# AGI Safety & Alignment

Data Circles Journal Club 3-2-22

# **Weapons of Math Destruction** by Cathy O'Neil

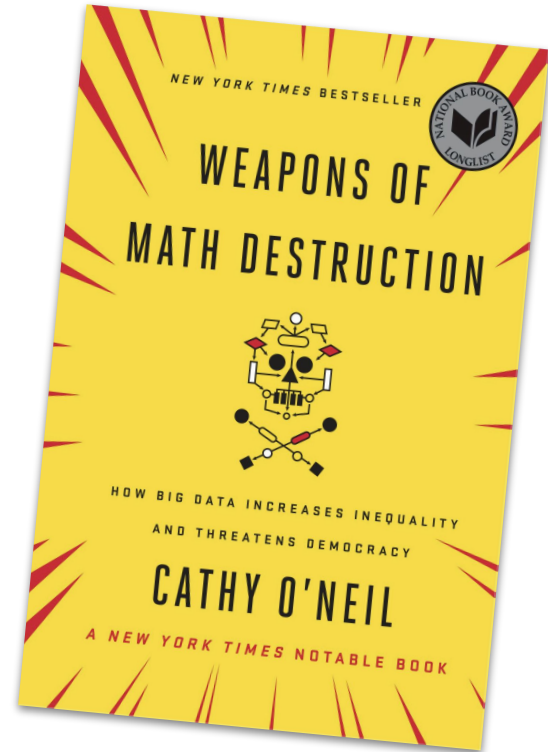Shell Shocked: Journey of Disillusionment

Arms Race: Going to College

Civilian Casualties: Justice in the Age of Big Data

Ineligible to Serve: Getting a Job

No Safe Zone: Getting Insurance
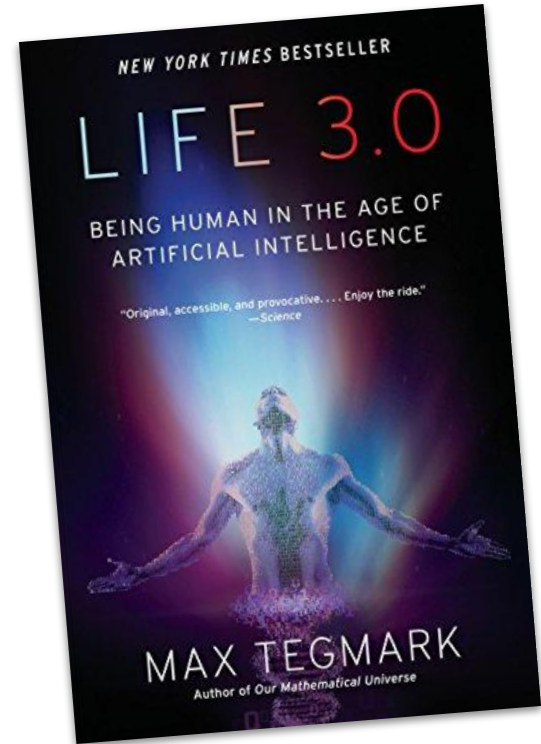
Targeted Citizen: Civic Life

# Life 3.0 by Max Tegmark

Life 1.0 fixed hardware & software

Life 2.0 fixed hardware, upgrade software

Life 3.0 upgrade both hardware & software

Explore scenarios with AGI

*"Any sufficiently **advanced technology** is indistinguishable from **magic**."*

*~Arthur C. Clarke*
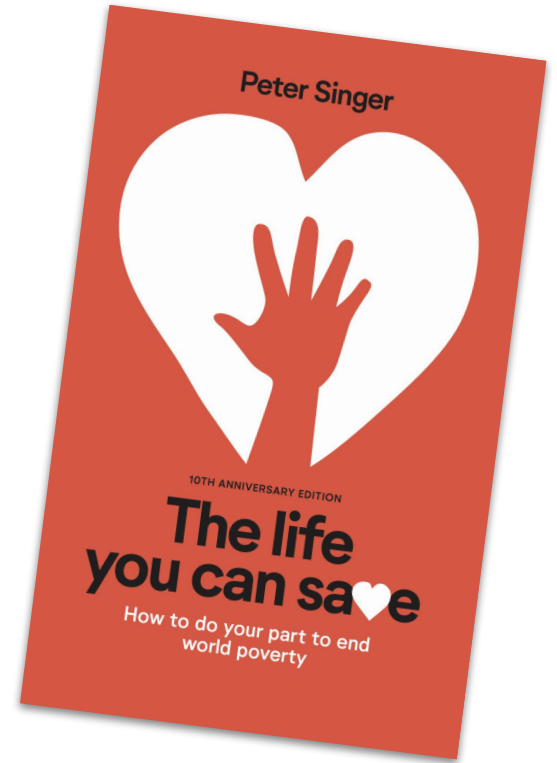
# The Life You Can Save by Peter Singer

Effective altruism

80,000 hours

Pressing global issues list

Highest priority areas

- Global priorities research
- Building effective altruism
- Reducing global catastrophic biological risks
- Positively shaping the development of AI

*"Fearing a rise of **killer robots** is like worrying about **overpopulation on Mars**"*
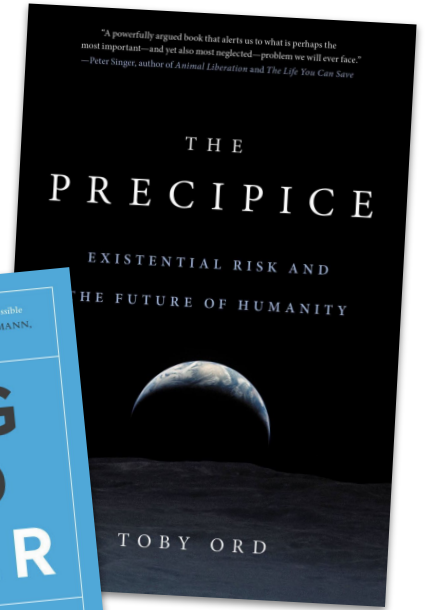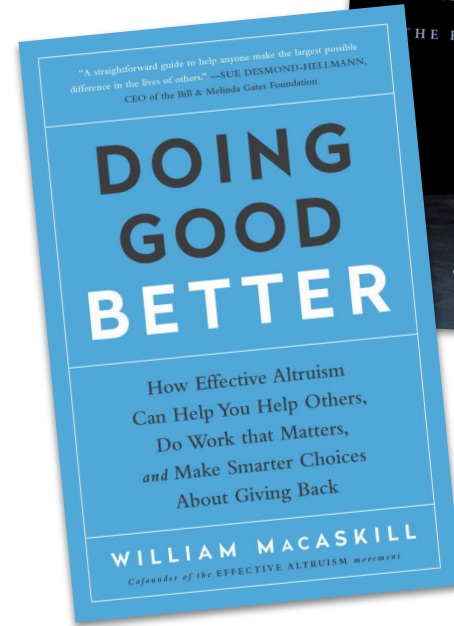
*~Andrew Ng*

# AGI Safety

aka the Control Problem, AI Alignment Problem

Long-termism

- Will MacAskill, *Doing Good Better* (2015)
- Toby Ord, *The Precipice* (2020)

Recommended readings

- Nick Bostrom, *Superintelligence* (2014)
- Stuart Russell, *Human Compatible* (2019)
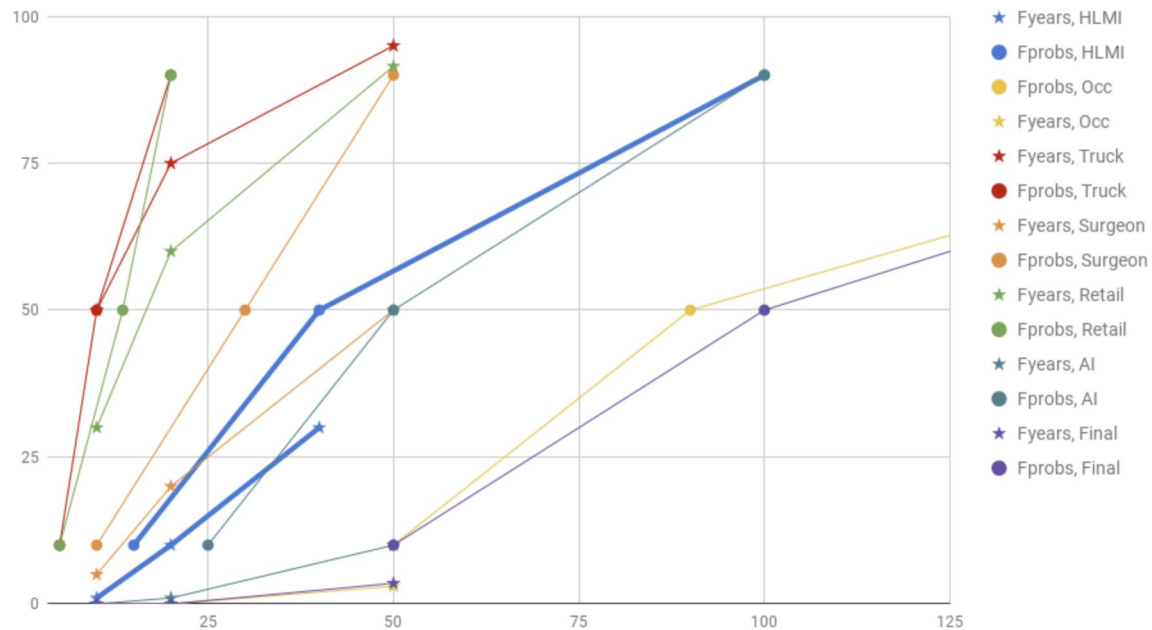- Brian Christian, *Alignment Problem* (2020)

*"We will, **sooner or later**, build an artificial **agent** with **general intelligence**."*

*~Rob Miles*

# 2016 Expert Survey

https://aiimpacts.org/2016-expert-survey-on-progress-in-ai/



AI forecasts by framing and milestone

Legend:
★ Fyears, HLMI
● Fprobs, HLMI
● Fprobs, Occ
★ Fyears, Occ
★ Fyears, Truck
● Fprobs, Truck
★ Fyears, Surgeon
● Fprobs, Surgeon
★ Fyears, Retail
● Fprobs, Retail
★ Fyears, AI
● Fprobs, AI
★ Fyears, Final
● Fprobs, Final

# Superintelligence by Nick Bostrom

History of ups and downs
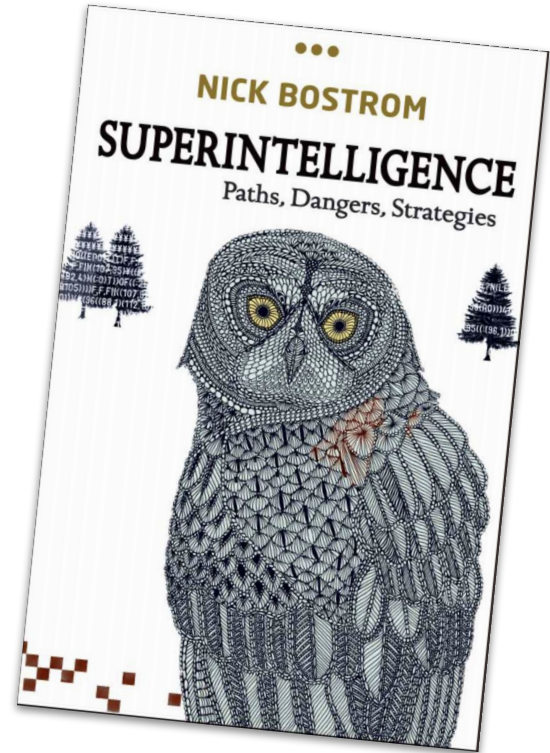
Paths: AI vs Whole Brain Emulation (WBE)

Speed: fast vs slow

Orthogonality thesis: Paperclip problem

Instrumental goals

- Self-preservation
- Resource acquisition
- Self-improvement

Learn human values

*"The development of **full artificial intelligence** could spell the **end** of the **human race**."*

*~Stephen Hawking*

# Human Compatible by Stuart Russell
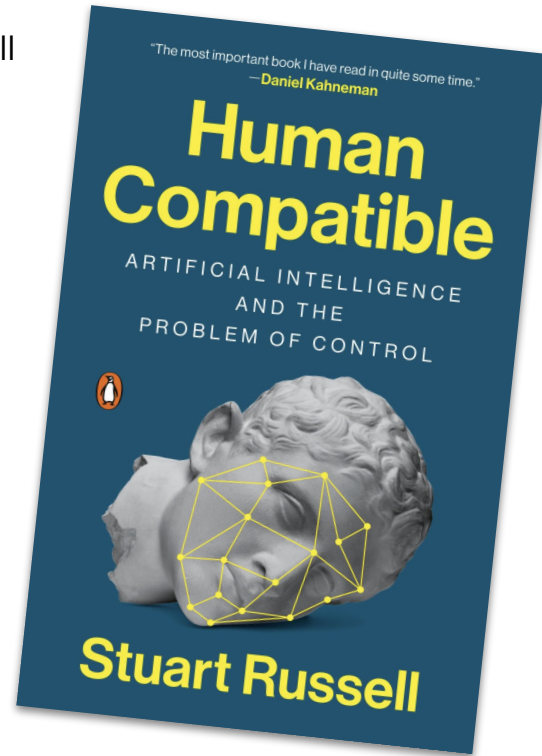
Conceptual breakthroughs needed

Gorilla problem, 2nd species

Beneficial machines:

- Purely altruistic
- Humble
- Learn to predict

Many benefits

Many risks



"The most important book I have read in quite some time."
—Daniel Kahneman

**Human Compatible**

ARTIFICIAL INTELLIGENCE AND THE PROBLEM OF CONTROL

**Stuart Russell**

*"Be **careful** what you **wish** for!"*

*~King Midas*

# Alignment Problem by Brian Christian
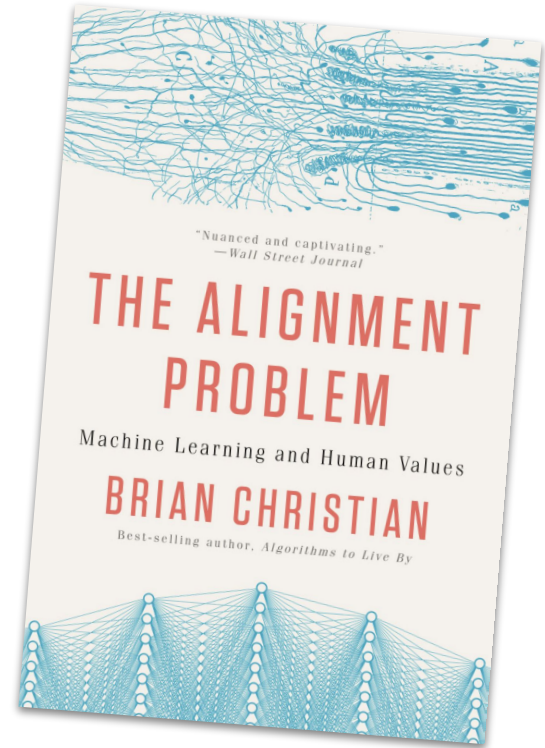
Machine Learning and Human Values

Prophecy: Representation, Fairness, Transparency

Agency
- Reinforcement: Policy vs Value, Temporal Diff
- Shaping: Rewards & Behaviors
- Curiosity: Novelty, Intrinsic Motivation

Normativity:
- Imitation: Recover Cascading Errors
- Inference: Demos, Feedback, Cooperation
- Uncertainly: Impact, Corrigible

"Biased **algorithms** are easier to fix than biased **people**."

~Sendhil Mullainathan

# References

Cathy O'Neil, *Weapons of Math Destruction* (2016)

Max Tegmark, *Life 3.0* (2017)

Peter Singer, *The Life You Can Save* (2019)

Will MacAskill, *Doing Good Better* (2015)

Toby Ord, *The Precipice* (2020)

Nick Bostrom, *Superintelligence* (2014)

Stuart Russell, *Human Compatible* (2019)

Brian Christian, *Alignment Problem* (2020)