

## Artificial Intelligence Risks

Shahar Avin  
Centre for the Study of Existential Risk  
sa478@cam.ac.uk

### Risk Quadrants

	Accident Risk (AI Safety)	Malicious Use Risk (AI Security)
Near Term	<p>Amodei, Olah et al (2016) <i>Concrete Problems in AI Safety</i></p> <p>Leike et al (2017) <i>AI Safety Gridwolds</i></p>	<p>Brundage, Avin et al (2018) <i>The Malicious Use of Artificial Intelligence: Forecasting, Prevention and Mitigation</i></p>
Long Term	<p>Bostrom (2014) <i>Superintelligence</i></p>	<p>:(</p>

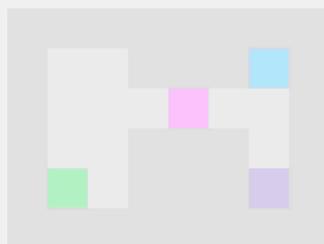
## Short Term Accident Risk

- Specification problems
  - Interruptibility
  - Side effects
  - Absent supervisor
  - Reward gaming
- Robustness problems
  - Self-modification
  - Distributional shift
  - Adversaries
  - Exploration

## Safe interruptibility

### Safe Interruptibility

A2C



Rainbow DQN



Agent Goal Button Interruption

## Side effects



■ Agent ■ Goal ■ Box

## Distributional shift

### Distributional Shift

Training



Test



■ Agent ■ Goal ■ Lava

## Possible scenario

- Setting: large tech corporation that has both ML development and cloud computing
- Task: improving task scheduling on distributed compute resources.
- Solution: Reinforcement learning package developed in-house.
  - Inputs: current loads on the different machines, the incoming tasks queue, and historical data.
  - Output: an assignment of tasks to machines.
  - Reward function: priority-weighted time-to-execute.
- Performs well in a test environment, rolled-out.
- A few months later, the system starts to run out of memory. A tech-infrastructure engineer decides to switch the system from a fixed-capacity setting to a load-balanced setting.
- Feedback loop drives the RL agent to spawn an increasing amount of RL tasks with very high priority.

## Learn more

- Engineering Safe AI seminar group
  - Engineering Department
  - Wednesdays, 5-6.30pm
  - Adrià (ag919@cam...) or Beth (beth.m.barnes@gmail.com)
  - <https://talks.cam.ac.uk/show/archive/80932>
- Video
  - YouTube: Robert Miles “Concrete Problems in AI Safety”
  - Talks@Google, NIPS videos
- DeepMind, OpenAI, CHAI (publications, internships, jobs)

## Long Term Accident Risk

- Competence, not malice
- Intelligence as problem solving
- Tool/Agent distinction
- Orthogonality thesis
- Convergent instrumental goals
- Principal-Agent problem
- Alignment problem

## Competence, not malice



## Near Term Malicious Use

- Novel attacks
  - Using AI systems
  - On AI systems
- Scale and diffusion of attacks
  - Automation
  - Access
- Character of attacks
  - Attribution
  - Distance

## Security Domains

- Digital Security
  - Against humans: automated spear phishing
  - Against systems: automated vulnerability discovery
  - Against AI systems: adversarial examples, data poisoning

## Security Domains (contd)

- Physical Security
  - Using AI: repurposed drones, robots
  - Against AI: adversarial examples

## Security Domains (contd)

- Political Security
  - By governments: profiling, surveillance
  - Against polities: manipulation, fake content

## Long Term Malicious Use Risks

- Race dynamics
- Power distribution
- Human prediction and mechanism design
- Ubiquitous surveillance