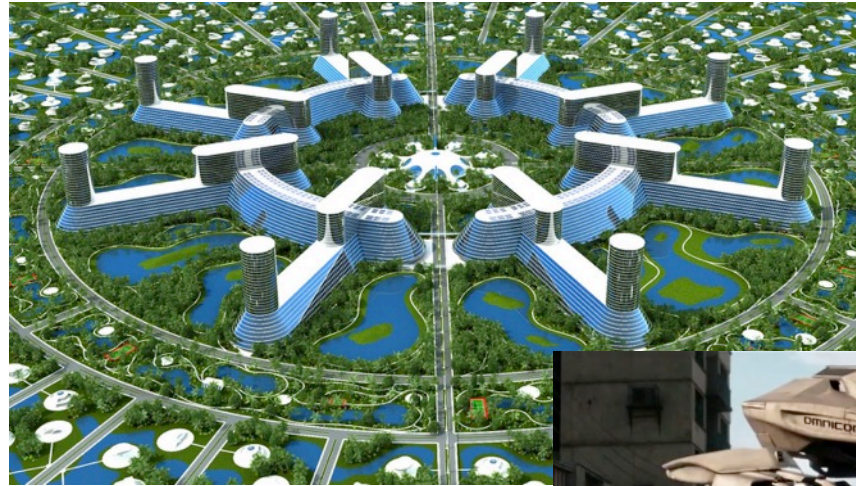# Exploring AGI Scenarios

Shahar Avin
sa478@cam.ac.uk

# AGI strategy

This is Bob.

Bob heads an AGI R&D lab.

What should Bob do?
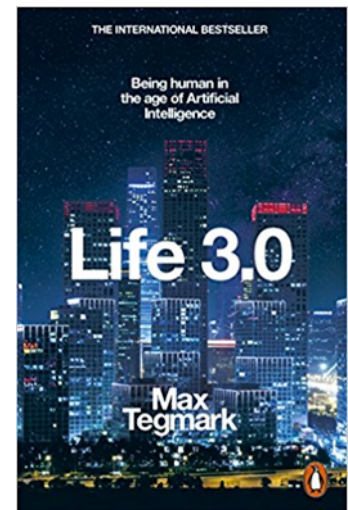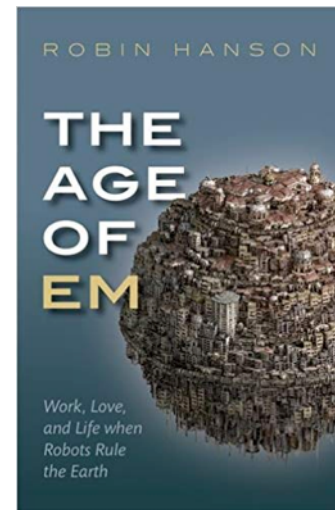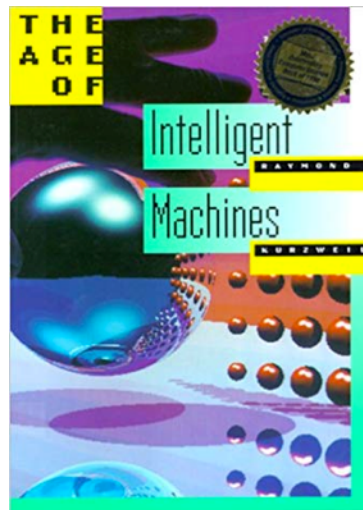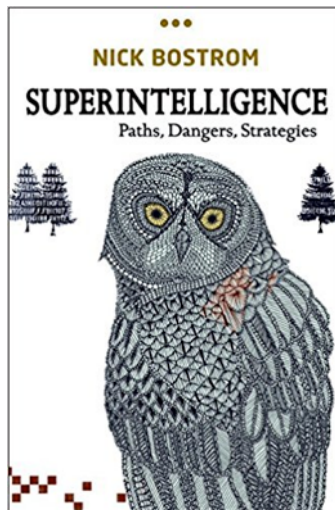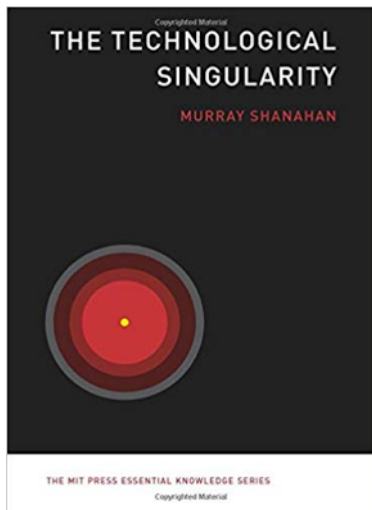
# AGI futures narratives

- Tech utopia

- Arms race
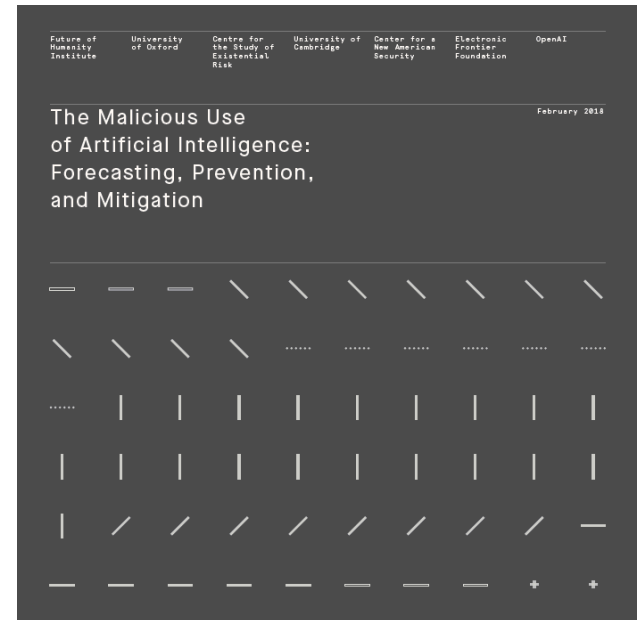
- Malicious use

- Existential risk

# How do we explore and communicate these futures?

Single author exploration

# How do we explore and communicate these futures?
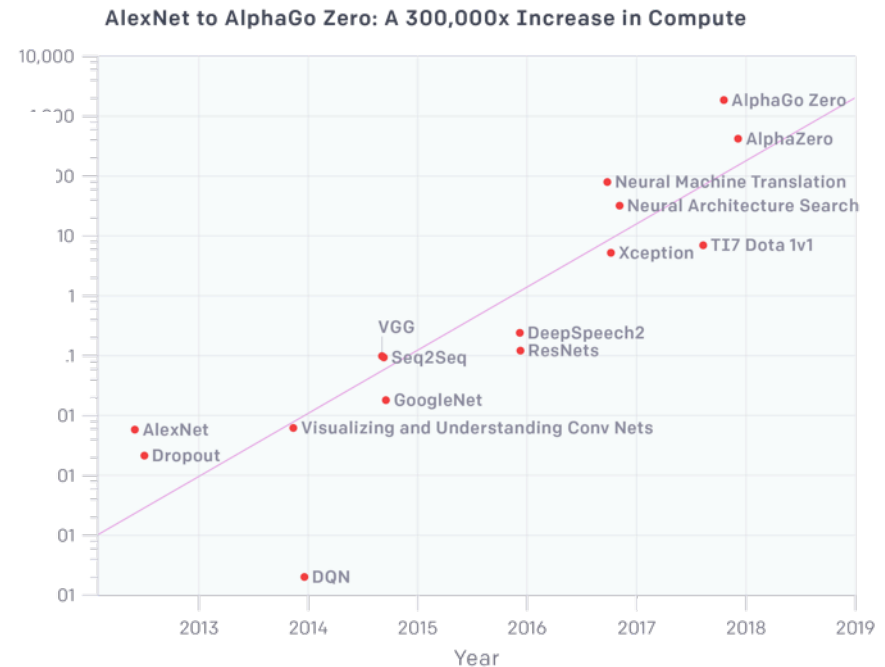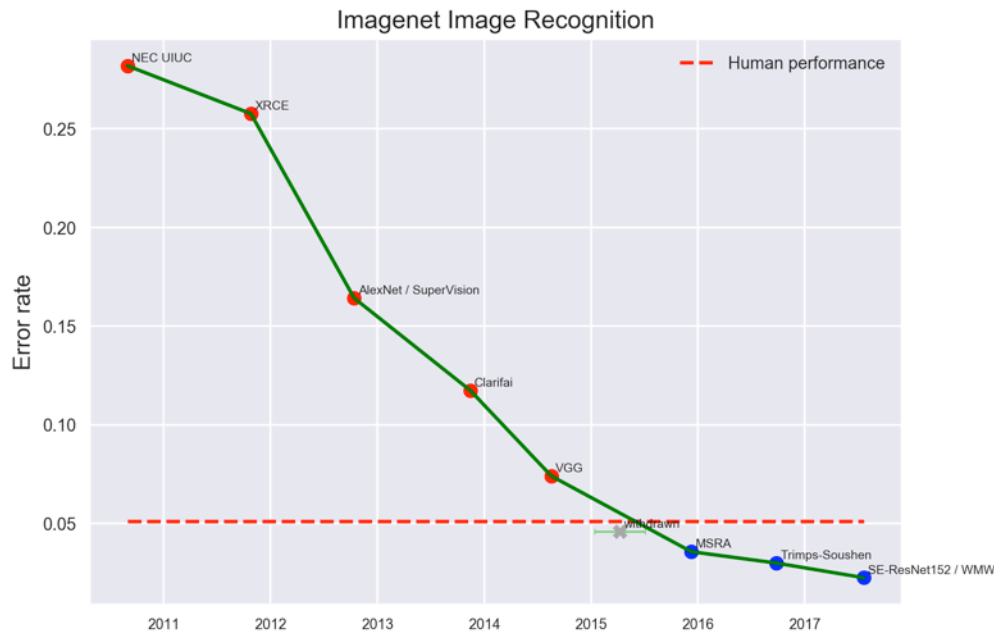
Expert workshops, multi-authored reports





http://maliciousaireport.com/

# How do we explore and communicate these futures?

## Data trends

https://www.eff.org/ai/metrics



https://blog.openai.com/

UNIVERSITY OF CAMBRIDGE

# How do we explore and communicate these futures?

## Aggregate probability estimates

https://www.getguesstimate.com





https://www.metaculus.com

UNIVERSITY OF CAMBRIDGE

# How do we explore and communicate these futures?

Video games



https://bit.ly/2uqXNUn



http://www.decisionproblem.com/paperclips/

# SPACECRAFT SCIENTIST/ENGINEER



What my friends think I do

What my parents think I do

What society thinks I do

What my boss thinks I do

What I think I do

What I really do

# What should we be looking at?

# Development factors

Inputs

# Development factors

## Nature of the problem

# Development factors

Control, incentives, openness



APRIL 9, 2018

**OpenAI Charter**

We're releasing a charter that describes the principles we use to execute on OpenAI's mission. This document reflects the strategy we've refined over the past two years, including feedback from many people internal and external to OpenAI. The timeline to AGI remains uncertain, but our charter will guide us in acting in the best interests of humanity throughout its development.



Solve intelligence. Use it to make the world a better place.

# Development factors

## Safety and Security

# Deployment factors

All of the above (I/O, Control, Safety & Security)!

Plus: generality, capability, domains of application

SIGHT

**Cloud Vision API**
Image recognition and classification.

**Cloud Video Intelligence API**
Scene-level video annotation.

**AutoML Vision** BETA
Custom image classification models.

LANGUAGE

**Cloud Translation API**
Language detection and translation.

**Cloud Natural Language API**
Text parsing and analysis.

**AutoML Transla**
Custom domain-spe

**AutoML Natural**
Custom text classifi

CONVERSATION

**Dialogflow Enterprise Edition**
Build conversational interfaces.

**Cloud Text-to-Speech API**
Convert text to speech.

The scale of intelligence:

mouse | village idiot

chimp | Einstein

recursively self-improved AI

# Landscape factors

Number and identity of actors

# Landscape factors

Inter-actor relationships

# Landscape factors

## International relations



### Artificial Intelligence Strategies

| Top row | |
|---|---|
| March: Pan-Canadian AI Strategy | Canada |
| May: AI Singapore Announced | Singapore |
| October: AI Strategy 2031 | UAE |
| December: Finland's AI Strategy | Finland |
| January: Budget for AI Taiwan | Taiwan |
| March: AI at the Service of Citizens | Italy |
| April: First Workshop for Strategy | Tunisia |
| April: UK AI Sector Deal | UK |
| May: White House Summit on AI | USA |
| May: Sweden's AI Strategy | Sweden |
| June: Towards an AI Strategy in Mexico | Mexico |
| Fall 2018: EU's AI Strategy | EU |

2017 — 2018

| Bottom row | |
|---|---|
| March: AI Technology Strategy | Japan |
| July: Next Generation AI Plan | China |
| December: Three-Year Action Plan | China |
| January: Blockchain and AI Task Force | Kenya |
| January: Strategy for Digital Growth | Denmark |
| March: France's AI Strategy | France |
| April: Communication on AI | EU |
| May: Australian Budget | Australia |
| May: AI R&D Strategy | South Korea |
| June: National Strategy for AI | India |
| Fall 2018: Germany's AI Strategy | Germany |

2018-07-13 | Politics + AI | Tim Dutton

UNIVERSITY OF CAMBRIDGE

# Landscape factors

## Society and culture



The New York Times

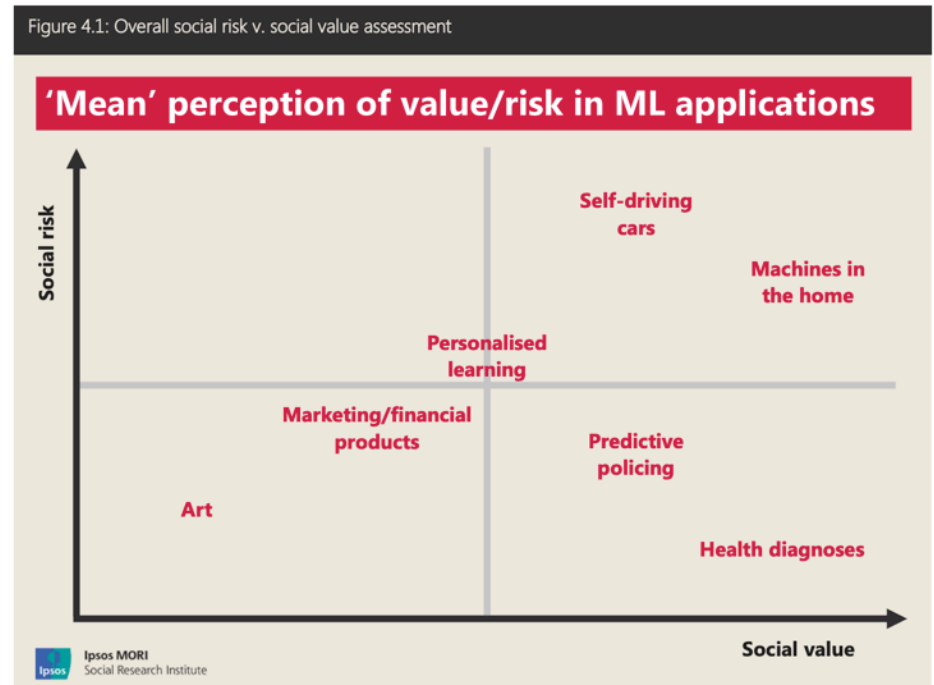**Wielding Rocks and Knives, Arizonans Attack Self-Driving Cars**



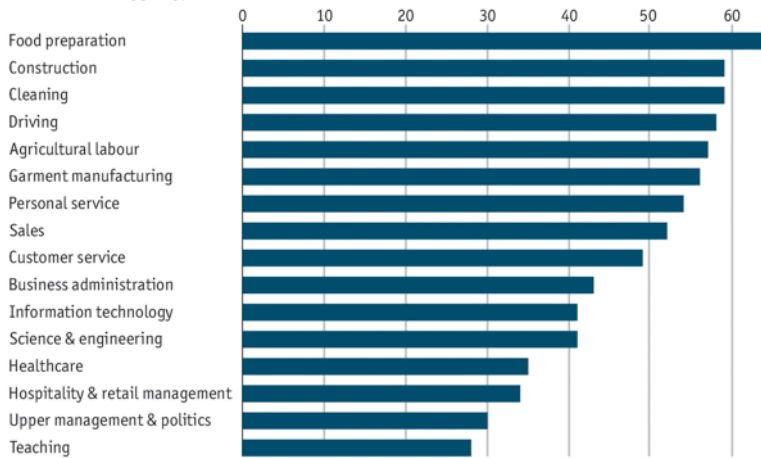Figure 4.1: Overall social risk v. social value assessment

'Mean' perception of value/risk in ML applications

- Self-driving cars
- Machines in the home
- Personalised learning
- Marketing/financial products
- Predictive policing
- Art
- Health diagnoses

Social risk / Social value

Ipsos MORI
Social Research Institute

UNIVERSITY OF CAMBRIDGE

# Landscape factors

## The economy



Productivity growth and hourly compensation growth, 1948–2017

**Automated for the people**
Automation risk by job type, %

Source: OECD

Economist.com

UNIVERSITY OF CAMBRIDGE

# Landscape factors

## The environment





AI for Earth

AI for Earth is a Microsoft program aimed at empowering people and organizations to solve global environmental challenges by increasing access to AI tools and educational opportunities, while accelerating innovation.

## Security

# How do we explore and communicate these futures?

## Scenario role-play







Peter Perla's
**The Art of Wargaming**
A guide for Professionals and Hobbyists

Edited By John Curry