



MIRI

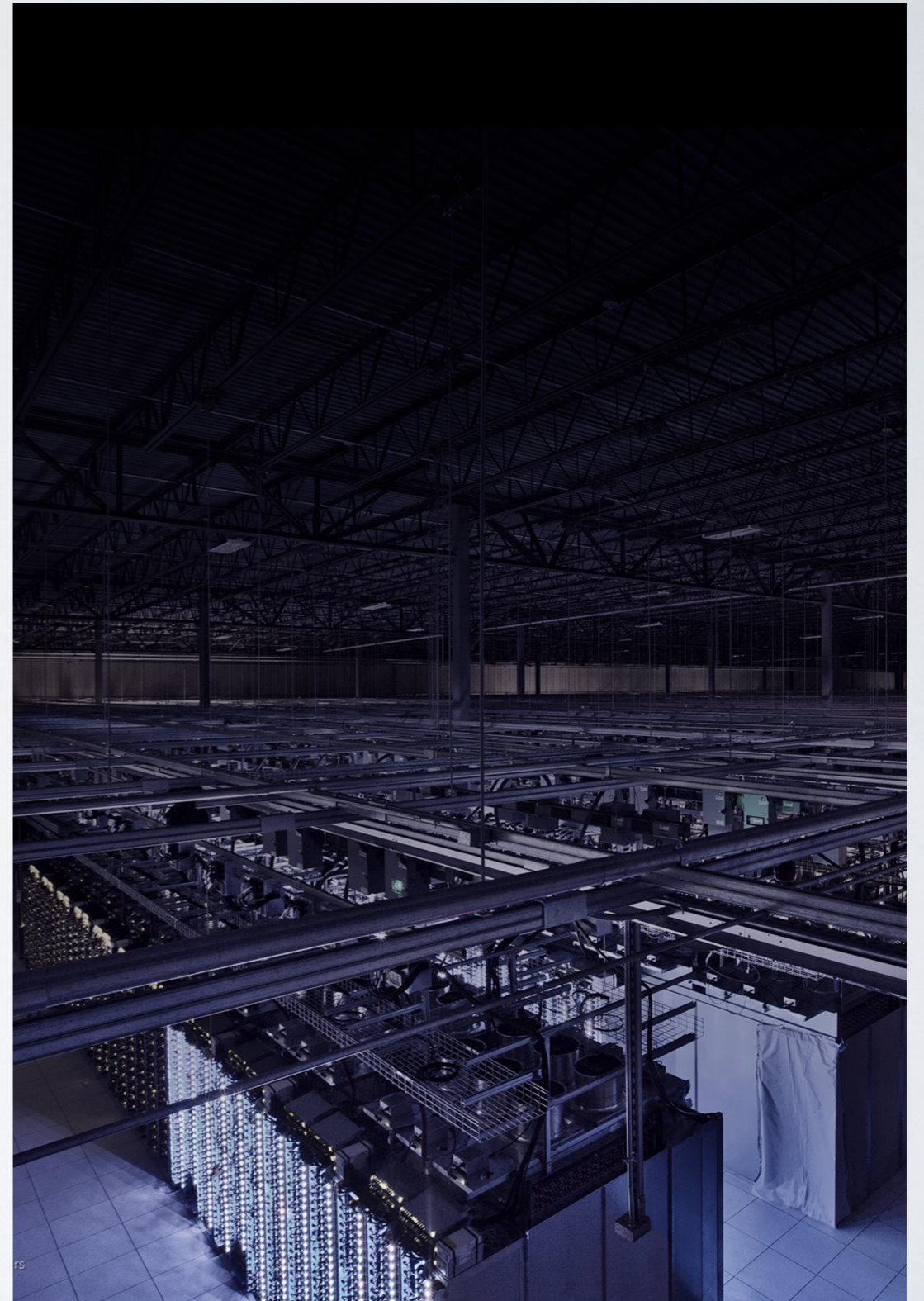
MACHINE INTELLIGENCE
RESEARCH INSTITUTE

Mission:

“to ensure that the creation of smarter-than-human intelligence has a positive impact”

Narrow AI

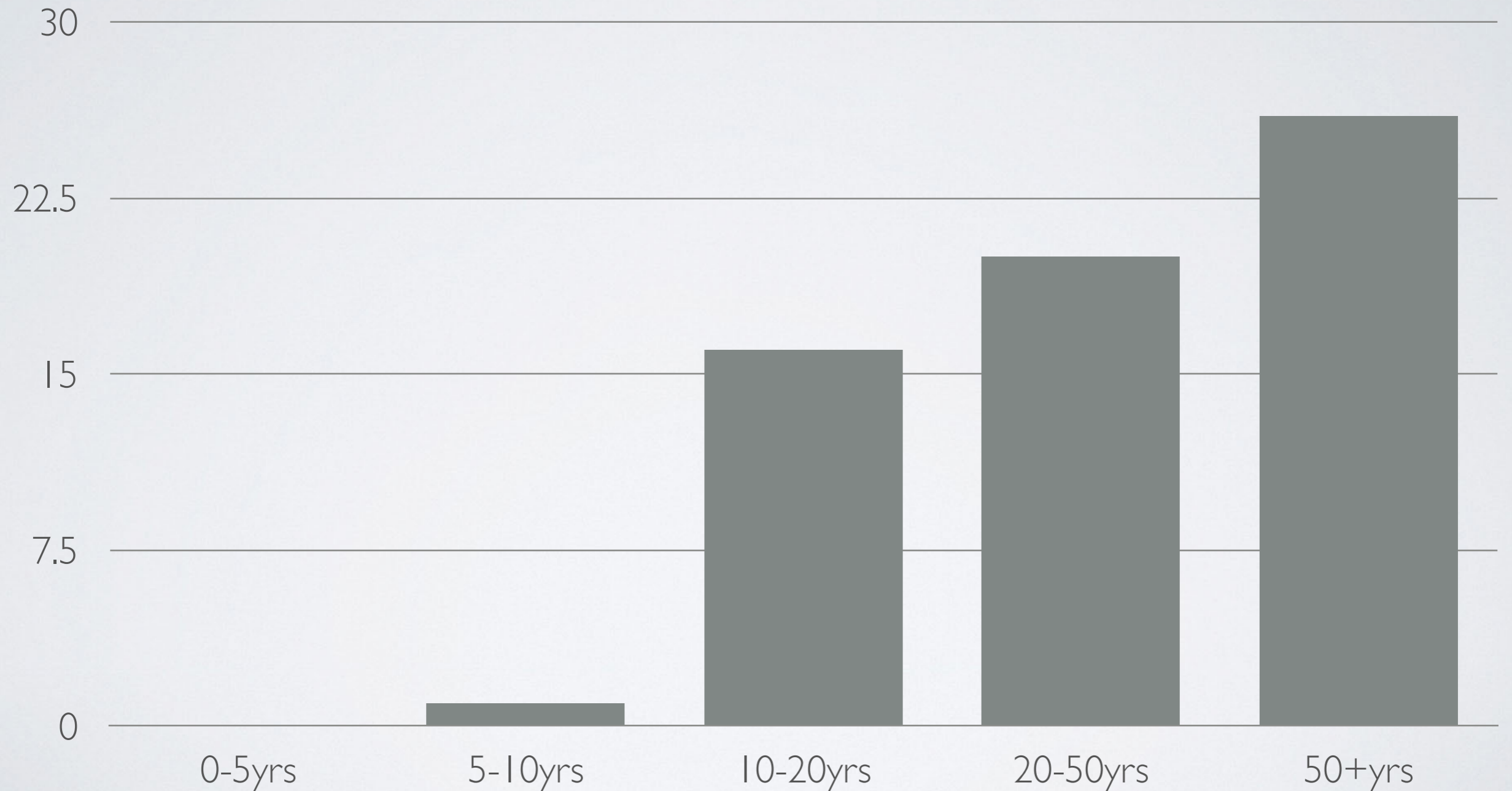
AGI



Expert opinion on AGI

Michie (1973), n=63:

Time to AGI



Expert opinion on AGI

“Assume that human scientific activity continues without major negative disruption. By what year would you see a (10% / 50% / 90%) probability for [AGI] to exist?”

Müller & Bostrom (2014), n=29:

	10%	50%	90%
median response	2024	2050	2070

Luke's estimate:

	10%	50%	90%
Luke	2030	2070	2140

See Muehlhauser, “When Will AI Be Created?” (2013)

Expert opinion on AGI

“Assume... that [AGI] will at some point exist. How likely do you then think it is that within (2 years / 30 years) thereafter there will be machine [superintelligence]?”

Müller & Bostrom (2014), n=29:

	2 years	30 years
median response	5%	50%

Luke's estimate:

	2 years	30 years
Luke	15%	85%

See Bostrom (2014) ch. 4 and Yudkowsky, “Intelligence Explosion Microeconomics”

Expert opinion on AGI

“Assume... that [AGI] will at some point exist. How positive or negative would be overall impact on humanity, in the long run?”

Müller & Bostrom (2014), n=29:

	Extremely good	On balance good	Neutral-ish	On balance bad	Extremely bad
mean	20%	40%	19%	13%	8%

Luke's estimate (on Tuesday and Thursdays, still volatile):

	Extremely good	On balance good	Neutral-ish	On balance bad	Extremely bad
Luke	19%	1%	~0%	5%	75%

See Bostrom (2014) and stuff that hasn't been published yet

What's the EA case for efforts toward Friendly AI?

- 1. Astronomical stakes:** Most QALYs are in the long-term future, so what matters most is that we shape the trajectory of the long-term future in a robustly positive direction. See Beckstead's PhD thesis, and Yudkowsky's talk at 10:40am tomorrow.
- 2. AI is the key lever on the long-term future:** See Bostrom's *Superintelligence* (2014).
- 3. Friendly AI work is urgent, tractable, and uncrowded:** Most AGIs do not stably optimize for desirable values, Friendly AI is strictly (much) harder than AGI, and today AGI progress is vastly outpacing Friendly AI progress. FAI work is also, it turns out, tractable and uncrowded. See my talk at 4pm tomorrow.

What's the EA case for efforts toward Friendly AI?

I. Astronomical stakes

QALYs we can produce in the long-term future

Eliezer: "If you occupy the incredibly rare and leverage-privileged position of being born into pre-AGI Earth, then the best thing you can do is to help make sure the reachable universe is converted into quality-adjusted life years."

QALYs we can produce in the next century

(Warning: dot may be too small to see on screens smaller than Saturn.)

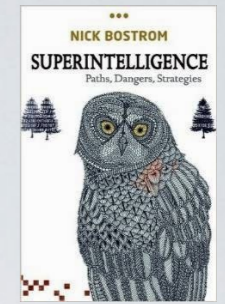
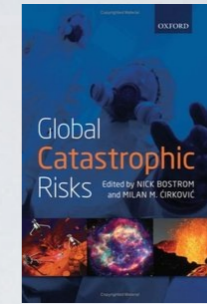
Yes, Friendly AI work is more speculative than, say, Deworm the World.

But if you take astronomical stakes seriously, then...

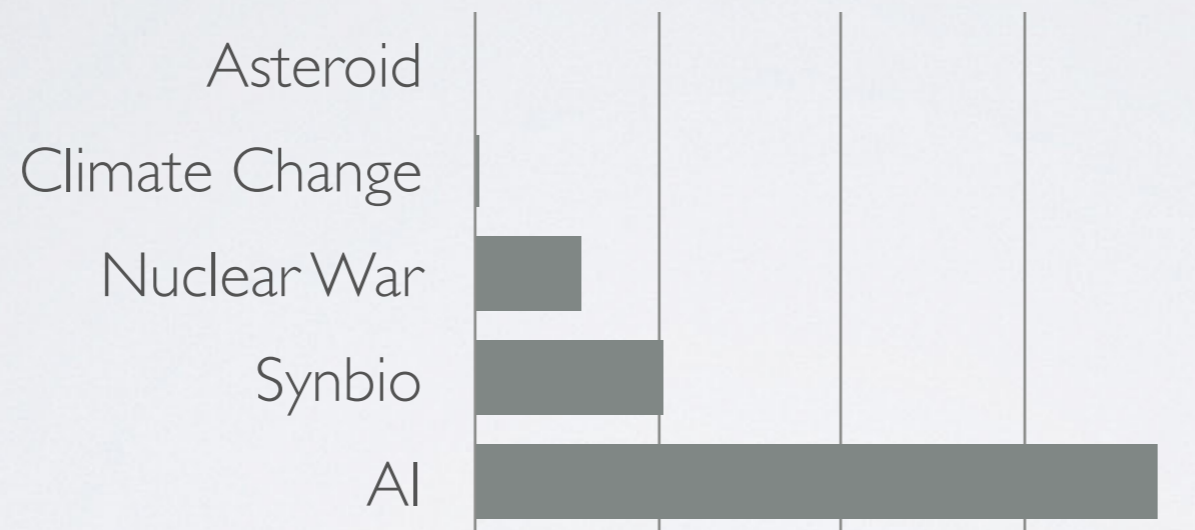
- FHI?
- MIRI?
- Another organization aimed at astronomical value?
- A new organization aimed at astronomical value?
- Some other thing that was never purposely aimed at astronomical value but happens to be optimal by accident?

What's the EA case for efforts toward Friendly AI?

2. AI as the key lever on the long-term future



Chance of being an x-risk in next century (my opinion)



Asymmetry #1:

FAI helps us mitigate other risks, but solving climate change, asteroids, etc. doesn't help us much with other risks.

Asymmetry #2:

FAI is the only technology that lets us convert the reachable universe into quality-adjusted life years.

What's the EA case for efforts toward Friendly AI?

3. Friendly AI work is urgent, tractable, and uncrowded

Uhhhhhhh... see my other talk. :)

So what does MIRI actually do?



Workshops

viewpoints

DOI:10.1145/2644257 Luke Muehlhauser and Bill Hibbard

Viewpoint

Exploratory Engineering in Artificial Intelligence

Using theoretical models to plan for AI safety.

WE REGULARLY SEE examples of new artificial intelligence (AI) capabilities. Google's self-driving car has safely traversed thousands of miles. IBM's Watson beat the "Jeopardy!" champions, and Deep Blue beat the chess champion. Boston Dynamics' Big Dog can walk over uneven terrain and fight itself when it falls over. From many angles, software can recognize faces as well as people can.

As their capabilities improve, AI systems will become increasingly independent of humans. We will be no more able to monitor their decisions than we are now able to check all the math done by today's computers. No doubt such automation will produce tremendous economic value, but will we be able to trust these advanced autonomous systems with so much capability?

The Procrastination Paradox (Brief technical note)

Eliesser Yudkowsky

This document is part of a collection of quick writings of results from the December 2013 MIRI research workshop, written during or directly after the workshop. It describes work done mainly by Alexander Hernandez, Josh Falkner, Benja Fallenstein and Stuart Armstrong, mostly at previous MIRI workshops.

Abstract

A theorem by Marcello Herreshoff, Benja Fallenstein, and Stuart Armstrong shows that if there exists an infinite series of theories T_i , extending T_{i-1} where each T_i proves the existence of T_{i+1} , then all the T_i must have only countable models. We call this the Procrastination Theorem. For reasons which will become apparent.

1 Review: Hierarchies of trust.

[This section primarily reviews material found in the paper "Eating Agents for Self-Scheduling AI, and the Liar's paradox"]

Notation:

T is an axiomatic system in classical logic whose consequences are recursively enumerable.

$T \vdash \phi$ is a syntactic consequence of the theory T .

$T \models \phi$ is semantic, i.e. ϕ is true in every model of T .

$\text{Pr}(T) = \{ \phi \mid T \vdash \phi \}$ is the Gödel number of a proof from the axioms of T whose conclusion is the theorem ϕ .

$\text{Pr}(\phi) = \{ p \mid \text{Pr}(T) \vdash \phi \}$ states that whenever T asserts ϕ , the order Peano arithmetic (PA) proves that there exists a T -proof of ϕ . (This follows from T being recursively enumerable and is true even if T is a more powerful system than PA.)

1. Trust. We say that T trusts \mathcal{U} , iff, for every formula ϕ with zero or one variable free, T asserts the uniform induction principle:

$$\forall x (P(x) \rightarrow (Q(x) \rightarrow \phi(x))) \rightarrow \phi(x)$$

That for every formula ϕ , T proves that for every x , if \mathcal{U} proves $\phi(x)$ then $T \vdash \phi(x)$. We can also say that T trusts \mathcal{U} , or that T proves \mathcal{U} sound.

Trust is transitive: if T trusts \mathcal{U} and \mathcal{U} trusts \mathcal{V} , then T trusts \mathcal{V} .

<http://rationalwiki.org/wiki/TrustAgency.pdf>

Robust Cooperation in the Prisoner's Dilemma: Program Equilibrium via Provability Logic

Mihaly Barasz, Paul Christiano, Benja Fallenstein, Marcello Herreshoff, Patrick LaViolette, and Eliesser Yudkowsky

January 23, 2014

Abstract

We consider the one-shot Prisoner's Dilemma between algorithms with read access to one another's source code, and we use the modal logic of provability to build agents that can achieve mutual cooperation in a manner that is robust, in that cooperation does not require exact equality of the agent's source code, and uncomputable, meaning that such an agent never cooperates when its opponent defects. We construct a general framework for such "modal agents", and study their properties.

1 Introduction

Can cooperation in a one-shot Prisoner's Dilemma be justified between rational agents? Hasegawa [18] argued in the 1960s that two agents with mutual knowledge of each other's rationality should be able to mutually cooperate. Howard [16] explains the argument thus:

Noncooperative arguments have been made in favor of playing C even in a single play of the PD. The one that interests us relies heavily on the usual assumption that both players are completely rational and know everything there is to know about the situation. (So for instance, Row knows that Column is rational, and Column knows that he knows it, and so on.) It can then be argued by Row that Column is an individual, very similar to himself, and in the same situation as himself. Hence whatever he eventually decides to do, Column will necessarily do the same (just as two good students given the same exam to calculate will necessarily arrive at the same answer). Hence if Row chooses D, so will Column, and each will get 1. However if Row chooses C, so will Column, and each will then get 3. Hence Row should choose C.

Hasegawa [18] described this line of reasoning as "superrationality", and held that knowledge of similar cognitive aptitudes should be enough to establish it, though the latter contention is (to say the least) controversial within decision theory. However, one may consider a stronger

Problems of self-reference in self-improving space-time embedded intelligence

Benja Fallenstein and Nate Soares

Machine Intelligence Research Institute
2000 Addison St. #300, Berkeley, CA 94704, USA
benja.fallenstein@miruscience.org

Abstract. By considering agents to be a part of their environment, Owen and Ring's previous embedded intelligence [16] is a better fit to the real world than the traditional agent framework. However, a self-modifying AGI that sees future versions of itself as an ordinary part of the environment may run into problems of self-reference. We show that in one particular model based on formal logic, some approaches either lead to inconsistency (meaning that unless an agent is part of an important sub-frame (the preservation period), or fail to allow the agent to justify even already self-revision (the Liar's paradox)). We argue that these problems have relevance beyond our particular formalism, and discuss partial solutions.

1 Introduction

Most formal models of artificial general intelligence (such as Hutter's AIXI [5] and the related formal notions of intelligence [6]) are based on the traditional agent framework, in which the agent interacts with an environment, but is not part of this environment. As Owen and Ring [16] point out, this is reminiscent of Cartesian dualism, the idea that the human mind is a non-physical substance external to the body [1]. A real-world AGI, on the other hand, will be part of the physical universe, and will need to deal with the possibility that external forces might observe or interfere with its internal operations.

The traditional separation of the agent from its environment seems even less attractive when one considers L.J. Good's idea that once AGI is sufficiently advanced, it may become better than any human at the task of making itself even smarter, leading to an "intelligence explosion" and leaving human intelligence far behind [7]. It seems plausible that an AGI undergoing an intelligence explosion may eventually want to adopt an architecture radically different from its initial one, such as one distributed over many different computers, where no single entity fulfills the agent's role from the traditional framework [8]. A formal model based on that framework cannot capture this.

How should one reason about such an agent? Owen and Ring [16] have proposed a formal model of open-time embedded intelligence to deal with this complexity. Their model consists of a set \mathcal{I} of policies, describing the state of the agent at a given point in time, an environment $\mathcal{E}(\tau_{i-1} | \tau_i)$, giving the

A Comparison of Decision Algorithms on Newcomblike Problems

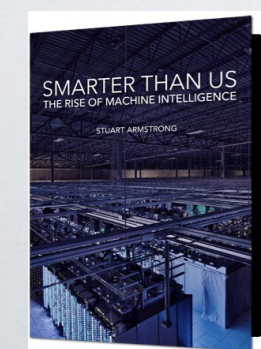
Alex Abriz
MIRI Research Fellow

Abstract

When formulated using Bayesian networks, two standard decision algorithms (Evidential Decision Theory and Causal Decision Theory) can be shown to fail systematically when faced with aspects of the prisoner's dilemma and so-called "Newcomblike" problems. We describe a new form of decision algorithm, called Tarskian Decision Theory, which consistently wins on these problems.

Alex Abriz, 2013. A Comparison of Decision Algorithms on Newcomblike Problems. Machine Intelligence Research Institute, Berkeley, CA.

Papers & reports



Scott Aaronson on Philosophical Progress

December 13, 2013 | Luke Muehlhauser | Comments

Scott Aaronson is an Associate Professor of Electrical Engineering and Computer Science at MIT. Before that, he did a PhD in computer science at UC Berkeley, as well as positions at the Institute for Advanced Study, Princeton, and the University of Waterloo. His research focuses on the qualitative and computational complexity and physics. Aaronson is known for his laid-back, as well as for founding the Complexity Zoo, an online encyclopedia of complexity classes. He's also written about quantum computing for Scientific American and the New York Times. His first book, Quantum Computing Since Democritus, was published this year by Cambridge University Press. He's received the Alan T. Waterman Award, the PECOS Award, and MIT's Junior Award for Excellence in Teaching.

Luke Muehlhauser: Though you're best known for your work in theoretical computer science, you've also produced some pretty interesting philosophical work, e.g. in Quantum Computing Since Democritus, "Why Philosophers Should Care About Computational Complexity," and "The Quest for the Quantum Turing Machine." You also taught a MIT 6.035 MIT class on Philosophy and Theoretical Computer Science.

Why are you so interested in philosophy? And what is the social value of philosophy, from your perspective?

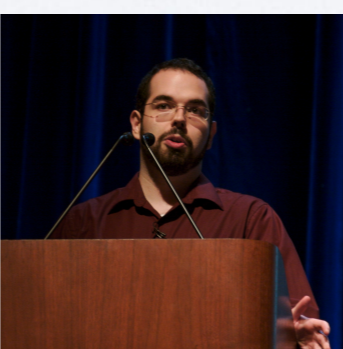
Scott Aaronson: I've always been offensively drawn to the biggest, most general questions that I can find myself to ask. You know, the old adage that if you could be a computer scientist or a philosopher, you should be a philosopher. I'm not sure if that's a good idea, but I do think it's important to have a good understanding of the world and the human condition. I'm not sure if that's a good idea, but I do think it's important to have a good understanding of the world and the human condition.

Can We Really Upload Johnny Depp's Brain?

A look at the science of Transcendental

By Luke Muehlhauser and Benja Fallenstein

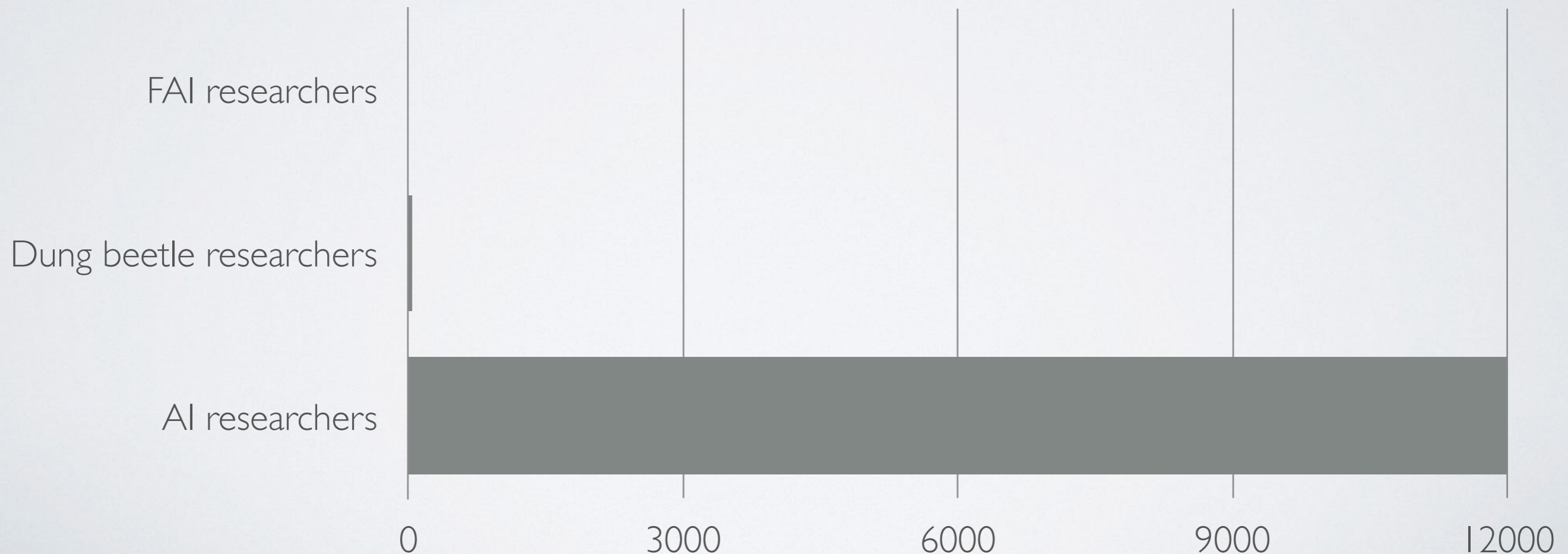
When Mike Pearson's Transcendental was released in April 2013, millions of people were captivated by the idea of uploading a human brain to a computer. The idea of uploading a human brain to a computer is a fascinating one, but it's also a very difficult one. The idea of uploading a human brain to a computer is a fascinating one, but it's also a very difficult one.



Some other stuff

Marginal dollars are mostly spent on:
finding/creating new FAI researchers.

Because:





MIRI

MACHINE INTELLIGENCE
RESEARCH INSTITUTE

For more information,
visit intelligence.org
or come talk to me.