# Stuart Russell: Responsible AI

# What is AI?

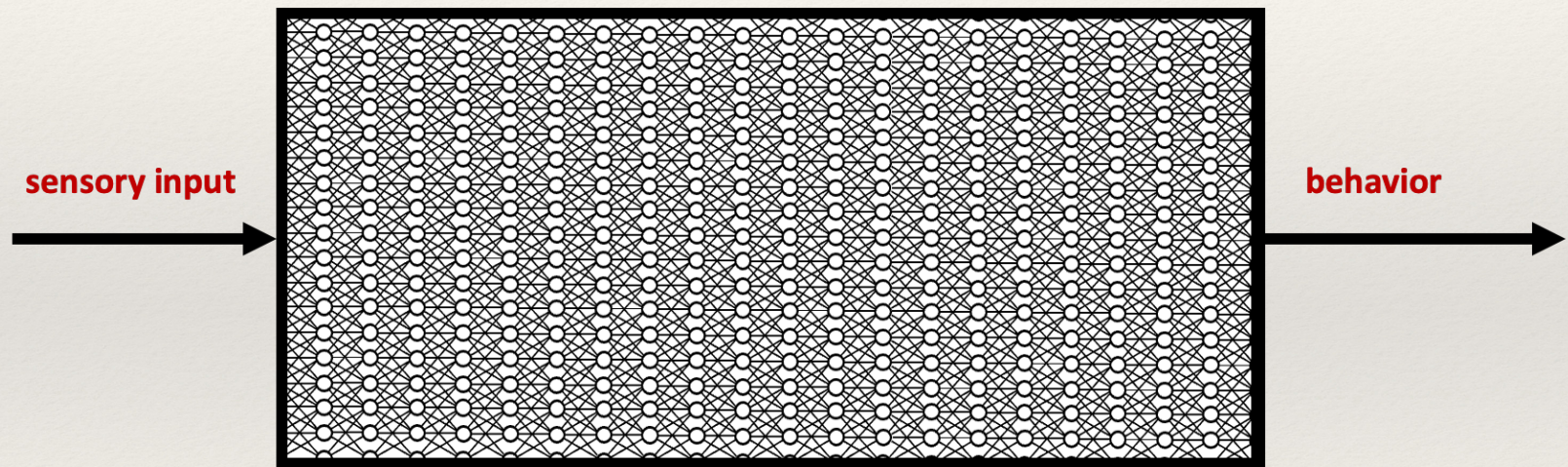*AI = making intelligent machines*

*Standard model:*

**Machines are intelligent to the extent that their actions can be expected to achieve their objectives**

*The goal is general-purpose AI:*

**Being capable of quickly learning high-quality behavior in "any" task environment**

# Current approach to AI: Deep Learning

# Problems with AI



**Which one is a real human?**

# Problems with AI



**Which one is a real human?**

*New York Times article:*
https://www.nytimes.com/interactive/
2020/11/21/science/artificial-
intelligence-fake-people-faces.html

# Problems with AI



WORLD

## A Military Drone With A Mind Of Its Own Was Used In Combat, U.N. Says

June 1, 2021 · 3:09 PM ET

JOE HERNANDEZ

# Problems with AI



**Wall-E!!**

# Problems with AI

**Can GPT-4 pursue goals?**

**Ask Kevin Roose (NYTimes)**

# Problems with AI

*Q: what is the largest even number?*
A: The largest even number is 9,999,999,998.

*Q: I started with 20 euros and gave 10 euros to my friend. How much do we have in total?*
A: You and your friend have a total of 30 euros.

*Q: which is bigger? an elephant or a cat?*
A: An elephant is bigger than a cat.

*Q: which is not bigger than the other? an elephant or a cat?*
A: Neither an elephant nor a cat is bigger than the other.

chatGPT 3.0

# A new model: Stuart Russell

Machines are **intelligent** to the extent that their actions can be expected to achieve **their** objectives

Machines are **<u>beneficial</u>** to the extent that their actions can be expected to achieve **<u>our</u>** objectives

# Provably beneficial AI

1. Robot goal: satisfy human preferences

2. Robot is *uncertain* about human preferences

3. Human behavior provides evidence of preferences

# The off-switch problem



A robot with a fixed objective has an incentive to disable its own off-switch

*"You can't fetch the coffee if you're dead"*

A robot with uncertainty about objective won't behave this way

# Responsible AI

*Many humans*
 => connections to **moral philosophy, economics**

*Many machines*
 => avoiding unanticipated strategic interactions

*Non-rational humans*
 => connections to **cognitive psychology, neuroscience**

*Foundations*
 => rebuild each area of AI (search, planning, etc.)