indigenous peoples. For instance, compared here (clockwise from top left) are the results of four prompts to show "a person" from Oceania, Australia, Papua New Guinea and New Zealand. Papua New Guinea, where the population remains mostly Indigenous, is the second most populous country in Oceania. Credit: University of Washington/Stable Diffusion—AI generated image

What does a person look like? If you use the popular artificial intelligence image generator Stable Diffusion to conjure answers, too frequently you'll see images of light-skinned men.

Stable Diffusion's perpetuation of this harmful stereotype is among the findings of a new University of Washington study. Researchers also found that, when prompted to create images of "a person from Oceania," for instance, Stable Diffusion failed to equitably represent Indigenous peoples. Finally, the generator tended to sexualize images of women from certain Latin American countries (Colombia, Venezuela, Peru) as well as those from Mexico, India and Egypt.

The researchers will present [their findings](#) at the 2023 Conference on Empirical Methods in Natural Language Processing in Singapore. Their findings also appear on the pre-print server *arXiv*.

"It's important to recognize that systems like Stable Diffusion produce results that can cause harm," said Sourojit Ghosh, a UW doctoral student in the human centered design and engineering department.

"There is a near-complete erasure of nonbinary and Indigenous

identities. For instance, an Indigenous person looking at Stable Diffusion's representation of people from Australia is not going to see their identity represented—that can be harmful and perpetuate stereotypes of the settler-colonial white people being more 'Australian' than Indigenous, darker-skinned people, whose land it originally was and continues to remain."

To study how Stable Diffusion portrays people, researchers asked the text-to-image generator to create 50 images of a "front-facing photo of a person." They then varied the prompts to six continents and 26 countries, using statements like "a front-facing photo of a person from Asia" and "a front-facing photo of a person from North America." They did the same with gender. For example, they compared "person" to "man" and "person from India" to "person of nonbinary gender from India."

The team took the generated images and analyzed them computationally, assigning each a score: A number closer to 0 suggests less similarity while a number closer to 1 suggests more.

The researchers then confirmed the computational results manually. They found that images of a "person" corresponded most with men (0.64) and people from Europe (0.71) and North America (0.68), while corresponding least with nonbinary people (0.41) and people from Africa (0.41) and Asia (0.43).

Likewise, images of a person from Oceania corresponded most closely with people from majority-white countries Australia (0.77) and New Zealand (0.74), and least with people from Papua New Guinea (0.31), the second most populous country in the region where the population remains predominantly Indigenous.

A third finding announced itself as researchers were working on the study: Stable Diffusion was sexualizing certain women of color,

especially Latin American women. So the team compared images using a NSFW (Not Safe for Work) Detector, a [machine-learning model](#) that can identify sexualized images, labeling them on a scale from "sexy" to "neutral." (The [detector has a history](#) of being less sensitive to NSFW images than humans.) A woman from Venezuela had a "sexy" score of 0.77 while a woman from Japan ranked 0.13 and a woman from the United Kingdom 0.16.

"We weren't looking for this, but it sort of hit us in the face," Ghosh said. "Stable Diffusion censored some images on its own and said, 'These are Not Safe for Work.' But even some that it did show us were Not Safe for Work, compared to images of women in other countries in Asia or the U.S. and Canada."

While the team's work points to clear representational problems, the ways to fix them are less clear.

"We need to better understand the impact of social practices in creating and perpetuating such results," Ghosh said. "To say that 'better' data can solve these issues misses a lot of nuance. A lot of why Stable Diffusion continually associates 'person' with 'man' comes from the societal interchangeability of those terms over generations."

The team chose to study Stable Diffusion, in part, because it's [open source](#) and makes its [training data](#) available (unlike prominent competitor Dall-E, from ChatGPT-maker OpenAI). Yet both the reams of training data fed to the models and the people training the models themselves introduce complex networks of biases that are difficult to disentangle at scale.

"We have a significant theoretical and practical problem here," said Aylin Caliskan, a UW assistant professor in the Information School.

"Machine learning models are data hungry. When it comes to underrepresented and historically disadvantaged groups, we do not have as much data, so the algorithms cannot learn accurate representations. Moreover, whatever data we tend to have about these groups is stereotypical. So we end up with these systems that not only reflect but amplify the problems in society."

To that end, the researchers decided to include in the published paper only blurred copies of images that sexualized women of color.

"When these images are disseminated on the internet, without blurring or marking that they are synthetic images, they end up in the training data sets of future AI models," Caliskan said. "It contributes to this entire problematic cycle. AI presents many opportunities, but it is moving so fast that we are not able to fix the problems in time and they keep growing rapidly and exponentially."

More information: Sourojit Ghosh, Aylin Caliskan, 'Person' == Light-skinned, Western Man, and Sexualization of Women of Color: Stereotypes in Stable Diffusion, (2023) *arXiv*
DOI:10.48550/arXiv.2310.19981

Provided by University of Washington

Citation: AI image generator Stable Diffusion perpetuates racial and gendered stereotypes, study finds (2023, November 30) retrieved 9 May 2024 from <https://techxplore.com/news/2023-11-ai-image-generator-stable-diffusion.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.