

## What is a black box? A computer scientist explains what it means when the inner workings of AIs are hidden

May 23 2023, by Saurabh Bagchi



Credit: Unsplash/CC0 Public Domain

For some people, the term "black box" brings to mind the recording devices in airplanes that are valuable for postmortem analyzes if the unthinkable happens. For others it evokes small, minimally outfitted theaters. But black box is also an important term in the world of artificial



intelligence.

AI <u>black boxes</u> refer to AI systems with internal workings that are invisible to the user. You can feed them input and get output, but you cannot examine the system's code or the logic that produced the output.

Machine learning is the dominant subset of artificial intelligence. It underlies generative AI systems like <u>ChatGPT</u> and <u>DALL-E 2</u>. There are three components to machine learning: an <u>algorithm</u> or a set of algorithms, <u>training data</u> and a model. An algorithm is a set of procedures. In machine learning, an algorithm learns to identify patterns after being trained on a large set of examples—the training data. Once a <u>machine-learning algorithm</u> has been trained, the result is a <u>machinelearning</u> model. The model is what people use.

For example, a machine-learning algorithm could be designed to identify patterns in images, and training data could be images of dogs. The resulting machine-learning model would be a dog spotter. You would feed it an image as input and get as output whether and where in the image a set of pixels represents a dog.

Any of the three components of a machine-learning system can be hidden, or in a black box. As is often the case, the algorithm is publicly known, which makes putting it in a black box less effective. So to protect their <u>intellectual property</u>, AI developers often put the model in a black box. Another approach software developers take is to obscure the data used to train the model—in other words, put the training data in a black box.

The opposite of a black box is sometimes referred to as a <u>glass box</u>. An AI glass box is a system whose algorithms, training data and model are all available for anyone to see. But researchers sometimes characterize aspects of even these as black box.



That's because researchers <u>don't fully understand</u> how machine-learning algorithms, particularly <u>deep-learning</u> algorithms, operate. The field of explainable AI is working to develop algorithms that, while not necessarily glass box, can be better understood by humans.

## Why AI black boxes matter

In many cases, there is good reason to be wary of black box machinelearning algorithms and models. Suppose a <u>machine-learning model</u> has made a diagnosis about your health. Would you want the model to be black box or glass box? What about the physician prescribing your course of treatment? Perhaps she would like to know how the model arrived at its decision.

What if a <u>machine-learning</u> model that determines whether you qualify for a business loan from a bank turns you down? Wouldn't you like to know why? If you did, you could more effectively appeal the decision, or change your situation to increase your chances of getting a loan the next time.

Black boxes also have important implications for software system security. For years, many people in the computing field thought that keeping software in a <u>black box</u> would prevent hackers from examining it and therefore it would be secure. This assumption has largely been proved wrong because hackers can <u>reverse-engineer</u> software—that is, build a facsimile by closely observing how a piece of software works—and discover vulnerabilities to exploit.

If software is in a glass box, then <u>software</u> testers and well-intentioned hackers can examine it and inform the creators of weaknesses, thereby minimizing cyberattacks.

This article is republished from <u>The Conversation</u> under a Creative



## Commons license. Read the original article.

## Provided by The Conversation

Citation: What is a black box? A computer scientist explains what it means when the inner workings of AIs are hidden (2023, May 23) retrieved 26 April 2024 from <u>https://techxplore.com/news/2023-05-black-scientist-ais-hidden.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.