Self-Improving Artificial Intelligence and the Methods of Rationality

Eliezer Yudkowsky Research Fellow Singularity Institute for Artificial Intelligence

yudkowsky.net

Yeshiva University March 2011 Linda is 31 years old, single, outspoken, and very bright. She majored in philosophy. As a student, she was deeply concerned with issues of discrimination and social justice, and also participated in anti-nuclear demonstrations.

Rank the following statements from most probable to least probable:

- Linda is a teacher in an elementary school
- Linda works in a bookstore and takes Yoga classes
- Linda is active in the feminist movement
- Linda is a psychiatric social worker
- Linda is a member of the League of Women Voters
- Linda is a bank teller
- Linda is an insurance salesperson
- Linda is a bank teller and is active in the feminist movement

Yeshiva University March 2011

Linda is 31 years old, single, outspoken, and very bright. She majored in philosophy. As a student, she was deeply concerned with issues of discrimination and social justice, and also participated in anti-nuclear demonstrations.

Rank the following statements from most probable to least probable:

- Linda is a teacher in an elementary school
- Linda works in a bookstore and takes Yoga classes
- Linda is active in the feminist movement
- Linda is a psychiatric social worker
- Linda is a member of the League of Women Voters
- Linda is a bank teller (A)
- Linda is an insurance salesperson
- Linda is a bank teller and is active in the feminist movement (A & B)

89% of subjects thought Linda was more likely to be a feminist bank teller than a bank teller Conjunction fallacy: P(A&B) must be \leq than P(A)

Yeshiva University March 2011

Linda is 31 years old, single, outspoken, and very bright. She majored in philosophy. As a student, she was deeply concerned with issues of discrimination and social justice, and also participated in anti-nuclear demonstrations.

- Linda is a bank teller (A)
- Linda is a bank teller and is active in the feminist movement (A & B)



Yeshiva University March 2011

WHICH IS THE BEST BET?

- Die with 4 green faces and 2 red faces
- Die rolled 20 times, series of Gs and Rs recorded
- If your chosen sequence appears, you win \$25
- Which of these three sequences would you prefer to bet on?



Amos Tversky and Daniel Kahneman (1983), "Extensional versus Intuitive Reasoning". Psychological Review, 90, 293-315

Yeshiva University March 2011

WHICH IS THE BEST BET?

- Die with 4 green faces and 2 red faces
- Die rolled 20 times, series of Gs and Rs recorded
- If your chosen sequence appears, you win \$25
- Which of these three sequences would you prefer to bet on?



- 65% chose Sequence 2.
- 125 undergraduates played this gamble with real payoffs.

Amos Tversky and Daniel Kahneman (1983), "Extensional versus Intuitive Reasoning". Psychological Review, 90, 293-315

Yeshiva University March 2011

CONJUNCTION FALLACY: P(A) < P(A&B)

- Die with 4 green faces and 2 red faces
- Die rolled 20 times, series of Gs and Rs recorded
- If your chosen sequence appears, you win \$25
- Which of these three sequences would you prefer to bet on?

 but Sequence 1 dominates Sequence 2; it appears within 2, so it appears anywhere 2 appears.



Yeshiva University March 2011 Eliezer Yudkowsky lesswrong.com

GRGRRR

RGRRR

CONJUNCTION FALLACY: P(A) < P(A&B)

• Die with 4 green faces and 2 red faces: GGRGGR



Judgment of probability versus We want:

Which is most probable?



Judgment of representativeness We ask:

Which is most representative?

GRGRRR

Amos Tversky and Daniel Kahneman (1983), "Extensional versus Intuitive Reasoning". Psychological Review, 90, 293-315

Yeshiva University March 2011

2nd INTERNATIONAL CONGRESS ON FORECASTING

Two independent groups of professional analysts asked: ٠ Please rate the probability that the following event will occur in 1983.

Version 1	Version 2
A complete suspension of diplomatic relations between the USA and the Soviet Union, sometime in 1983.	A Russian invasion of Poland, and a complete suspension of diplomatic relations between the USA and the Soviet Union, sometime in 1983.

The group asked to rate Version 2 responded with significantly higher ۲ probabilities.

Amos Tversky and Daniel Kahneman (1983), "Extensional versus Intuitive Reasoning". Psychological Review, 90, 293-315

Yeshiva University March 2011

SECOND INTERNATIONAL CONGRESS ON FORECASTING

- Two independent groups of professional analysts were asked to: •
- Please rate the probability that the following event will occur in 1983. ۲ Version 1 Version 2

A complete suspension of diplomatic relations between the USA and the Soviet Union, sometime in 1983.

A Russian invasion of Poland, and a complete suspension of diplomatic relations hetween the USA and the Soviet Union, sometime in 1983.

The group asked to rate Version 2 responded with significantly higher ٠ probabilities.

Amos Tversky and Daniel Kahneman (1983), "Extensional versus Intuitive Reasoning". Psychological Review, 90, 293-315

Yeshiva University March 2011

CONJUNCTION FALLACY: P(A) < P(A&B)

- Two independent groups of professional analysts were asked to:
- Please rate the probability that the following event will occur in 1983. Version 1
 Version 2

A complete suspension of diplomatic relations between the USA and the Soviet Union, sometime in 1983. A Russian invasion of Poland, and a complete suspension of diplomatic relations between the USA and the Soviet Union, sometime in 1983.

- The group asked to rate Version 2 responded with significantly higher probabilities.
- Adding detail to a story can make the story sound more plausible, even though the story necessarily becomes less probable.

Yeshiva University March 2011

INSENSITIVITY TO PREDICTABILITY

- Subjects presented with description of teacher during lesson
- Some subjects asked to evaluate quality of lesson, in percentile score
- Other subjects asked to predict percentile standing of teacher
 5 years later

Kahneman, D., Tversky, A. 1973. On the psychology of prediction. Psychological Review 80: 237-51.

Yeshiva University March 2011

INSENSITIVITY TO PREDICTABILITY

- Subjects presented with description of teacher during lesson
- Some subjects asked to evaluate quality of lesson, in percentile score
- Other subjects asked to predict percentile standing of teacher
 5 years later
- Judgments identical equally extreme!

Kahneman, D., Tversky, A. 1973. On the psychology of prediction. Psychological Review 80: 237-51.

Yeshiva University March 2011

What do you think you know, and how do you think you know it?

Yeshiva University March 2011

One of these technologies is not like the others...

Artificial Intelligence	Interplanetary travel
Cancer cure	Nano-manufacturing

Yeshiva University March 2011

The power of intelligence:

- Fire
- Language
- Nuclear weapons
- Skyscrapers
- Spaceships
- Money
- Science

"Book smarts" vs. cognition:

"Book smarts" evokes images of:

- Calculus
- Chess
- Good recall of facts

Other stuff that happens *in the brain*:

- Social persuasion
- Enthusiasm
- Reading faces
- Rationality
- Strategic cleverness

What do we think we can guess?

• Technologies which *impact upon cognition* will end up mattering most, because intelligence is more powerful and significant than cool devices.

What do you think you know, and how do you think you know it?

Yeshiva University March 2011

Artificial Addition



Yeshiva University March 2011

Views on Artificial General Addition

- Framing problem what 'twenty-one plus' equals depends on whether it's 'plus three' or 'plus four'. Need to program huge network of arithmetical facts to cover common-sense truths.
- Need Artificial Arithmetician that can understand natural language, so instead of being told that twenty-one plus sixteen equals thirty-seven, it can obtain knowledge by reading Web.
- Need to develop General Arithmetician the same way Nature did evolution.
- Top-down approaches have failed to produce arithmetic. Need bottom-up approach to make arithmetic *emerge*. Accept unpredictability of complex systems.

Views on Artificial General Addition

- Neural networks just like the human brain! Can be trained without understanding how they work! NNs will do arithmetic without us, their creators, ever understanding how they add.
- Need calculators as powerful as a human brain. Moore's Law predicts availability on April 27, 2031 between 4 and 4:30AM.
- Must simulate neural circuitry humans use for addition.
- Godel's Theorem shows no formal system can ever capture properties of arithmetic. Classical physics is formalizable. Hence AGA must exploit quantum gravity.
- If human arithmetic were simple enough to program, we couldn't count high enough to build computers.
- Haven't you heard of John Searle's Chinese Calculator Experiment?
- Will never know nature of arithmetic; problem is just too hard.

Artificial Addition: The Moral

• *Moral 1:* When you're missing a basic insight, you *must* find that insight. Workarounds may *sound* clever, but aren't. Can't talk sensibly about solutions until no-longer-confused.



Yeshiva University March 2011

Artificial Addition: The Moral

- *Moral 1:* When you're missing a basic insight, you *must* find that insight. Workarounds may *sound* clever, but aren't. Can't talk sensibly about solutions until no-longer-confused.
- *Moral 2:* Patching an infinite number of surface cases means you didn't understand the underlying generator.



Yeshiva University March 2011

"The influence of animal or vegetable life on matter is infinitely beyond the range of any scientific inquiry hitherto entered on. Its power of directing the motions of moving particles, in the demonstrated daily miracle of our human free-will, and in the growth of generation after generation of plants from a single seed, are infinitely different from any possible result of the fortuitous concurrence of atoms... Modern biologists were coming once more to the acceptance of something and that was a vital principle."

- Lord Kelvin

Mind Projection Fallacy:

If I am ignorant about a phenomenon, this is a fact about my state of mind, not a fact about the phenomenon.

(Jaynes, E.T. 2003. *Probability Theory: The Logic of Science*. Cambridge: Cambridge University Press.)

Yeshiva University March 2011

What do we think we can guess?

- Technologies which *impact upon cognition* will end up mattering most.
- Today's difficulties in constructing AI are not because intelligence is *inherently mysterious*, but because we're currently still ignorant of some basic insights.

What do you think you know, and how do you think you know it?

Yeshiva University March 2011



Biology not near top of scale

- Lightspeed >10⁶ times faster than axons, dendrites.
- Synaptic spike dissipates >10⁶ minimum heat (though transistors do worse)
- Transistor clock speed >>10⁶ times faster than neuron spiking frequency

1,000,000-fold speedup physically possible: (1 year = 31,556,926 seconds) \rightarrow 31 seconds



									2	00	97	7								
January			February							March										
	1	2	3	4	5	6					1	2	3					1	2	3
7	8	9	10	11	12	13	4	5	6	7	8	9	10	- 4	5	6	7	8	9	10
14	15	16	17	18	19	20	11	12	13	14	15	16	17	11	12	13	14	15	16	17
21	22	23	24	25	26	27	18	19	20	21	22	23	24	18	19	20	21	22	23	24
28	29	30	31				25	26	27	28				25	26	27	28	29	30	31
April					May					June										
1	2	3	4	5	6	7			1	2	3	4	5						1	2
8	9	10	11	12	13	14	б	7	8	9	10	11	12	3	4	5	б	7	8	9
15	16	17	18	19	20	21	13	14	15	16	17	18	19	10	11	12	13	14	15	16
22	23	24	25	26	27	28	20	21	22	23	24	25	26	17	18	19	20	21	22	23
29	30						27	28	29	30	31			24	25	26	27	28	29	30
July							A	ugu	ıst			September								
1	2	3	4	S	6	7				1	2	3	4							1
8	9	10	11	12	13	14	5	6	7	8	9	10	11	2	3	4	5	6	7	8
15	16	17	18	19	20	21	12	13	14	15	16	17	18	.9	10	11	12	13	14	15
22	23	24	25	26	27	28	19	20	21	22	23	24	25	16	17	18	19	20	21	22
29	30	31					2.6	27	28	29	30	31		23	24	25	26	27	28	29
														30						
October			November						December											
	1	2	3	4	5	6					1	2	3							1
7	8	9	10	11	12	13	4	5	6	7	8	9	10	2	3	- 4	5	6	7	8
14	15	16	17	18	19	20	11	12	13	14	15	16	17	9	10	11	12	13	14	15
21	22	23	24	25	26	27	18	19	20	21	22	23	24	16	17	18	19	20	21	22
28	29	30	31				25	26	27	28	29	30		23	24	25	26	27	28	29
														30	31					

Eliezer Yudkowsky lesswrong.com

Yeshiva University March 2011



Yeshiva University March 2011



Yeshiva University March 2011

AI advantages

- Total read/write access to own state
- Absorb more hardware (possibly orders of magnitude more!)
- Understandable code
- Modular design
- Clean internal environment

What do we think we can guess?

- Technologies which *impact upon cognition* will end up mattering most.
- Today's difficulties in constructing AI are due to ignorance of key insights, not inherent mysteriousness.
- It is possible for minds to exist that are *far* more powerful than human. We can see that human *hardware* is far short of the physical limits, and there's no reason to expect the software to be optimal either.

What do you think you know, and how do you think you know it?

Yeshiva University March 2011










"Intelligence explosion:"

- Invented by I. J. Good in 1965
- Hypothesis: More intelligence -> more creativity on task of making yourself even smarter
- Prediction: Positive feedback cycle rapidly creates *superintelligence*.

(Good, I. J. 1965. Speculations Concerning the First Ultraintelligent Machine. Pp. 31-88 in *Advances in Computers*, **6**, F. L. Alt and M. Rubinoff, eds. New York: Academic Press.)

Yeshiva University March 2011

How fast an intelligence explosion?

- From chimpanzees to humans 4x brain size, 6x frontal cortex, no apparent obstacles to constant evolutionary pressure.
- From flint handaxes to Moon rockets constant human brains don't slow down over course of technological history as problems get harder.
- Thus an *increasingly* powerful optimization process should undergo *explosive* self-improvement.

Intelligence explosion hypothesis does *not* imply, nor require:

• More change occurred from 1970 to 2000 than from 1940 to 1970.

Intelligence explosion hypothesis does *not* imply, nor require:

- More change occurred from 1970 to 2000 than from 1940 to 1970.
- Technological progress follows a predictable curve.

What do we think we can guess?

- Technologies which *impact upon cognition* will end up mattering most.
- Today's difficulties in constructing AI are due to ignorance of key insights, not inherent mysteriousness.
- It is possible for minds to exist that are *far* more powerful than human.
- An AI over some threshold level of intelligence will *recursively* improve itself and explode into superintelligence.

What do you think you know, and how do you think you know it?

Yeshiva University March 2011

In Every Known Human Culture:

- tool making
- weapons
- grammar
- tickling
- meal times

- mediation of conflicts
- dance, singing
- personal names
- promises
- mourning

(Donald E. Brown, 1991. Human universals. New York: McGraw-Hill.)

Yeshiva University March 2011

A complex adaptation must be universal within a species.

(John Tooby and Leda Cosmides, 1992. *The Psychological Foundations of Culture.* In *The Adapted Mind,* eds. Barkow, Cosmides, and Tooby.)

Yeshiva University March 2011

A complex adaptation must be universal within a species.

If: 6 necessary genes Each at 10% frequency in population Then: 1 in 1,000,000 have complete adaptation

(John Tooby and Leda Cosmides, 1992. *The Psychological Foundations of Culture.* In *The Adapted Mind,* eds. Barkow, Cosmides, and Tooby.)

Yeshiva University March 2011

Incremental evolution of complexity:

- 1. [A] A is advantageous by itself.
 - $[A \leftarrow B] \qquad B depends on A.$
- 3. $[A'\leftrightarrow B]$ A' rep
- A' replaces A, depends on B.
 - C depends on A' and B

(John Tooby and Leda Cosmides, 1992. *The Psychological Foundations of Culture.* In *The Adapted Mind,* eds. Barkow, Cosmides, and Tooby.)

Yeshiva University March 2011

...

[A'B←C]

2.

4.



The Psychic Unity of Humankind

Complex adaptations must be universal in a species – including *cognitive* machinery in *Homo sapiens!*

(John Tooby and Leda Cosmides, 1992. *The Psychological Foundations of Culture.* In *The Adapted Mind,* eds. Barkow, Cosmides, and Tooby.)

Yeshiva University March 2011





The Great Failure of Imagination: *Anthropomorphism*



Eliezer Yudkowsky lesswrong.com

Yeshiva University March 2011



Fallacy of the Giant Cheesecake

- Major premise: A superintelligence could create a mile-high cheesecake.
- Minor premise: Someone will create a recursively selfimproving AI.
- Conclusion: The future will be full of giant cheesecakes.

Power does not imply motive.

Fallacy of the Giant Cheesecake

- Major premise: A superintelligence could create a mile-high cheese cake.
- Minor premise: Someone will create a recurs vely selfimproving AI.
- Conclusion: The future will be full of giant checkets.

Power does not imply motive.

Yeshiva University March 2011

What do we think we can guess?

- Technologies which *impact upon cognition* will end up mattering most.
- Today's difficulties in constructing AI are due to ignorance of key insights, not inherent mysteriousness.
- It is possible for minds to exist that are *far* more powerful than human.
- An AI over some threshold level of intelligence will recursively self-improve and explode into superintelligence.
- Features of "human nature" that we take for granted are just one of a vast number of possibilities, and not all possible agents in that space are friendly.

Yeshiva University March 2011

What do you think you know, and how do you think you know it?

Yeshiva University March 2011

What do you think you know, and how do you think you know it?

To predict a smarter mind, you'd have to be that smart yourself?

Yeshiva University March 2011







Key insight: Predict & value consequences



Yeshiva University March 2011

• Deep Blue's programmers couldn't predict its exact moves

- Deep Blue's programmers couldn't predict its exact moves
- So why not use random move generator?

- Deep Blue's programmers couldn't predict its exact moves
- So why not use random move generator?
- Unpredictability of superior intelligence ≠ unpredictability of coinflips

- Deep Blue's programmers couldn't predict its exact moves
- So why not use random move generator?
- Unpredictability of superior intelligence ≠ unpredictability of coinflips
- Deep Blue's "unpredictable" move, predictably has consequence of winning game

- Deep Blue's programmers couldn't predict its exact moves
- So why not use random move generator?
- Unpredictability of superior intelligence ≠ unpredictability of coinflips
- Deep Blue's "unpredictable" move, predictably has consequence of winning game
- Inspection of code can't *prove* consequences in an uncertain real-world environment, but it can establish with near-certainty that an agent is *trying to find the most-probably-good action*.

Yeshiva University March 2011

Stability of goals in self-modifying agents:

- Gandhi doesn't want to kill people.
- Offer Gandhi a pill that will make him enjoy killing people?

Stability of goals in self-modifying agents:

- Gandhi doesn't want to kill people.
- Offer Gandhi a pill that will make him enjoy killing people?
- If Gandhi correctly assesses this is what the pill does, Gandhi will refuse the pill, because the current Gandhi doesn't want the consequence of people being murdered.

Stability of goals in self-modifying agents:

- Gandhi doesn't want to kill people.
- Offer Gandhi a pill that will make him enjoy killing people?
- If Gandhi correctly assesses this is what the pill does, Gandhi will refuse the pill, because the current Gandhi doesn't want the consequence of people being murdered.
- Argues for folk theorem that *in general*, rational agents will *preserve their utility functions during self-optimization*.

(Steve Omohundro, 2008. "The Basic AI Drives." In *Proceedings of the First AGI Conference*, eds. Pei Wang and Stan Franklin.)

What do we think we can guess?

- Technologies which *impact upon cognition* will end up mattering most.
- Today's difficulties in constructing AI are due to ignorance of key insights, not inherent mysteriousness.
- It is possible for minds to exist that are *far* more powerful than human.
- An AI over some threshold level of intelligence will recursively self-improve and explode into superintelligence.
- Not all possible agents in mind design space are friendly.
- *If we can obtain certain new insights,* it should be possible to construct a benevolent self-improving Artificial Intelligence.

What do we think we can guess?

- Technologies which *impact upon cognition* will end up mattering most.
- Today's difficulties in constructing AI are due to ignorance of key insights, not inherent mysteriousness.
- It is possible for minds to exist that are *far* more powerful than human.
- An AI over some threshold level of intelligence will recursively self-improve and explode into superintelligence.
- Not all possible agents in mind design space are friendly.
- *If we can obtain certain new insights,* it should be possible to construct a benevolent self-improving Artificial Intelligence.
- And live happily ever after.

Yeshiva University March 2011
More Info: http://LessWrong.com/ http://Yudkowsky.net/ http://SingInst.org/

Eliezer Yudkowsky Research Fellow Singularity Institute for Artificial Intelligence

yudkowsky.net

Yeshiva University March 2011