

Forget dystopian scenarios—AI is pervasive today, and the risks are often hidden

November 22 2023, by Anjana Susarla



Credit: CC0 Public Domain

The turmoil at ChatGPT-maker OpenAI, sparked by the board of directors <u>firing high-profile CEO Sam Altman</u> on Nov. 17, 2023, has put a spotlight on artificial intelligence safety and concerns about the rapid



development of artificial general intelligence, or AGI. AGI is loosely defined as <u>human-level intelligence across a range of tasks</u>.

The OpenAI board stated that <u>Altman's termination was for lack of</u> <u>candor</u>, but speculation has centered on a rift between Altman and members of the board over concerns that OpenAI's remarkable growth—products such as ChatGPT and Dall-E have <u>acquired hundreds</u> <u>of millions of users worldwide</u>—has <u>hindered the company's ability</u> to focus on <u>catastrophic risks</u> posed by AGI.

OpenAI's goal of developing AGI has become entwined with the idea of <u>AI acquiring superintelligent capabilities</u> and the need to safeguard against the technology being misused or going rogue. But for now, AGI and its attendant risks are speculative. Task-specific forms of AI, meanwhile, are very real, have become widespread and often fly under the radar.

As a <u>researcher of information systems and responsible AI</u>, I study how these everyday algorithms work—and how they can harm people.

AI is pervasive

AI plays a visible part in many people's daily lives, from face recognition unlocking your phone to speech recognition powering your digital assistant. It also plays roles you might be vaguely aware of—for example, shaping your social media and online shopping sessions, guiding your video-watching choices and <u>matching you with a driver</u> in a ride-sharing service.

AI also affects your life in ways that might completely escape your notice. If you're applying for a job, <u>many employers use AI in the hiring process</u>. Your bosses might be using it to identify employees <u>who are likely to quit</u>. If you're applying for a loan, odds are your bank is using



AI to decide whether to grant it. If you're being treated for a medical condition, your <u>health care providers</u> might use it to <u>assess your medical</u> <u>images</u>. And if you know someone caught up in the <u>criminal justice</u> <u>system</u>, AI could well play a role in <u>determining the course of their life</u>.

Algorithmic harms

Many of the AI systems that fly under the radar have biases that can cause harm. For example, machine learning methods use <u>inductive logic</u>, which starts with a set of premises to generalize patterns from <u>training data</u>. A machine learning-based <u>resume screening tool was found to be biased against women</u> because the training data reflected past practices when most resumes were submitted by men.

The use of predictive methods in areas ranging from health care to child welfare could exhibit <u>biases such as cohort bias</u> that lead to unequal risk assessments across different groups in society. Even when legal practices prohibit discrimination based on attributes such as race and gender—for example, in consumer lending—proxy discrimination can still occur. This happens when algorithmic decision-making models do not use characteristics that are legally protected, such as race, and instead use characteristics that are highly correlated or connected with the legally protected characteristic, like neighborhood. Studies have found that risk-equivalent Black and Latino borrowers pay significantly higher interest rates on government-sponsored enterprise securitized and Federal Housing Authority-insured loans than white borrowers.

Another form of bias occurs when decision-makers use an algorithm differently from how the algorithm's designers intended. In a well-known example, a <u>neural network</u> learned to <u>associate asthma with a</u> <u>lower risk of death from pneumonia</u>. This was because asthmatics with pneumonia are traditionally given more aggressive treatment that lowers their mortality risk compared to the overall population. However, if the



outcome from such a neural network is used in hospital bed allocation, then those with asthma and admitted with pneumonia would be dangerously deprioritized.

Biases from algorithms can also result from complex societal feedback loops. For example, when predicting recidivism, authorities attempt to predict which people convicted of crimes are <u>likely to commit crimes</u> again. But the data used to train predictive algorithms is actually about who is likely to get re-arrested.

AI safety in the here and now

The Biden administration's recent executive order and enforcement efforts by <u>federal agencies</u> such as the Federal Trade Commission are the first steps in recognizing and safeguarding against algorithmic harms.

And though <u>large language models</u>, such as GPT-3 that powers ChatGPT, and <u>multimodal large language models</u>, such as GPT-4, are steps on the road toward artificial general intelligence, they are also algorithms people are increasingly using in school, work and daily life. It's important to consider the biases that result from the widespread use of <u>large language models</u>.

For example, these models could exhibit biases resulting from <u>negative</u> stereotyping involving gender, race or religion, as well as biases in the representation of <u>minorities and disabled people</u>. As these models demonstrate the ability to outperform <u>humans on tests such as the bar</u> <u>exam</u>, I believe that they require greater scrutiny to ensure that AIaugmented work conforms to <u>standards of transparency</u>, accuracy and <u>source crediting</u>, and <u>that stakeholders have the authority</u> to enforce such standards.

Ultimately, who wins and loses from large-scale deployment of AI may



not be about rogue superintelligence but about understanding who is vulnerable when algorithmic decision-making is ubiquitous.

This article is republished from <u>The Conversation</u> under a Creative Commons license. Read the <u>original article</u>.

Provided by The Conversation

Citation: Forget dystopian scenarios—AI is pervasive today, and the risks are often hidden (2023, November 22) retrieved 9 May 2024 from https://techxplore.com/news/2023-11-dystopian-scenariosai-pervasive-today-hidden.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.