

Opinion: Why humans can't trust AI—we don't know how it works or whether it'll serve our interests

September 14 2023, by Mark Bailey



In neural networks, the strength of the connections between 'neurons' changes as data passes from the input layer through hidden layers to the output layer, enabling the network to 'learn' patterns. Credit: <u>Wiso via Wikimedia Commons</u>



There are alien minds among us. Not the little green men of science fiction, but the alien minds that power the facial recognition in your smartphone, <u>determine your creditworthiness</u> and write <u>poetry</u> and <u>computer code</u>. These alien minds are artificial intelligence systems, the ghost in the machine that you encounter daily.

But AI systems have a significant limitation: Many of their inner workings are <u>impenetrable</u>, <u>making them fundamentally unexplainable</u> and unpredictable. Furthermore, constructing AI systems that behave in ways that people expect is a significant challenge.

If you fundamentally don't understand something as unpredictable as AI, how can you trust it?

Why AI is unpredictable

<u>Trust</u> is grounded in predictability. It depends on your ability to anticipate the behavior of others. If you trust someone and they don't do what you expect, then your perception of their trustworthiness diminishes.

Many AI systems are built on <u>deep learning</u> neural networks, which in some ways emulate the <u>human brain</u>. These networks contain interconnected "neurons" with variables or "parameters" that affect the strength of connections between the neurons. As a naïve network is presented with <u>training data</u>, it <u>"learns" how to classify the data</u> by adjusting these parameters. In this way, the AI system learns to classify data it hasn't seen before. It doesn't memorize what each data point is, but instead predicts what a data point might be.

Many of the most powerful AI systems contain <u>trillions of parameters</u>. Because of this, the reasons AI systems make the decisions that they do are often opaque. This is the <u>AI explainability problem</u>—the



impenetrable black box of AI decision-making.

Consider a variation of the <u>"Trolley Problem</u>." Imagine that you are a passenger in a <u>self-driving vehicle</u>, controlled by an AI. A small child runs into the road, and the AI must now decide: run over the child or swerve and crash, potentially injuring its passengers. This choice would be difficult for a human to make, but a human has the benefit of being able to explain their decision. Their rationalization—shaped by ethical norms, the perceptions of others and expected behavior—supports trust.

In contrast, an AI can't rationalize its decision-making. You can't look under the hood of the self-driving vehicle at its trillions of parameters to explain why it made the decision that it did. AI fails the predictive requirement for trust.

AI behavior and human expectations

Trust relies not only on predictability, but also on <u>normative or ethical</u> motivations. You typically expect people to act not only as you assume they will, but also as they should. Human values are influenced by <u>common experience</u>, and <u>moral reasoning</u> is a <u>dynamic process</u>, shaped by ethical standards and others' perceptions.

Unlike humans, AI doesn't adjust its behavior based on how it is perceived by others or by adhering to ethical norms. AI's internal representation of the world is largely static, set by its training data. Its <u>decision-making process</u> is grounded in an unchanging model of the world, unfazed by the dynamic, nuanced social interactions constantly influencing <u>human behavior</u>. Researchers are working on programming AI to include ethics, but that's <u>proving challenging</u>.

The self-driving car scenario illustrates this issue. How can you ensure that the car's AI makes decisions that align with human expectations?



For example, the car could decide that hitting the child is the optimal course of action, something most human drivers would instinctively avoid. This issue is the <u>AI alignment problem</u>, and it's another source of uncertainty that erects barriers to trust.

Critical systems and trusting AI

One way to reduce uncertainty and boost trust is to ensure people are in on the decisions AI systems make. This is the <u>approach taken by the U.S.</u> <u>Department of Defense</u>, which requires that for all AI decision-making, a human must be either in the loop or <u>on the loop</u>. In the loop means the AI system makes a recommendation but a human is required to initiate an action. On the loop means that while an AI system can initiate an action on its own, a human monitor can interrupt or alter it.

While keeping humans involved is a great first step, I am not convinced that this will be sustainable long term. As companies and governments continue to adopt AI, the future will likely include nested AI systems, where rapid decision-making limits the opportunities for people to intervene. It is important to resolve the explainability and alignment issues before the <u>critical point</u> is reached where human intervention becomes impossible. At that point, there will be no option other than to trust AI.

Avoiding that threshold is especially important because AI is increasingly being integrated into <u>critical systems</u>, which include things such as electric grids, the internet and <u>military systems</u>. In <u>critical</u> <u>systems</u>, trust is paramount, and undesirable behavior could have deadly consequences. As AI integration becomes more complex, it becomes even more important to resolve issues that limit trustworthiness.

Can people ever trust AI?



AI is alien—an intelligent system into which people have little insight. Humans are largely predictable to other humans because we share the same human experience, but this doesn't extend to artificial intelligence, even though humans created it.

If trustworthiness has inherently predictable and normative elements, AI fundamentally lacks the qualities that would make it worthy of trust. More research in this area will hopefully shed light on this issue, ensuring that AI systems of the future are worthy of our <u>trust</u>.

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: Opinion: Why humans can't trust AI—we don't know how it works or whether it'll serve our interests (2023, September 14) retrieved 12 May 2024 from <u>https://techxplore.com/news/2023-09-opinion-humans-aiwe-dont-itll.html</u>

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.