# Unintended consequences of AI

May 2024
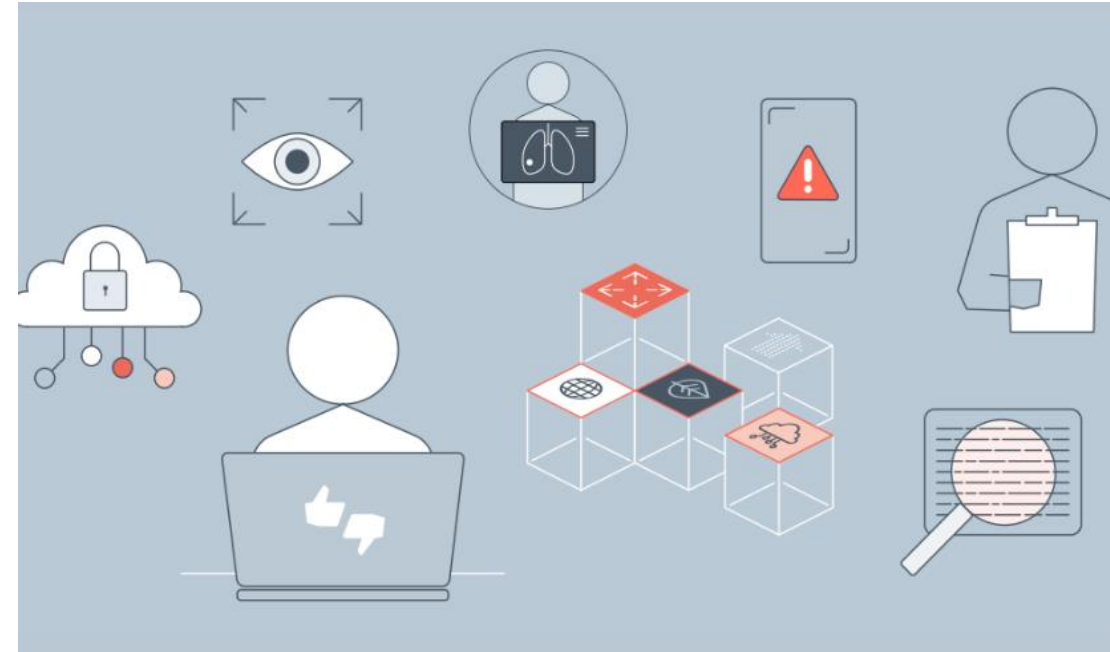
Fran Bennett

Ada Lovelace Institute

# About the Ada Lovelace Institute

The Ada Lovelace Institute is an independent research institute with a mission to make data and AI work for people and society.

Tina Ricky

# Sam Altman Calls AI 'Most Important Step Yet' for Humans, Tech

## OpenAI CEO Sam Altman Says AI Could End Poverty

BY **PYMNTS** | JUNE 23, 2023
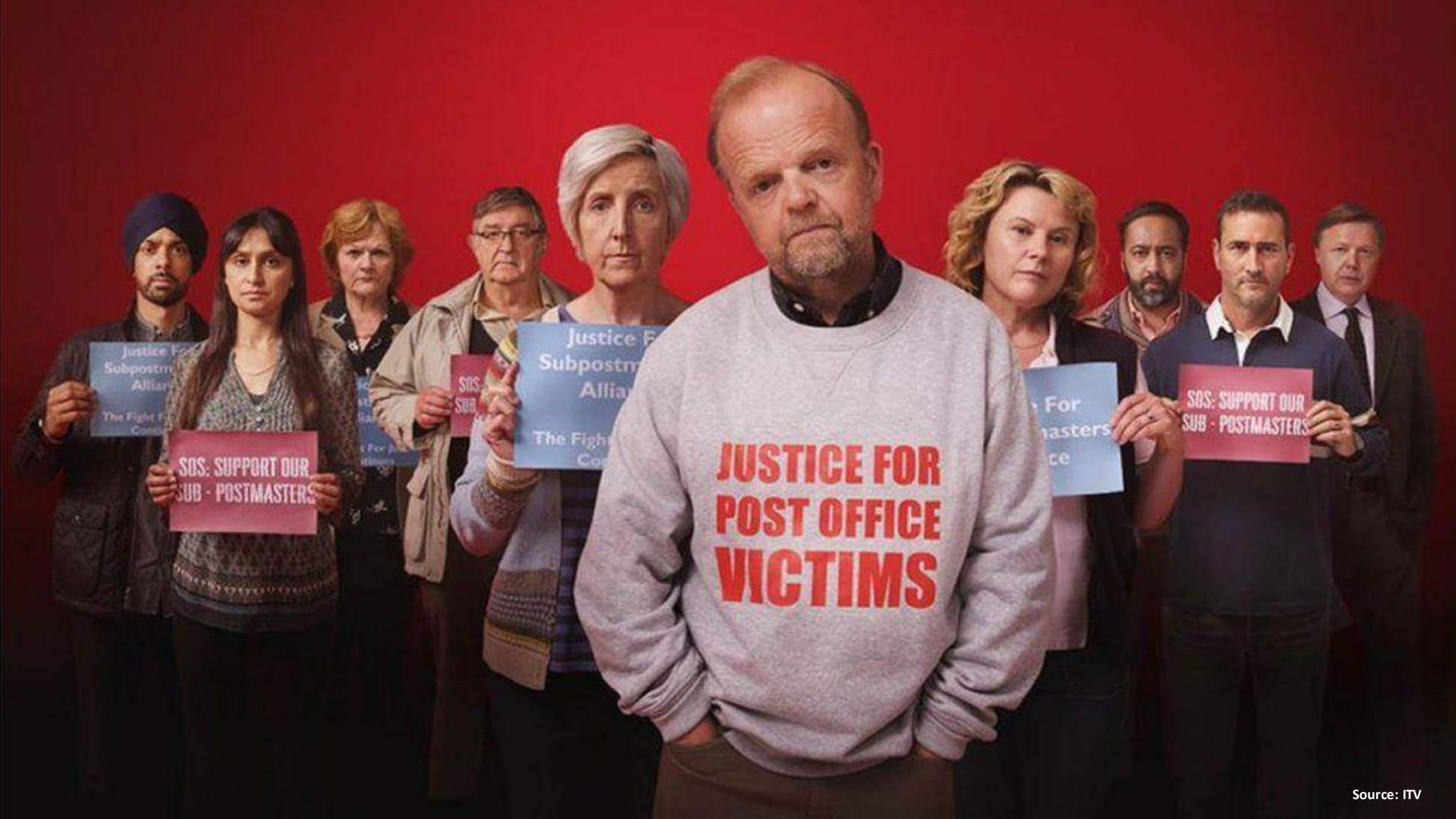
UNIVERSALLY BAD HEALTHCARE

**OPENAI CEO SAYS AI WILL GIVE MEDICAL ADVICE TO PEOPLE TOO POOR TO AFFORD DOCTORS**

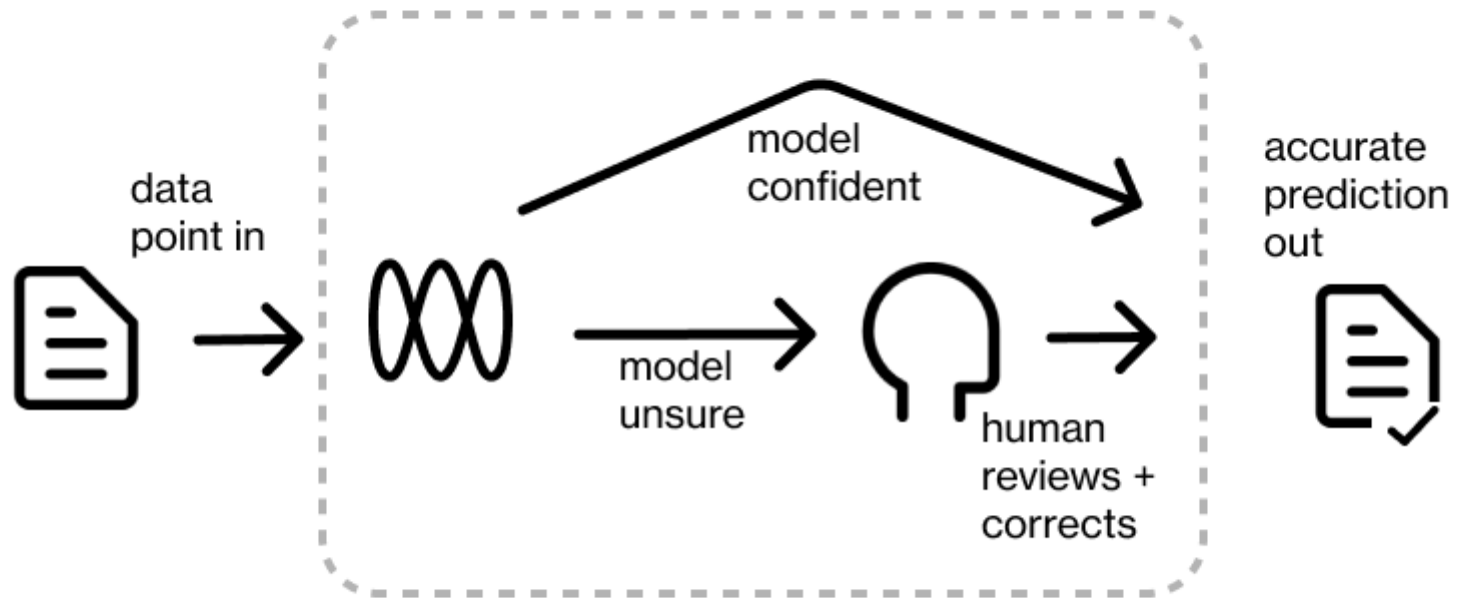| | |
|---|---|
| Impossible Tasks | Conceptually Impossible |
| | Practically Impossible |
| Engineering Failures | Design Failures |
| | Implementation Failures |
| | Missing Safety Features |
| Post-Deployment Failures | Robustness Issues |
| | Failure under Adversarial Attacks |
| | Unanticipated Interactions |
| Communication Failures | Falsified or Overstated Capabilities |
| | Misrepresented Capabilities |

data point in

model confident

accurate prediction out

model unsure

human reviews + corrects

**Workers-in-the-Loop AI deployment**

〽 Humanloop

# Figure 2: FDA-style oversight for foundation models



UPSTREAM

DATA
LAYER

COMPUTE
LAYER

**HOW (e.g.):** Pre-notifications
data documentation

**HOW (e.g.):** Third-party
evaluations

FOUNDATION MODEL
DEVELOPER LAYER

HOST LAYER

Cross-sector
pre-approval gate

**HOW (e.g.):** Algorithmic
impact assessments

APPLICATION
LAYER

DOWNSTREAM

APPLICATION
USER

Sector specific
pre-approval gate

**HOW (e.g.):** post-market monitoring,
incident reporting

**Source:** Ada Lovelace Institute

- **Biases** from the training data or feedback process

- **Black box** and **stochastic** models – can't investigate why they do what they do, and outputs are a little different each time

- **Non-transparent supply chain –** what was the input data? What was the feedback process? Are there extra control layers?

- **Test safety in a crash** – as we do with cars

- **Assign liability for harm** – as we do with food poisoning

- **Hear voices of those affected** – as we do with planning permission

- **Right to contest decisions and request explanations** – as we do with human bureaucracies

# Thank you

Fran Bennett, Ada Lovelace Institute

*@AdaLovelaceInstitute*

Ada
Lovelace
Institute