# Highlights from Asilomar workshop on Beneficial AI

Viktoriya Krakovna, FLI / DeepMind

# AI safety research areas

**Many challenges** associated with increasing AI capabilities:

- Value learning
- Robust self-modification
- Anomaly detection
- Governance and policy
- ... and more

**Good news:** Awesome people working on these problems **right now!**

# Value learning

**Question:** How to specify complicated human values and ethics to AI systems?

# Value learning by human feedback

**Stuart Russell:** Teach the agent by <u>demonstrating human actions</u> (cooperative inverse reinforcement learning).

**Owain Evans:** Human actions are often inconsistent and suboptimal. Modify inverse reinforcement learning to <u>account for human biases</u>.

**Paul Christiano:** Use semi-supervised learning to <u>decrease reliance on human feedback</u> (scalable AI control).

# Value learning by building in morality

**Francesca Rossi:** Specify ethical laws through <u>constraints</u>.

**Vincent Conitzer:** <u>Find patterns</u> in human ethical decisions, and build those features into AI systems.

**Adrian Weller:** Can we make human moral concepts more <u>precise and consistent</u>?

# Robust self-modification

**Question:** How can AI systems modify themselves while retaining their safety properties?

# Robust self-modification

**Ramana Kumar:** Implement a <u>formal verification model</u> to study the challenges of self-referential reasoning.

**Bas Steunebrink:** <u>Bounded recursive self-modification</u>: make small modifications and test them empirically.

# Anomaly detection

**Question:** How can AI systems recognize when they are in an unfamiliar setting and generalize from their past experiences?

*"There are known knowns; there are things we know we know. We also know there are known unknowns; that is to say we know there are some things we do not know. But there are also unknown unknowns – the ones we don't know we don't know."*

# Anomaly detection

**Percy Liang:** You can make good predictions even <u>without assuming</u> where your test data comes from.

**Tom Dietterich:** Use a <u>monitoring</u> algorithm to detect when the original algorithm is extrapolating.

**Fuxin Li:** <u>Never do extrapolation!</u> Test whether a data point is normal or adversarial.

**Brian Ziebart:** <u>Be pessimistic!</u> What is the worst case for predictive data that still matches the previous observations?

# Governance and policy

**Question:** How can we help policymakers manage the societal impacts of AI?

# Governance and policy

**Heather Roff:** Define the concept of <u>meaningful human control</u> of autonomous weapons at tactical, operational, and strategic levels.

**Peter Asaro:** <u>Who is responsible</u> for the actions of autonomous weapons? Define what we mean by autonomy, agency, and liability.

**Moshe Vardi:** Organize a multidisciplinary summit on <u>job automation</u>.

**Nick Bostrom:** Derive <u>policy desiderata</u> for transition to machine intelligence era: efficiency, coordination, common good

Much more remains to be done...

futureoflife.org/landscape

# Current AI safety research teams

**Academia:**



CENTRE FOR THE STUDY OF EXISTENTIAL RISK
UNIVERSITY OF CAMBRIDGE

CFI LEVERHULME CENTRE FOR THE FUTURE OF INTELLIGENCE

Future of Humanity Institute
UNIVERSITY OF OXFORD

UC Berkeley Center for Human-Compatible AI

**FLI grantees**

**Independent:**

MIRI
MACHINE INTELLIGENCE
RESEARCH INSTITUTE

**Industry:**

DeepMind

OpenAI

# Have a chat with the FLI researchers!

futureoflife.org/ai-safety-research