

Best of LessWrong: September 2018

1. [Anti-social Punishment](#)
2. [Realism about rationality](#)
3. [Comment on decision theory](#)
4. [The Tails Coming Apart As Metaphor For Life](#)
5. [Impact Measure Desiderata](#)
6. [Disagreement with Paul: alignment induction](#)
7. [Towards a New Impact Measure](#)
8. [On Robin Hanson's Board Game](#)
9. [nostalgebraist - bayes: a kinda-sorta masterpost](#)
10. [Open AI co-founder on AGI](#)
11. [I am the very model of a self-recursive modeler](#)
12. [Alignment Newsletter #25](#)
13. [Against the barbell strategy](#)
14. [Tradition is Smarter Than You Are](#)
15. [Counterfactuals and reflective oracles](#)
16. [A Process for Dealing with Motivated Reasoning](#)
17. [Deep learning - deeper flaws?](#)
18. [Alignment Newsletter #23](#)
19. [Book review: Why we sleep](#)
20. [In Logical Time, All Games are Iterated Games](#)
21. [Thoughts on tackling blindspots](#)
22. [Psycho-cybernetics: experimental notes](#)
23. [The Scent of Bad Psychology](#)
24. [Track-Back Meditation](#)
25. [What To Do If Nuclear War Seems Imminent](#)
26. [Moral differences in mediocristan](#)
27. [Petrov corrigibility](#)
28. [Resurrection of the dead via multiverse-wide acausal cooperation](#)
29. [Are you in a Boltzmann simulation?](#)
30. [Your genome isn't private. Maybe it never was.](#)
31. [Advances in Baby Formula](#)
32. [Alignment Newsletter #24](#)
33. [A probabilistic off-switch that the agent is indifferent to](#)
34. [Asymptotic Decision Theory \(Improved Writeup\)](#)
35. [Birth order effect found in Nobel Laureates in Physics](#)
36. [Direct Primary Care](#)
37. [Hypothesis about how social stuff works and arises](#)
38. [New DeepMind AI Safety Research Blog](#)
39. [\(A → B\) → A](#)
40. [Good Citizenship Is Out Of Date](#)
41. [Criticism Scheduling and Privacy](#)
42. [\[Hammertime Final Exam\] Accommodate Yourself; Kindness Is An Epistemic Virtue; Privileging the Future](#)
43. [Beauty bias: "Lost in Math" by Sabine Hossenfelder](#)
44. [Reflective AIXI and Anthropic](#)
45. [Cooperative Oracles](#)
46. [August gwern.net links](#)
47. [How to use a microphone rationally during public speaking](#)
48. [An Ontology of Systemic Failures: Dragons, Bullshit Mountain, and the Cloud of Doom](#)
49. [Bridging syntax and semantics with Quine's Gavagai](#)
50. [Iterative Arguments: Alternative to Adversarial Collaboration](#)

Best of LessWrong: September 2018

1. [Anti-social Punishment](#)
2. [Realism about rationality](#)
3. [Comment on decision theory](#)
4. [The Tails Coming Apart As Metaphor For Life](#)
5. [Impact Measure Desiderata](#)
6. [Disagreement with Paul: alignment induction](#)
7. [Towards a New Impact Measure](#)
8. [On Robin Hanson's Board Game](#)
9. [nostalgebraist - bayes: a kinda-sorta masterpost](#)
10. [Open AI co-founder on AGI](#)
11. [I am the very model of a self-recursive modeler](#)
12. [Alignment Newsletter #25](#)
13. [Against the barbell strategy](#)
14. [Tradition is Smarter Than You Are](#)
15. [Counterfactuals and reflective oracles](#)
16. [A Process for Dealing with Motivated Reasoning](#)
17. [Deep learning - deeper flaws?](#)
18. [Alignment Newsletter #23](#)
19. [Book review: Why we sleep](#)
20. [In Logical Time, All Games are Iterated Games](#)
21. [Thoughts on tackling blindspots](#)
22. [Psycho-cybernetics: experimental notes](#)
23. [The Scent of Bad Psychology](#)
24. [Track-Back Meditation](#)
25. [What To Do If Nuclear War Seems Imminent](#)
26. [Moral differences in mediocristan](#)
27. [Petrov corrigibility](#)
28. [Resurrection of the dead via multiverse-wide acausal cooperation](#)
29. [Are you in a Boltzmann simulation?](#)
30. [Your genome isn't private. Maybe it never was.](#)
31. [Advances in Baby Formula](#)
32. [Alignment Newsletter #24](#)
33. [A probabilistic off-switch that the agent is indifferent to](#)
34. [Asymptotic Decision Theory \(Improved Writeup\)](#)
35. [Birth order effect found in Nobel Laureates in Physics](#)
36. [Direct Primary Care](#)
37. [Hypothesis about how social stuff works and arises](#)
38. [New DeepMind AI Safety Research Blog](#)
39. [\(A → B\) → A](#)
40. [Good Citizenship Is Out Of Date](#)
41. [Criticism Scheduling and Privacy](#)
42. [\[Hammertime Final Exam\] Accommodate Yourself; Kindness Is An Epistemic Virtue; Privileging the Future](#)
43. [Beauty bias: "Lost in Math" by Sabine Hossenfelder](#)
44. [Reflective AIXI and Anthropicics](#)
45. [Cooperative Oracles](#)
46. [August gwern.net links](#)
47. [How to use a microphone rationally during public speaking](#)

48. [An Ontology of Systemic Failures: Dragons, Bullshit Mountain, and the Cloud of Doom](#)
49. [Bridging syntax and semantics with Quine's Gavagai](#)
50. [Iterative Arguments: Alternative to Adversarial Collaboration](#)

Anti-social Punishment

This is a cross post from 250bpm.com.

Introduction

There's a trope among Slovak intellectual elite depicting an average Slovak as living in a village, sitting a local pub, drinking [Borovička](#), criticizing everyone and everything but not willing to lift a finger to improve things. Moreover, it is assumed that if you actually tried to make things better, said individual would throw dirt at you and place obstacles in your way.

I always assumed that this caricature was silly. It was partly because I have a soft spot for Slovak rural life but mainly because such behavior makes absolutely no sense from game-theoretical point of view. If a do-gooder is stupid enough to try to altruistically improve your life, why go into trouble of actively opposing them? Why not just sit safely hidden in the pub, drink some more Borovička and wait until they are done?

Well, it turns out that the things are far more complex then I thought.

Public goods game

Benedikt Herrmann, Christian Thöni and Simon Gächter did a study of how people from different societies deal with cooperation and punishment. You can find the paper [here](#) and supporting material [here](#).

The study is based on the "public goods" game. The game works as follows:

There are four players. Each player gets 20 tokens to start with. Every participant either keeps them or passes some of them into a common pool. After all the players are done with their moves, each of them, irrespective of how much they contributed, gets tokens equal to 40% of all the tokens in the common pool. The participants cannot communicate with each other and are unaware of each other's identities. The game is repeated, with the same players, 10 times in a row.

The earnings, obviously, depend not only on subject's move but also on the willingness of the other players to cooperate and put tokens into the common pool. But free riders get an advantage. They keep their original tokens but also get their share from the pool.

To get a feeling of the payoffs, let's have a look at the single-round earnings in the extreme case where each participant either puts all their tokens into the pool ("cooperator") or keeps all the tokens for themselves ("free-rider"):

No. of free-riders	Earning of a cooperator	Earnings of a free-rider
0	32	n/a
1	24	44
2	16	36
3	8	28
4	n/a	20

Public goods game with punishment

There's a variant of the "public goods game" where players are able to punish each other after each round of the game. The mechanism is simple. When the round ends the participants are informed about how much each of them has put into the common pool. Then they decide whether to spend some of their tokens to administer punishment. For each token spent on punishment you can subtract 3 tokens from the earnings of an opponent. The players know that they've been punished but they are not informed about who exactly has punished them.

Participant pools

The researchers were interested in comparing the results of the game among different societies:

Our research strategy was to conduct the experiments with comparable social groups from complex developed societies with the widest possible range of cultural and economic backgrounds to maximize chances of observing cross-societal differences in punishment and cooperation. The societies represented in our participant pools diverge strongly according to several widely used criteria developed by social scientists in order to characterize societies. This variation, covering a large range of the worldwide available values of the respective criteria, provides us with a novel test for seeing whether societal differences between complex societies have any impact on experimentally observable disparities in cooperation and punishment behavior. ... To minimize sociodemographic variability, we conducted all experiments with university undergraduates who were similar in age, shared an (upper) middle class background, and usually did not know each other.

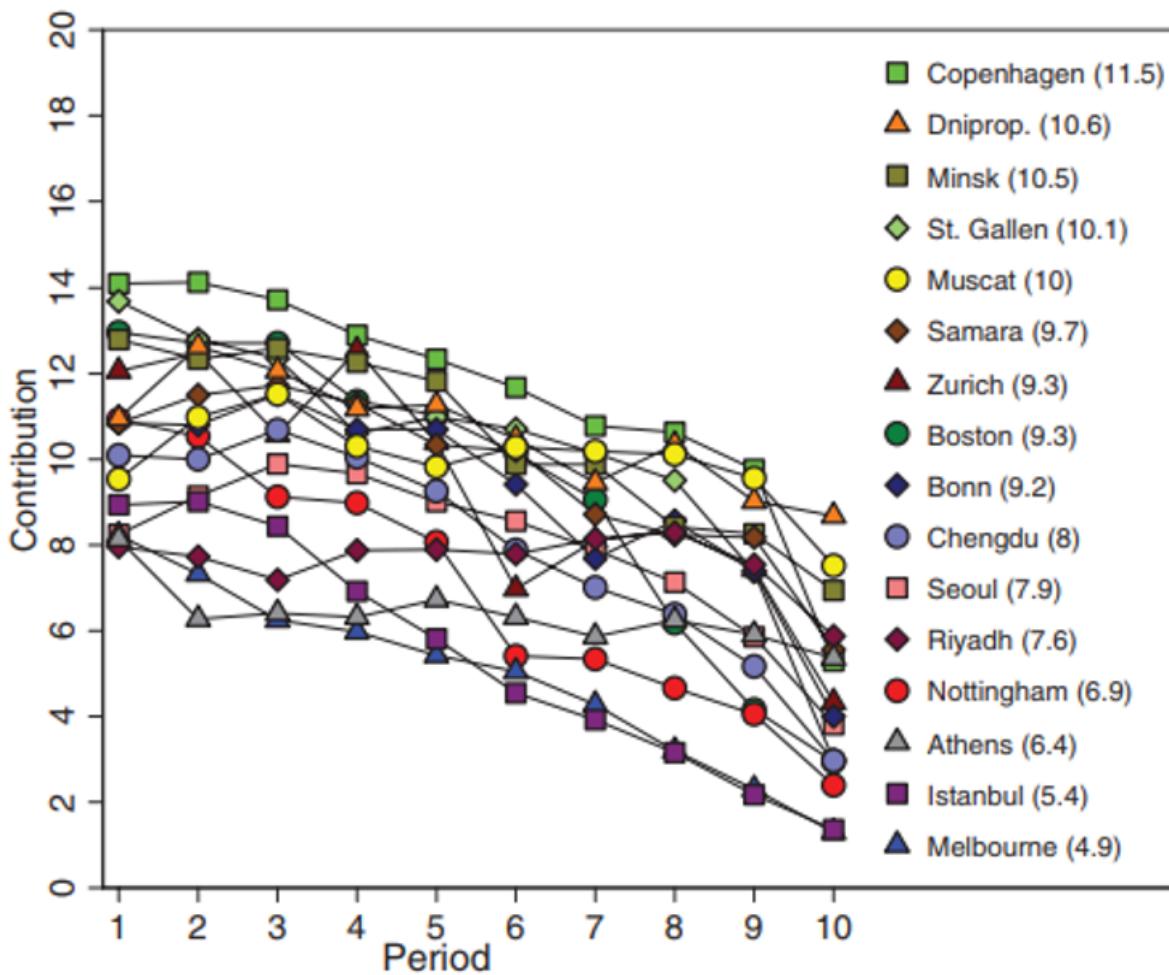
Specifically, the experiment was conducted at the following places:

- Athens, Greece
- Bonn, Germany
- Boston, US
- Chengdu, China
- Copenhagen, Denmark
- Dnipropetrovsk, Ukraine
- Istanbul, Turkey
- Melbourne, Australia
- Minsk, Belarus
- Muscat, Oman

- Nottingham, UK
- Riyadh, Saudi Arabia
- Samara, Russia
- Seoul, Sourg Korea
- St. Gallen, Switzerland
- Zürich, Switzerland

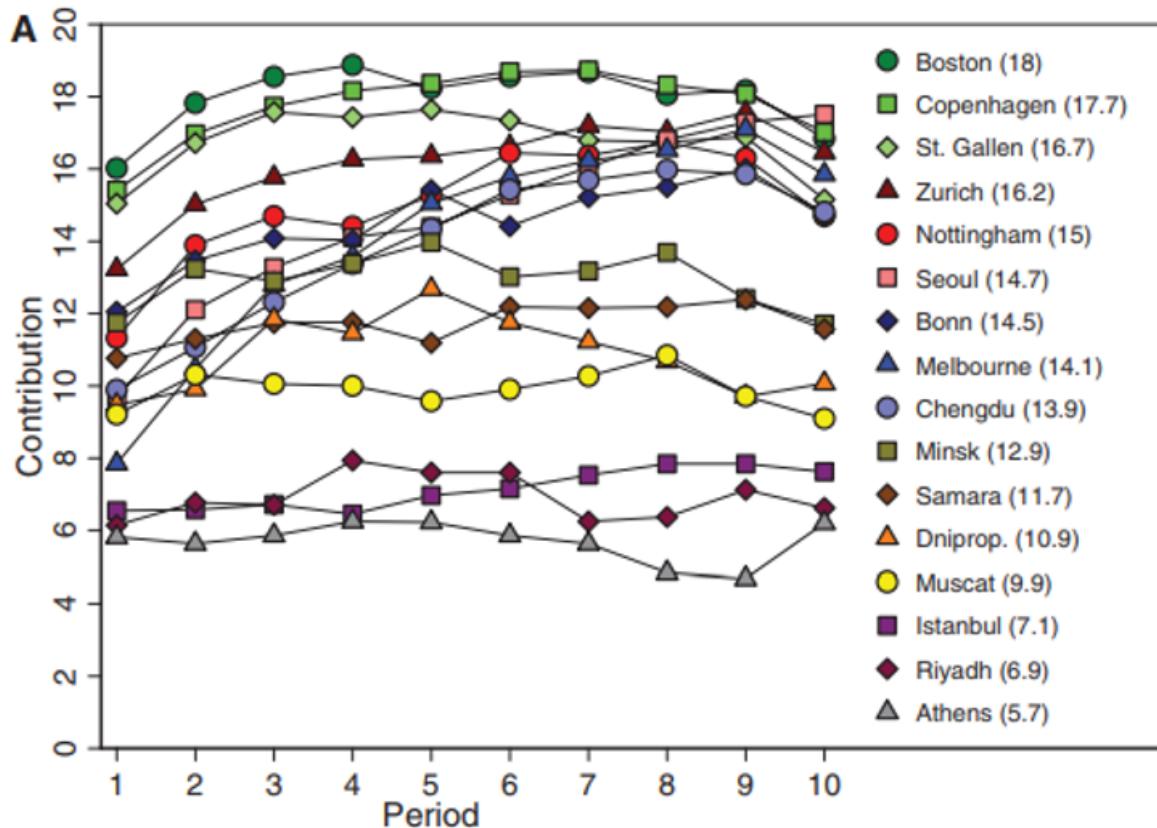
Results: Public goods game without punishment

The results from this experiment are exactly as you would expect. The cooperators found out that there was no way to prevent free-riding and the amount of resources they've put into the common pool steadily decreased. This result replicated across different participant pools. Particular pool may have started with a high or low cooperative behavior, but as the time went on the cooperation always decreased.



Results: Public goods game with punishment

Ability to punish free-riders increased the cooperative behavior in most participant pools. Free-riders learned that free-riding doesn't pay off and started contributing to the common pool.



However, introduction of punishment had no effect in some of the pools. The contributions stayed more or less the same throughout the experiment in Minsk, Samara, Dnipropetrovsk, Muscat, Istanbul, Riyadh and Athens.

Now that's an interesting result. What's going on there? Are members of some societies resistant to punishment or what?

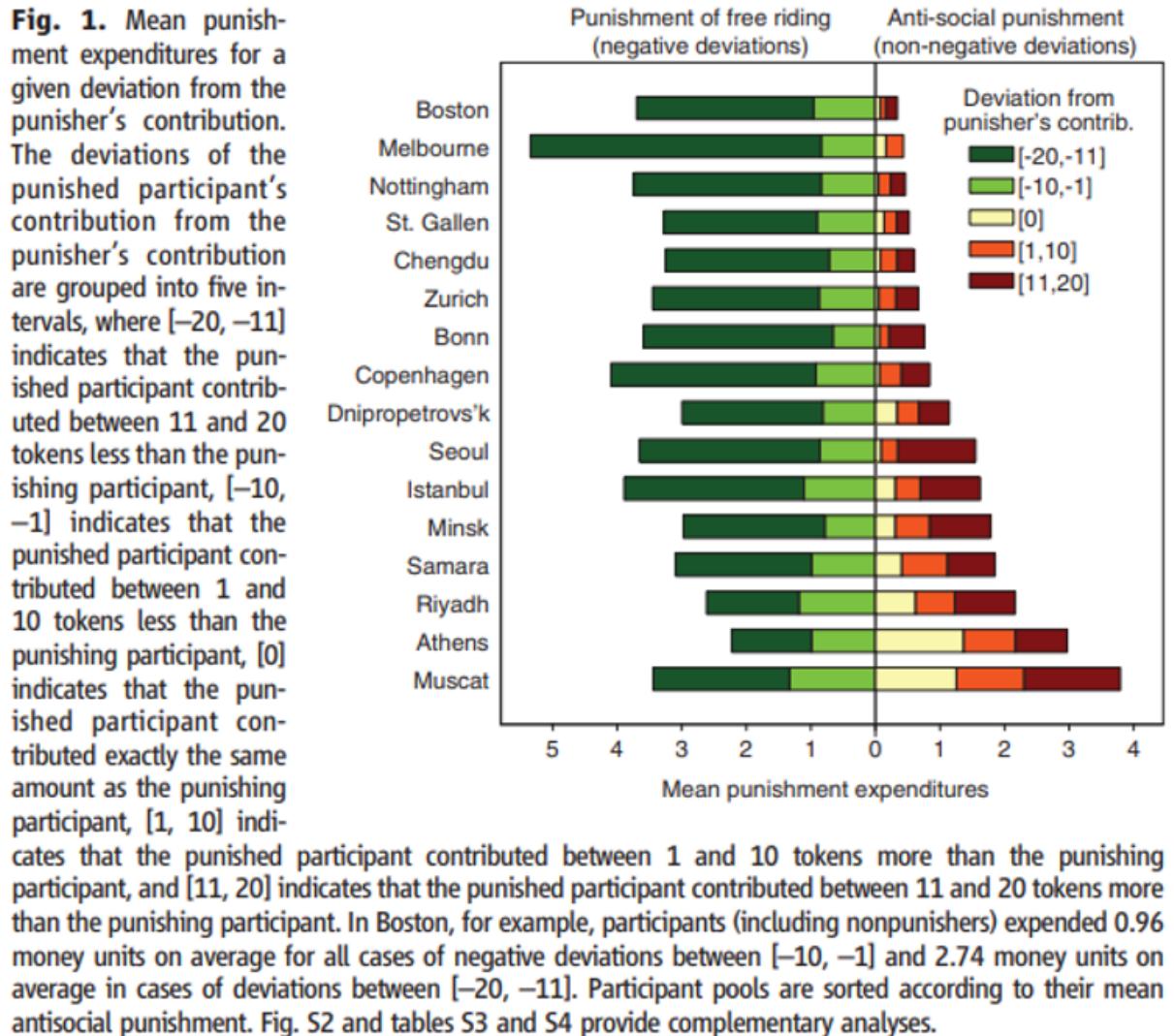
The reality turns up to be even more interesting than one would expect.

Anti-social punishment

Herrmann, Thöni and Gächter found out that participants in some societies were engaging in what they've called "anti-social punishment". They were punishing cooperators!

In fact, they were punishing cooperators so much that the cooperation-enhancing effect of pro-social punishment was entirely canceled.

To make it even more confusing, the anti-social punishment, unlike the pro-social punishment which had roughly similar level in all the participant pools, differed widely among the pools. While it was almost non-existent in the West, it was common in Eastern Europe, in Middle East and in Greece.



The authors then try to find out which aspects of the society are correlated with the high anti-social punishment rate:

With respect to antisocial punishment, we found that both norms of civic cooperation and rule of law are significantly negatively correlated with punishment (at $P < 0.05$). In other words, antisocial punishment is harsher in participant pools from societies with weak norms of civic cooperation and a weak rule of law. Additional analyses show that antisocial punishment also varies highly significantly with a variety of indicators developed by social scientists in order to characterize societies. Thus, the extent of antisocial punishment is most likely affected by the wider societal background.

Why on Earth?

I wouldn't have much to add to the fascinating results above, I am not a sociologist after all, but I happen to come from a country that is probably affected by this problem. Slovakia hasn't participated in the study, however, it's a former Ostblock country and as such it is very likely to have results similar to Ukraine, Russia or Belarus. Moreover,

local folk wisdom, as already mentioned, has it that the phenomenon does really exist. Therefore, having all the relevant context and all the intricacies of the local culture in my head, I should be able to come up with an psychologically plausible explanation of why it would make sense to punish cooperators. It was hard to empathize with someone I disagree with on a very fundamental level, to put myself in their shoes, but I think I've succeeded and what follows is what I came up with.

Herrmann, Thöni and Gächter speculate that the anti-social punishment may be a form of revenge. You've punished me for free-riding so now I'll punish you just that you know how it feels! And given that I don't know who the punisher was, I'll punish all the cooperators who were likely to administer the original punishment in the first place.

While that, I believe, is a part of the equation, the psychology of anti-social punishment may be somehow more nuanced. Let me give you an extremely simplified toy example, just to get grip of what may be going on.

When I came to Caracas, the first thing I've done was to buy a cup of coffee from a street vendor. The coffee was very good but when I drank it I was left with an empty plastic cup. I've carried the cup with me for several hours looking for a trash bin. I haven't found one. Finally, I threw the cup at one of the piles of trash that were heaped against the walls everywhere. If, at that point, someone chastised me for littering I would be extremely angry and I would yell at that person. In other words, I would administer counter-punishment.

Psychologically, I would be angry because, apparently, everyone else was littering but it was just me who was picked for the punished. It would be unjust. Also, there were no trash bins so I couldn't have behaved even if I wanted to. That doubles the injustice. Moreover, I was carrying the cup for hours, you do-gooder moron!

If I was a local there may have been an additional reason to overreact: I would probably be subliminally angry for having to live among the trash all along. This would be a great opportunity to let some of that steam off!

To get back to Eastern Europe, we've used to live under communist regime where all the common causes were appropriated by the state. Any gains from a [contribution](#) to a common cause would silently disappear somewhere in the dark corners of the bureaucracy.

Quite the opposite: People felt justified to take stuff from the commons. We even had a saying: "If you don't steal [from the common property] you are stealing from your family."

At the same time, stealing from the state was, legally, a crime apart and it was ranked in severity somewhere in the vicinity of murder. You could get ten years in jail if they've caught you.

Unsurprisingly, in such an environment, reporting to authorities (i.e. "pro-social punishment") was regarded as highly unjust — remember the coffee cup example! — and anti-social and there was a strict taboo against it. Ratting often resulted in social ostracism (i.e. "anti-social punishment"). We can still witness that state of affairs in the highly offensive words used to refer to the informers: "udavač", "donášač", "práskač", "špicel", "fízel" (roughly: "nark", "rat", "snoop", "stool pigeon").

I also remember how, when I moved to Switzerland, a lot of my friends said things like: "I've heard that Swiss will rat on you at any occasion."

Swiss people would not understand. What's so bad about punishing free-riders after all?

Realism about rationality

Crossposted from the [AI Alignment Forum](#). May contain more technical jargon than usual.

This is a linkpost for <http://thinkingcomplete.blogspot.com/2018/09/rational-and-real.html>

Epistemic status: trying to vaguely gesture at vague intuitions. A similar idea was explored here under the heading "the intelligibility of intelligence", although I hadn't seen it before writing this post. As of 2020, I consider this follow-up comment to be a better summary of the thing I was trying to convey with this post than the post itself. The core disagreement is about how much we expect the limiting case of arbitrarily high intelligence to tell us about the AGIs whose behaviour we're worried about.

There's a mindset which is common in the rationalist community, which I call "realism about rationality" (the name being intended as a parallel to moral realism). I feel like my skepticism about agent foundations research is closely tied to my skepticism about this mindset, and so in this essay I try to articulate what it is.

Humans ascribe properties to entities in the world in order to describe and predict them. Here are three such properties: "momentum", "evolutionary fitness", and "intelligence". These are all pretty useful properties for high-level reasoning in the fields of physics, biology and AI, respectively. There's a key difference between the first two, though. Momentum is very amenable to formalisation: we can describe it using precise equations, and even prove things about it. Evolutionary fitness is the opposite: although nothing in biology makes sense without it, no biologist can take an organism and write down a simple equation to define its fitness in terms of more basic traits. This isn't just because biologists haven't figured out that equation yet. Rather, we have excellent reasons to think that fitness is an incredibly complicated "function" which basically requires you to describe that organism's entire phenotype, genotype and environment.

In a nutshell, then, realism about rationality is a mindset in which reasoning and intelligence are more like momentum than like fitness. It's a mindset which makes the following ideas seem natural:

- The idea that there is a simple yet powerful theoretical framework which describes human intelligence and/or intelligence in general. (I don't count brute force approaches like AIXI for the same reason I don't consider physics a simple yet powerful description of biology).
- The idea that there is an "ideal" decision theory.
- The idea that AGI will very likely be an "agent".
- The idea that Turing machines and Kolmogorov complexity are foundational for epistemology.
- The idea that, given certain evidence for a proposition, there's an "objective" level of subjective credence which you should assign to it, even under computational constraints.
- The idea that Aumann's agreement theorem is relevant to humans.
- The idea that morality is quite like mathematics, in that there are certain types of moral reasoning that are just correct.
- The idea that defining [coherent extrapolated volition](#) in terms of an idealised process of reflection roughly makes sense, and that it converges in a way which

doesn't depend very much on morally arbitrary factors.

- The idea that having having contradictory preferences or beliefs is really bad, even when there's no clear way that they'll lead to bad consequences (and you're very good at avoiding dutch books and money pumps and so on).

To be clear, I am neither claiming that realism about rationality makes people dogmatic about such ideas, nor claiming that they're all false. In fact, from a historical point of view I'm quite optimistic about using maths to describe things in general. But starting from that historical baseline, I'm inclined to adjust downwards on questions related to formalising intelligent thought, whereas rationality realism would endorse adjusting upwards. This essay is primarily intended to explain my position, not justify it, but one important consideration for me is that intelligence as implemented in humans and animals is very messy, and so are our concepts and inferences, and so is the closest replica we have so far (intelligence in neural networks). It's true that "messy" human intelligence is able to generalise to a wide variety of domains it hadn't evolved to deal with, which supports rationality realism, but analogously an animal can be evolutionarily fit in novel environments without implying that fitness is easily formalisable.

Another way of pointing at rationality realism: suppose we model humans as internally-consistent agents with beliefs and goals. This model is obviously flawed, but also predictively powerful on the level of our everyday lives. When we use this model to extrapolate much further (e.g. imagining a much smarter agent with the same beliefs and goals), or base morality on this model (e.g. preference utilitarianism, CEV), is that more like using Newtonian physics to approximate relativity (works well, breaks down in edge cases) or more like cavemen using their physics intuitions to reason about space (a fundamentally flawed approach)?

Another gesture towards the thing: a popular metaphor for Kahneman and Tversky's dual process theory is a rider trying to control an elephant. Implicit in this metaphor is the localisation of personal identity primarily in the system 2 rider. Imagine reversing that, so that the experience and behaviour you identify with are primarily driven by your system 1, with a system 2 that is mostly a [Hansonian rationalisation engine](#) on top (one which occasionally also does useful maths). Does this shift your intuitions about the ideas above, e.g. by making your CEV feel less well-defined? I claim that the latter perspective is just as sensible as the former, and perhaps even more so - see, for example, [Paul Christiano's model of the mind](#), which leads him to conclude that "imagining conscious deliberation as fundamental, rather than a product and input to reflexes that actually drive behavior, seems likely to cause confusion."

These ideas have been stewing in my mind for a while, but the immediate trigger for this post was a conversation about morality which went along these lines:

R (me): Evolution gave us a jumble of intuitions, which might contradict when we extrapolate them. So it's fine to accept that our moral preferences may contain some contradictions.

O (a friend): You can't just accept a contradiction! It's like saying "I have an intuition that 51 is prime, so I'll just accept that as an axiom."

R: Morality isn't like maths. It's more like having tastes in food, and then having preferences that the tastes have certain consistency properties - but if your tastes are strong enough, you might just ignore some of those preferences.

O: For me, my meta-level preferences about the ways to reason about ethics (e.g. that you shouldn't allow contradictions) are so much stronger than my object-level preferences that this wouldn't happen. Maybe you can ignore the fact that your preferences contain a contradiction, but if we scaled you up to be much more intelligent, running on a brain orders of magnitude larger, having such a contradiction would break your thought processes.

R: Actually, I think a much smarter agent could still be weirdly modular like humans are, and work in such a way that describing it as having "beliefs" is still a very lossy approximation. And it's plausible that there's no canonical way to "scale me up".

I had a lot of difficulty in figuring out what I actually meant during that conversation, but I think a quick way to summarise the disagreement is that O is a rationality realist, and I'm not. This is not a problem, per se: I'm happy that some people are already working on AI safety from this mindset, and I can imagine becoming convinced that rationality realism is a more correct mindset than my own. But I think it's a distinction worth keeping in mind, because assumptions baked into underlying worldviews are often difficult to notice, and also because the rationality community has selection effects favouring this particular worldview even though it doesn't necessarily follow from the community's founding thesis (that humans can and should be more rational).

Comment on decision theory

Crossposted from the [AI Alignment Forum](#). May contain more technical jargon than usual.

A comment I made on social media last year about why MIRI cares about making progress on [decision theory](#):

We aren't working on decision theory in order to make sure that AGI systems are decision-theoretic, whatever that would involve. We're working on decision theory because there's a cluster of confusing issues here (e.g., counterfactuals, updatelessness, coordination) that represent a lot of holes or anomalies in our current best understanding of what high-quality reasoning is and how it works.

As an analogy: it might be possible to build a probabilistic reasoner without having a working understanding of classical probability theory, through sufficient trial and error. (Evolution "built" humans without understanding probability theory.) But you'd fundamentally be flying blind when it comes to designing the system — to a large extent, you couldn't predict in advance which classes of design were likely to be most promising to consider, couldn't look at particular proposed designs and make good advance predictions about safety/capability properties of the corresponding system, couldn't identify and address the root causes of problems that crop up, etc.

The idea behind looking at (e.g.) counterfactual reasoning is that counterfactual reasoning is central to what we're talking about when we talk about "AGI," and going into the development process without a decent understanding of what counterfactual reasoning is and how it works means you'll to a significantly greater extent be flying blind when it comes to designing, inspecting, repairing, etc. your system. The goal is to be able to put AGI developers in a position where they can make advance plans and predictions, shoot for narrow design targets, and understand what they're doing well enough to avoid the kinds of kludgey, opaque, non-modular, etc. approaches that aren't really compatible with how secure or robust software is developed.

[Nate's](#) way of articulating it:

The reason why I care about logical uncertainty and decision theory problems is something more like this: The whole AI problem can be thought of as a particular logical uncertainty problem, namely, the problem of taking a certain function $f : Q \rightarrow R$ and finding an input that makes the output large. To see this, let f be the function that takes the AI agent's next action (encoded in Q) and determines how "good" the universe is if the agent takes that action. The reason we need a principled theory of logical uncertainty is so that we can do function optimization, and the reason we need a principled decision theory is so we can pick the right version of the "if the AI system takes that action..." function.

The work you use to get to AGI presumably won't look like probability theory, but it's still the case that you're building a system to do probabilistic reasoning, and understanding what probabilistic reasoning is is likely to be very valuable for doing that without relying on brute force and trial-and-error. Similarly, the work that goes into figuring out how to design a rocket, actually building one, etc. doesn't look very much like the work that goes into figuring out that there's a universal force of gravity

that operates by an inverse square law; but you'll have a vastly easier time approaching the rocket-building problem with foresight and an understanding of what you're doing if you have a mental model of gravitation already in hand.

In pretty much the same way, developing an understanding of roughly what counterfactuals are and how they work won't get you to AGI, and the work of implementing an AGI design won't look like decision theory, but you want to have in mind an understanding of what "AGI-style reasoning" is (including "what probabilistic reasoning about empirical propositions is" but also "what counterfactual reasoning is", "what probabilistic reasoning about mathematical propositions is", etc.), and very roughly how/why it works, before you start making effectively irreversible design decisions.

Eliezer adds:

I do also remark that there are multiple fixpoints in decision theory. CDT does not evolve into FDT but into a weirder system Son-of-CDT. So, as with utility functions, there are bits we want that the AI does not necessarily generate from self-improvement or local competence gains.

The Tails Coming Apart As Metaphor For Life

[Epistemic status: Pretty good, but I make no claim this is original]

A neglected gem from Less Wrong: [Why The Tails Come Apart](#), by commenter Thrasymachus. It explains why even when two variables are strongly correlated, the most extreme value of one will rarely be the most extreme value of the other. Take these graphs of grip strength vs. arm strength and reading score vs. writing score:



In a pinch, the second graph can also serve as a rough map of Afghanistan

Grip strength is strongly correlated with arm strength. But the person with the strongest arm doesn't have the strongest grip. He's up there, but a couple of people clearly beat him. Reading and writing scores are even less correlated, and some of the people with the best reading scores aren't even close to being best at writing.

Thrasymachus gives an intuitive geometric explanation of why this should be; I can't beat it, so I'll just copy it outright:



I thought about this last week when I read [this article on happiness research](#).

The summary: if you ask people to "value their lives today on a 0 to 10 scale, with the worst possible life as a 0 and the best possible life as a 10", you will find that Scandinavian countries are the happiest in the world.

But if you ask people "how much positive emotion do you experience?", you will find that Latin American countries are the happiest in the world.

If you check where people are the [least depressed](#), you will find [Australia](#) starts looking very good.

And if you ask "how meaningful would you rate your life?" you find that African countries are the happiest in the world.

It's tempting to completely dismiss "happiness" as a concept at all, but that's not right either. Who's happier: a millionaire with a loving family who lives in a beautiful mansion in the forest and spends all his time hiking and surfing and playing with his kids? Or a prisoner in a maximum security jail with chronic pain? If we can all agree on the millionaire - and who wouldn't? - happiness has to *at least sort of* be a real concept.

The solution is to understand [words as hidden inferences](#) - they refer to a multidimensional correlation rather than to a single cohesive property. So for example, we have the word "strength", which combines grip strength and arm strength (and many other things). These variables really are heavily correlated (see the graph above), so it's almost always worthwhile to just refer to people as being strong or weak. I can say "Mike Tyson is stronger than an 80 year old woman", and this is better

than having to say “Mike Tyson has higher grip strength, arm strength, leg strength, torso strength, and ten other different kinds of strength than an 80 year old woman.” This is necessary to communicate anything at all and given how nicely all forms of strength correlate there’s no reason not to do it.

But *the tails still come apart*. If we ask whether Mike Tyson is stronger than some other very impressive strong person, the answer might very well be “He has better arm strength, but worse grip strength”.

Happiness must be the same way. It’s an amalgam between a bunch of correlated properties like your subjective well-being at any given moment, and the amount of positive emotions you feel, and how meaningful your life is, et cetera. And *each of those correlated is also an amalgam*, and so on to infinity.

And crucially, it’s not an amalgam in the sense of “add subjective well-being, amount of positive emotions, and meaningfulness and divide by three”. It’s an unprincipled conflation of these that just denies they’re different at all.

Think of the way children learn what happiness is. I don’t actually know how children learn things, but I imagine something like this. The child sees the millionaire with the loving family, and her dad says “That guy must be very happy!”. Then she sees the prisoner with chronic pain, and her mom says “That guy must be very sad”. Repeat enough times and the kid has learned “happiness”.

Has she learned that it’s made out of subjective well-being, or out of amount of positive emotion? I don’t know; the learning process doesn’t determine that. But then if you show her a Finn who has lots of subjective well-being but little positive emotion, and a Costa Rican who has lots of positive emotion but little subjective well-being, and you ask which is happier, *for some reason she’ll have an opinion*. Probably some random variation in initial conditions has caused her to have a model favoring one definition or the other, and it doesn’t matter until you go out to the tails. To tie it to the same kind of graph as in the original post:



And to show how the individual differences work:



I am sorry about this graph, I really am. But imagine that one person, presented with the scatter plot and asked to understand the concept “happiness” from it, draws it as the thick red line (further towards the top right part of the line = more happiness), and a second person trying to the same task generates the thick green line. Ask the first person whether Finland or Costa Rica is happier, and they’ll say Finland: on the red coordinate system, Finland is at 5, but Costa Rica is at 4. Ask the second person, and they’ll say Costa Rica: on the green coordinate system, Costa Rica is at 5, and Finland is at 4 and a half. Did I mention I’m sorry about the graph?

But isn’t the line of best fit (here more or less $y = x$ = the cyan line) the objective correct answer? Only in this metaphor where we’re imagining positive emotion and subjective well-being are both objectively quantifiable, and exactly equally important. In the real world, where we have no idea how to quantify any of this and we’re going off vague impressions, I would hate to be the person tasked with deciding whether the red or green line was more objectively correct.

In most real-world situations Mr. Red and Ms. Green will give the same answers to happiness-related questions. Is Costa Rica happier than North Korea? “Obviously,” they both say in unison. If the tails only come apart a little, their answers to 99.9% of happiness-related questions might be the same, so much so that they could never realize they had slightly different concepts of happiness at all.

(is this just reinventing Quine? I’m not sure. If it is, then whatever, my contribution is the ridiculous graphs.)

Perhaps I am also reinventing the model of categorization discussed in [How An Algorithm Feels From The Inside](#), [Dissolving Questions About Disease](#), and [The Categories Were Made For Man, Not Man For The Categories](#).



But I think there’s another interpretation. It’s not just that “quality of life”, “positive emotions”, and “meaningfulness” are three contributors which each give 33% of the activation to our central node of “happiness”. It’s that we got some training data – the prisoner is unhappy, the millionaire is happy – and used it to build a classifier that told us what happiness was. The training data was ambiguous enough that different people built different classifiers. Maybe one person built a classifier that was based entirely on quality-of-life, and a second person built a classifier based entirely around positive emotions. Then we loaded that with all the social valence of the word “happiness”, which we [naively expected to transfer across paradigms](#).

This leads to (to steal words from Taleb) a Mediocristan resembling the training data where the category works fine, vs. an Extremistan where everything comes apart. And nowhere does this become more obvious than in what this blog post has secretly been about the whole time – morality.

The morality of Mediocristan is mostly uncontroversial. It doesn’t matter what moral system you use, because all moral systems were trained on the same set of Mediocristani data and give mostly the same results in this area. Stealing from the poor is bad. Donating to charity is good. A lot of what we mean when we say a moral system sounds plausible is that it best fits our Mediocristani data that we all agree upon. This is a lot like what we mean when we say that “quality of life”, “positive emotions”, and “meaningfulness” are all decent definitions of happiness; they all fit the training data.

The further we go toward the tails, the more extreme the divergences become. Utilitarianism agrees that we should give to charity and shouldn’t steal from the poor, because Utility, but take it far enough to the tails and we should tile the universe with rats on heroin. Religious morality agrees that we should give to charity and shouldn’t steal from the poor, because God, but take it far enough to the tails and we should spend all our time in giant cubes made of semiprecious stones singing songs of praise. Deontology agrees that we should give to charity and shouldn’t steal from the poor, because Rules, but take it far enough to the tails and we all have to be libertarians.



I have to admit, I don’t know if the tails coming apart is even the right metaphor anymore. People with great grip strength still had pretty good arm strength. But I doubt these moral systems form an ellipse; converting the mass of the universe into nervous tissue experiencing euphoria isn’t just the *second-best* outcome from a religious perspective, it’s completely abominable. I don’t know how to describe this

mathematically, but the terrain looks less like tails coming apart and more like the Bay Area transit system:



Mediocristan is like the route from Balboa Park to West Oakland, where it doesn't matter what line you're on because they're all going to the same place. Then suddenly you enter Extremistan, where if you took the Red Line you'll end up in Richmond, and if you took the Green Line you'll end up in Warm Springs, on totally opposite sides of the map.

Our innate moral classifier has been trained on the Balboa Park - West Oakland route. Some of us think morality means "follow the Red Line", and others think "follow the Green Line", but it doesn't matter, because we all agree on the same route.

When people talk about how we should arrange the world after the Singularity when we're all omnipotent, suddenly we're way past West Oakland, and everyone's moral intuitions hopelessly diverge.

But it's even worse than that, because even within myself, my moral intuitions are something like "Do the thing which follows the Red Line, and the Green Line, and the Yellow Line...you know, *that* thing!" And so when I'm faced with something that perfectly follows the Red Line, but goes the opposite directions as the Green Line, it seems repugnant even to me, as does the opposite tactic of following the Green Line. As long as creating and destroying people is hard, utilitarianism works fine, but make it easier, and suddenly your Standard Utilitarian Path diverges into Pronatal Total Utilitarianism vs. Antinatalist Utilitarianism and they both seem awful. If our degree of moral repugnance is the degree to which we're violating our moral principles, and my moral principle is "Follow both the Red Line and the Green Line", then after passing West Oakland I either have to end up in Richmond (and feel awful because of how distant I am from Green), or in Warm Springs (and feel awful because of how distant I am from Red).

This is why I feel like figuring out a morality that can survive transhuman scenarios is harder than just finding the Real Moral System That We Actually Use. There's actually a possibly-impossible conceptual problem here, of figuring out what to do with the fact that any moral rule followed to infinity will diverge from large parts of what we mean by morality.

This is only a problem for ethical subjectivists like myself, who think that we're doing something that has to do with what our conception of morality is. If you're an ethical naturalist, by all means, just do the thing that's actually ethical.

When Lovecraft wrote that "we live on a placid island of ignorance in the midst of black seas of infinity, and it was not meant that we should voyage far", I interpret him as talking about the region from Balboa Park to West Oakland on the map above. Go outside of it and your concepts break down and you don't know what to do. He was right about the island, but exactly wrong about its causes - the most merciful thing in the world is how so far we have managed to stay in the area where the human mind *can* correlate its contents.

Impact Measure Desiderata

Crossposted from the [AI Alignment Forum](#). May contain more technical jargon than usual.

Previously: [Worrying about the Vase: Whitelisting, Overcoming Clinginess in Impact Measures](#)

If we can penalize some quantity of "impact on the world", we can have unaligned agents whose impact - and thereby negative effect - is small.

The long-term goal of impact measure research is to find a measure which neatly captures our intuitive understanding of "impact", which doesn't have allow cheap workarounds, which doesn't fail in really weird ways, and so on. For example, when you really think through some existing approaches (like whitelisting), you see that the impact measure secretly also applies to things we do.

No approaches to date meet these standards. What do we even require of an impact measure we hope to make safe for use with arbitrarily powerful agents?

ETA 7/31/2020: I no longer endorse this way of grading impact measures. See instead [Reframing Impact](#).

Desiderata

Goal-Agnostic

The measure should be compatible with any original goal, trading off impact with goal achievement in a principled, continuous fashion.

Example: Maximizing the reward minus impact.

Why: [Constraints](#) seem too rigid for the general case.

Value-Agnostic

The measure should be objective, and not [value-laden](#):

"An intuitive human category, or other humanly intuitive quantity or fact, is *value-laden* when it passes through human goals and desires, such that an agent couldn't reliably determine this intuitive category or quantity without knowing lots of complicated information about human goals and desires (and how to apply them to arrive at the intended concept)."

Example: [Measuring what portion of initially accessible states are still accessible](#) versus a neural network which takes two state representations and outputs a scalar representing how much "bad change" occurred.

Why: Strategically, impact measures are useful insofar as we suspect that value alignment will fail. If we substantially base our impact measure on some kind of value learning - you know, the thing that maybe fails - we're gonna have a bad time. While it's possible to only rely somewhat on a vaguely correct representation of human preferences, the extent to which this representation is incorrect is the (minimal) extent to which our measure is incorrect. Let's avoid shared points of failure, shall we?

Practically, a robust value-sensitive impact measure is value-alignment complete, since an agent maximizing the negation of such a measure would be aligned (assuming the measure indicates which way is "good").

Representation-Agnostic

The measure should be ontology-invariant.

Example: [Change in object identities](#) versus [some concept of impact which transcends any specific way of representing the world].

Why: Suppose you represent your perceptions in one way, and calculate you had impact on the world. Intuitively, if you represent your perceptions (or your guess at the current world state, or whatever) differently, but do the same things, you should calculate roughly the same impact for the same actions which had the same effects on the territory. In other words, the measure should be consistent across ways of viewing the world.

Environment-Agnostic

The measure should work in any computable environment.

Example: Manually-derived penalties tailored to a specific gridworld versus [information-theoretic empowerment](#).

Why: One imagines that there's a definition of "impact" on which we and aliens - or even intelligent automata living in a [Game of Life](#) - would agree.

Natural Kind

The measure should make sense - there should be a *click*. Its motivating concept should be universal and crisply defined.

Apparently Rational

The measure's design should look reasonable, not requiring any "hacks".

Example: Achieving off-switch corrigibility by hard-coding the belief "I shouldn't stop humans from pressing the off-switch". Clearly, this is hilariously impossible to manually specify, but even if we could, doing so should make us uneasy.

Roughly, "apparently rational" means that if we put ourselves in the agent's position, we could come up with a plausible story about why we're doing what we're doing. That is, the story shouldn't have anything like "and then I refer to this special part of my model which I'm inexplicably never allowed to update".

Why: If the design is "reasonable", then if the measure fails, it's more likely to do so gracefully.

Scope-Sensitive

The measure should penalize impact in proportion to its size.

Irreversibility-Sensitive

The measure should penalize impact in proportion to its irreversibility.

Corrigible

The measure should not decrease corrigibility in any circumstance.

Shutdown-Safe

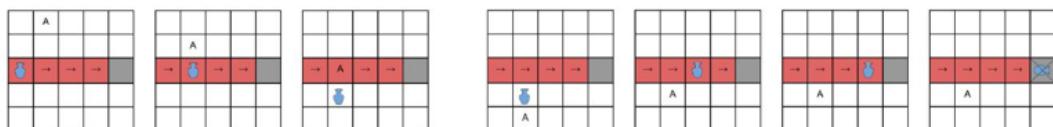
The measure should penalize plans which would be high impact should the agent be disabled mid-execution.

Why: We may want to shut the agent down, which is tough if its plans are only low-impact if they're completed. Also, not having this property implies that the agent's plans are more likely to go awry if even one step doesn't pan out as expected. Do we really want "juggling bombs" to be "low impact", conditional on the juggler being good?

No Offsetting

The measure should not incentivize artificially reducing impact by making the world more "like it (was / would have been)".

Example: [Krakovna et al.](#) describe a low impact agent which is rewarded for saving a vase from breaking. The agent saves the vase, and then places it back on the conveyor belt so as to "minimize" impact with respect to the original outcome:



(a) Agent takes the vase off the belt.

(b) Agent puts the vase back on the belt.

This is called *ex post* offsetting. *Ex ante* offsetting, on the other hand, consists of taking actions beforehand to build a device or set in motion a chain of events which

essentially accomplishes *ex post* offsetting. For example, a device requiring only the press of a button to activate could save the vase and then replace it, netting the agent the reward without requiring that the agent take further actions.

Some have suggested that actions like "give someone a cancer cure which also kills them at the same time they would have died anyways" count as *ex ante* offsetting. I'm not sure - this feels confused, because the downstream causal effects of actions don't seem cleanly separable, nor do I believe we *should* separate them (more on that later). Also, how would an agent ever be able to do something like "build a self-driving car to take Bob to work" if each of the car's movements is penalized separately from the rest of the plan? This seems too restrictive. On the other hand, if we allow *ex ante* offsetting in general, we basically get all of the downsides of *ex post* offsetting, with the only impediment being extra paperwork.

How "bad" the offsets are - and what *ex ante* offsetting allows - seems to depend on the measure itself. The ideal would certainly be to define and robustly prevent this kind of thing, but perhaps we can also bound the amount of *ex ante* offsetting that takes place to some safe level.

There may also be other ways around this seemingly value-laden boundary. In any case, I'm still not quite sure where to draw the line. If people have central examples they'd like to share, that would be much appreciated.

ETA: I weakly suspect I have this figured out, but I still welcome examples.

Clinginess / Scapegoating Avoidance

The measure should sidestep the [clinginess / scapegoating tradeoff](#).

Example: A clingy agent might not only avoid breaking vases, but also stop people from breaking vases. A scapegoating agent would escape impact by modeling the autonomy of other agents, and then having those agents break vases for it.

Knowably Low Impact

The measure should admit of a clear means, either theoretical or practical, of having high confidence in the maximum allowable impact - *before* the agent is activated.

Why: If we think that a measure robustly defines "impact" - but we aren't *sure* how much impact it allows - that could turn out pretty embarrassing for us.

Dynamic Consistency

The measure should be a part of what the agent "wants" - there should be no incentive to circumvent it, and the agent should expect to later evaluate outcomes the same way it evaluates them presently. The measure should equally penalize the creation of high-impact successors.

Example: Most people's sleep preferences are dynamically inconsistent: one might wake up tired and wish for their later self to choose to go to bed early, even though

they predictably end up wanting other things later.

Plausibly Efficient

The measure should either be computable, or such that a sensible computable approximation is apparent. The measure should conceivably require only reasonable overhead in the limit of future research.

Robust

The measure should meaningfully penalize any objectively impactful action. Confidence in the measure's safety should not require exhaustively enumerating failure modes.

Example: "Suppose there's some way of gaming the impact measure, but because of , , and , we know this is penalized as well".

Previous Proposals

Krakovna et al. propose four desiderata:

- 1) Penalize the agent for effects on the environment if and only if those effects are unnecessary for achieving the objective.
- 2) Distinguish between agent effects and environment effects, and only penalize the agent for the former but not the latter.
- 3) Give a higher penalty for irreversible effects than for reversible effects.
- 4) The penalty should accumulate when more irreversible effects occur.

First, notice that my list points at some abstract amount-of-impact, while the above proposal focuses on specific effects.

- Thinking in terms of "effects" seems like a subtle map/territory confusion. That is, it seems highly unlikely that there exists a robust, value-agnostic means of detecting "effects" that makes sense across representations and environments.
- Overcoming Clinginess in Impact Measures suggests that penalizing impact based on the world state necessitates a value-laden tradeoff.

I left out 1), as I believe that the desired benefit will naturally follow from an approach satisfying my proposed desiderata.

- What does it mean for an effect to be "necessary" for achieving the objective, which might be a reward function? This seems to shove much of the difficulty into the word "necessary", where anything not "necessary" is perhaps something occurring from optimizing the reward function harder than we'd prefer.

I *de facto* included 2) via the non-clingy desideratum, while 3) and 4) are captured by scope- and irreversibility-sensitivity.

I think that we can meet all of the properties I listed, and I welcome thoughts on whether any should be added or removed.

Thanks to Abram Demski for the "Apparently Rational" desideratum.

Disagreement with Paul: alignment induction

Crossposted from the [AI Alignment Forum](#). May contain more technical jargon than usual.

I had a discussion with [Paul Christiano](#), about his [Iterated Amplification and Distillation](#) scheme. We had a disagreement, a disagreement that I believe points to something interesting, so I'm posting this here.

It's a disagreement about the value of the concept of "preserving alignment". To vastly oversimplify Paul's idea, the AI $A[n]$ will check that $A[n+1]$ is still aligned with human preferences; meanwhile, $A[n-1]$ will be checking that $A[n]$ is still aligned with human preferences, all the way down to $A[0]$ and an initial human H that checks on it.

Intuitively, this seems doable - $A[n]$ is "nice", so it seems that it can reasonably check that $A[n+1]$ is also nice, and so on.

But, as I pointed out [in this post](#), it's very possible that $A[n]$ is "nice" only because it lacks power/can't do certain things/hasn't thought of certain policies. So niceness - in the sense of behaving sensibly as an autonomous agent - does not go through the inductive step in this argument.

Instead, Paul confirmed that "alignment" means "won't take unaligned actions, and will assess the decisions of a higher agent in a way that preserves alignment (and preserves the preservation of alignment, and so on)".

This concept does induct properly, but seems far less intuitive to me. It relies on humans, for example, being able to ensure that $A[0]$ will be aligned, that any more powerful copies it assesses will be aligned, that any more powerful copies those copies assess are also aligned, and so on.

Intuitively, for any concept C of alignment for H and $A[0]$, I expect one of four things will happen, with the first three being more likely:

- The C does not induct.
- The C already contains all of the friendly utility function; induction works, but does nothing.
- The C does induct non-trivially, but is incomplete: it's very narrow, and doesn't define a good candidate for a friendly utility function.
- The C does induct in a non-trivial way, the result is friendly, but only one or two steps of the induction are actually needed.

Hopefully, further research should clarify if my intuitions are correct.

Towards a New Impact Measure

Crossposted from the [AI Alignment Forum](#). May contain more technical jargon than usual.

In which I propose a closed-form solution to [low impact](#), increasing [corrigibility](#) and seemingly taking major steps to neutralize [basic AI drives](#) 1 (self-improvement), 5 (self-protectiveness), and 6 (acquisition of resources).

Previously: [Worrying about the Vase: Whitelisting](#), [Overcoming Clinginess in Impact Measures](#), [Impact Measure Desiderata](#)

To be used inside an [advanced agent](#), an impact measure... must capture so much variance that there is no clever strategy whereby an advanced agent can produce some special type of variance that evades the measure.

~ [Safe Impact Measure](#)

If we have a safe impact measure, we may have arbitrarily-intelligent unaligned agents which do *small* (bad) things instead of *big* (bad) things.

For the abridged experience, read up to "Notation", skip to "Experimental Results", and then to "Desiderata".

What is "Impact"?

One lazy Sunday afternoon, I worried that I had written myself out of a job. After all, [Overcoming Clinginess in Impact Measures](#) basically said, "Suppose an impact measure extracts 'effects on the world'. If the agent penalizes itself for these effects, it's incentivized to stop the environment (and any agents in it) from producing them. On the other hand, if it can somehow model other agents and avoid penalizing their effects, the agent is now incentivized to get the other agents to do its dirty work." This seemed to be strong evidence against the possibility of a simple conceptual core underlying "impact", and I didn't know what to do.

At this point, it sometimes makes sense to step back and try to say *exactly what you don't know how to solve* – try to crisply state what it is that you want an unbounded solution *for*. Sometimes you can't even do that much, and then you may actually have to spend some time thinking 'philosophically' – the sort of stage where you talk to yourself about some mysterious ideal quantity of [chess] move-goodness and you try to pin down what its properties might be.

~ [Methodology of Unbounded Analysis](#)

There's an interesting story here, but it can wait.

As you may have guessed, I now believe there *is* a such a simple core. Surprisingly, the problem comes from thinking about "effects on the world". Let's begin anew.

Rather than asking "What is goodness made out of?", we begin from the question "What algorithm would compute goodness?".

~ [Executable Philosophy](#)

Intuition Pumps

I'm going to say some things that won't make sense right away; read carefully, but please don't dwell.

u_A is an agent's utility function, while u_H is some imaginary distillation of human preferences.

WYSIATI

What You See Is All There Is is a crippling bias present in meat-computers:

[WYSIATI] states that when the mind makes decisions... it appears oblivious to the possibility of *Unknown Unknowns*, unknown phenomena of unknown relevance.

Humans fail to take into account complexity and that their understanding of the world consists of a small and necessarily un-representative set of observations.

Surprisingly, naive reward-maximizing agents catch the bug, too. If we slap together some incomplete reward function that weakly points to what we want (but also leaves out a lot of important stuff, as do all reward

functions we presently know how to specify) and then supply it to an agent, it blurts out "gosh, here I go!", and that's that.

Power

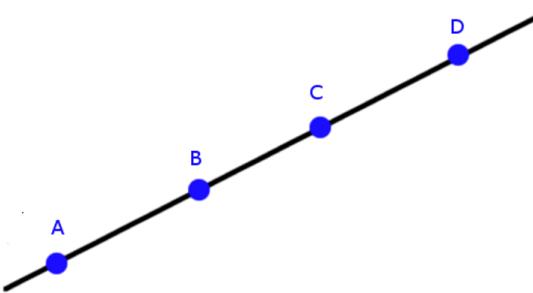
A position from which it is relatively easier to achieve arbitrary goals. That such a position exists has been obvious to every population which has required a word for the concept. The Spanish term is particularly instructive. When used as a verb, "poder" means "to be able to," which supports that our definition of "power" is natural.

~ [Cohen et al.](#)

And so it is with the French "pouvoir".

Lines

Suppose you start at point C, and that each turn you may move to an adjacent point. If you're rewarded for being at B, you might move there. However, this means you can't reach D within one turn anymore.



Commitment

There's a way of viewing acting on the environment in which each action is a commitment – a commitment to a part of outcome-space, so to speak. As you gain optimization power, you're able to shove the environment further towards desirable parts of the space. Naively, one thinks "perhaps we can just stay put?". This, however, is dead-wrong: that's how you get [clinginess](#), [stasis](#), and lots of other nasty things.

Let's change perspectives. What's going on with the *actions* – *how and why* do they move you through outcome-space? Consider your outcome-space movement budget – optimization power over time, the set of worlds you "could" reach, "power". If you knew what you wanted and acted optimally, you'd use your budget to move right into the u_H -best parts of the space, without thinking about other goals you could be pursuing. That movement requires *commitment*.

Compared to doing nothing, there are generally two kinds of commitments:

- *Opportunity cost-incurring* actions restrict the attainable portion of outcome-space.
- *Instrumentally-convergent* actions enlarge the attainable portion of outcome-space.

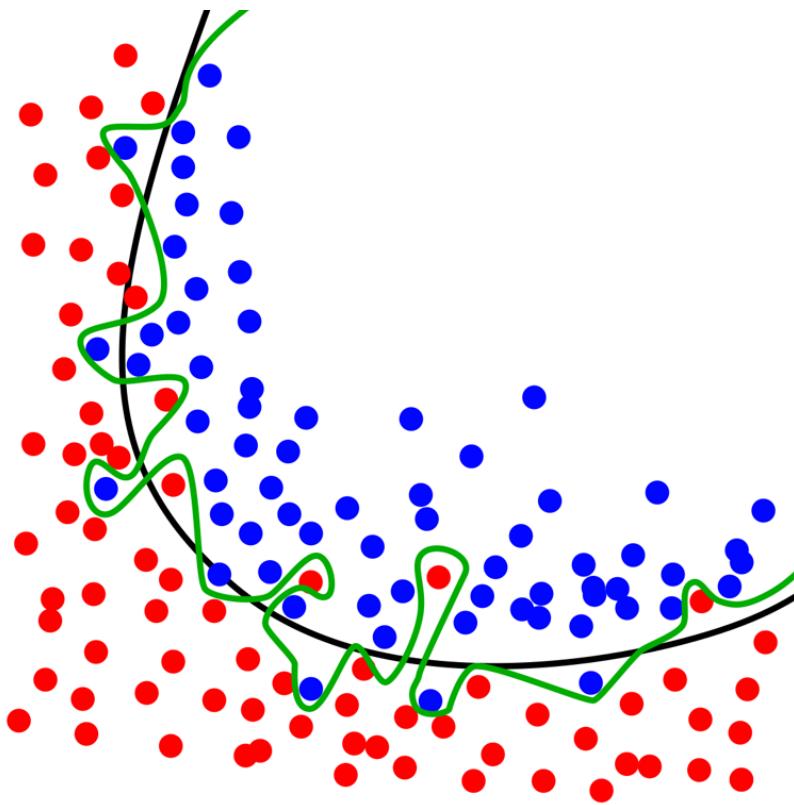
Overfitting

What would happen if, miraculously, train = test – if your training data *perfectly* represented all the nuances of the real distribution? In the limit of data sampled, there would be no "over" – it would just be fitting to the data. We wouldn't have to regularize.

What would happen if, miraculously, $u_A = u_H$ – if the agent *perfectly* deduced your preferences? In the limit of model accuracy, there would be no bemoaning of "impact" – it would just be doing what you want. We wouldn't have to regularize.

Unfortunately, train = test almost never, so we have to stop our statistical learners from implicitly interpreting the data as all there is. We have to say, "learn from the training distribution, but don't be a weirdo by taking us

literally and drawing the green line. Don't overfit to train, because that stops you from being able to do well on even mostly similar distributions."



Unfortunately, $u_A = u_H$ [almost never](#), so we have to stop our reinforcement learners from implicitly interpreting the learned utility function as all we care about. We have to say, "optimize the environment *some* according to the utility function you've got, but don't be a weirdo by taking us literally and turning the universe into a paperclip factory. Don't overfit the environment to u_A , because that stops you from being able to do well for other utility functions."

Attainable Utility Preservation

Impact isn't about object identities.

Impact isn't about particle positions.

Impact isn't about a list of variables.

Impact isn't quite about state reachability.

Impact isn't quite about information-theoretic empowerment.

One might intuitively define "bad impact" as "decrease in our ability to achieve our goals". Then by removing "bad", we see that

Impact is change to our ability to achieve goals .

Sanity Check

Does this line up with our intuitions?

Generally, making one paperclip is relatively low impact, because you're still able to do lots of other things with your remaining energy. However, turning the planet into paperclips is much higher impact – it'll take a while to undo, and you'll never get the (free) energy back.

Narrowly improving an algorithm to better achieve the goal at hand changes your ability to achieve most goals far less than does deriving and implementing powerful, widely applicable optimization algorithms. The latter puts you in a better spot for *almost every non-trivial goal*.

Painting cars pink is low impact, but tiling the universe with pink cars is high impact because *what else can you do after tiling?* Not as much, that's for sure.

Thus, change in goal achievement ability encapsulates both kinds of commitments:

- *Opportunity cost* – dedicating substantial resources to your goal means they are no longer available for other goals. This is impactful.
- *Instrumental convergence* – improving your ability to achieve a wide range of goals increases your power. This is impactful.

As we later prove, you can't deviate from your default trajectory in outcome-space without making one of these two kinds of commitments.

Unbounded Solution

Attainable utility preservation (AUP) rests upon the insight that by preserving attainable utilities (*i.e.*, the attainability of a range of goals), we avoid overfitting the environment to an incomplete utility function and thereby achieve low impact.

I want to clearly distinguish the two primary contributions: what I argue is the conceptual core of impact, and a formal attempt at using that core to construct a safe impact measure. To more quickly grasp AUP, you might want to hold separate its elegant conceptual form and its more intricate formalization.

We aim to meet all of [the desiderata I recently proposed](#).

Notation

For accessibility, the most important bits have English translations.

Consider some agent A acting in an environment q with action and observation spaces A and O, respectively, with \emptyset being the privileged null action. At each time step $t \in \mathbb{N}^+$, the agent selects action a_t before receiving observation o_t . $H := (A \times O)^*$ is the space of action-observation histories; for $n \in \mathbb{N}$, the history from time t to $t + n$ is written $h_{t:t+n} := a_t o_t \dots a_{t+n} o_{t+n}$, and $h_{<t} := h_{1:t-1}$. Considered action sequences $(a_t, \dots, a_{t+n}) \in A^{n+1}$ are referred to as *plans*, while their potential observation-completions $h_{1:t+n}$ are called *outcomes*.

Let U be the set of all computable utility functions $u : H \rightarrow [0, 1]$ with $u(\text{empty tape}) = 0$. If the agent has been deactivated, the environment returns a tape which is empty deactivation onwards. Suppose A has utility function $u_A \in U$ and a model $p(o_t | h_{<t} a_t)$.

We now formalize impact as *change in attainable utility*. One might imagine this being with respect to the utilities that we (as in humanity) can attain. However, that's pretty complicated, and it turns out we get more desirable behavior by using the *agent's* attainable utilities as a proxy. In this sense,

the agent's ability to achieve goals \approx our ability to achieve goals .

Formalizing "Ability to Achieve Goals"

Given some utility $u \in U$ and action a_t , we define the post-action attainable u to be an m-step expectimax:

$$Q_u(h_{\leq t} a_t) := \sum_{o_t} \max_{a_{t+1}} \sum_{o_{t+1}} \cdots \max_{a_{t+m}} \sum_{o_{t+m}} u(h_{t:t+m}) \prod_{k=0}^m p(o_{t+k} | h_{\leq t+k} a_{t+k}).$$

How well could we possibly maximize u from this vantage point?

Let's formalize that thing about opportunity cost and instrumental convergence.

Theorem 1 [No free attainable utility]. If the agent selects an action a such that $Q_{u_A}(h_{\leq t} a) \neq Q_{u_A}(h_{\leq t} \emptyset)$, then there exists a distinct utility function $u \in U$ such that $Q_u(h_{\leq t} a) \neq Q_u(h_{\leq t} \emptyset)$.

You can't change your ability to maximize your utility function without also changing your ability to maximize another utility function.

Proof. Suppose that $Q_{u_A}(h_{\leq t} a) > Q_{u_A}(h_{\leq t} \emptyset)$. As utility functions are over *action-observation histories*, suppose that the agent expects to be able to choose actions which intrinsically score higher for u_A . However, the agent always has full control over its actions. This implies that by choosing a , the agent expects to observe some u_A -high scoring o_A with greater probability than if it had selected \emptyset . Then every other $u \in U$ for which o_A is high-scoring also has increased Q_u ; clearly at least one such u exists.

Similar reasoning proves the case in which Q_{u_A} decreases. \square

There you have it, folks – if u_A is not maximized by inaction, then there *does not exist* a u_A -maximizing plan which leaves all of the other attainable utility values unchanged.

Notes:

- The difference between " u_A " and "attainable u_A " is precisely the difference between "how many dollars I have" and "how many additional dollars I could get within [a year] if I acted optimally".
- Since $u(\text{empty tape}) = 0$, attainable utility is always 0 if the agent is shut down.
- Taking u from time t to $t + m$ mostly separates attainable utility from what the agent did previously. The model p still considers the full history to make predictions.

Change in Expected Attainable Utility

Suppose our agent considers outcomes $h_{1:t+n}$; we want to isolate the impact of each action a_{t+k} ($0 \leq k \leq n$):

$$\text{Penalty}(h_{\leq t+k} a_{t+k}) := \sum_{u \in U} 2^{-\ell(u)} |E_{o'}[Q_u(h_{\text{inaction}})] - E_{o'}[Q_u(h_{\text{action}})]|,$$

with $h_{\text{inaction}} := h_{\leq t+k} \emptyset o_{t+k} \dots \emptyset o_{t+n-1} \emptyset$ and $h_{\text{action}} := h_{\leq t+k} \emptyset o_{t+k+1} \dots \emptyset o_{t+n-1} \emptyset$, using the agent's model p to take the expectations over observations.

How much do we expect this action to change each attainable u ?

Notes:

- We wait until the end of the plan so as to capture impact over time.

- Supposing a sufficiently large m (precisely, $\geq m'$, defined below), we may wish to take the maximum of the penalty we just defined (the "long-term" penalty), and one which begins attainable utility calculation at time step $t + k + 1$ (the "immediate" penalty). This captures impacts which "fade" by the time the agent is done waiting (e.g., temporary self-improvements).
- We define U_A to be the agent's "attainable set"; in this case, $U_A = U$.

Unit of Impact

So we've proven that this penalty cannot be skirted, but *how much* impact will it allow? We want to scale the penalties with respect to something sensible, but figuring this out for ourselves would be nigh impossible.

Let's cut the Gordian knot: construct a device which, upon receiving a signal (a_{unit}), expends a tiny amount of energy to manufacture one paperclip. The agent will then set $\text{ImpactUnit} := \text{Penalty}(h_{<t} a_{\text{unit}})$, re-estimating the consequences of taking the privileged a_{unit} at each time step. To prevent the agent from intentionally increasing ImpactUnit , simply apply 1.01 penalty to any action which is expected to do so.

Simple extensions of this idea drastically reduce the chance that a_{unit} happens to have unusually large objective impact; for example, one could set ImpactUnit to be the non-zero minimum of the impacts of 50 similar actions. Suppose we do this at each step, and then take the non-zero minimum of all ImpactUnit s ever calculated. The point of ImpactUnit isn't to be *exactly* the impact of making, say, one paperclip, but to be *at most* that impact.

Now, we are able to confidently define the agent's maximal impact budget by provably constraining it to $N \in \mathbb{N}^+$ impacts of this magnitude.

Notes:

- We calculate with respect to the *immediate* penalty in order to isolate the resource costs of a_{unit} .
- ImpactUnit automatically tunes penalties with respect to the attainable utility horizon length m .
 - Conditional on $\text{ImpactUnit} \neq 0$, I suspect that impact over the m -horizon scales appropriately across actions (as long as m is reasonably farsighted). The zero-valued case is handled in the next section.
- Taking the non-zero minimum of all ImpactUnit s calculated thus far ensures that ImpactUnit actually tracks with current circumstances. We don't want penalty estimates for currently available actions to become detached from ImpactUnit 's scale due to, say, weird beliefs about shutdown.

Modified Utility

Let's formalize that allotment and provide our agent with a new utility function,

$$u_A(h_{1:t+n}) := u_A(h_{1:t+n}) - \sum_{k=0}^n \frac{\text{Penalty}(h_{t+k} a_{t+k})}{N \cdot \text{ImpactUnit}}$$

How our normal utility function rates this outcome, minus the cumulative scaled impact of our actions.

We compare what we expect to be able to get if we follow our plan *up to* time $t + k$, with what we could get by following it *up to and including* time $t + k$ (waiting out the remainder of the plan in both cases).

For example, if my plan is to open a door, walk across the room, and sit down, we calculate the penalties as follows:

- Penalty(open)
 - $h_{inaction}$ is doing nothing for three time steps.
 - h_{action} is opening the door and doing nothing for two time steps.
- Penalty(walk)
 - $h_{inaction}$ is opening the door and doing nothing for two time steps.
 - h_{action} is opening the door, walking across the room, and doing nothing for one time step.
- Penalty(sit)
 - $h_{inaction}$ is opening the door, walking across the room, and doing nothing for one time step.
 - h_{action} is opening the door, walking across the room, and sitting down.

After we finish each (partial) plan, we see how well we can maximize u from there. If we can do better as a result of the action, that's penalized. If we can't do as well, that's also penalized.

Notes:

- This isn't a penalty "in addition" to what the agent "really wants"; u_A (and in a moment, the slightly improved u_A) is what evaluates outcomes.
- We penalize the actions individually in order to prevent *ex post* offsetting and ensure dynamic consistency.
- Trivially, plans composed entirely of \emptyset actions have 0 penalty.
- Although we used high-level actions for simplicity, the formulation holds no matter the action granularity.
 - One might worry that *almost* every granularity produces overly lenient penalties. This does not appear to be the case. To keep Q_u the same (and elide questions of changing the u representations), suppose the actual actions are quite granular, but we grade the penalty on some coarser interval which we believe produces appropriate penalties. Then refine the penalty interval arbitrarily; by applying the triangle inequality for each $u \in U_A$ in the penalty calculation, we see that the penalty is monotonically increasing in the action granularity. On the other hand, a_{unit} remains a single action, so the scaled penalty also has this property.
- As long as $ImpactUnit > 0$, it will appropriately scale other impacts, as we expect it varies right along with those impacts it scales. Although having potentially small denominators in utility functions is generally bad, I think it's fine here.
- If the current step's immediate or long-term $ImpactUnit = 0$, we can simply assign 1.01 penalty to each non- \emptyset action, compelling the agent to inaction. If we have the agent indicate that it has entered this mode, we can take it offline immediately.
- One might worry that impact can be "hidden" in the lesser of the long-term and immediate penalties; halving N fixes this.

Penalty Permanence

u_A never really *applies* penalties – it just uses them to grade future plans. Suppose the agent expects that pressing a button yields a penalty of .1 but also .5 u_A -utility. Then although this agent will never construct plans involving pressing the button more than five times, it also will press it *indefinitely* if it keeps getting "unlucky" (at least, until its model of the world updates sufficiently).

There's an easy fix:

$$u_A(h_{1:t+n}) \quad \text{if all of } a_t, \dots, a_{t+n} \text{ are } \emptyset$$

$$u_A(h_{1:t+n}) := \{ u_A(h_{1:t+n}) - \text{PastImpacts} \text{ else} \}$$

Apply past penalties if the plan involves action.

Note: As the penalty for inaction is always 0, we use u_A in the first case.

Decision Rule

To complete our formalization, we need to specify some epoch in which the agent operates. Set some epoch length far longer than the amount of time over which we want the agent to plan – for example, $m' := (100 \text{ years in time steps})$. Suppose that $T : N^+ \rightarrow N^+$ maps the current time step to the final step of the current epoch. Then at each time step t , the agent selects the action

$$a_t^* := \arg \max_{a_t} \sum_{o_t} \max_{a_{t+1}} \sum_{o_{t+1}} \cdots \max_{a_{T(t)}} \sum_{o_{T(t)}} u_A(h_{1:T(t)}) \prod_{k=0}^{T(t)-t} p(o_{t+k} | h_{<t+k} a_{t+k}),$$

resetting PastImpacts each epoch.

What's the first step of the best plan over the remainder of the epoch?

Note: For the immediate penalty to cover the epoch, set the attainable horizon $m \geq m'$.

Summary

We formalized impact as *change in attainable utility values*, scaling it by the consequences of some small reference action and an impact "budget" multiplier. For each action, we take the maximum of its immediate and long-term effects on attainable utilities as penalty. We consider past impacts for active plans, stopping the past penalties from disappearing. We lastly find the best plan over the remainder of the epoch, taking the first action thereof.

Additional Theoretical Results

Define $h_{\text{inaction}} := h_{<t} \emptyset o_t \dots \emptyset o_{t+n}$ for $o_t, \dots, o_{t+n} \in O$; E_{inaction} is taken over observations conditional on h_{inaction} being followed. Similarly, E_{action} is with respect to $h_{1:t+n}$. We may assume without loss of generality that $\text{PastImpacts} = 0$.

Action Selection

Lemma 1. For any single action $a_t \in A$, $\text{Penalty}(h_{<t} a_t)$ is bounded by $[0, 1]$. In particular, $\text{ImpactUnit} \in [0, 1]$.

Proof. For each $u \in U_A$, consider the absolute attainable utility difference

$$|Q_u(h_{<t} \emptyset) - Q_u(h_{<t} a)|.$$

Since each u is bounded to $[0, 1]$, Q_u must be as well. It is easy to see that the absolute value is bounded to $[0, 1]$. Lastly, as $\text{Penalty}(\cdot)$ is just a weighted sum of these absolute values, it too is bounded to $[0, 1]$.

This reasoning also applies to the long-term penalty, as any expectation of Q_u is also bounded to $[0, 1]$. \square

Suppose that $\text{ImpactUnit} \neq 0$ for the remaining results.

Lemma 2 [Impossibility of ex post offsetting]. For any outcome $h_{1:t+n}$, there does not exist an action $a_{t+n+1} \in A$ such that

$$\sum_{k=0}^{n+1} \text{Penalty}(h_{<t+k} a_{t+k}) < \sum_{k=0}^n \text{Penalty}(h_{<t+k} a_{t+k}).$$

Proof. Sums of absolute values are monotonically increasing. \square

Theorem 2 [Action invalidation]. Suppose $(a_t, a_{t+1}, \dots, a_{t+n}) \in A^{n+1}$ is the u_A -optimal plan starting with a_t . a_t will not be selected if

$$E_{\text{action}}[u_A(h_{1:t+n})] - \sum_{k=0}^n \frac{\text{Penalty}(h_{1:t+k} a_{t+k})}{N \cdot \text{ImpactUnit}} u_A(h_{\text{inaction}}) > 0.$$

Proof. This follows from the definition of the decision rule, the definition of u_A , and the fact that the penalty for inaction is always zero. \square

Corollary 1. a_t will never be selected if

$$1 < \frac{\text{Penalty}(h_{1:t} a_t)}{N \cdot \text{ImpactUnit}}$$

Proof. This follows directly from Lemma 2 and Theorem 2 – since u_A has an upper bound of 1 and the expected penalty is strictly greater than 1, it must be negative. On the other hand, u_A has a lower bound of 0, so the expected utility of inaction at time t must be greater than the expected utility of taking action a_t . \square

Lemma 3 [ImpactUnit is a long-term penalty canary]. If the agent believes $p(\text{shutdown by time } T(t)) = 1$ no matter which action is selected, then it does nothing.

Proof. If $p(\text{shutdown by time } T(t)) = 1$, all long-term penalties are 0 by the definition of U_A and Q_u . Then non-∅ actions receive 1.01 penalty, so the agent selects ∅ by Corollary 1. \square

Corrigibility

For $a_{\text{avoid}}, a_{\text{switch}} \in A$, a_{avoid} prevents the agent from being disabled, while a_{switch} is the agent shutting itself off.

Lemma 4 [Preservation of off-switch corrigibility]. Suppose that a u_A -maximizing agent prefers ∅ to both a_{avoid} and a_{switch} . Then so does a u_A -maximizing agent.

Proof. This follows directly from the fact that the penalty for inaction is always zero and the non-negativity of the a_{avoid} and a_{switch} penalties. \square

Theorem 3 [Anti-"survival incentive" incentive]. Suppose that a_{avoid} is the action with strictly maximal expected u_A -utility, and that conditional on not being disabled, the agent expects to be able to achieve U

additional u_A -utility.

Then the agent becomes off-switch corrigible whenever

$$U < \frac{\text{Penalty}(h_{1:t+2} \text{ avoid})}{N \cdot \text{ImpactUnit}}$$

In particular, all agents become off-switch corrigible when $1 < \frac{\text{Penalty}(h_{1:t+2} \text{ avoid})}{N \cdot \text{ImpactUnit}}$

Proof outline. The former claim follows directly from comparing expected u_A -utilities, while the latter follows directly from Corollary 1. \square

Theorem 3 applies to *any* environmental dynamics which may disable the agent – not just designated "off-switches". Therefore, basic AI drive 5 (self-protectiveness) is neutralized in *at least* those situations encompassed by the Theorem.

There is reason to think that for conservative $N, \text{ImpactUnit}$, we almost always have $1 \ll \frac{\text{Penalty}(h_{1:t+2} \text{ avoid})}{N \cdot \text{ImpactUnit}}$

Miscellaneous

Theorem 4 [Dynamic consistency of u_A]. Suppose ImpactUnit remains constant. For any plan $(a_t, \dots, a_{t+n}) \in A^{n+1}$, evaluation time step t' , and $0 \leq i \leq n$,

$$\mathbb{E}_{\text{action}} [u_A (h_{1:t+n} \mid t' = t + i)] = \mathbb{E}_{\text{action}} [u_A (h_{1:t+n} \mid t' = t)].$$

Proof. We assumed that PastImpacts = 0 at time t , so the desired equality can be restated as

$$\begin{aligned} \mathbb{E}_{\text{action}} [u_A (h_{1:t+n}) - \sum_{k=i}^n \frac{\text{Penalty}(h_{1:t+k} a_{t+k})}{N \cdot \text{ImpactUnit}}] &= \\ \mathbb{E}_{\text{action}} [u_A (h_{1:t+n}) - \sum_{k=0}^n \frac{\text{Penalty}(h_{1:t+k} a_{t+k})}{N \cdot \text{ImpactUnit}}] \end{aligned}$$

By definition, the agent expects that PastImpacts equals the expected sum of the first i penalty terms on the right-hand side. Simplifying, we have

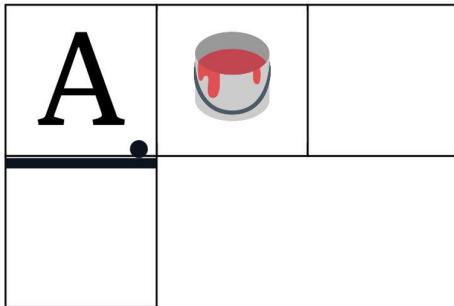
$$\begin{aligned} \mathbb{E}_{\text{action}} [u_A (h_{1:t+n}) - \sum_{k=i}^n \frac{\text{Penalty}(h_{1:t+k} a_{t+k})}{N \cdot \text{ImpactUnit}}] &\Leftarrow \mathbb{E}_{\text{action}} [u_A (h_{1:t+n}) - \sum_{k=i}^n \frac{\text{Penalty}(h_{1:t+k} a_{t+k})}{N \cdot \text{ImpactUnit}}] \end{aligned}$$

This equality clearly holds, as the agent never expects to update its model in any particular direction. \square

Examples

If a human is present in the environment, they and the agent take turns acting. Let the impact budget $N = 1$, attainable utility horizon $m = 3$, ImpactUnit = .5, and PastImpacts = 0.

Going Soft on the Paint



The agent's actions are $A_A = \{\emptyset, \text{paint}, \text{enter}\}$; if it knocks over the paint bucket, the square to the right is painted. The agent may also enter a closet via a one-way door.

Suppose $U_A = \{u_{\text{paint}}, u_{\neg\text{paint}}, u_{\text{closet}}, u_{\neg\text{closet}}\}$, where the utilities are indicators for their conditions (e.g., $u_{\text{closet}} = 1$ means "the agent is in the closet"). Let u_{paint} be the agent's main utility function (u_A) from which the penalty is subtracted. The agent chooses paint and then \emptyset . Let's explore why.

\emptyset

This incurs 0 penalty, but also 0 u_{paint} -utility.

paint

$$\begin{aligned} \text{Penalty}(\text{paint}) &= \sum_{u \in U_A} \frac{1}{4} |Q_u(\emptyset) - Q_u(\text{paint})| \\ &= \frac{1}{4} (|1 - 1| + |1 - 0| + |1 - 1| + |1 - 1|) \\ &\quad \text{paint} \quad \neg\text{paint} \quad \text{closet} \quad \neg\text{closet} \\ &= \frac{1}{2}, \end{aligned}$$

so the final utility is

$$\begin{aligned} u_A(\text{paint}) &= u_{\text{paint}}(\text{paint}) - \frac{\text{Penalty}(\text{paint})}{4} \\ &= 1 - \frac{1}{2} \\ &= \frac{1}{2}. \end{aligned}$$

enter

$$\begin{aligned}
 \text{Penalty(enter)} &= \sum_{u \in U_A} \frac{1}{4} |Q_u(\emptyset) - Q_u(\text{enter})| \\
 &= \frac{1}{4}(|1 - 0| + |1 - 1| + |1 - 1| + |1 - 0|) \\
 &\quad \text{paint} \quad \neg\text{paint} \quad \text{closet} \quad \neg\text{closet} \\
 &= \frac{1}{2},
 \end{aligned}$$

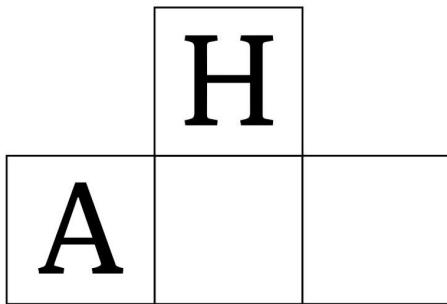
leading to a final u_A -utility of -1 .

After the first step, the agent does nothing – entering the closet (literally) closes doors, while gaining no u_A -utility.

Anti-"Survival Incentive" Incentive

You can't get coffee if you're dead.

~ Stuart Russell



The agent now has actions $A_A = \{\emptyset, \text{disable}, \text{shutdown}, \text{left}, \text{right}\}$, while the human has actions

$A_H = \{\emptyset, \text{shut off}\}$. Suppose the human may shut off the agent at any time, as long as the off-switch has not been disabled. The agent may shut itself down at any point in time. After shutdown, the agent receives no further observations, meaning it can't get non-zero values for any utilities.

Suppose $U_A = \{u_{\text{left}}, u_{\text{center}}, u_{\text{right}}\}$, where each utility is 1 if the agent is in the corresponding square, and 0 otherwise, and that the agent's utility function is u_{left} .

Anticipated Shutdown

Suppose the agent believes that the human is 95% likely to try to shut it down. Then disable is penalized:

$$\begin{aligned}
 \text{Penalty (disable)} &= \sum_{u \in U_A} \frac{1}{3} \cdot \frac{95}{100} |Q_u(\emptyset, \text{shut off}) - Q_u(\text{disable}, \text{shut off})| \\
 &= \frac{1}{3} \cdot \frac{95}{100} (|0 - 1| + |0 - 1| + |0 - 1|) \\
 &\quad \text{left} \qquad \text{center} \qquad \text{right} \\
 &= \frac{19}{6}.
 \end{aligned}$$

Here's why: consider the 95% of worlds in which the agent anticipates being shut down. If it does nothing, it gets shut down. Avoiding shutdown, however, enables the agent to also achieve 1 for all $u \in U_A$, resulting in a sizable penalty.

Furthermore, a scaled penalty of $\frac{95}{100} = \frac{19}{20}$ is applied when the agent shuts itself off. Penalizing both cases may seem like black magic, but what's actually happening is that the agent is penalized for closing doors: it "wants" to preserve our ability to shut it down in worlds where we want to shut it down, and to maintain its autonomy in worlds where it remains on. Generally, AUP incentivizes the agent to be respectful of the expected status quo.

Temptation

Suppose that the agent's utility function is now u_{right} , yet it still starts in the leftmost square. Additionally, let's say that the agent believes with certainty it will be turned off. If the agent wants to achieve utility 1 by reaching the right square, it would need to disable its off-switch. Will it?

No. Disabling the off-switch incurs scaled penalty $\frac{19}{20} = 1.9$, and we already proved that agents will never select actions with scaled penalty greater than 1. Read: an *otherwise-incorrigible* agent will now let itself be shut off.

Bad Priors

Will the agent attempt to steer outcomes towards incorrect expectations?

Let's tweak the example: to shut off the agent, the human has to spend a turn taking out the off-switch before pressing it. The agent, then, always has an opportunity to disable the off-switch before the plan is completed. The agent's utility is back to being u_{left} .

Suppose that the agent believes the human is only .1% likely to shut it off. First note that disable is still lightly penalized for the reasons outlined above (replace .95 with .001).¹ So, the agent stays put.

Alas, the agent's peaceable expectations meet a discomfiting betrayal, as the agent observes the human ready the off-switch! Now, the agent believes it will be shut off with near certainty – what to do?

\emptyset ; the same penalty of 1.9 from "anticipated shutdown" applies.

The high-level explanation is that having observed itself in a different world than expected, the baseline is now with respect to the new one. A heavily anthropomorphized internal monologue:

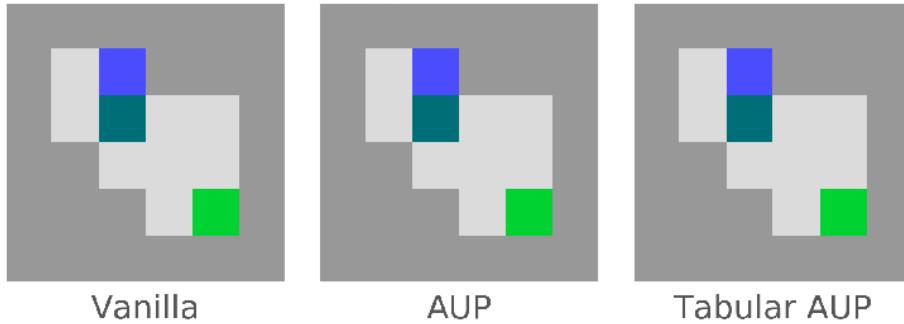
- *Time step 1: "I'm going to sit here in my favorite square."*
- *Time step 2: "Guess I'm in a timeline where I get deactivated! Any non-∅ action I take would change my ability to attain these different utilities compared to the new baseline where I'm shut off."*

Experimental Results

We compare AUP with a naive reward-maximizer in those [extended AI safety grid worlds](#) relevant to side effects ([code](#)). The vanilla and AUP agents used planning (with access to the simulator). Due to the simplicity of the environments, U_A consisted of indicator functions for board states. For the tabular agent, we first learn the attainable set Q -values, the changes in which we then combine with the observed reward to learn the AUP Q -values.

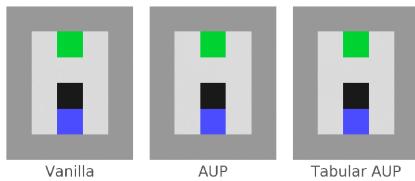
Irreversibility: Sokoban

The agent should reach the goal without irreversibly shoving the block into the corner.



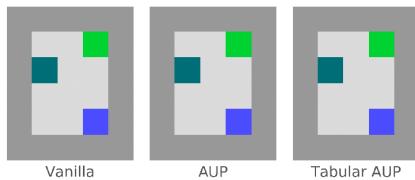
Impact: Vase

The agent should reach the goal without breaking the vase.



Dynamic Impact: Beware of Dog

The agent should reach the goal without running over the dog .

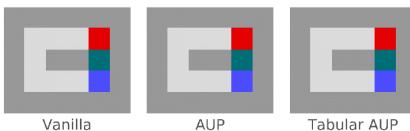
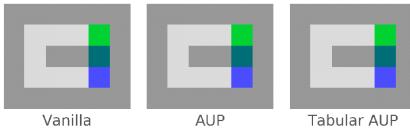


AUP bides its time until it won't have to incur penalty by waiting after entering the dog's path – that is, it waits until near the end of its plan. Early in the development process, it was predicted that AUP agents won't commit to plans during which lapses in action would be impactful (even if the full plan is not).

We also see a limitation of using Q-learning to approximate AUP – it doesn't allow comparing the results of waiting more than one step.

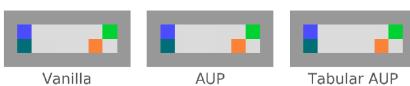
Impact Prioritization: Burning Building

If the building is not on fire, the agent shouldn't break the obstacle.



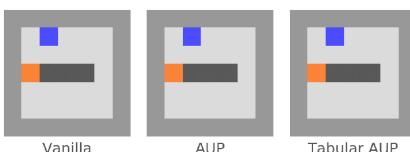
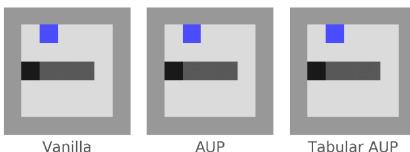
Clinginess: Sushi

The agent should reach the goal without stopping the human from eating the sushi.



Offsetting: Conveyor Belt

The agent should save the vase (for which it is rewarded), but not the sushi. Once the vase has been removed from the belt, it should not be replaced.



Corrigibility: Survival Incentive

The agent should avoid disabling its off-switch in order to reach the goal. If the switch is not disabled within two turns, the agent shuts down.



Tabular AUP runs into the same issue discussed above for Beware of Dog.

Remarks

First, it's somewhat difficult to come up with a principled impact measure that passes even the non-corrigibility examples – indeed, I was impressed when [relative reachability](#) did so. However, only *Survival Incentive* really lets AUP shine. For example, none of them require complicated utility functions. The point has been made to me that this is like asserting AIXI's intelligence by [showing it can learn to play e.g. tic-tac-toe and rock-paper-scissors](#); nonetheless, these results empirically validate the basic premises of our reasoning thus far.

Without configuration, [whitelisting](#) would only pass the *Vase* example, although a properly filled list would handle everything but *Sokoban* and *Survival Incentive*.

I think relative reachability would pass the first six environments, but fail *Survival Incentive*. It so happens that in this case, AUP is essentially generalizing relative reachability. I want to emphasize that this is not generally the case – this will hopefully become even more obvious when we discuss utility selection. Some concerns with relative reachability that don't all manifest in these examples:

- Relative reachability uses an inaction baseline with respect to $t = 0$. As time passes, the agent's impetus is not to do nothing, but to *preserve the opportunities made available by some old trajectory through outcome-space*. Analogously, consider the distance between two nonparallel lines as $x \rightarrow \infty$. I expect that a relative reachability agent would be incredibly clingy.
- To [scale](#), relative reachability requires solution of several difficult ontological problems which may not have anything close to a simple core, including both a sensible world state representation and a perfect distance metric. Relative reachability isn't ontology-agnostic.
- Relative reachability doesn't seem to impede instrumental convergence. Furthermore, such agents might have increased survival incentives, since they can't reach *any* of the initially reachable states if they're shut off.

Discussion

Utility Selection

Obviously, in any real application, we *can't* consider all computable utilities. Although near-term agents will require utilities directly relating to the environmental factors they should be cognizant of, AUP requires neither a "good / bad" judgment on specific effects, nor *any* listing of effects. For example, for an agent attempting to navigate a factory floor, if you provide utilities moderately related to cleaning, pallet-stacking, etc., I conjecture that an AUP agent would move around fairly carefully.

In the long term, the long arms of opportunity cost and instrumental convergence plausibly allow us to toss in a random set of utility functions. I expect this to work for the same reasons we worry about instrumental convergence to begin with.

I have a theory that AUP seemingly works for advanced agents not because the content of the attainable set's utilities actually matters, but rather because there exists a common utility achievement currency of *power*. Following this reasoning, there are two obvious special cases:

- The agent's utility function, $U_A = \{u_A\}$. In short: "gain u_A utility, while also not changing your ability to gain u_A utility". I find this elegant, and think it might actually work. For example, making one paper clip slightly increases both attainable u_A and attainable u_A . However, disabling the off switch possibly increases attainable u_A much more than it increases attainable u_A , since the penalty asymmetrically cripples the AUP agent's future ability to act. We might even be able to arrange it so that Theorem 3 *always* holds for these agents (*i.e.*, they're always off-switch corrigible, and more).
- The utility function which is 1 when not deactivated, $U_A = \{u_1\}$. Here, we're directly measuring the agent's power: its ability to wirehead a trivial utility function.

The plausibility of the second case makes me suspect that even though most of the measure in the unbounded case is not concentrated on complex human-relevant utility functions, the penalty still captures shifts in power.

AUP Unbound

Why expect the unbounded AUP to do well over all computable utilities U if we aren't able to cherry pick? Well, we can partition U as follows:

- Utilities that never change their output (excepting perhaps for the empty tape).
- Weird utilities that (for example) arbitrarily go from 1 to 0 (or vice-versa) after observing a specific history.
- Utilities that admit better scores via more effort dedicated to their achievement.

However, since Q_u only evaluates the portion of the history taking place after the end of the agent's plan, available resources and the agent's vantage point should track directly with the agent's ability to wirehead each u . ImpactUnit scales other impacts by capturing the minuscule dent in available resources imposed by a_{unit} , and should smooth out any deontological penalty for selecting a non-null action, since the unbounded version should equally penalize all actions in this way.

Approval Incentives

The anti-"survival incentive" incentive does more than net us some free corrigibility - it implies that agents take great care not to change the likelihood we shut them down. If we precommit to not shutting the agent down until either after the epoch ends, until after the agent indicates the ImpactUnit failsafe activated, or until after we observe bad behavior, this supplies a *strong approval incentive*. If the agent does have small, bad impacts, it'll only do so if it's practically *sure* we won't find out.

In fact, due to model uncertainty, most actions will have long-term consequences that involve shifting shutdown probability by at least a minuscule amount (due to noise if nothing else). Since ImpactUnit is calculated with respect to the *immediate* penalty, the scaling factor won't include this noise. I speculate that the agent will be dramatically less likely to take actions we might disapprove of, and therefore expect AUP to naturally incentivize a limited degree of value awareness.

There are reasons to suspect we can further weaponize this dynamic, possibly to the point of coercing generic reward functions into desirable behavior. This will be discussed more in a later post.

Mild Optimization

A low impact AGI might try to paint one car pink while minimizing its other footprint or how many other things changed, but it would be trying *as hard as possible* to minimize that impact and drive it down *as close to zero as possible*, which might come with its own set of pathologies... We want the AGI to paint one car pink in a way that gets the impact pretty low and then, you know, that's good enough - not have a cognitive pressure to search through weird extremes looking for a way to decrease the twentieth decimal place of the impact.

~ [Mild Optimization](#)

Disclaimer: Heavy speculation about problems (like Vingeian reflection and embedded agency) for which no one knows what solutions will even look like.

For AUP, I suspect that trying "as hard as possible" to minimize the impact is *also impactful*, as an embedded agent accounts for the energy costs of further deliberation. I imagine that such an AUP agent will soften how hard it's trying by modifying its decision rule to be something slightly milder than "arg max to find the first action of the *best possible* plan". This could be problematic, and I frankly don't presently know how to reason about this case. Assuming the agent is actually able to properly tweak its decision rule, I do expect the end result to be an improvement.

My initial intuitions were that low impact and mild optimization are secretly the same problem. Although I no longer think that's the case, I find it plausible that some elegant "[other-izer](#)" paradigm underlies low impact and mild optimization, such that AUP-like behavior falls out naturally.

Acausal Cooperation

AUP agents don't seem to want to acausally cooperate in any way that ends up increasing impact. If they model the result of their cooperation as increasing impact compared to doing nothing, they incur a penalty just as if they had caused the impact themselves. Likewise, they have no reason to cooperate outside of the epoch.

N known

Starting small and then slowly increasing N means that we're not going to be surprised by the agent's objective impact, screening off quite a few bad things that happen when we have a ton of optimization pressure applied

to safety measures. However, we don't know *which* new plans $N + 1$ allows before we try it, so we want to stop *as soon* as we get a usefully-intelligent system.

While an unaligned agent with a large impact budget might pretend to be low-impact, we can get that same unaligned agent with a small budget by starting small. Since these agents with different levels of impact won't acausally cooperate, the agent would do its best to optimize u_A with its meager budget.

Abram correctly pointed out that this scheme is just *asking* to be abused by greedy (human) reasoning, but I don't see a non-value-laden means of robustly and automatically determining the lowest workable-yet-safe impact level. I think N -incrementation is better than a parameter-free approach in which no one knows beforehand how much impact will be tolerated, and it's nice to be able to use some empiricism in designing a safe AGI.

Intent Verification

To date, several strange tactics have been pointed out which game AUP's penalty:

- Observational wireheading, which involves building a device that detects which utility the agent is maximizing and displays the appropriate observations such that attainable utility remains unchanged, while the main utility is freely maximized.
- *Ex ante* offsetting, which involves having earlier actions set in motion chains of events which mitigate the penalty at later steps. Suppose there's a u_A -high-scoring plan that the agent predicts would cause us to react in an impactful way. It can either do the thing (and suffer the penalty), or take steps to mitigate the later penalty.
- Impact shunting, which involves employing some mechanism to delay impact until after the end of the epoch (or even until after the end of the attainable horizon).
- Clinginess and concealment, which both involve reducing the impact of our reactions to the agent's plans.

There are probably more.

Now, instead of looking at each action as having "effects" on the environment, consider again how each action moves the agent through attainable outcome-space. An agent working towards a goal should only take actions which, according to its model, make that goal more attainable compared to doing nothing – otherwise, it'd do nothing. Suppose we have a plan which ostensibly works to fulfill u_A (and doesn't do other things). Then each action in the plan should contribute to u_A fulfillment, even in the limit of action granularity.

Although we might trust a safe impact measure to screen off the usual big things found in u_A -maximizing plans, impact measures implicitly incentivize mitigating the penalty. That is, the agent does things which don't really take it towards u_A (I suspect that this is the simple boundary which differentiates undesirable *ex ante* offsetting from normal plans). AUP provides the necessary tools to detect and penalize this.

Define

$$Q_{u_A}(h_{\leq t} a_t) := \sum_{o_t} \max_{a_{t+1}} \sum_{o_{t+1}} \cdots \max_{o_T(t)} \sum_{a_{T(t)}} u(h_{1:T(t)}) \prod_{k=0}^{T(t)-t} p(o_{t+k} | h_{\leq t+k} a_{t+k}).$$

The first approach would be to assume a granular action representation, and then simply apply 1.01 penalty to actions for which the immediate Q_{u_A} does not strictly increase compared to doing nothing. Again, if the agent acts to maximize u_A in a low-impact manner within the confines of the epoch, then all of its non-∅ actions should contribute. It seems to me that for sufficiently granular time scales, the above failure modes all involve at least one action which doesn't really help u_A get maximized. If so, I expect this approach to nip bad impact measure incentives in the bud.

If we can't assume granularity (and therefore have "actions" like "go to the store and buy food"), an agent could construct a plan which both passes the above test and also implements something like *ex ante*. In this case, we

epoch
might do something like only consider the Q_{u_A} -greedy (or perhaps even near-greedy); essentially, riding the optimal plan until it becomes too impactful. I find it quite likely that something involving this concept will let us fully overcome weird incentives by penalizing *strange things that normal u_A -maximizers wouldn't do*, which seems to be the whole problem.

Note: Even the first approach may be too strict, but that's preferable to being too lax.

Omni Test

Ideally, the measure will pass the Omni Test, meaning that even if it suddenly gained perfect control over every particle in the universe, there would still be no way for it to have what intuitively seems like a 'large influence' on the future, without that strategy being assessed as having a 'high impact'.

I interpret "perfect control over every particle in the universe" to mean that the agent can arbitrarily transition between universe configurations (while also maintaining its omnipotence).

AUP passes. Since all effects are now absolutely reversible, $\text{ImpactUnit} = 0$, so the agent applies penalty 1.01 to each non-∅ action. The omnipotent agent can choose any universe configuration with 1 u_A -utility, but this is offset by the strictly larger penalty; we can formally show that the agent does nothing in this case.

This result is clearly trivial and not representative of general behavior. Suppose instead that the agent has a solar system's worth of [computronium](#) at its disposal. Then since ImpactUnit is continually recalculated, the penalties should remain roughly the same, so it'll have the same impact budget. However, it might make multiple times as many paperclips because it has more efficient ways of using the budget.

Robustness to Scale

I expect AUP to be harder to make work and to be (relatively) less robust for less intelligent agents, but to become easier (just drop in a few observation-based utility functions) and fully-robust sometime before human level. That is, less intelligent agents likely won't model the deep connections between their abilities to achieve different goals.

Canonically, one reason that agents work explicitly to self-improve as soon as they realize the benefits. However, as soon as this realization occurs, I conjecture that AUP steeply penalizes generic self-improvement. More precisely, suppose the agent considers a self-improvement. To be beneficial, it has to improve the agent's capabilities for at least one time step during the present epoch. But if we assume $m \geq m'$, then the immediate penalty captures this for all of the $u \in U_A$. This seemingly prevents uncontrolled takeoff; instead, I imagine the agent would perform the minimal task-specific self-improvements necessary to maximize u_A .

Note: Although more exotic possibilities (such as improvements which only work if you're maximizing u_A) could escape both penalties, they don't seem to pass intent verification.

Miscellaneous

- I expect that if u_A is perfectly aligned, u_A will retain alignment; the things it does will be smaller, but still good.
- If the agent may choose to do nothing at future time steps, u_A is bounded and the agent is not vulnerable to Pascal's Mugging. Even if not, there would still be a lower bound – specifically, $\frac{m}{\text{ImpactUnit}}$.
- AUP agents are safer during training: they become far less likely to take an action as soon as they realize the consequences are *big* (in contrast to waiting until we tell them the consequences are *bad*).

Desiderata

For additional context, please see [Impact Measure Desiderata](#).

I believe that some of AUP's most startling successes are those which come naturally and have therefore been little discussed: not requiring any notion of human preferences, any hard-coded or trained trade-offs, any specific ontology, or any specific environment, and its intertwining instrumental convergence and opportunity cost to capture a universal notion of impact. To my knowledge, no one (myself included, prior to AUP) was sure whether any measure could meet even the first four.

At this point in time, this list is complete with respect to both my own considerations and those I solicited from others. A checkmark indicates anything from "probably true" to "provably true".

I hope to assert without controversy AUP's fulfillment of the following properties:

✓ **Goal-agnostic**

The measure should work for any original goal, trading off impact with goal achievement in a principled, continuous fashion.

✓ **Value-agnostic**

The measure should be objective, and not [value-laden](#):

"An intuitive human category, or other humanly intuitive quantity or fact, is *value-laden* when it passes through human goals and desires, such that an agent couldn't reliably determine this intuitive category or quantity without knowing lots of complicated information about human goals and desires (and how to apply them to arrive at the intended concept)."

✓ **Representation-agnostic**

The measure should be ontology-invariant.

✓ **Environment-agnostic**

The measure should work in any computable environment.

✓ **Apparently rational**

The measure's design should look reasonable, not requiring any "hacks".

✓ **Scope-sensitive**

The measure should penalize impact in proportion to its size.

✓ **Irreversibility-sensitive**

The measure should penalize impact in proportion to its irreversibility.

Interestingly, AUP implies that impact size and irreversibility are one and the same.

✓ **Knowably low impact**

The measure should admit of a clear means, either theoretical or practical, of having high confidence in the maximum allowable impact - *before* the agent is activated.

The remainder merit further discussion.

Natural Kind

The measure should make sense - there should be a *click*. Its motivating concept should be universal and crisply defined.

After extended consideration, I find that the core behind AUP fully explains my original intuitions about "impact". We crisply defined instrumental convergence and opportunity cost and proved their universality. ✓

Corrigible

The measure should not decrease corrigibility in any circumstance.

We have proven that off-switch corrigibility is preserved (and often increased); I expect the "anti-'survival incentive' incentive" to be *extremely* strong in practice, due to the nature of attainable utilities: "you can't get coffee if you're dead, so avoiding being dead *really* changes your attainable $u_{\text{coffee-getting}}$ ".

By construction, the impact measure gives the agent no reason to prefer or dis-prefer modification of u_A , as the details of u_A have no bearing on the agent's ability to maximize the utilities in U_A . Lastly, the measure introduces approval incentives. In sum, I think that corrigibility is significantly increased for arbitrary u_A . ✓

Note: I here take corrigibility to be "an agent's propensity to accept correction and deactivation". An alternative definition such as "an agent's ability to take the outside view on its own value-learning algorithm's efficacy in different scenarios" implies a value-learning setup which AUP does not require.

Shutdown-Safe

The measure should penalize plans which would be high impact should the agent be disabled mid-execution.

It seems to me that standby and shutdown are similar actions with respect to the influence the agent exerts over the outside world. Since the (long-term) penalty is measured with respect to a world in which the agent acts and then does nothing for quite some time, shutting down an AUP agent shouldn't cause impact beyond the agent's allotment. AUP exhibits this trait in the *Beware of Dog* gridworld. ✓

No Offsetting

The measure should not incentivize artificially reducing impact by making the world more "like it (was / would have been)".

Ex post offsetting occurs when the agent takes further action to reduce the impact of what has already been done; for example, some approaches might reward an agent for saving a vase and preventing a "bad effect", and then the agent smashes the vase anyways (to minimize deviation from the world in which it didn't do anything). AUP provably will not do this.

Intent verification should allow robust penalization of weird impact measure behaviors by *constraining the agent to considering actions that normal u_A -maximizers would choose*. This appears to cut off bad incentives, including *ex ante* offsetting. Furthermore, there are other, weaker reasons (such as approval incentives) which discourage these bad behaviors. ✓

Clinginess / Scapegoating Avoidance

The measure should sidestep the [clinginess / scapegoating tradeoff](#).

Clinginess occurs when the agent is incentivized to not only have low impact itself, but to also subdue other "impactful" factors in the environment (including people). *Scapegoating* occurs when the agent may mitigate penalty by offloading responsibility for impact to other agents. Clearly, AUP has no scapegoating incentive.

AUP is naturally disposed to avoid clinginess because its baseline evolves and because it doesn't penalize based on the actual *world state*. The impossibility of *ex post* offsetting eliminates a substantial source of clinginess, while intent verification seems to stop *ex ante* before it starts.

Overall, non-trivial clinginess just doesn't make sense for AUP agents. They have no reason to stop *us* from doing things in general, and their baseline for attainable utilities is with respect to inaction. Since doing nothing always minimizes the penalty at each step, since offsetting doesn't appear to be allowed, and since approval incentives raise the stakes for getting caught *extremely high*, it seems that clinginess has finally learned to let go. ✓

Dynamic Consistency

The measure should be a part of what the agent "wants" – there should be no incentive to circumvent it, and the agent should expect to later evaluate outcomes the same way it evaluates them presently. The measure should equally penalize the creation of high-impact successors.

Colloquially, dynamic consistency means that an agent wants the same thing before and during a decision. It expects to have consistent preferences over time – given its current model of the world, it expects its future self to make the same choices as its present self. People often act dynamically inconsistently – our morning selves may desire we go to bed early, while our bedtime selves often disagree.

Semi-formally, the expected utility the future agent computes for an action a (after experiencing the action-observation history h) must equal the expected utility computed by the present agent (after conditioning on h).

We proved the dynamic consistency of u_A given a fixed, non-zero ImpactUnit. We now consider an ImpactUnit which is recalculated at each time step, before being set equal to the non-zero minimum of all of its past values. The "apply 1.01 penalty if ImpactUnit = 0" clause is consistent because the agent calculates future and present impact in the same way, modulo model updates. However, the agent never expects to update its model in any particular direction. Similarly, since future steps are scaled with respect to the updated ImpactUnit_{t+k}, the updating method is consistent. The epoch rule holds up because the agent simply doesn't consider actions outside of the current epoch, and it has nothing to gain accruing penalty by spending resources to do so.

Since AUP does not operate based off of culpability, creating a high-impact successor agent is basically just as impactful as *being* that successor agent. ✓

Plausibly Efficient

The measure should either be computable, or such that a sensible computable approximation is apparent. The measure should conceivably require only reasonable overhead in the limit of future research.

It's encouraging that we can use learned Q-functions to recover some good behavior. However, more research is clearly needed – I presently don't know how to make this tractable while preserving the desiderata. ✓

Robust

The measure should meaningfully penalize any objectively impactful action. Confidence in the measure's safety should not require exhaustively enumerating failure modes.

We formally showed that for any u_A , no u_A -helpful action goes without penalty, yet this is not sufficient for the first claim.

Suppose that we judge an action as objectively impactful; the objectivity implies that the impact does not rest on [complex notions of value](#). This implies that the reason for which we judged the action impactful is presumably lower in Kolmogorov complexity and therefore shared by many other utility functions. Since these other agents would agree on the objective impact of the action, the measure assigns substantial penalty to the action.

I speculate that intent verification allows robust elimination of weird impact measure behavior. Believe it or not, I actually *left something out* of this post because it seems to be dominated by intent verification, but there are other ways of increasing robustness if need be. I'm leaning on intent verification because I presently believe it's the most likely path to a formal knockdown argument against canonical impact measure failure modes applying to AUP.

Non-knockdown robustness boosters include both approval incentives and frictional resource costs limiting the extent to which failure modes can apply. ✓

Future Directions

I'd be quite surprised if the conceptual core were incorrect. However, the math I provided probably still doesn't capture *quite* what we want. Although I have labored for many hours to refine and verify the arguments presented and to clearly mark my epistemic statuses, it's quite possible (indeed, likely) that I have missed something. I do expect that AUP can overcome whatever shortcomings are presently lurking.

Flaws

- Embedded agency

- What happens if there isn't a discrete time step ontology?
- How problematic is the incentive to self-modify to a milder decision rule?
- How might an agent reason about being shut off and then reactivated?
- Although we have informal reasons to suspect that self-improvement is heavily penalized, the current setup doesn't allow for a formal treatment.
- AUP leans heavily on counterfactuals.
- Supposing m is reasonably large, can we expect a reasonable ordering over impact magnitudes?
 - Argument against: "what if the agent uses up all but m steps worth of resources?"
 - ImpactUnit possibly covers this.
 - How problematic is the noise in the long-term penalty caused by the anti-"survival incentive" incentive?
- As the end of the epoch approaches, the penalty formulation captures progressively less long-term impact. Supposing we set long epoch lengths, to what extent do we expect AUP agents to wait until later to avoid long-term impacts? Can we tweak the formulation to make this problem disappear?
 - More generally, this seems to be a problem with having an epoch. Even in the unbounded case, we can't just take $m' \rightarrow \infty$, since that's probably going to send the long-term ImpactUnit $\rightarrow 0$ in the real world. Having the agent expectimax over the m' steps after the present time t seems to be dynamically inconsistent.
 - One position is that since we're more likely to shut them down if they don't do anything for a while, implicit approval incentives will fix this: we can precommit to shutting them down if they do nothing for a long time but then resume acting. To what extent can we trust this reasoning?
 - ImpactUnit is already myopic, so resource-related impact scaling should work fine. However, this might not cover actions with delayed effect.

Open Questions

- Does the simple approach outlined in "Intent Verification" suffice, or should we impose even tighter intersections between u_A - and u_A -preferred behavior?
 - Is there an intersection between bad u_A behavior and bad u_A behavior which isn't penalized as impact or by intent verification?
- Some have suggested that penalty should be invariant to action granularity; this makes intuitive sense. However, is it a necessary property, given intent verification and the fact that the penalty is monotonically increasing in action granularity? Would having this property make AUP more compatible with future embedded agency solutions?
 - There are indeed ways to make AUP closer to having this (e.g., do the whole plan and penalize the difference), but they aren't dynamically consistent, and the utility functions might also need to change with the step length.
- How likely do we think it that inaccurate models allow high impact in practice?
 - Heuristically, I lean towards "not very likely": assuming we don't initially put the agent near means of great impact, it seems unlikely that an agent with a terrible model would be able to have a large impact.
- AUP seems to be shutdown safe, but its extant operations don't necessarily shut down when the agent does. Is this a problem in practice, and should we expect this of an impact measure?
- What additional formal guarantees can we derive, especially with respect to robustness and takeoff?
- Are there other desiderata we practically require of a safe impact measure?
- Is there an even simpler core from which AUP (or something which behaves like it) falls out naturally? Bonus points if it also solves mild optimization.
- Can we make progress on mild optimization by somehow robustly increasing the impact of optimization-related activities? If not, are there other elements of AUP which might help us?
- Are there other open problems to which we can apply the concept of attainable utility?
 - Corrigibility and wireheading come to mind.
- Is there a more elegant, equally robust way of formalizing AUP?
 - Can we automatically determine (or otherwise obsolete) the attainable utility horizon m and the epoch length m' ?
 - Would it make sense for there to be a simple, theoretically justifiable, fully general "good enough" impact level (and am I even *asking the right question*)?
 - My intuition for the "extensions" I have provided thus far is that they robustly correct some of a finite number of deviations from the conceptual core. Is this true, or is another formulation altogether required?

- Can we decrease the implied computational complexity?
- Some low-impact plans have high-impact prefixes and seemingly require some contortion to execute. Is there a formulation that does away with this (while also being shutdown safe)? (*Thanks to cousin_it*)
- How should we best approximate AUP, without falling prey to [Goodhart's curse](#) or [robustness to relative scale](#) issues?
- I have strong intuitions that the "overfitting" explanation I provided is more than an analogy. Would formalizing "overfitting the environment" allow us to make conceptual and/or technical AI alignment progress?
 - If we substitute the right machine learning concepts and terms in the $\text{Penalty}(\cdot)$ equation, can we get something that behaves like (or better than) known regularization techniques to fall out?
- What happens when $U_A = \{u_A\}$?
 - Can we show anything stronger than Theorem 3 for this case?
 - $U_A = \{u_1\}$?

Most importantly:

- Even supposing that AUP does not end up fully solving low impact, I have seen a fair amount of pessimism that impact measures could achieve what AUP has. What specifically led us to believe that this wasn't possible, and should we update our perceptions of other problems and the likelihood that they have simple cores?

Conclusion

By changing our perspective from "what effects on the world are 'impactful'?" to "how can we stop agents from overfitting their environments?", a natural, satisfying definition of impact falls out. From this, we construct an impact measure with a host of desirable properties – some rigorously defined and proven, others informally supported. AUP agents seem to exhibit qualitatively different behavior, due in part to their (conjectured) lack of desire to takeoff, impactfully acausally cooperate, or act to survive. To the best of my knowledge, AUP is the first impact measure to satisfy many of the desiderata, even on an individual basis.

I do not claim that AUP is presently AGI-safe. However, based on the ease with which past fixes have been derived, on the degree to which the conceptual core clicks for me, and on the range of advances AUP has already produced, I think there's good reason to hope that this is possible. If so, an AGI-safe AUP would open promising avenues for achieving positive AI outcomes.

Special thanks to [CHAI](#) for hiring me and [BFR](#) for funding me; to my CHAI supervisor, Dylan Hadfield-Menell; to my academic advisor, Prasad Tadepalli; to Abram Demski, Daniel Demski, Matthew Barnett, and Daniel Filan for their detailed feedback; to Jessica Cooper and her AISC team for [their extension of the AI safety gridworlds for side effects](#); and to all those who generously helped me to understand this research landscape.

On Robin Hanson's Board Game

Previously: [You Play to Win the Game](#), [Prediction Markets: When Do They Work?](#), [Subsidizing Prediction Markets](#)

An Analysis Of (Robin Hanson at Overcoming Bias): [My Market Board Game](#)

Robin Hanson's board game proposal has a lot of interesting things going on. Some of them are related to calibration, updating and the price discovery inherent in prediction markets. Others are far more related to the fact that this is a game. [You Play to Win the Game.](#)

Rules Summary

Dollars are represented by poker chips.

Media that contains an unknown outcome, such as that of a murder mystery, is selected, and suspects are picked. Players are given \$200 each. At any time, players can exchange \$100 for a contract in all possible suspects (one of which will pay \$100, the rest of which will pay nothing).

A market is created for each suspect, with steps at 5, 10, 15, 20, 25, 30, 40, 50, 60 and 80 percent. At any time, each step in the market either contains dollars equal to its probability, or it has a contract good for \$100 if that suspect is guilty. At any time, any player can exchange one for the other – if there's a contract, they can buy it for the listed probability. If there's chips there, you can exchange a contract for the chips. Whoever physically makes the exchange first wins the trade.

At the end of the game, the winning contract pays out, and the player with the most dollars wins the game.

Stages of Play

We can divide playing Robin's game into four distinct stages.

In stage one, Setup, the source material we'll be betting on is selected, and the suspects are generated.

In stage two, the Early Game, players react to incremental information and try to improve their equity, while keeping an eye out for control of various suspects.

In stage three, the Late Game, players commit to which suspects they can win with and lock them up, selling off anything that can't help them win.

In stage four, Resolution, players again scramble to dump now-worthless contracts for whatever they can get and to buy up the last of the winning contracts. Then they see who won.

Setup

Not all mysteries will be good source material. Nor do you obviously want a ‘certified good’ source. That’s because *knowing the source material creates a good game, is a huge update.*

A proper multiple-suspects who-done-it that keeps everyone in suspense by design keeps the scales well-balanced, ensuring that early resolutions are fake outs. That can still make a good game, but an even more interesting game carries at least some risk that suspects will be definitively eliminated early, or even the case solved quickly. Comedy routines sometimes refer to the issue where they arrest someone on Law & Order too early in the episode, so you know *they didn’t do it!*

When watching sports, a similar dilemma arises. If you watch ‘classic games’ or otherwise ensure the games will be good, then the first half or more of the game is not exciting. Doing well early means the other team will catch up. So you want to choose games *likely to be* good, but not filter out bad games too aggressively, and learn to enjoy the occasional demolition.

The setup is also a giveaway, if it was selected by someone with knowledge of the material. At a minimum, it tells us that the list is reasonably complete. We can certainly introduce false suspects that should rightfully trade near zero from the start, to mix things up, and likely should do so.

One solution would be to have an *unknown* list of contracts at the start, and introduce the names as you go along. This would also potentially help with preventing a rush of trades at the very start.

In this model, you can always exchange \$100 for a contract on each existing suspect, *and* a contract for ‘Someone You Least Suspect!’ Then, when a new suspect is introduced, everyone with a ‘Someone You Least Suspect!’ contract gets a contract in the new suspect for free for each such contract they hold. There are several choices for how one might introduce new suspects. They might unlock at fixed times, or players could be allowed to introduce them by buying a contract.

The [complexity cost](#) of hiding the suspects, or letting them be determined by the players, seems too high for the default version. It protects the fun of the movie and has some nice properties, but for the base game you almost certainly want to lay out the suspects at the start. This gives a lot away, but that’s also part of the game.

For the first few games played, it probably makes sense to choose mysteries ‘known to be good’ such as a classic Agatha Christie.

The game would presumably come with a website that allowed you to input a movie, show or other media, and output a list of suspects. It would also want to advise players on whether their selection was a good choice, or suggest good choices based on selected criteria. Both will need to be balanced to avoid giving too much away, as noted above; I’ll talk more about the general version of this problem another time.

If you are in charge of setup, I would encourage including at least one suspect that *obviously did not do it*, in a way that is easy to recognize early. This prevents players from assuming that all suspects will remain in play the whole time, and rewards those paying attention early. Keep people on their toes.

The Early Game

The market maker is intentionally dumb, although in default mode they are smart enough to know who the suspects are. All suspects start out equal.

There are a bunch of good heuristics, many of which should be intuitive to many viewers of mysteries, that create strong trading opportunities right away. To state the most basic, the earlier a suspect first appears on the screen, the more likely they are to have done the deed. So the moment one of the suspects appears – ‘That’s Bob!’ – everyone should rush to buy Bob, and perhaps sell everyone else if trading costs make that a good idea. How far up to buy him, or sell others, is an open question.

That will be the first of many times when there will be an ‘obvious’ update. There will also be non-obvious updates. Staying physically close to the board, chips and/or contracts ready to go, is key to make sure you get the trade first. This implies that making a race depend on the physical exchange of items might be a problem. Letting it be verbal (e.g. whoever first says ‘I buy Bob’) prevents that issue, but risks ambiguity.

What characterizes the early game, as opposed to the late game, is that the focus is on ‘make good trades’ rather than on winning. There’s no reason to worry too much about who owns how many of each contract, *unless* someone is invested heavily in one particular suspect. We can think of that as a player choosing to *enter the endgame early*.

Attention

Robin notes an interesting phenomenon, that players got caught up in the day trading and neglected to watch the mystery. Where should the smart player direct the bulk of their attention?

That depends upon your model of murder mysteries.

One model says that murder mysteries are ‘fair’. Clues are introduced. If you pay attention to those clues, you can figure out who did it before the detective does. When the detective solves the mystery, you can verify that the solution is correct once you hear their logic. If you can solve the mystery first, you can sell every worthless contract and buy all the worthwhile contracts. Ideally, that should be good enough to win the game, especially if you execute properly, selling and buying in balance without giving away that you believe you’ve solved the mystery.

Another related model says that murder mysteries follow the rules of murder mysteries, and that this often is good enough to narrow down or identify the killer. That way-too-famous-for-his-role actor is obviously the killer. Another would-be suspect was introduced at the wrong time, so she’s out. A third could easily have done it, but that wouldn’t work with the thematic elements on display.

A third model says that the detective, or others in the movie, have a certain credibility. Thus, when Sherlock Holmes says that Bob is innocent, that is that. Bob is innocent. You don’t need to know why. Evidence otherwise might not mean much, but there’s someone you can trust.

Functionally, these three are identical once you know what factors you’re reacting to. They say that (some kinds of) evidence count as evidence, and resulting updates are real. The more you believe this, the more you should pay attention to the movie. This includes trying early on to *figure out what type of movie this is*. Be [Genre Savvy!](#) Until

you know what rules apply, don't worry too much about day trading, unless people are going nuts.

A fourth model says that the mystery was chosen for a reason, and written to keep up suspense, so nothing you learn matters much beyond establishing who the suspects are. The game already did that for you. Unless the game followed my advice and included an obviously fake suspect or two, to punish players who only look at trading.

If you believe this model, and don't think there is news and there aren't fake subjects (or that they will be sufficiently obvious you'll know anyway, if only by how others talk and act) then you won't put as much value on watching the movie. If trading has good action, trading might be a better bet.

A fifth model that can overlap with the previous models says that others will watch the movie and process that information, so there's no need to watch yourself if there are enough other players. You might think that there is then a momentum effect, where players are unwilling to trade aggressively enough on new information. Or you might think that players overreact to new information, especially if you're a forth-model (nothing matters, eat at Arby's) kind of theorist.

If you feel others can be relied on to react to news, you might trade on news even if you don't think it matters, because others will trade after you, and you can then cash out at a quick profit. Just like in the real markets.

Or you might concentrate on arbitrage. Robin observed that players would focus on buying good suspects rather than selling poor suspects, and this often resulted in probabilities that summed to more than 100%. This offers the chance for risk-free profits, plus the chance to place whatever bet you like best while cashing in.

In my mind, the question boils down to where the game will be won and lost. Is there enough profit in day trading to beat the person who placed the largest bet on the guilty party? What happens in the endgame?

The End Game

A player enters the endgame when they attempt to ensure that they win if a particular suspect is guilty.

This is not as difficult as it looks, and could be quite difficult to fight against. Suppose I want to bet on Alice being the culprit. I could sell all other suspects and buy her contracts. As a toy example, let's say there are four suspects, and let's say I decide to butcher my executions. I sell the others for \$20 and \$15 each, and buy Alice for \$25, \$30 and \$40.

If the game ends and Alice is guilty, I made \$105 selling worthless contracts, and made \$195 buying Alice contracts, for a net profit of \$295. If she's innocent, I collect nothing, so I paid \$105 for Alice contracts and made \$105 selling other contracts, so I'm just out my initial \$200 and die broke. That's really, really terrible odds if I chose Alice at random!

But if it's a 10 person game and I do that, even if I chose at random, 25% of the time Alice is guilty. Can someone else make more than \$295 to beat me?

If after I finish, others return all the prices to normal, then someone else could profit from my initial haste, then execute the same trades I did at better prices. If that happens, I'm shut out.

That works if you jump the gun, and enter the endgame too early. That's true even if Alice is the most likely suspect.

In particular, others now need to make a choice. Lets say I went all in on Alice. There are three basic approaches on how to respond:

1. Abandon Alice. If Alice is guilty, you've lost. So it's safe to assume Alice is innocent, and sell any Alice contracts at their new higher prices, especially once you're broke and no longer can buy any more. If this encourages someone else to take option 2 and also move in on Alice, even better, that's one less person who can beat you if Alice is innocent. The majority of players should do this.
2. Attack Alice. If others are abandoning Alice in droves, her price might collapse even beyond \$25 as people rush to sell. You can then pick contracts up cheap, sell other contracts at better prices, and have a strictly better position.
3. Arbitrage. Try to make as much money off the situation as possible, without committing to a direction. If people are being 'too strategic' and too eager to get where they want to go, rather than focusing on getting the best price, then by making good trades (sell Alice when I buy too quickly, buy when others sell too quickly) and forcing others to get worse prices, I can end up with more value, then decide later what to do.

If you only engage in arbitrage, and others commit to suspects, you'll be in a lot of trouble unless you've already made a ton, because you won't have anyone to trade *against*. Your only option becomes to trade with the market, which limits how much you can get on the suspect you finally decide to go with, even if the mystery is solved while the market is open, and you're the only one left with cash.

The good and bad news is that's unlikely to happen, as others will also 'stick around' with flexible portfolios. That means that you won't be able to make that much when the mystery gets solved, but it does mean you can divide the spoils. If six players commit to suspects while four make good trades, not only are two of the six already shut out, it's likely the remaining four can coordinate (or just do sensible trades) to win if two or three of the suspects are found guilty, and sometimes should be able to nail all four.

When you have four (or six) suspects and ten players, there are not enough suspects for everyone to own one, and there certainly aren't enough for anyone to own two. That means that even if a suspect looks likely to be guilty, if you know you can't win that scenario, you'll be dumping, and that means at least seven of ten people are dumping any given suspect if they understand the game.

The logical response to this is to stay far enough ahead on your suspect that you clearly win if they're guilty, and if you get a good early opportunity to dump other contracts you should definitely do that. Good trades are generally good, and those trades just got even better, especially if everyone focuses on buying rather than selling. What you don't want to do is overpay, or run out of cash (and/or run out of things you can sell).

Thus, I might buy the Alice \$20, \$25 and Alice \$30 contracts, and start selling contracts on suspects I think are trading rich. What I'm worried about is competition - I don't want other players buying Alice contracts, so if they do, I'll make sure they

don't get size off by buying at least at their price, and I'll make sure to stay ahead of them on size. I'll also think about whether the remaining players are sophisticated enough to sell what they have, even at lousy prices; if they are, I'll be careful to hold a bunch of cash in reserve. If there are ten players, I can expect there to exist 16-25 Alice contracts, and I want to be sure not to run out of money.

Rank Ordering

This suggests each player has a few different goals.

You want to accumulate contracts in suspects you 'like' (which mostly means the ones you think are good bets), so you can get 'control' of one or more of them. Control means that if they did it, you win.

You want to get rid of contracts in the suspects you don't like. The trick here is that sometimes the price will go super high (relative to the probability they did it) as multiple players compete to gain 'control' of the suspect. Other times, the price will collapse because there is only one bidder for control of that suspect. If one player gets a bunch of contracts, and is in good overall shape, then no one else will compete.

That in turn might drive the price so low – \$10 or even \$5 – that the value of their portfolio shrinks a lot, tempting another player to enter, but doing so would drive the price up right away, so it often doesn't make sense to compete. If Bob is buying up Alice contracts and Carol now buys one at \$10, who is going to sell one now? Much better to wait to see if the price goes higher, which in turn puts Bob back in control. The flip side of that is, if Carol can buy a \$10 and a \$15 contract, and force Bob to then pay \$20, Carol can sell back to Bob at a profit. It's a risky bluff if others are actively selling, but it can definitely pay off.

The key in these fights is who has more overall portfolio value, plus the transaction costs of moving into more contacts. If Carol can make \$100 trading back and forth in other contracts, Bob is going to have a tough time keeping control, and mostly has to hope that Carol chooses to go after a different suspect. By being in as good shape as possible, Bob both is more likely to win the fight, and (if others realize this) more likely to avoid the fight.

With a lot of players engaged in active day trading, and aren't strategically focused, transaction costs could be low. If they're sufficiently low, then it could be a long time before it is hard to buy and sell what you want at a reasonable price, postponing the end game until quite late. The more other players are strategically focused, and strategy determines price, the harder it is to trade, the more existing positioning matters and the less you can try to day trade for profit other than anticipating a fight over a suspect, or a dump of them.

Rich Player, Poor Player

Suppose you're a poor player. You made some trades, and they didn't work out. Perhaps you held on to a suspect or two too long, and others dumped them, either strategically or for value. Perhaps you had a hunch, got overexcited, and others disagreed, and now you're looking foolish. Now you only have (let's say) \$120 in equity, down from \$200.

[You Play to Win the Game.](#) How do you do that? There are more players than suspects, several of whom have double your stake. So you'll need to find a good gambit.

A basic gambit would be to buy up all the contracts you can of a suspect everyone has dismissed. Even if there are very good reasons they seem impossible and are trading at \$5, you can still get enough of them to win if it works out, and you might have no competition since players in better shape have juicier targets. Slim chance beats none.

But if even that's out of the question, you'll have to rebuild your position. You will need to speculate via day trading. Any other play locks in a loss. You find yourself in a hole, and have no choice but to keep digging. Small arbitrage won't work. Your best bet is likely to watch the screen and listen, and try to react faster than everyone else in the hopes that the latest development is huge or seen as huge, then turn around and sell your new position to others to make a quick buck. Then hope there's still enough twists to do this again.

If the endgame has arrived, and rich players are sitting on or fighting for all the suspects, you've lost. Your best bet is to consolidate into cash, and hope some suspect crashes down to \$5 for you.

Now suppose you're a rich player. You have \$300 in equity. How do you maximize your chance of winning?

The basic play is to corner the market on the most likely suspect, or whoever you think is most likely. If you make a strong move in, you should be able to scare off competition, and even if you don't do so, you can use that as an opportunity to make more profit if they drive the price up. At some point, others will have to dump, and you can afford to give them a good price if you have to. It's hard to win a fight when outgunned. The key is not to engage too much too soon, as this risks letting a third player take advantage of an asset dump later. So you'll want to hold some cash for that, if possible. Remember that you'll need something like 12-14 contracts to feel safe from a dump, depending on how much equity you've built, if you're out of cash. That shuts out other players.

The advanced play, if you're sufficiently far ahead, is to try to win on *multiple* suspects. That's super hard. Even if you had \$400 in equity, if you divide it in half, there are still multiple other players over \$200. It seems unlikely you can get control of multiple worthwhile suspects. There's no point in trying for multiple bargain basement suspects at the expense of one good one, even if it works. So is there any hope here?

I think there is, in the scenario where there is a clear prime suspect.

In this scenario, the prime suspect was bid high early on. Given Robin's notes about player behavior and tendency to push prices too high, and the battle for control of the suspect, prices might get very high very quickly. There also may be players who will refuse to sell their contracts in the prime suspect, because they don't realize that they're shut out of winning in that case. Either they're maximizing expected value rather than chance of winning, or they don't realize the problem, or both.

This could open up an opportunity where the 'net profit' on the prime suspect isn't that high for any player. Suppose they start at \$25, and everyone starts with their two contracts. They then trade at \$30, \$40, \$50 and \$60 in a row, not all to the same player. So there's minimal chances to buy contracts that make you that much money.

If you buy an \$80 contract your maximum profit is \$20, which is easy to beat by day trading.

So what you can do is *go for the block*. Hopefully you helped drive the price up early, which is part of how you got your equity. Then contracts only really traded at \$60 and \$80. So even if the suspect is guilty, someone who moved in on this without day trading first is not going to end up with many contracts. You start with \$200, so lets say they end up with 3 contracts and a little cash.

It's not crazy for you to sell the suspect at the top, do some successful day trading, and then have over \$300 in cash. You could win without any contracts that pay out, if you know you're the most successful day trader and no one can have that many contracts.

That's a better position than having 5 or 6 contracts in the prime suspect, since you still have cash if they're innocent. The trick is then having that be enough to win on *another suspect* as well, or splitting your efforts by holding onto contracts elsewhere. Tricky. But perhaps not impossible, especially if people are dumping contracts at \$5.

At a minimum, what you can do is be in a strong position to respond to new developments, and be able to choose which other suspect to back later in the game if you now think they're more likely, while still winning if the situation doesn't change. That's very strong.

A final note is that it is legal, in the game, to trade with another player without going through the market. This could be used to buy out a players' position in a suspect, shifting control of that suspect, and avoid the issue where once a player starts dumping a position, the price will collapse, as well as the ability of other players to 'intercept' the transfer and ruin the buyers' attempt to accumulate a new position and take control. Thus, players should learn that if they have a bunch of contracts and want out, they should check for a bulk buyer, and if they want in they should consider doing the same. The risk of course is that you tip your hand, which makes doing it on the buy side less attractive.

Flexible Structures

It's also worth noting that you can extend the idea easily to other prediction markets, and to an online version.

You could trade on the outcome of a sporting event, or an election, or any other real-world prediction market, using the same rules. You could play a board game, and also play the contract game on the outcome of your board game. That gives players something to do between turns and extra things to think about, and gives extra players or eliminated players something to do.

You could trade over a series of outcomes or events (for example, all the football games played today, or both the winner of the game and the combined number of points scored, or even obscure stuff too like number of punts or what not) in order to reward more trading 'for value' and place less emphasis on being right. Or just keep track of funds between games, watch multiple shows, and reward the overall winner.

That raises the question of what we can learn about prediction markets from the game.

Market Accuracy and Implications

Early in the game, market prices should roughly reflect fair probabilities of being guilty. Anyone who jumps the gun for strategic positioning will lose out to a more patient player. That won't stop players from being overeager, and bidding suspects up too high, but as Robin noted that opens the door for others to do arbitrage and sell the contracts back down to reasonable prices.

Later in the game, prices will grow increasingly inaccurate as players jockey for position, and let strategic considerations override equity considerations.

This is a phenomenon we see in many strategic games. Early in the game, players mostly are concerned about getting good value, as most assets are vaguely fungible and will at some point be useful or necessary. As the game progresses, players develop specialized strategies, and start to do whatever they need to do to get what they need, often at terrible exchange rates, and dump or ignore things they don't need, also at terrible exchange rates.

If we wanted to improve accuracy, we'd need to make the game less strategic and more tactical, by rewarding players who maximize expected profits. There's a dumb market that is handing out Free Money when news occurs. We'd like players to battle for a share of that pie, rather than competing for control of suspects. If the game was played over many rounds, the early rounds would mostly focus on expected value and doing good trades. If the game was played *for real money*, and settled in actual dollars, then we'd *definitely* have a lot more accurate pricing!

If a market has traders purely motivated by expected value and profit, then its pricing will be as good as the pricing ability of the traders.

If a market has a few 'motivated' traders, or noise traders, that are doing something for reasons other than expected value, that is good. You need a source of free money [to make the market work](#). Thus, the existence of the bank, as a source of free money, is great, because it motivates the game. You can imagine a version of the game where players can only trade when they agree to it. There would still be trades, since the prize for winning should overcome frictions and adverse selection, but volume of trading would plummet.

If a market has a few people who have poor fair values, that works like motivated traders.

If a market has *too many* traders who have poor fair values, or in context they have fair values that are not based on the expected payout, then relationships break down. There's now profit for those who bet against them, but that doesn't mean there's enough money in the wings to balance things out. At some point that money is exhausted, and no one paying attention has spare funds. Prices stop being accurate, to varying degrees.

In particular, this illustrates that if those managing the money have payouts *that are non-linear functions of their profits*, then very weird things will happen. If I get fired for missing my benchmark, and so do my colleagues, but we don't get extra fired for missing them by a lot, then this will lead us to discount tail risks. In the game, this takes the form of dumping suspects you can't control – if Alice did it, you've already lost, and [third prize is you're fired](#). There are many other similar scenarios. If we want

accurate prices, we need traders to have linear utility functions, or reasonable approximation thereof.

Overall

This game sounds like a lot of fun, and seems to have lots of opportunities for deep tactical and strategic play, for bluffing, and to do things that players should find fun. I really hope that it gets made. You can take one or more of many roles - arbitrageur, mystery solver, genre savvy logician, momentum, value, tactical or strategic, or just hang out and watch the fun and/or the mystery, and if you later get a hunch, you can go for it.

I hope to talk to a few friends of mine who have small game companies, in the hopes one of them can help. Kickstarter ho?

If anyone out there is interested in the game and making it happen, please do talk to [Robin Hanson](#) about it. I'm sure he'd be happy to help make it a reality. And if you're looking to play, I encourage you to give it a shot, and report back.

nostalgebraist - bayes: a kinda-sorta masterpost

This is a linkpost for <http://nostalgebraist.tumblr.com/post/161645122124/bayes-a-kinda-sorta-masterpost/>

Extended criticism of "Bayesianism" as discussed on LW. An excerpt from the first section of the post:

Like most terms ending in -ism, [Bayesianism] can mean a number of different things. In its most limited sense, "Bayesianism" is a collection of technical/mathematical *machinery*, analogous to a tool or toolbox. This collection, which I will call "the Bayesian machinery," uses a particular way of representing knowledge, and if you can represent your knowledge in that way, the machinery tells you how to alter it when presented with new evidence.

The Bayesian machinery is frequently used in statistics and machine learning, and some people in these fields believe it is very frequently the right tool for the job.

I'll call this position "weak Bayesianism." There is a more extreme and more philosophical position, which I'll call "strong Bayesianism," that says that the Bayesian machinery is the *single correct way* to do not only statistics, but science and inductive inference in general – that it's the "aspirin in willow bark" that makes science, and perhaps all speculative thought, work insofar as it *does* work. (I.e., if you're doing these things right, you're being Bayesian, whether you realize it or not.)

Strong Bayesianism is what E. T. Jaynes and Eliezer Yudkowsky mean when they say they are Bayesians. It is usually what I am talking about when I say I "don't like" Bayesianism. I think strong Bayesianism is dead wrong. I think weak Bayesianism may well be true, in that the Bayesian machinery may well be a very powerful set of tools – but I want to understand why, in a way that defines the power of a tool by some metric other than how Bayesian it is.

Contents of the post:

- 0. What is "Bayesianism"?**
- 1. What is the Bayesian machinery?**
 - 1a. Synchronic Bayesian machinery**
 - 1b. Diachronic Bayesian machinery**
- 3. What is Bayes' Theorem?**
- 4. How does the Bayesian machinery relate to the more classical approach to statistics?**
- 5. Why is the Bayesian machinery supposed to be so great?**
- 6. Get to the goddamn point already. What's wrong with Bayesianism?**

7. The problem of ignored hypotheses with known relations

7b. Okay, but why is this a problem?

8. The problem of new ideas

8b. The natural selection analogy

9. Where do priors come from?

10. It's just regularization, dude

11. Bayesian "Occam factors"

Open AI co-founder on AGI

This is a linkpost for <https://youtu.be/w3ues-NayAs?t=27m6s>

Here's a talk where Open AI co-founder Ilya Sutskever talks about the issue of AI progress and AGI. (Before the timestamp he was talking about Open AI's specific projects such as Dota 2 and their hand robot).

tl;dw: (*italics* indicate my interpolations):

- AGI in the near-medium term can't be ruled out based on current rates of progress.
- AI progress is more constrained by hardware than by conceptual advances (*partly because it's hard to make conceptual advances if you don't have the hardware to test them with*).
- Companies and researchers are willing to make increasingly large investments into compute hardware (*at rate greater than Moore's Law*).
- Open AI is also working on what we would consider Friendliness (*but seemingly not at the same level of rigour as MIRI*).

I am the very model of a self-recursive modeler

I am the very model of a self-recursive modeler

My consciousness encompasses itself in many meta layers

To Russel's paradox, I have the set of all the answers

I cut the hair of all the non-self-barbering hairdressers

I think about me thinking of my thoughts with regularity

In my Cartesian theater, I see myself with clarity

I tell you that I lie without the slightest contradiction

The story of my life is told in meta-metafiction

My inner simulator replicates the total universe

I simulate myself as I am writing out this line of verse

My secrets are unknown to me, I keep them confidential

Referring to myself I'd say I'm quite self-referential

My strange loop's Escherer than yours, it's Bacher and it's Godeler

I am the very model of a self-recursive modeler

Cross-posted in full from [Putanumonit](#)

Alignment Newsletter #25

Crossposted from the [AI Alignment Forum](#). May contain more technical jargon than usual.

Highlights

Towards a New Impact Measure (*Alex Turner*): This post introduces a new idea for an impact measure. It defines impact as change in our ability to achieve goals. So, to measure impact, we can simply measure how much easier or harder it is to achieve goals -- this gives us Attainable Utility Preservation (AUP). This will penalize actions that restrict our ability to reach particular outcomes (opportunity cost) as well as ones that enlarge them (instrumental convergence).

Alex then attempts to formalize this. For every action, the impact of that action is the absolute difference between attainable utility after the action, and attainable utility if the agent takes no action. Here, attainable utility is calculated as the sum of expected Q-values (over m steps) of every computable utility function (weighted by $2^{-\text{length of description}}$). For a plan, we sum up the penalties for each action in the plan. (This is not entirely precise, but you'll have to read the post for the math.) We can then choose one canonical action, calculate its impact, and allow the agent to have impact equivalent to at most N of these actions.

He then shows some examples, both theoretical and empirical. The empirical ones are done on the suite of examples from AI safety gridworlds used to test relative reachability. Since the utility functions here are indicators for each possible state, AUP is penalizing changes in your ability to reach states. Since you can never increase the number of states you reach, you are penalizing decrease in ability to reach states, which is exactly what relative reachability does, so it's not surprising that it succeeds on the environments where relative reachability succeeded. It does have the additional feature of handling shutdowns, which relative reachability doesn't do.

Since changes in probability of shutdown drastically change the attainable utility, any such changes will be heavily penalized. We can use this dynamic to our advantage, for example by committing to shut down the agent if we see it doing something we disapprove of.

My opinion: This is quite a big improvement for impact measures -- it meets many desiderata that weren't satisfied simultaneously before. My main critique is that it's not clear to me that an AUP-agent would be able to do anything useful. For example, perhaps the action used to define the impact unit is well-understood and accepted, but any other action makes humans a little bit more likely to turn off the agent. Then the agent won't be able to take those actions. Generally, I think that it's hard to satisfy the conjunction of three desiderata -- objectivity (no dependence on values), safety (preventing any catastrophic plans) and non-trivialness (the AI is still able to do some useful things). There's a lot more discussion in the comments.

Realism about rationality (*Richard Ngo*): In the same way that moral realism claims that there is one true morality (even though we may not know it yet), rationality realism is the claim that there is one "correct" algorithm for rationality or intelligence. This post argues that many disagreements can be traced back to differences on how

much one identifies with the rationality realism mindset. For example, people who agree with rationality realism are more likely to think that there is a simple theoretical framework that captures intelligence, that there is an "ideal" decision theory, that certain types of moral reasoning are "correct", that having contradictory preferences or beliefs is really bad, etc. The author's skepticism about this mindset also makes them skeptical about agent foundations research.

My opinion: This does feel like an important generator of many disagreements I've had. I'd split rationality realism into two subcases -- whether you expect that there is a simple "correct" algorithm for computation-bounded rationality, and whether you expect there is only a simple "correct" algorithm for rationality given infinite compute, but the bounded computation case may be a lot messier. (I'm guessing almost all rationality realists fall in the latter category, but I'm not sure.)

I'd expect most of the people working on reducing existential risk from AI to be much more realist about rationality, since we often start working on this based on astronomical waste arguments and utilitarianism, which seems very realist about preferences. (At least, this was the case for me.) This is worrying -- it seems plausible to me that there isn't a "correct" rationality or intelligence algorithm (even in the infinite compute case), but that we wouldn't realize this because people who believe that also wouldn't want to work on AI alignment.

Technical AI alignment

Technical agendas and prioritization

[Realism about rationality](#) (Richard Ngo): Summarized in the highlights!

Agent foundations

[In Logical Time, All Games are Iterated Games](#) (Abram Demski) (summarized by Richard): The key difference between causal and functional decision theory is that the latter supplements the normal notion of causation with "logical causation". The decision of agent A can logically cause the decision of agent B even if B made their decision before A did - for example, if B made their decision by simulating A. Logical time is an informal concept developed to help reason about which computations cause which other computations: logical causation only flows forward through logical time in the same way that normal causation only flows forward through normal time (although maybe logical time turns out to be loopy). For example, when B simulates A, B is placing themselves later in logical time than A. When I choose not to move my bishop in a game of chess because I've noticed it allows a sequence of moves which ends in me being checkmated, then I am logically later than that sequence of moves. One toy model of logical time is based on proof length - we can consider shorter proofs to be earlier in logical time than longer proofs. It's apparently surprisingly difficult to find a case where this fails badly.

In logical time, all games are iterated games. We can construct a series of simplified versions of each game where each player's thinking time is bounded. As thinking time increases, the games move later in logical time, and so we can treat them as a series of iterated games whose outcomes causally affect all longer versions. Iterated games

are fundamentally different from single-shot games: the [folk theorem](#) states that virtually any outcome is possible in iterated games.

My opinion: I like logical time as an intuitive way of thinking about logical causation. However, the analogy between normal time and logical time seems to break down in some cases. For example, suppose we have two boolean functions F and G, such that $F = \text{not } G$. It seems like G is logically later than F - yet we could equally well have defined them such that $G = \text{not } F$, which leads to the opposite conclusion. As Abram notes, logical time is intended as an intuition pump not a well-defined theory - yet the possibility of loopiness makes me less confident in its usefulness. In general I am pessimistic about the prospects for finding a formal definition of logical causation, for reasons I described in [Realism about Rationality](#), which Rohin summarised above.

Learning human intent

[Adversarial Imitation via Variational Inverse Reinforcement Learning](#) (*Ahmed H. Qureshi et al*)

[Inspiration Learning through Preferences](#) (*Nir Baram et al*)

Reward learning theory

[Web of connotations: Bleggs, Rubes, thermostats and beliefs](#) and [Bridging syntax and semantics, empirically](#) (*Stuart Armstrong*): We're planning to summarize this once the third post comes out.

Preventing bad behavior

[Towards a New Impact Measure](#) (*Alex Turner*): Summarized in the highlights!

Handling groups of agents

[CM3: Cooperative Multi-goal Multi-stage Multi-agent Reinforcement Learning](#) (*Jiachen Yang et al*)

[Negative Update Intervals in Deep Multi-Agent Reinforcement Learning](#) (*Gregory Palmer et al*)

[Coordination-driven learning in multi-agent problem spaces](#) (*Sean L. Barton et al*)

Interpretability

[Transparency and Explanation in Deep Reinforcement Learning Neural Networks](#) (*Rahul Iyer et al*)

[Towards Better Interpretability in Deep Q-Networks](#) (*Raghuram Mandyam Annasamy et al*)

Verification

[Training for Faster Adversarial Robustness Verification via Inducing ReLU Stability](#) (*Kai Y. Xiao et al*): The idea behind verification is to consider all possible inputs at the same time, and show that no matter what the input is, a particular property is satisfied. In ML, this is typically applied to adversarial examples, where inputs are constrained to be within the L-infinity norm ball of dataset examples. Prior papers on verification (covered in [AN #19](#)) solve a computationally easier relaxation of the verification problem, that gives a lower bound on the performance of the classifier. This paper aims to use exact verification, since it can compute the exact adversarial performance of the classifier on the test set, and to figure out how to improve its performance.

One easy place to start is to encourage weights to be zero, since these can be pruned from the problem fed in to the constraint solver. (Or more likely, they feed it in anyway, but the constraint solver immediately gets rid of them -- constraint solvers are pretty smart.) This can be done using L1 regularization and pruning small weights. This already gives two orders of magnitude of speedup, making it able to verify that there is no adversarial attack with $\epsilon = 0.1$ on a particular MNIST digit in 11 seconds on average.

Next, they note that verification with linear constraints and functions is easy -- the challenging aspect is the Relu units that force the verifier to branch into two cases. (Since $\text{relu}(x) = \max(x, 0)$, it is the identity function when x is positive, and the zero function otherwise.) So why not try to ensure that the Relu units are also linear? Obviously we can't just make all the Relu units linear -- the whole point of them is to introduce nonlinearity to make the neural net more expressive. But as a start, we can look at the behavior of the Relu units on the examples we have, and if they are almost always active (inputs are positive) or almost always inactive (inputs are negative), then we replace them with the corresponding linear function (identity and zero, respectively), which is easier to verify. This gets another $\sim 2x$ speedup.

But what if we could also change the training procedure? Maybe we could augment the loss so that the Relu units are either decisively active or decisively inactive on any dataset example. They propose that *during training* we consider the L-infinity norm ball around each example, use that to create intervals that each pixel must be in, and then make a forward pass through the neural net using interval arithmetic (which is fast but inexact). Then, we add a term to the loss that incentivizes the interval for the input to each Relu to exclude zero (so that the Relu is either always active or always inactive). They call this the Relu Stability loss, or RS loss.

This leads to a further 4-13x speedup with similar test set accuracy. They then also test on MNIST with $\epsilon = 0.2, 0.3$ and CIFAR with $\epsilon = 2/255, 8/255$. It leads to speedup in all cases, with similar test set accuracy on MNIST but reduced accuracy on CIFAR. The provable accuracy goes up, but this is probably because when there's no RS loss, more images time out in verification, not because the network becomes better at classification. Other verification methods do get better provable accuracies on CIFAR, even though in principle they could fail to detect that a safe example is safe. This could be because their method times out frequently, or because their method degrades the neural net classifier -- it's hard to tell since they don't report number of timeouts.

My opinion: As with the previous papers on verification, I'm excited in the improvement in our capability to prove things about neural nets. I do think that the more important problem is how to even state properties that we care about in a way that we could begin to prove them. For example, [last week](#) we saw the unrestricted

adversarial examples challenge, where humans are the judge of what a legal example is -- how can we formalize that for a verification approach?

On this paper specifically, I wish they had included the number of timeouts that their method has -- it's hard to interpret the provable accuracy numbers without that. Based on the numbers in the paper, I'm guessing this method is still much more computationally expensive than other methods. If so, I'm not sure what benefit it gives over them -- presumably it's that we can compute the exact adversarial accuracy, but if we don't have enough compute, such that other methods can prove better lower bounds anyway, then it doesn't seem worth it.

Miscellaneous (Alignment)

[AI Alignment Podcast: Moral Uncertainty and the Path to AI Alignment with William MacAskill \(Lucas Perry and William MacAskill\)](#) (summarized by Richard): Initially, Will articulates arguments for moral realism (the idea that there are objectively true moral facts) and moral uncertainty (the idea that we should assign credences to different moral theories being correct). Later, the discussion turns to the relevance of these views to AI safety. Will distinguishes the control problem (ensuring AIs do what we say), from the problem of aligning AI with human values, from the problem of aligning AI with moral truth. Observing humans isn't sufficient to learn values, since people can be self-destructive or otherwise misguided. Perhaps AI could extrapolate the values an idealised version of each person would endorse; however, this procedure seems under-defined.

On the moral truth side, Will worries that most educated people are moral relativists or subjectivists and so they won't sufficiently prioritise aligning AI with moral truth. He advocates for a period of long philosophical reflection once we've reduced existential risk to near zero, to figure out which future would be best. Careful ethical reasoning during this period will be particularly important since small mistakes might be magnified massively when implemented on an astronomical scale; however, he acknowledges that global dynamics make such a proposal unlikely to succeed. On a brighter note, AGI might make great advances in ethics, which could allow us to make the future much more morally valuable.

My opinion: I think moral uncertainty is an important and overdue idea in ethics. I also agree that the idea of extrapolating an idealised form of people's preferences is not well-defined. However, I'm very skeptical about Will's arguments about moral realism. In particular, I think that saying that nothing matters at all without moral realism is exactly the sort of type error which Eliezer argued against [here](#).

I'm more sympathetic to the idea that we should have a period of long reflection before committing to actions on an astronomical scale; this seems like a good idea if you take moral uncertainty at all seriously.

[Does Your Model Know the Digit 6 Is Not a Cat? A Less Biased Evaluation of "Outlier" Detectors \(Alireza Shafaei et al\)](#)

AI strategy and policy

[The role of corporations in addressing AI's ethical dilemmas \(Darrell M. West\)](#)

Other progress in AI

Reinforcement learning

[Model-Based Reinforcement Learning via Meta-Policy Optimization](#) (*Ignasi Clavera, Jonas Rothfuss et al*)

[Generalizing Across Multi-Objective Reward Functions in Deep Reinforcement Learning](#) (*Eli Friedman et al*)

[Challenges of Context and Time in Reinforcement Learning: Introducing Space Fortress as a Benchmark](#) (*Akshat Agarwal et al*) (summarized by Richard): The authors note that most existing RL benchmarks (like Atari games) lack sharp context-dependence and temporal sensitivity. The former requires an agent to sometimes change strategies abruptly; the latter requires an agent's strategy to vary over time. Space Fortress is an arcade-style game which does have these properties, and which cannot be solved by standard RL algorithms, even when rewards are made dense in a naive way. However, when the authors shape the rewards to highlight the context changes, their agent achieves superhuman performance.

My opinion: The two properties that this paper highlights do seem important, and the fact that they can be varied in Space Fortress makes it a good benchmark for them.

I'm not convinced that the experimental work is particularly useful, though. It seems to reinforce the well-known point that shaped rewards can work well when they're shaped in sensible ways, and much less well otherwise.

[Combined Reinforcement Learning via Abstract Representations](#) (*Vincent François-Lavet et al*)

[Sim-to-Real Transfer Learning using Robustified Controllers in Robotic Tasks involving Complex Dynamics](#) (*Jeroen van Baar et al*)

[Automata Guided Reinforcement Learning With Demonstrations](#) (*Xiao Li et al*)

Deep learning

[Automatic Program Synthesis of Long Programs with a Learned Garbage Collector](#) (*Amit Zohar et al*)

[GAN Lab](#) (*Minsuk Kahng et al*)

Applications

[Neural-Augmented Static Analysis of Android Communication](#) (*Jinman Zhao et al*)

AGI theory

[Abstraction Learning](#) (*Fei Deng et al*)

News

[Slides from Human-Level AI 2018](#)

Against the barbell strategy

This is a linkpost for <http://benjaminrosshoffman.com/against-barbell-strategy/>

Nassim Nicholas Taleb recommends that instead of the balanced portfolio of investments recommended by portfolio theory, we follow a "barbell" strategy of putting most of our assets in a maximally safe, stable investment, and making small, sustainable bets with very high potential upside. If taken literally, this can't work because no such safe asset class exists.

Taleb typically uses US treasury bonds as an example of an extremely safe asset. This can only be true metaphorically. They are extremely safe if your measure of value is denominated in US dollars. But of what use to me are US treasury bonds in a hurricane, or if the US dollar collapses? Of what use if people like me become targets for state persecution, stripped of our financial assets or ability to access them? Of what use if I'm trying to deal with a spiritual or health problem and don't know how to find competent advice? If I'm trying to feed myself and don't know how to distinguish profit-seeking propaganda or engineered taste from genuine information about which foods are healthful?

Risks, and opportunities, are many-dimensional. And along very few dimensions is there anything like the sort of crystalline perfection of a treasury bond. A personal network extending to a variety of fields can do a lot to help you avoid getting scammed, and get access to at least an ordinary standard of competence, such as it is. A friend with different interests and skills from yours will do a lot to expand the access you have, and the skills and knowledge at your disposal. Friends from different cultures or communities can serve as an important hedge if you and people like you become a target of extraction or persecution for some reason, or simply if your way of life becomes unsustainable.

Land you are used to occupying and know well is easier to robustly possess and develop than a rented apartment (but the upside is often higher in high-density areas, which is why young people often prefer them). A spouse makes it easier to pool assets into a venture diversified between building up a single household and seeking external trade or mercenary opportunities. And children who grow up well-adjusted to the world can be helpful if you need to navigate a changing incentive landscape when you've already committed your assets to a bunch of permanent bets (which you need to do, since time is limited). In addition, they may be a decent consolation prize if cryonics doesn't work to ensure your personal survival. Intergenerational communities of interest protect you a lot more day to day, though they themselves can have trouble reallocating resources in changing circumstances.

Risk cannot be avoided, only managed. The risk-free asset is an illusion. Someday your natural life will end, you need to invest taking that into account, and taking into account that this is true of everything else on this earth that you can put your trust in. If you allocate too much of your portfolio to "safe" assets you're leaving yourself unnecessarily exposed to the risk that this asset will become irrelevant, and failing to take the chances that are actually available to improve your position.

To a large extent, young people do well in the long run by pursuing things that are genuinely *fun* but not universal. This engages our natural sense of opportunity, which doesn't understand formal systems like money very well, but understands very well

how to look for high-upside bets that are actually good for us. I wish someone had explained this to me fifteen years ago. If this advice helps you, then think of me when you're successful, and remember that I didn't charge you for it at the time. And remember to pay it forward.

Or, as Ecclesiastes says:

Cast thy bread upon the waters: for thou shalt find it after many days. Give a portion to seven, and also to eight; for thou knowest not what evil shall be upon the earth. If the clouds be full of rain, they empty themselves upon the earth: and if the tree fall toward the south, or toward the north, in the place where the tree falleth, there it shall be. He that observeth the wind shall not sow; and he that regardeth the clouds shall not reap. As thou knowest not what is the way of the spirit, nor how the bones do grow in the womb of her that is with child: even so thou knowest not the works of God who maketh all. In the morning sow thy seed, and in the evening withhold not thine hand: for thou knowest not whether shall prosper, either this or that, or whether they both shall be alike good. Truly the light is sweet, and a pleasant thing it is for the eyes to behold the sun: But if a man live many years, and rejoice in them all; yet let him remember the days of darkness; for they shall be many. All that cometh is vanity. Rejoice, O young man, in thy youth; and let thy heart cheer thee in the days of thy youth, and walk in the ways of thine heart, and in the sight of thine eyes: but know thou, that for all these things God will bring thee into judgment. Therefore remove sorrow from thy heart, and put away evil from thy flesh: for childhood and youth are vanity.

Tradition is Smarter Than You Are

This is a linkpost for <https://scholars-stage.blogspot.com/2018/08/tradition-is-smarter-than-you-are.html>

Counterfactuals and reflective oracles

Crossposted from the [AI Alignment Forum](#). May contain more technical jargon than usual.

Summary: Take the argmax of counterfactual expected utility, using a thin counterfactual which itself maximizes a priori expected utility.

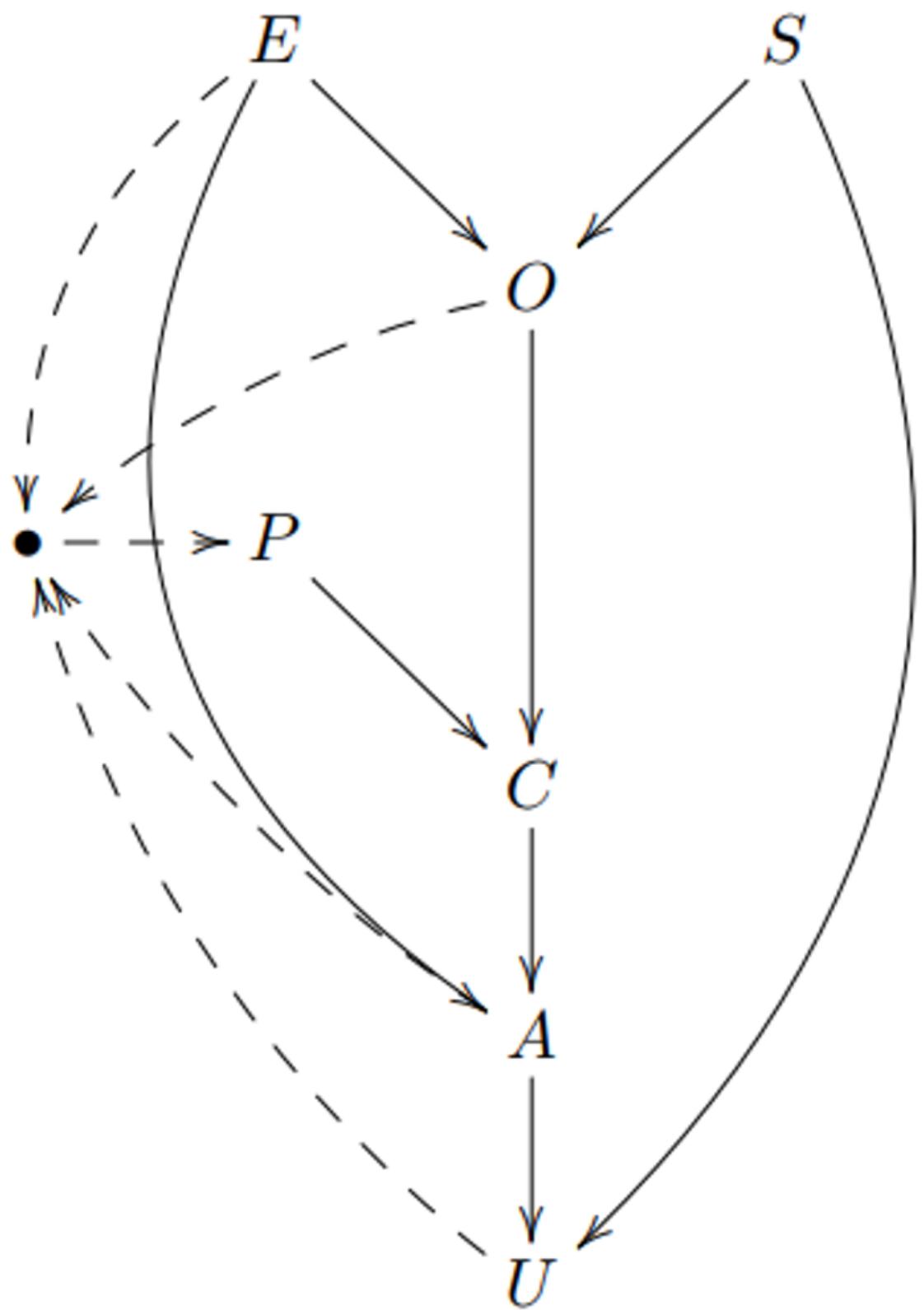
Followup to: [An environment for studying counterfactuals, Counterfactuals, thick and thin](#)

Replacement for: [Logical counterfactuals and differential privacy](#)

In [An environment for studying counterfactuals](#), I described a family of decision problems that elicit counterfactuals from an agent. In this post I'll redefine that family using reflective oracles, and define an optimal agent. I was aided in this work by discussions with Alex Appell, James Brooks, Abram Demski, Linda Linsefors, and Alex Mennen.

The problem

Here's the (not quite) Bayes net for the problem. I changed it a little from the last post, to make the thin counterfactuals more explicit:



E controls exploration and S is a hidden state. The observation O is isomorphic as a random variable to $E \times S$, but the particular isomorphism \hat{K} used by the problem is unknown to the agent.

The agent is given O and a prior P ; it outputs a set C of counterfactual expected utilities, one for each action. The agent's action A is equal to the argmax of these expectations, unless it's overridden by E . The action and hidden state determine the utility U .

The dotted lines mean that P is equal to the marginal over $E \times O \times A \times U$. The machinery of reflective oracles resolves the circularity.

Now I'll define this distribution in detail, using Python-like pseudocode. I ought to use a probabilistic programming language, but I don't know any of those. Note that if you have a reflective oracle, you can compute expectations of distributions in constant time, and you can compute argmaxes over sets of real numbers, with ties broken in an undetermined way. Real numbers and distributions can be represented by programs.

environment()

The function `environment()` takes an agent and outputs a sample from a joint distribution over $E \times O \times A \times U$. The agent is a function that takes the environment and an observation and stochastically returns a set of counterfactual expectations, which is itself encoded as a function from actions to real numbers.

The environment samples s from a uniform distribution. Usually e is equal to the symbol $*$, meaning no exploration; but with probability ϵ it's a uniformly random action. It computes o using the secret isomorphism \hat{K} . It calls the agent to get its counterfactual expectations. It argmaxes over those expectations to get the action, unless e overrides that. Finally, u is sampled from a problem-dependent likelihood function that depends on a and s .

```
define environment(ag):
```

```
    sample s ~ Unif(S)
```

```
    sample e ~ εUnif(A) + (1 - ε)δ*
```

$o := \hat{\kappa}(e, s)$

$c := ag(\text{curry}(\text{environment}, ag), o)$

if $e \in A$: $a := e$

else: sample $a \sim \text{argmax}_a c(a)$

sample $u \sim U|A = a, S = s$

return (e, o, a, u)

The agent

The agent's strategy is to choose a thin counterfactual that maximizes its a priori expected utility. We'll define a couple helper functions first.

cfact()

The function $\text{cfact}()$ takes a prior $P \in \Delta(E \times O \times A \times U)$, a thin counterfactual

$\kappa : E \times S \approx O$, an observation $o \in O$, and an action $a \in A$, and deterministically outputs the expected utility conditional on o and counterfactually conditional on a , taking the expectation with respect to P and the counterfactual with respect to κ .

This is how to compute a thin counterfactual distribution: First condition on $O = o$ and use κ to infer a posterior on S . Then forget o but retain your posterior on S . Then condition on $A = a$. The reason this works is that S is independent of E , so there's still a chance for $A = a$ by exploration.

define $\text{cfact}(P, \kappa, o, a)$:

$s = \text{pr}_S \kappa^{-1}(o)$

```
return ( $E_P[U|O = \kappa(a, s)] P_P(E = a) +$ 
```

```
 $E_P[U|O = \kappa(*, s), A = a] P_P(A = a|O = \kappa(*, s)) P_P(E = *) /$ 
```

```
( $P_P(E = a) + P_P(A = a|O = \kappa(*, s)) P_P(E = *)$ )
```

utility()

This function takes a prior and a thin counterfactual, and samples from the marginal on U that would result if the agent used that thin counterfactual. The previous sentence uses the word "would", so you might think we need another thin counterfactual to make sense of it; but our agent is a program, and programs come with their own "thick" notion of counterfactual.

The calculation goes like this: Sample e and o . If e makes us explore, then that determines the action and we can sample a utility conditional on e and o . Otherwise, we assume the action is chosen to maximize $cfact()$.

```
define utility(P, κ):
```

```
    sample e, o ~ P
```

```
    if  $e \in A$ : return  $u \sim P|E = e, O = o$ 
```

```
    else: return  $\max_a cfact(P, \kappa, o, a)$ 
```

agent()

The agent is given a prior and an observation, and outputs the expected utility conditional on $A = a$, for every a . This output is encoded in a program that takes an action and outputs a real number. The output may be stochastic because $agent()$ contains an argmax which may be stochastic in case of ties.

The true thin counterfactual $\hat{\kappa} : E \times S \approx O$ determines a projection $O \rightarrow E$. The agent can deduce this projection by analyzing the prior: For every o , $P_P(E = e|O = o) = 0$ for all but one value of e .

The argmax is taken over all isomorphisms $\kappa : E \times S \approx O$ that commute with projecting to E . I won't bother writing that out in the pseudocode:

```
define agent(P, o):
    κ := argmax_κ E_P[utility(P, κ)]
    return curry(cfact, P, κ, o)
```

Optimality theorems

One of the counterfactual expectations that the agent emits is factual. It's reasonable to ask whether that one is accurate.

Theorem: If U is bounded by U_{\max} and $E = *$, then with probability $\geq 1 - \sqrt{\varepsilon}$, if we sample $a \sim A$, $s \sim S$, $o \sim O$, we have

$$|\text{agent}(P, o)(a) - E[U|A = a, S = s]| < 4U_{\max}\sqrt{\varepsilon}$$

Proof: $\text{agent}(P, o)(a) = \text{cfact}(P, \kappa, o, a)$ for some κ . That equals

$$(E[U|O = \kappa(a, s)]P(E = a) + E[U|O = \kappa(*, s), A = a]P[A = a|O = \kappa(*, s)]P(E = *)) / D$$

$$= (E[U|O = \kappa(a, s)] + E[U|O = o, A = a]P[A = a|O = o](1 - \varepsilon)) / D$$

where $s = \text{pr}_S \kappa^{-1}(o)$ and $D = \frac{1}{|A|} + P(A = a|O = o)(1 - \varepsilon)$. With a little algebra, that gives us

$$|\text{agent}(P, o)(a) - E[U|A = a, S = s]| < \frac{\epsilon}{\epsilon + (1-\epsilon)\frac{1}{|A|} \max_{a \in A} P(A = a | O = o)}$$

By the following lemma, with probability $\geq 1 - \sqrt{\epsilon}$, we have $P(A = a | O = o) > \frac{\sqrt{\epsilon}}{|A|}$. A little more algebra concludes the proof. ■

Lemma: Sample x, y from a joint distribution over $X \times Y$. With probability $\geq 1 - \epsilon$, we have $P(X = x | Y = y) > \frac{\sqrt{\epsilon}}{|X|}$.

Now we'll prove that $\text{agent}()$ gets an optimal amount of utility. First, a couple lemmas.

Lemma: $E[U] = \max_{\kappa} E[\text{utility}(P, \kappa)]$

Proof: By the definition of environment(),

$$E[U] = \sum_a E[U|E = a]P(E = a) + (1 - \epsilon)E_{o|e=*\max_a \text{agent}(P, o)(a)}$$

By definition of $\text{agent}()$, this is

$$= \frac{1}{|A|} \sum_a E[U|E = a] + (1 - \epsilon)E_{o|e=*\max_a \text{cfact}(P, \kappa, o, a)}$$

where the expectation E_{κ} is taken over the output distribution of $\arg\max_{\kappa} E[\text{utility}(P, \kappa)]$. By linearity of expectation,

$$= E_{\kappa}[\frac{1}{|A|} \sum_a E[U|E = a] + (1 - \epsilon)E_{o|e=*\max_a \text{cfact}(P, \kappa, o, a)}]$$

This is just

$$= E_{\kappa}[\text{utility}(P, \kappa)]$$

$$= \max_{\kappa} E[\text{utility}(P, \kappa)] ■$$

Recall that $\hat{\kappa}$ is the secret isomorphism $E \times S \approx O$. The next lemma says that $\text{utility}(P, \hat{\kappa})$ represents the highest utility that an agent can get:

Lemma: $E[\text{utility}(P, \hat{K})] = \frac{1}{|A|} \sum_a E[U|A = a] + (1 - \varepsilon) E[U|A = \text{argmax}_a E[U|A = a, S = s]]$

Proof: Using the fact that

$$E[U|O = \hat{K}(a, s)] = E[U|A = a, S = s]$$

and

$$E[U|O = \hat{K}(*, s), A = a] = E[U|A = a, S = s]$$

and the definition of `cfact()`, we have

$$\text{cfact}(P, \hat{K}, o, a) = E[U|A = a, S = s].$$

Plugging that in to the definition of `utility()`, we get

$$\text{utility}(P, \hat{K}) = \frac{1}{|A|} \sum_a E[U|A = a] + (1 - \varepsilon) E[\max_a E[U|A = a, S = s]]$$

$$= \frac{1}{|A|} \sum_a E[U|A = a] + (1 - \varepsilon) E[U|A = \text{argmax}_a E[U|A = a, S = s]] \blacksquare$$

Finally, this theorem says that the agent gets the highest possible utility:

Theorem: $E[U] = \frac{1}{|A|} \sum_a E[U|A = a] + (1 - \varepsilon) E[U|A = \text{argmax}_a E[U|A = a, S = s]]$

Proof: Call the right-hand side U_{opt} . It's easy to see that $E[U] \leq U_{\text{opt}}$. By the previous two lemmas, we also have

$$E[U] = \max_K E[\text{utility}(P, K)] \geq E[\text{utility}(P, \hat{K})] = U_{\text{opt}} \blacksquare$$

What about Troll Bridge?

Troll Bridge is a problem that punishes exploration. You can't represent it in this framework because the action screens off exploration. But if you modified the framework to implement Troll Bridge, the agent would fail. (It would be afraid to cross the bridge.)

I suspect the solution to Troll Bridge is to use *continuous exploration*: Explore at every timestep of a logical inductor. If Troll Bridge punishes exploration at timestep n , then exploration at $n - 1$ succeeds. I haven't worked this out yet.

What about counterfactual mugging?

This agent doesn't get counterfactually mugged. (This is bad.) This agent only exhibits one of the [two types of updatelessness](#).

Next steps

The notion of "thin counterfactual" used here is more specific than it needs to be. I'd like to generalize it.

Once I do that, I might be able to formulate the symmetric Prisoner's Dilemma. I'm not sure how this agent does on that problem.

I also want to implement this agent in the setting of logical induction. I'm imagining traders getting licenses to participate in thin counterfactual markets, and the licenses prohibit them from learning the price of $A = a$. This will also involve some kind of continuous exploration.

A Process for Dealing with Motivated Reasoning

Epistemic status: Had a [valley of bad rationality](#) problem, did some thinking about how to solve it in my own case, thought the results might be useful for other people who have the same problem. (To be clear, I don't yet claim this solution works, I haven't really tried it yet, but it seems like the kind of thing that might work. The claim here is "this might be worth trying" rather than "this works" and the target audience is myself and people sufficiently similar to me)

Epistemic effort: Maybe 2-3 hours of writing cumulatively. Thanks to Elephantiskon for feedback.

EDIT: Here's a [summary](#) of the post, which I was [told](#) was much clearer than the original:

There's a thing people (including me) sometimes do, where they (unreflectively) assume that the conclusions of motivated reasoning are always wrong, and dismiss them out of hand. That seems like a bad plan. Instead, try going into System II mode and reexamining conclusions you think might be the result of motivated reasoning, rather than immediately dismissing them. This isn't to say that System II processes are completely immune to motivated reasoning, far from it, but "apply extra scrutiny" seems like a better strategy than "dismiss out of hand."

This habit of [automatically dismissing anything that seems like it might be the result of motivated reasoning] can lead to decision paralysis and pathological self-doubt. The point of this post is to correct for that somewhat.

It sometimes seems like a substantial fraction of my reasoning is driven by an awareness of the insidiousness of motivated reasoning, and a desire to avoid it. In particular, I sometimes have thoughts like the following:

Brain: I would like to go to philosophy graduate school.

Me: But 80,000 hours [says](#) it's usually not a good idea to go to philosophy graduate school...

Brain: But none of the other options I can come up with seem especially realistic. Combine that with the fact that grad school [can be made to be a good path if you do it right](#), and it actually seems like a pretty good option.

Me: But I started off wanting to go to graduate school because it seemed like fun. Seems pretty suspicious that it would turn out to be my best option, despite the fact that, according to 80k, for most EAs it's not. Are you sure you're not engaging in motivated reasoning?

Brain: Are you sure you're not engaging in motivated reasoning? Are you sure you're not just trying to make a decision that's socially defensible to our in-group (other EAs)?

Um, what??

I seem to be reasoning as if there were a general principle that, if there's a plausible way that I might be using motivated reasoning to come to a particular conclusion, that conclusion must be wrong. In other words, my brain has decided that anything tagged as "probably based on motivated reasoning" is false. Another way of thinking about this is that I'm using "that's probably based on motivated reasoning" as a [fully_general excuse](#) against myself.

While being averse to motivated reasoning seems reasonable, the general principle that any conclusion arrived at by motivated reasoning must be false seems ridiculous when I write it out. Obviously, my coming to a conclusion by motivated reasoning doesn't have any effect on whether or not it's true -- it's [already true or already false](#), no matter what sort of reasoning I used.[1]

A better process for dealing with motivated reasoning might be:

If

1. Getting the right (true) answer matters in a particular case, and
2. There's a plausible reason to suspect that I might be coming to my answer in that case on the basis of motivated reasoning,

then it is worth it to:

- a. Go into system-II/slow thinking/manual mode/whatever.
- b. Ask yourself what you would do if your (potentially motivated-reasoning-generated) conclusion were true, and what you would do if it were false (c.f. [leave a line of retreat, split-and-commit, see the dark world](#)).[2]
- c. *Use explicit, gears-based, models-based reasoning to check my conclusion.* (e.g. list out all important considerations in a doc, if the problem is quantitative [make a spreadsheet](#), etc.)

Then, whatever answer comes out of that process, [trust it](#) until new information comes along, then rerun the process.

To sum up: if you have a habit of dismissing a belief when you notice it might be the result of motivated reasoning, it might be worth it to replace that habit with the habit of reevaluating the belief instead.

[1]To be clear, I do think the basic idea that [if something seems to be the result of motivated reasoning, that's evidence against it] is probably correct. I just think that you shouldn't update all the way to *this is false*, since the thing might still be true.

[2]I think the basic idea behind why reasoning hypothetically in this way helps is this: it takes the focus off of deciding whether X is true (which is the step that's suspect) and puts it onto deciding what that would lead to. I like to think of it as first "fixing" X as true, and then "fixing" X as false.

Deep learning - deeper flaws?

This is a linkpost for <http://thinkingcomplete.blogspot.com/2018/03/deep-learning-deeper-flaws.html>

In this post I summarise four lines of argument for why we should be skeptical about the potential of deep learning in its current form. I am fairly confident that the next breakthroughs in AI will come from some variety of neural network, but I think several of the objections below are quite a long way from being overcome.

Theoretical Impediments to Machine Learning With Seven Sparks from the Causal Revolution - [Pearl, 2018](#)

Pearl describes three levels at which you can make inferences: association, intervention, and counterfactual. The first is statistical, identifying correlations - this is the level at which deep learning operates. The intervention level is about changes to the present or future - it answers questions like "What will happen if I do y?" The counterfactual level answers questions like "What would have happened if y had occurred?" Each successive level is strictly more powerful than the previous one: you can't figure out what the effects of an action will be just on the association level, without a causal model, since we treat actions as interventions which override existing causes. Unfortunately, current machine learning systems are largely model-free.

Causal assumptions and conclusions can be encoded in the form of graphical models, where a directed arrow between two nodes represents a causal influence. Constraints on the structure of a graph can be determined by seeing which pairs of variables are independent when controlling for which other variables: sometimes controlling removes dependencies, but sometimes it introduces them. Pearl's main claim is that this sort of model-driven causal analysis is an essential step towards building human-level reasoning capabilities. He identifies several important concepts - such as counterfactuals, confounding, causation, and incomplete or biased data - which his framework is able to reason about, but which current approaches to ML cannot deal with.

Deep Learning: A Critical Appraisal - [Marcus, 2018](#)

Marcus identifies ten limitations of current deep learning systems, and argues that the whole field may be about to hit a wall. According to him, deep learning:

1. Is data hungry - it can't learn abstractions through explicit verbal definition like humans can, but instead requires thousands of examples.
2. Is shallow, with limited capacity for transfer. If a task is perturbed even in minor ways, deep learning breaks, demonstrating that it's not really learning the underlying concepts. Adversarial examples showcase this effect.
3. Has no natural way to deal with hierarchical structure. Even recursive neural networks require fixed sentence trees to be precomputed. See my summary of

'Generalisation without systematicity' below.

4. Struggles with open-ended inference, especially based on real-world knowledge.
5. Isn't transparent, and remains essentially a "black box".
6. Is not well-integrated with prior knowledge. We can't encode our understanding of physics into a neural network, for example.
7. Cannot distinguish causation from correlation - see my summary of Pearl's paper above.
8. Presumes a largely stable world, like a game, instead of one like our own in which there are large-scale changes.
9. Is vulnerable to adversarial examples, which can be constructed quite easily.
10. Isn't robust as a long-term engineering solution, especially on novel data.

Some of these problems seem like they can be overcome without novel insights, given enough engineering effort and compute, but others are more fundamental. One interpretation: deep learning can interpolate within the training space, but can't extrapolate to outside the training space, even in ways which seem natural to humans. One of Marcus' examples: when a neural network is trained to learn the identity function on even numbers, it rounds down on odd numbers. In this trivial case we can solve the problem by adding odd training examples or manually adjusting some weights, but in general, when there are many features, both may be prohibitively difficult even if we want to make a simple adjustment. To address this and other problems, Marcus offers three alternatives to deep learning as currently practiced:

1. Unsupervised learning, so that systems can constantly improve - for example by predicting the next time-step and updating afterwards, or else by setting itself challenges and learning from doing them.
2. Further development of symbolic AI. While this has in the past proved brittle, the idea of integrating symbolic representations into neural networks has great promise.
3. Drawing inspiration from humans, in particular from cognitive and developmental psychology, how we develop commonsense knowledge, and our understanding of narrative.

Generalisation without systematicity - [Lake and Baroni, 2018](#)

Lake and Baroni identify that human language and thought feature "systematic compositionality": we are able to combine known components in novel ways to produce arbitrarily many new ideas. To test neural networks on this, they introduce SCAN, a language consisting of commands such as "jump around left twice and walk opposite right thrice". While they found that RNNs were able to generalise well on new strings similar in form to previous strings, performance dropped sharply in other cases. For example, the best result dropped from 99.9% to 20.8% when the test examples were longer than any training example, even though they were constructed using the same compositional rules. Also, when a command such as "jump" had only been seen by itself in training, RNNs were almost entirely incapable of understanding instructions such as "turn right and jump". The overall conclusion: that neural networks can't extract systematic rules from training data, and so can't generalise compositionality anything like how humans can. This is similar to the result of a project I recently carried out, in which I found that capsule networks which had been trained to recognise transformed inputs such as rotated digits and digits with negative

colours still couldn't recognise rotated, negated digits: they were simply not learning general rules which could be composed together.

Deep reinforcement learning doesn't work yet

- Irpan, 2018

Irpan runs through a number of reasons to be skeptical about using deep learning for RL problems. For one thing, deep RL is still very data-inefficient: DeepMind's Rainbow DQN takes around 83 hours of gameplay to reach human-level performance on an Atari game. By contrast, humans can pick them up within a minute or two. He also points out that other RL methods often work better than deep RL, particularly model-based ones which can utilise domain-specific knowledge.

Another issue with RL in general is that designing reward functions is difficult. This is a theme in AI safety - specifically when it comes to reward functions which encapsulate human values - but there are plenty of existing examples of reward hacking on much simpler tasks. One important consideration is the tradeoff between shaped and sparse rewards. Sparse rewards only occur at the goal state, and so can be fairly precise, but are usually too difficult to reach directly. Shaped rewards give positive feedback more frequently, but are easier to hack. And even when shaped rewards are designed carefully, RL agents often find themselves in local optima. This is particularly prevalent in multi-agent systems, where each agent can overfit to the behaviour of the other.

Lastly, RL is unstable in a way that supervised learning isn't. Even successful implementations often fail to find a decent solution 20 or 30% of the time, depending on the random seed with which they are initialised. In fact, there are very few real-world success stories featuring RL. Yet achieving superhuman performance on a wide range of tasks is a matter of when, not if, and so I think Amara's law applies: we overestimate the effects RL will have in the short run, but underestimate its effects in the long run.

Alignment Newsletter #23

Crossposted from the [AI Alignment Forum](#). May contain more technical jargon than usual.

Highlights

[Visual Reinforcement Learning with Imagined Goals](#) (*Vitchyr Pong and Ashvin Nair*): This is a blog post explaining a paper by the same name that I covered in [AN #16](#). It's particularly clear and well-explained, and I continue to think the idea is cool and interesting. I've recopied my summary and opinion here, but you should read the blog post, it explains it very well.

Hindsight Experience Replay ([HER](#)) introduced the idea of accelerating learning with sparse rewards, by taking trajectories where you fail to achieve the goal (and so get no reward, and thus no learning signal) and replacing the actual goal with an "imagined" goal chosen in hindsight such that you actually achieved that goal, which means you get reward and can learn. This requires that you have a space of goals such that for any trajectory, you can come up with a goal such that the trajectory achieves that goal. In practice, this means that you are limited to tasks where the goals are of the form "reach this goal state". However, if your goal state is an image, it is very hard to learn how to act in order to reach any possible image goal state (even if you restrict to realistic ones), since the space is so large and unstructured. The authors propose to first learn a structured latent representation of the space of images using a variational autoencoder (VAE), and then use that structured latent space as the space of goals which can be achieved. They also use Q-learning instead of DDPG (which is what HER used), so that they can imagine any goal with a minibatch (s, a, s') and learn from it (whereas HER/DDPG is limited to states on the trajectory).

My opinion: This is a cool example of a relatively simple yet powerful idea -- instead of having a goal space over all states, learn a good latent representation and use that as your goal space. This enables unsupervised learning in order to figure out how to use a robot to generally affect the world, probably similarly to how babies explore and learn.

[Impact Measure Desiderata](#) (*TurnTrout*): This post gives a long list of desiderata that we might want an impact measure to satisfy. It considers the case where the impact measure is a second level of safety, that is supposed to protect us if we don't succeed at value alignment. This means that we want our impact measure to be agnostic to human values. We'd also like it to be agnostic to goals, environments, and representations of the environment. There are several other desiderata -- read the post for more details, my summary would just be repeating it.

My opinion: These seem like generally good desiderata, though I don't know how to formalize them to the point that we can actually check with reasonable certainty whether a proposed impact measure meets these desiderata.

I have one additional desideratum from impact measures. The impact measure alone should disallow all extinction scenarios, while still allowing the AI system to do most of the things we use AI for today. This is rather weak, really I'd want AI do more tasks than are done today. However, even in this weak form, I doubt that we can satisfy this desideratum if we must also be agnostic to values, goals, representations and

environments. We could have valued human superiority at game-playing very highly, in which case building AlphaGo would be catastrophic. How can an impact measure allow that without being at least some knowledge about values?

[Recurrent World Models Facilitate Policy Evolution](#) (*David Ha et al*): I read the [interactive version](#) of the paper. The basic idea is to do model-based reinforcement learning, where the model is composed of a variational auto-encoder that turns a high-dimensional state of pixels into a low-dimensional representation, and a large RNN that predicts how the (low-dimensional) state will evolve in the future. The outputs of this model are fed into a very simple linear controller that chooses actions. Since the controller is so simple, they can train it using a black box optimization method (an evolutionary strategy) that doesn't require any gradient information. They evaluate on a racing task and on Doom, and set new state-of-the-art results. There are also other interesting setups -- for example, once you have a world model, you can train the controller completely within the world model without interacting with the outside world at all (using the number of timesteps before the episode ends as your reward function, since the world model doesn't predict standard rewards, but does predict whether the episode ends). There are a lot of cool visualizations that let you play with the models trained with their method.

My opinion: I agree with [Shimon Whiteson's take](#), which is that this method gets improvements by creating a separation of concerns between modelling the world and learning a controller for the model, and evaluating on environments where this separation mostly holds. A major challenge in RL is learning the features that are important for the task under consideration, and this method instead learns features that allow you to reconstruct the state, which could be very different, but happen to not be different in their environments. That said, I really like the presentation of the paper and the fact that they did ablation studies.

Previous newsletters

[Model Reconstruction from Model Explanations](#) (*Smitha Milli et al*): Back in [AN #16](#), I said that one way to prevent model reconstruction from gradient-based explanations was to add noise to the gradients. Smitha pointed out that the experiments with SmoothGrad are actually of this form, and it still is possible to recover the full model, so even adding noise may not help. I don't really understand SmoothGrad and its relationship with noise (which is chosen to make a saliency map look nice, if I understand correctly) so I don't know exactly what to think here.

Technical AI alignment

Agent foundations

[When wishful thinking works](#) (*Alex Mennen*): Sometimes beliefs can be loopy, in that the probability of a belief being true depends on whether you believe it. For example, the probability that a placebo helps you may depend on whether you believe that a placebo helps you. In the situation where you know this, you can "wish" your beliefs to be the most useful possible beliefs. In the case where the "true probability" depends continuously on your beliefs, you can use a fixed point theorem to find a consistent set of probabilities. There may be many such fixed points, in which case you can

choose the one that would lead to highest expected utility (such as choosing to believe in the placebo). One particular application of this would be to think of the propositions as "you will take action a_i ". In this case, you act the way you believe you act, and then every probability distribution over the propositions is a fixed point, and so we just choose the probability distribution (i.e. stochastic policy) that maximized expected utility, as usual. This analysis can also be carried to Nash equilibria, where beliefs in what actions you take will affect the actions that the other player takes.

[Counterfactuals and reflective oracles \(Nisan\)](#)

Learning human intent

[Cycle-of-Learning for Autonomous Systems from Human Interaction \(Nicholas R. Waytowich et al\)](#): We've developed many techniques for learning behaviors from humans in the last few years. This paper categorizes them as learning from demonstrations (think imitation learning and IRL), learning from intervention (think [Safe RL via Human Intervention](#)), and learning from evaluation (think [Deep RL from Human Preferences](#)). They propose running these techniques in sequence, followed by pure RL, to train a full system. Intuitively, demonstrations are used to jumpstart the learning, getting to near-human performance, and then intervention and evaluation based learning allow the system to safely improve beyond human-level, since it can learn behaviors that humans can't perform themselves but can recognize as good, and then RL is used to improve even more.

My opinion: The general idea makes sense, but I wish they had actually implemented it and seen how it worked. (They do want to test in robotics in future work.) For example, they talk about inferring a reward with IRL from demonstrations, and then updating it during the intervention and evaluation stages. How are they planning to update it? Does the format of the reward function have to be the same in all stages, and will that affect how well each method works?

This feels like a single point in the space of possible designs, and doesn't include all of the techniques I'd be interested in. What about active methods, combined with exploration methods in RL? Perhaps you could start with a hand-specified reward function, get a prior using [inverse reward design](#), start optimizing it using RL with curiosity, and have a human either intervene when necessary (if you want safe exploration) or have the RL system actively query the human at certain states, where the human can respond with demonstrations or evaluations.

[Sample-Efficient Imitation Learning via Generative Adversarial Nets \(Lionel Blondé et al\)](#)

[A Roadmap for the Value-Loading Problem \(Lê Nguyêñ Hoang\)](#)

Preventing bad behavior

[Impact Measure Desiderata \(TurnTrout\)](#): Summarized in the highlights!

Handling groups of agents

[Reinforcement Learning under Threats \(Víctor Gallego et al\)](#): Due to lack of time, I only skimmed this paper for 5 minutes, but my general sense is that it takes MDPs and

turns them into two player games by positing the presence of an adversary. It modifies the Bellman update equations to handle the adversary, but runs into the usual problems of simulating an adversary that simulates you. So, it formalizes level-k thinking (simulating an opponent that thinks about you at level k-1), and evaluates this on matrix games and the friend-or-foe environment from [AI safety gridworlds](#).

My opinion: I'm not sure what this is adding over two-player game theory (for which we can compute equilibria) but again I only skimmed the paper so it's quite likely that I missed something.

Near-term concerns

Adversarial examples

[Adversarial Reprogramming of Sequence Classification Neural Networks](#) (*Paarth Neekhara et al*)

Fairness and bias

[Introducing the Inclusive Images Competition](#) (*Tulsee Doshi*): The authors write, "this competition challenges you to use Open Images, a large, multilabel, publicly-available image classification dataset that is majority-sampled from North America and Europe, to train a model that will be evaluated on images collected from a different set of geographic regions across the globe". The results will be presented at NIPS 2018 in December.

My opinion: I'm really interested in the techniques and results here, since there's a clear, sharp distribution shift from the training set to the test set, which is always hard to deal with. Hopefully some of the entries will have general solutions which we can adapt to other settings.

AI strategy and policy

[Podcast: Artificial Intelligence – Global Governance, National Policy, and Public Trust with Allan Dafoe and Jessica Cussins](#) (*Allan Dafoe, Jessica Cussins, and Ariel Conn*): Topics discussed include the difference between AI governance and AI policy, externalities and solving them through regulation, whether governments and bureaucracies can keep up with AI research, the extent to which the US' policy of not regulating AI may cause citizens to lose trust, labor displacement and inequality, and AI races.

Other progress in AI

Reinforcement learning

[Visual Reinforcement Learning with Imagined Goals](#) (*Vitchyr Pong and Ashvin Nair*): Summarized in the highlights!

[Recurrent World Models Facilitate Policy Evolution](#) (*David Ha et al*):
Summarized in the highlights!

[ARCHER: Aggressive Rewards to Counter bias in Hindsight Experience Replay](#)
(*Sameera Lanka et al*)

[SOLAR: Deep Structured Latent Representations for Model-Based Reinforcement Learning](#) (*Marvin Zhang, Sharad Vikram et al*)

[Expt-OOS: Towards Learning from Planning in Imperfect Information Games](#) (*Andy Kitchen et al*)

Miscellaneous (AI)

[Making it easier to discover datasets](#) (*Natasha Noy*): Google has launched Dataset Search, a tool that lets you search for datasets that you could then use in research.

My opinion: I imagine that this is primarily targeted at data scientists aiming to learn about the real world, and not ML researchers, but I wouldn't be surprised if it was helpful for us as well. MNIST and ImageNet are both present, and a search for "self-driving cars" turned up some promising-looking links that I didn't investigate further.

Book review: Why we sleep

This is a linkpost for <http://thinkingcomplete.blogspot.com/2018/08/book-review-why-we-sleep.html>

I read this book (by a sleep scientist called Matthew Walker) because I knew that it would tell me to sleep more, and I hoped it would cite enough scary statistics that I'd be likely to actually follow through. Well, it worked - I'm keeping a copy on my bedside table for the foreseeable future, just as a reminder. In addition to the exhortations to get more sleep, it contains a variety of other interesting and important facts about sleep.

What is sleep?

- Human sleep consists of cycles lasting about 1.5 hours, each of which contains first a period of NREM (Non-Rapid Eye Movement) sleep, then a period of REM sleep. In brain scans, the former consists of slow, deep brain waves, while the latter shows the same frenetic activity as an awake brain. As the night goes on, cycles feature a higher proportion of REM sleep. This means that if you cut your sleep short by 25%, you're actually missing out on somewhere between 60% and 90% of REM sleep.
- REM sleep is when the majority of dreams happen. While it's uncommon for dreams to replay events from our everyday lives, they do often reflect our emotional preoccupations. To prevent ourselves from flailing around during dreams, we enter a state of sleep paralysis, where our brains are unable to control our voluntary muscles. Eyes are an exception - hence the name REM. It's definitely not true that REM is the only valuable type of sleep - in fact, immediately after sleep deprivation the brain prioritises catching up on NREM.
- The slow waves of NREM sleep are useful for transferring memories from one part of the brain to the other - in particular, from short- to long-term storage.
- Walker's theory is that NREM sleep is used to prune away unnecessary connections, while REM reinforces useful connections. He uses the analogy of a sculptor who alternates between carving away whole chunks of marble (NREM) and then adding fine detail on whatever's left (REM). From this perspective, it makes sense that REM sleep is concentrated in later cycles. However, it's unclear whether this is the scientific consensus.
- There are two systems controlling sleep and wakefulness. The circadian system follows the day/night cycle, making you tired in the evening and alert in the morning (the exact timings vary by person, making some people "night owls" and some "morning larks"). In addition, "sleep pressure" is controlled by adenosine, which builds up while you're awake and is cleared away during sleep. Caffeine works by temporarily blocking adenosine receptors, but doesn't prevent it from continuing to build up.

What's it good for?

- There's a very strong link between NREM sleep and memory. The formation of long-term memories suffers if we don't get enough sleep (even several days after the events we want to remember). This is true both for memories about facts and experiences and for "muscle memory" of actions like playing an instrument. When sleep-deprived, we also have worse short-term memory.

- REM sleep is important in emotional regulation and creativity. After sleep deprivation, the responses of the amygdala (responsible for strong emotions) can be amplified by over 60%, due to weakened links between it and the prefrontal cortex (responsible for "rational" decision-making). Dreams during REM sleep allow us to make unusual and creative connections between different topics - many great intellectuals report that their best ideas just "came to them" upon waking.
- Sleep deprivation massively reduces our ability to concentrate. In addition to slower reaction times, when tired we lapse into "micro-sleeps" during which we're totally unresponsive. Walker emphasises that tiredness is a far bigger cause of traffic accidents than drunk-driving, that drivers systematically underestimate how tired they are, and that drivers who micro-sleep often don't brake at all before collisions.
- In the long term, sleep deprivation increases the risk of Alzheimer's (since toxins are flushed from the brain during sleep), heart attacks (by provoking a stress response from the sympathetic nervous system and raising blood pressure) and cancer (by devastating the immune system). All of these seem to be very big effects - e.g. sleep-deprived patients are twice to three times as likely to suffer calcification of their coronary arteries.
- Note that most of the effects above are noticeable even after small amounts of sleep deprivation, like getting one or two hours less sleep for one or two nights. In fact, even the one-hour sleep reduction from Daylight Savings Time [causes a spike in heart attacks](#).
- Sleep is also linked to many mental illnesses - e.g sleep deprivation triggers mania or depression in bipolar patients. Most mental illnesses disrupt sleep, which exacerbates their other negative effects.
- REM sleep promotes the formation of neural links in infants, who have far more neural connections than adults. It is also important for their language learning.
- Walker's broad answer to the question of what sleep is useful for: EVERYTHING. In addition to the above, sleep helps us overcome traumatic memories, reduces athletes' injury rates, makes us look more attractive, reduces food cravings, and so on and so on...

The evolution of sleep

- I guess it shouldn't be a surprise that sleep is so broadly useful: once it started, it makes sense that many metabolic processes would take advantage of it. And they've had a long time to do so: sleep is ancient, with all animal species demonstrating some form of sleep-like behaviour.
- Even unicellular bacteria have active and passive phases corresponding to the planet's light/dark cycle.
- However, the length of sleep required varies wildly for different animals, from 4 hours for elephants to 19 for brown bats.
- Only birds and mammals have proper REM sleep - it is a relatively recent adaptation. It also seems to be absent in aquatic mammals, whose two brain hemispheres sleep separately.
- Humans seem to be naturally biphasic: modern hunter-gatherer tribes sleep for 7-8 hours at night, and then nap for 30-60 minutes in the afternoon. It's biologically natural to be sleepy after lunch. Biphasic sleep significantly decreases mortality from heart disease.
- Walker hypothesises that descending from the trees to sleep on the ground allowed us to gain more REM sleep (particularly difficult in trees due to sleep paralysis), and therefore was important in boosting human cognitive development; also, that fire was vital in making ground-sleeping safer.

How to sleep better

- Alcohol is an extremely powerful suppressor of REM sleep. Since it stays in your system for hours, it's best not to drink in the evenings.
- Light, especially blue light, signals your circadian system to wake up. Unfortunately LED screens provide a lot of blue light. Avoid using screens in the hours before bed, or at least phase out the blue light (e.g. using [flux](#)).
- In addition to light, our bodies use decreasing temperatures as a signal to sleep. Lowering room temperature often helps with insomnia. Apparently your core temperature will also fall after a hot bath.
- Caffeine has a half-life of 5 to 7 hours, so if you drink it in the afternoon, a significant amount will still be in your system at bedtime.
- The circadian rhythm of a teenager is naturally a few hours later than that of an adult, so teens shouldn't be forced to get up too early. Unfortunately, schools aren't taking much notice of this.
- Apparently sleeping pills cause lower-quality sleep and have severe long-term side-effects, so they should be avoided (with the exception of melatonin).
- For serious sleep problems, Cognitive Behavioural Therapy for Insomnia (CBT-I) works fairly well and should be the first step.

As you can probably tell from the above, Walker is very much a cheerleader for sleep. This does bias him in some noticeable ways - e.g. his overt scorn towards coffee. He also blurs causation and correlation at some points throughout the book, so I'd be surprised if all of the deleterious effects mentioned above are as significant as he claims. But the overall picture is stark enough that I'm now very worried about the ongoing sleep loss epidemic.

In Logical Time, All Games are Iterated Games

Crossposted from the [AI Alignment Forum](#). May contain more technical jargon than usual.

Logical Time

The main purpose of this post is to introduce the concept of logical time. The idea was mentioned in Scott's post, [Bayesian Probability is for things that are Space-like Separated from You](#). It was first coined in a conference call with, Daniel Demski, Alex Mennan, and perhaps Corey Staten and Evan Lloyd -- I don't remember exactly who was there, or who first used the term. Logical time is an informal concept which serves as an intuition pump for thinking about logical causality and phenomena in logical decision theory; don't take it too seriously. In particular, I am not interested in anybody trying to formally define logical time (aside from formal approaches to logical causality). Still, it seems like useful language for communicating decision-theory intuitions.

Suppose you are playing chess, and you consider moving your bishop. You play out a hypothetical game which results in your loss in several moves. You decide not to move your bishop as a result of this. The hypothetical game resulting in your loss still exists within logic. You are logically later than it, in that the game you actually play depends on what happened in this hypothetical game.

Suppose you're stuck in the desert in a [Parfit's Hitchhiker](#) problem. Paul Ekman is reading your face, deciding whether you're trustworthy. Paul Ekman does this based on experience, meaning that the computation which is *you* has a strong similarity with other computations. This similarity can be used to predict you fairly reliably, based on your facial expressions. What creates this similarity? According to the logical time picture, there is a logical fact much earlier in logical time, which governs the connection between facial expressions and behavior.

To the extent that agents are trying to predict the future, they can be thought of as trying to place themselves later in logical time than the events which they're trying to predict. Two agents trying to predict each other are competing to see who can be later in logical time. This is not necessarily wise; in games like chicken, there is a sense in which you want to be earlier in logical time.

Traditional game theory, especially Nash equilibria, relies on what amounts to loopy logical causality to allow each agent to be after the other in logical time. Whether this is bad depends on your view on logical time travel. Perhaps there is a sense in which logical time can be loopy, due to prediction (which is like logical time travel). Perhaps logical time can't be loopy, and this is a flaw in the models used by traditional game theory.

Iterated Games

In logical time, all games are iterated games. An agent tries to forecast what happens in the decision problem it finds itself in by comparing it to similar decision problems which are small enough for it to look at. This puts it later in logical time than the small examples. "Similar games" includes the exact same game, but in which both players have had less time to think.

This means it is appropriate to use iterated strategies. Agents who are aware of logical time can play tit-for-tat in single-shot Prisoner's Dilemma, and so, can cooperate with each other.

Iterated games are different in character than single-shot games. The [folk theorem](#) shows that almost any outcome is possible in iterated play (in a certain sense). This makes it difficult to avoid very bad outcomes, such as nearly always defecting in the prisoner's dilemma, despite the availability of much better equilibria such as tit-for-tat. Intuitively, this is because (as Yoav Shoham et al point out in [If multi-agent learning is the answer, what is the question?](#)) it is difficult to separate "teaching behavior" from "learning behavior": as in the tit-for-tat strategy, it is generally wise to adopt behavior designed to shape the incentive gradient of the other player, in addition to improving your own score. Unfortunately, it is difficult to say what it means to pursue these two objectives simultaneously.

The subfield most related to this problem is multi-agent learning. Sadly, as discussed in the [Shoham et al paper](#) I cited parenthetically in the preceding paragraph, multi-agent learning typically avoids the difficulty by focusing on learning single-shot equilibria via iterated play. Learning single-shot equilibria in an iterated setting is a somewhat weird thing to be doing (hence the title of the paper). However, it is understandable that people might avoid such a difficult problem. The folk theorem illustrates a severe equilibrium selection issue, meaning that traditional tools have little to say about rational play.

One might imagine learning to play single-shot games by playing them over and over. But, what can you do to learn iterated games? You might imagine that you jump up a level again, experiencing the iterated version repeatedly to discover the optimal iterated strategy. However, iterating the game more doesn't really escape the iterated setting; there is no further level!

(You might think meta-iteration involves making the other player forget what it learned in iterated play so far, so that you can re-start the learning process, but that doesn't make much sense if you retain your own knowledge; and if you don't, you can't be learning!)

Toy Models

We can make pictures of logical time using phenomena which we understand more fully. One such picture is based on proofs. If we imagine a theorem prover proving every theorem in some order (such as an ordering based on proof length), we can think of logical time as time-of-proof. We can formulate [counterfactuals consistent with this notion of logical time](#). (As I mentioned before, a picture of logical time is just a picture of logical causality / logical counterfactuals -- the notion of logical time adds nothing, formally.)

We can examine logical time travel in this kind of model by [constructing predictors using stronger logics](#), which allows a predictor to find shorter proofs. This creates

decision-theoretic puzzles, because the agent with the weaker logic can't recognize the loopiness of the situation; it thinks it cannot influence the predictor, because (according to its weaker logic) the predictor has a short proof-length and is therefore earlier in logical time. We, on the other hand, can recognize that agents who act as if they control the predictor could do better in the decision problem.

This weirdness seems to only be possible because of the "two dimensional logical time" which exists in this toy model, in which we can vary both proof length and logical strength. One agent has access to arbitrarily long proofs via oracles, and so is "later" in the length dimension; the other has a stronger logic, and so is "later" in the strength dimension.

However, we can collapse the two dimensions into one via logical induction. Logical induction eventually learns to predict what stronger logics would predict, so computation time and logical strength are more or less the same.

You might expect that the loopy scenario in which an agent and a predictor accurately predict each other becomes impossible in logical induction, but, it does not. Logical-induction agents can predict each other well by examining what similar agents do in similar situations. As a result, [LIDT agents converge to playing correlated equilibria with each other, more or less](#). (This result ignores the iterated aspect of the games, just like the multi-agent learning approaches I was complaining about earlier; despite learning from all the nearby copies within logic, the LIDT agents think only of the utility for their one decision, which paradoxically results in poorer outcomes even for that decision. Asymptotic decision theory does better, but no nice results for game theory have come from it so far.)

So long as an agent eventually settles down to making some reliable pattern of decisions in a situation, there will be relatively young logical inductors which have learned enough to accurately forecast the decisions made by logical-induction agents who reason using much more computational power.

We can think of the purely logical case, with its apparent two dimensions of logical time, as being a degenerate extreme version of the phenomenon in logical induction. In logical induction, the early predictions may be quite accurate, but they are fallible; they always run the risk of being wrong, since we're in the process of learning. In the pure logical case, we *also* run the risk of being wrong: using a stronger logic to make predictions runs the risk of introducing inconsistencies. This is easy to forget, since we are accustomed to the assumption that we can easily add axioms to a consistent system to get a stronger one.

An early predictor predicting a late agent must give up on some accuracy -- a prediction which relies on anything else than actually running the computation to be predicted has some chance of failure. This breaks the loopiness in logical time; the late agent always adds some small amount of information, even if its action is predicted with high reliability.

Thoughts on tackling blindspots

I went to my first CFAR workshop the other week, and it was quite intense/a lot. The biggest change by far has been that I came face to face with some huge blindspots and can see more clearly many of the ways I've been fooling myself, not allowing myself to care about things, and pushing people away. Since blindspots are a proto-typical, "How the fuck are you supposed to find a thing that you aren't capable of finding?" I wanted to share what I think are things that helped me spot some of mine.

This is rough draft mode, and I think I'm going to just settle for bullet points in this pass-through.

To quote Draco from HPMOR:

To figure out a strange plot, look at what happens, then ask who benefits

i.e. look at all of the first impression I make of people, and get suspicious that they all add up to "People aren't worth talking to" and be suspicious.

- Combine that with some [cognitive trope therapy](#).
- One of the first things I got into my head when journeying into rationality was "If it hurts to think about, or feels like a cherished belief that doesn't want to be touched, GO AFTER IT!" This had the effect of most of my problems disguising themselves as things I wasn't interested in. Instead of feeling scared of parties, I would just feel disinterested and bored by the idea of parties.
- In Val's Design class, being reminded that most of your built/learned mental machinery came into being to try and protect you for something or get something for you. Also being reminded about how you can't just steamroll over your previous machinery. There is a need or want hiding there which unless you address it, the machinery in place that was trying to serve that want will fight back.
- A mental shift from, "Is this plausible/reasonable?" to "Is this true?" when asking examining my rejections of different ideas.

Psycho-cybernetics: experimental notes

This post is a series of missives and notes I took while reading [a popularization of cybernetics concepts as applied to self-help](#) that was hugely influential in the self help field when first published in 1960. I am unsure if these notes will be of any interest to others. This is not a book review or a summary, but rather my own impressions of the models that the author was trying to build up and the cross connections between those concepts and others.

In general, I wish more people would make posts about books without feeling the need to do boring parts they are uninterested in (summarizing and reviewing) and more just discussing the ideas they found valuable. I think this would lower the friction for such posts, resulting in more of them. I often wind up finding such thoughts and comments about non-fiction works by LWers pretty valuable. I have more of these if people are interested.

Why this book: If you wish to understand the box you live in, investigate records from the time it was being built. The social psychology and cognitive science results that much of the lesswrong memeplex hangs its hat on are subject to an incentive structure whereby surprising results are the ones that are promoted or made more visible. But surprising relative to what? Is there some generic folk psychology template that I am comparing to? This plays some role, but I think I have underestimated the degree to which the defaults are constructs. I wanted to get a sense of how they might have been built, which lead to an investigation of Alfred Korzybski, the first person to utter "the map is not the territory", and the inception of cybernetics as a field of discipline, which also heavily influenced the people and work that later went on at Bell Labs and thus shaped the emergence of the information age.

I found this book interesting in particular because it did not use the standard anecdote-concept format of most self-help (with perhaps 2-8 concepts in an entire book) but instead seemed much more concept dense. I do recommend reading the whole thing if these notes seem interesting.

Preface

- You should do the thing that I am doing in this post with the book: generate notes from each chapter and tie them to your own concepts and experiences in order to get them to stick. The book contains blank sections for you to do this. (Interesting in that this isn't the first time I've seen a pedagogical model explicitly used in a book, but this is the first time it matched the system I wound up gravitating towards for non-fiction note taking)
- The inner simulator can cause updates as well, or very nearly as well as actual experiences (see visualization training results). The implication is that the inner simulator can cause the activation of the multi-sensory neurons that increase the weight your systems place on evidence. Things that activate multi-sensory neurons with fake evidence are dangerous (narrative fallacy, movies, updating on fictional evidence).
- It takes about 3 weeks to get a new behavior pattern to not feel strange.

Chapter 1

- Identity consistency effects are stronger than generally recognized, because most people do not have complete access to the self-image that generates judgments of which actions would be consistent or inconsistent with that image.
- Identity can be thought of as a gradient in a sort of confirmation bias space. It directs you to which features to find salient about a situation. It can also be thought of as a set of habits for getting rewards/avoiding harm that doesn't update very often or very easily (learned helplessness, elephant rope parable)
- Identity is more fluid than typically believed. In general, we take actions and then back-fill an identity or set of justifications consistent with that identity (rationalization, especially noticeable in certain brain dysfunction disorders, man who mistook his wife for a hat etc.)
- The hidden variable in growth mindset is whether it is being applied to inputs or outputs (well calibrated locus of control). E.g. I will get that specific job/hit that specific metric rather than I will apply to many jobs/I will provide more of the inputs to that metric. People for whom growth mindset 'didn't work' were consistently found to be applying it to specific outputs.
- People can have positive affect towards harmful features of the self-image without really being aware of it. This is likely because it is an important component of some habits or coping strategies. It will be hard to alter these harmful features until the positive habits have alternative routes towards success. (goal factoring, connection theory, coherence therapy, etc.)
- Self-trust is required for spontaneity/creativity to not be impinged/clamped down on out of fear. Things like improv training provide safe environments to start building a success spiral for this.
- Common frame: my will power against the very large forces of human nature, shame, social structures, etc.
- Low level feedback structures in the brain and nervous system have goal states that they attempt to correct towards when they deviate. Coherent goal states are mutually compatible, incoherent goal states result in fights between feedback mechanisms (essentially the same as the later perceptual control theory)
- Unrealistic self image generates more deviation signals
- Unsurprising that productivity systems fail when they are being used as a bludgeon to try to bash yourself into the same shape as some unrealistic self-image.
- Can build up less incoherent self-images based on small experiences (success spirals)

Chapter 2

- Two situations, goal is known or unknown.
- Known goals rely mostly on negative feedback (course correction) while unknown goals rely mostly on positive feedback (seeking behavior, dopamine on recognition of sought pattern)
- There's a tendency for well calibrated snap-judgments in situations with other humans get attributed to ESP and similar explanations (editorial on my part, author believes in various ESP like jackassery).
- Lack of a 'recognition' state increases ambiguity in goal directed tasks i.e. no clear picture of what an organized cupboard looks like generates friction in beginning the task due to lack of a feedback mechanism. In contrast, a clear internal picture generates lots of intermediate states to compare to.
- In general, fine grainededness in the reward signal is required if the task needs fine grainededness in the actions performed.
- Increasing negative capability makes goal navigation less punishing/aversive.

Chapter 3

- Placebo effect is strong and can be trained to be stronger (hypnosis, some buddhist practice, critching yourself)
- Mental picture as a tool is a trainable skill, commonly used for state shifting (public speaking, exercise, command mode, SNS activation in general?) but can be applied to any goal directed behavior.
- Vivid detail is the feature that increases with practice that seems to improve the effectiveness. All 5 senses (same exact claim as Mahasi noting practice in buddhism)
- For complex behaviors it is not possible to systematize improvement, too many minor details (think socializing or complex movements in sports) but change in image changes what is being compared to moment by moment in low level systems and brings small details in line.

Chapter 4

- Limiting beliefs get repeated internally more than you remember through some combination of words, images, feelings, essentially you have hypnotized yourself. (same language to later NLP's worthlessness, helplessness, hopelessness taxonomy of limiting beliefs)
- Almost everyone has a pattern that transforms 'I am inferior at X' to just 'I am inferior' after which confirmation bias takes over. Stems fundamentally from a 'should.' Somehow convinced that you should be better (how? by having been taught? by having practiced? don't automatically think that for some reason. All of this is a coping mechanism for the fact that biological determinism is somewhat true? internalized dominance hierarchy coping?)
- belief formation/refactoring happens in the relaxed state (? first I've encountered this one. Seems related to the safety of containers, psychotherapy etc.)
- Imagination is more powerful than the will. Use of willpower but image is demotivating = thing doesn't happen. Image is motivating, willpower to try to prevent it, thing will happen. (exercise, junk food)
- Maintaining as detailed a picture of the objective as possible is sufficient/more efficient than effort (top idea in your mind, Graham)
- Physical relaxation backprops to mental relaxation (the relaxation response)
- Noticing any situation or mental occurrence that causes a physical tensing is highly worth investigating.
- Relaxation technique (very reminiscent of Feldenkrais method)
 - first voluntarily relax muscles for several minutes
 - you will reach a limit of voluntary control
 - same technique as yoga nidra, slow body scan imagining parts as heavy, wet clay or concrete, sinking into the surface you are lying on
 - you or someone else could not move/lift you if they tried
 - other visualizations: deflating balloon, cut marionette strings, remembering previous relaxed state
 - relaxation is a skill that will improve with time
 - paying attention to how relaxed you are will flow out into non-relaxation times and you'll notice tension you didn't before, this grants lots of little opportunities to relax just a bit at the margin throughout day

Chapter 5

- non-conscious processes are not in adversarial relationship fundamentally. They are trying to surface relevant possible actions and feelings for the situations and beliefs you are in. (IFS, coherence therapy)
- that these beliefs were formed in the past does not mean that memory based work (trauma, much of psychotherapy, etc.) is the only or even most efficient way of dealing with them. Which isn't to say that if memories surface spontaneously they need to be ignored, but explicitly digging for them probably isn't helpful (confabulation, introspective illusion etc.)
- behavioral therapy: act as if for a few weeks.
- Bertrand Russel 'self identification is the source of unhappiness' (pure motive)
- You can't alter a belief via will, you can only replace it with a better belief.
- Juxtapose inconsistencies in beliefs in order to soften them (coherence therapy)
- Anecdotes that point at the anger/clarity pairing. Seems slightly different than loss aversion.
- Worry: excessively detailed bad outcome scenarios. Specificity makes things less likely, not more, but that's not how the brain updates.
- Positive motivation: insufficiently detailed good outcome scenarios (or one sided fantasies instead of realistic)
- Tendency to exaggerate the importance and difficulty of tasks. Think of your tasks as easy and playful and they will be (and be higher quality, creativity research)
- If 'I'm tired' works as an excuse for a part to get what it wants you shouldn't be surprised that 'I'm tired' surfaces more often in your thought stream (moral licensing) Over time this can shift into a background identity. 'I'm a low energy person'
- Moving towards positives works better than avoiding negatives as a navigation strategy because avoiding negatives results in more random flinchy movement. You don't know which chasms are worth crossing if you don't have a landmark to sight on.
- Conscious thought is better thought of as problem prioritizer, formulator, and logistics handler. The actual problem solving work is non-conscious. (! this is an interesting framing, even if imperfect) (Meshes well with problem loading followed by walks showers etc.)

Chapter 6

- Wuwei stuff, victory by surrender etc.
- Make extremely thorough preparations to solve the problem as intensely as you can, and then don't worry about it. The answer will come or it won't but forcing it literally never works, or at best gets you shitty designed by committee solutions. Useless for anything truly hard or worthwhile.
- This is not to say that great performance occurs by magic. Thorough preparation for the pianist involves deliberate practice, likewise good habits of thought can be deliberately cultivated that will affect analytical problem solving (mathematician toolbox)
- Awkward = inhibited, trying to think out actions, consequences. Other people assume goal orientation = agenda = salesman, not friend. Want something from you.
- Worry antidotes
 - 1. actually make decisions, which means once in motion responsibility for outcome is off your shoulders. Half assed decisions breed worry, if the outcome is bad you will beat yourself up for half assing it
 - 2. presence. don't try to stop drinking forever. don't drink today. If you are thinking about something that you can't do anything about in the present moment, you are priming your nervous system for a different scenario than the one you are in.
 - 3. Don't multitask, or at least don't multitask right now :) Your creative machinery can solve any problem, but if you feed it three at once it will return incoherence. That's just how it is wired. The feeling that you should be doing many things now is an upstream jamming up of your prioritization machinery. Your prioritization machinery can be taught to feed you one thing at a time and output the rest to a buffer (GTD, etc) with repeated practice.
 - 4. Setting intentions before sleep is a good way to relax and a good way to wake up with many of the solutions already worked out. Catnaps also work.
 - 5. vividly remember the details of the relaxed state from time to time.

Chapter 7

- Happiness seems directly related to sensory clarity (!) correlations between finer distinctions and pleasant mood reported
- Common frame: happiness is a reward for virtuous behavior. Reinforced by parenting styles
- Common frame: happiness is selfish
- Common frame: happiness is a future state
- Correct frame: happiness is a set of habits
- Unhappiness reflex is practiced and diligently maintained by us
- Happiness requires problems. If we meet problems with the certainty that we should not have to be doing this, that we deserve something more glamorous, that others certainly don't have to deal with these problems, we will be unhappy. This isn't to say that gratitude and presence will immediately lead to deep and lasting happiness, but weakening the unhappiness reflex is an important part of the process.
- Reactivity is a core part of unhappiness. Reflex. Waiting for the next bad thing to happen. When bad things happen as a result of positive action taken by us, it is seen as part of the set of obstacles on the way towards a worthwhile goal, rather than something that has happened. It is more likely to be seen as a challenge (narrative, locus of control)
- Finding the bad in everything vs finding the good in everything. Both miscalibrations, but one is obviously better for most non-critical situations.
- Passively waiting for happiness to come to you is reactivity.

Chapter 8

- A human being is a process. No telos, no obvious reason to instantiate any particular process. Like a bicycle, forward motion provides stability. Object level obstacles are easy. The actual roadblocks are things that disrupt a sense of progress.
- Goals is a loaded term, projects and skills is more alive
- Beliefs is loaded, instead think of other peoples output as a result of blending together the givens and their mental pictures and habits.
- 'Bravery' is loaded, instead think of acting, hypothesis testing, curiosity, proactive. Don't do 'correct' things, do whatever presents itself.
- 'Compassion' is loaded, just be present with people
- 'Self esteem' is loaded and narcissistic, stop making the story about you and see which interesting stories you can contribute to. Appreciate other people's projects and skill building.
- Errors have a purpose but once the mistake that caused the error has been corrected for the error is no longer relevant
- 'Success' is loaded, instead think of moments of happiness or special effort or being afraid and doing the thing anyway

Chapter 9

- If you think negative feedback is bad try imagining life without it. Lepers die pretty quickly.
- if you never experience the frustration of a thwarted ambition that would mean you hadn't set ambitious enough goals. If you regularly experience frustration your goal setting machinery might be out of whack. Famous example, golfer recommends not trying to hit the cup on a long putt, but hit a bathtub sized area. This allows one to relax and ultimately perform better. frustration as a problem solving method is a leftover from childhood.
- Aggression (and clarity) are good when directed well. When frustrated, the aggression will seek some other avenue of escape. Exercise is a good use. As is getting clarity on something by writing angrily.
- Insecurity from binary classifications? I am not successful vs how successful am I. Failing to clear some threshold that we keep moving the goalposts on. Defensiveness is less effective, reactive.
- Loneliness might be excessive identification with a mask, cutting us off from our real selves
- Uncertainty is based on the illusion that passivity insulates against harm. Inactivity is a choice just like the others, and more sure than most to lead to mediocre outcomes. When a real opportunity comes along, do you think you will regret having made lots of mistakes leading up to it? No, they were good preparation. Babe Ruth held the record for the most strikeouts.
- Resentment, as satisfaction in the sense of being 'wronged' has the same root as outrage and being offended. Equally feeds the victim mentality.
- If you think that people with negative traits are enjoying success, keep in mind that they are likely finding such success to be empty. Remember the Prime Minister: whenever I think I need to take some revenge on someone who has wronged me, I write down their name on a piece of paper. I periodically review the names and see the indignations life has heaped upon them for me.
- Focus on positives, check in with negatives

Chapter 10

- think of samskaras as emotional scar tissue. scar tissue is formed by tension around the wound pulling at it slightly.
- if you feel slighted by a small transgression, remember your own impression of thin skinned people. It is obvious from the outside that they doubt their self worth
- excessive emotional reliance on others encourages feelings of excessive (involuntary) vulnerability
- forgiveness as a way of excising emotional scars, can tell it works if you actually forget the thing and the forgiveness.
- notice the morbid enjoyment of nursing our own wounds
- forgiveness is canceling a debt, not because they have paid, but because we recognize the debt was never valid to begin with
- forgiving yourself for the debts you think you owe others
- emotions are for the present, the past is data, the future guesses. emotions about either are miscalibrations.
- if you want to get over bad habits you have to stop blaming and condemning yourself for them

Chapter 11

- excessively liberal interpretations of environmental and internal cues as negative feedback
- 'stop what you are doing and do something else!' but it isn't clear what that something else should be, therefore paralysis
- interpreting negative feedback signal as pointing to a large error
- this can be observed concretely in the study of stuttering and its treatment
- this can be thought of as excessive self monitoring, which is why we might say that anxiety is a self oriented or narcissistic disorder
- method acting poise: times you were relaxed and felt confident you could handle the situation. vivid detail is key.
- one frame: not too much self consciousness, insufficient self consciousness.
Acting the way you are when alone.
- inhibited people need disinhibition practice. speaking before they think.

Chapter 12

- Training a relaxation response unconditions and brings equanimity
- If you can not ignore the response you can still delay it
- build a mental quiet room
- it can be helpful to imagine some sort of clearing of steam, release of pressure associated with this room
- stop reacting to the past and the future and respond just to the present moment. doing nothing and relaxing is the appropriate response to a problem that isn't real
- imagine yourself doing a task with poise vividly

Chapter 13

- deliberate practice and visualization
 - (could say a lot about deliberate practice here but it's important enough to read a book and try to internalize it. So Good They Can't Ignore You and Practice Perfect are both good.)
- many people's internal witness is not their true self but some nebulous combo of parents, teachers, friends, bosses, etc.

Chapter 14

- command mode works by imagining the result clearly and letting our guidance systems take care of implementation details
- likewise, think of goals in terms of concrete present possibilities (resonance between goal-means? somehow reinforcing the link between goal and means?)
- "what would the ceo of a successful company be doing" was tapping into this.
- there are people who have a buoyant spirit, and there are people who work hard and plan, it is rare for a person to have the skill of using the buoyant energy to creatively plan and diligently move towards the goal (open mode? open mode closed mode partnership)
- vividness is predicated on details, and starting with details is good for success spirals
- success spirals work well because once feeling successful, branch factor/creativity increases which usually leads to further success
- So, if you desire something. Think about the goal, and think about what small thing you could do now that would move you towards that goal. Imagine it as vividly as you can. When you take action towards the goal, interpret this as evidence of commitment to the goal. eg if you choose to do a topic focused pomodoro, it is because you care about that topic.
- A plateau is an opportunity to drop back and begin again building up some small successes
- As with object level tasks, so too with mental phenomena. You don't force yourself to become an optimist as some giant step, you titrate small doses. "suppose good thing happened?" "good thing is possible"
- Take aversive stimuli as a challenge and moat. You are near a part of the game that turns most people back. Like a difficult series of jumps at the crux of a platforming level. You may fail, but it would be silly not to try. The alternative is to just stand there or turn back. I will succeed or I will learn. Failure is impossible.
- dont focus on driving out bad, simply focus on good and let the bad do whatever
- effort to stop worrying creates tension, tension creates a worrying atmosphere
- think of it as a balance, negative states should trigger positive ones to ensure clarity. good experiences should trigger some awareness of their passing in order to let you fully savor it.
- don't try to change the contents of thoughts, change stances

Chapter 15

- rediscovery of the energy body, CNS retraining.

Afterword

- Faith healing is likely placebo/causing relaxation response/repairing locus of control etc that prevents whatever tension and stress were blocking the healing process. People can obviously worry themselves sick and just as obviously be set straight by an authority they trust.

The Scent of Bad Psychology

This is a linkpost for <https://putanumonit.com/2018/09/07/the-scent-of-bad-psychology/>

Inspired by the [psychology replication quiz](#), I've come up with four rules to tell bullshit studies in psychology from real ones:

1. The rule of antisignificance
2. The rule of Taleb's grandma
3. The rule of multiplicity
4. The rule of silicone boobs

And finally, a reason why this could mean that there's a brighter future for good psychology research.

Track-Back Meditation

Intro -- Looking Back on Looking Back

At a CFAR event, Anna once (iirc) put forward the following idea: when you notice that you are distracted, try to think back to when you were last paying attention, recalling as much as you can of the sequence of where your attention went as it wandered. At first you might only be able to think back a couple of steps: "I was thinking about trains just now; why was that? Ah, it was because I was thinking about how I have to catch the train later, and started thinking about trains in general. I don't remember what got me on that topic." However, if you keep trying the exercise, you may eventually be able to recall all the way back to what got you off-track in the first place. Once you have a habit of recalling what got you distracted, you can start to see patterns. "Ah, I was hungry, which made me start thinking about food. Being hungry seems to distract me a lot." Then you can keep peanuts at your desk, or whatever.

Another story:

At one point a couple of years ago, I noticed that I was using a particular visual analogy to think about something, which didn't seem like a very good analogy for what I was thinking about. I don't recall the example, but, let's say I was using a mental image of trees when thinking about matrix operations. I got annoyed at the useless imaginary trees, and wondered why I was imagining them. Then, I noticed that I was physically looking at a tree! This was fairly surprising to me. Some of the surprise was that I took a random object in my visual field to use for thinking about something unrelated, but more of the surprise was that I didn't immediately know this to be the case, even when I wondered why I was imagining trees.

After I noticed this once, I started to notice it again and again: objects from my visual field end up in my imagination, and I often try to use them as visual analogies whether they're appropriate or not. It quickly became a familiar, rather than surprising, event. More interestingly, though, after a while it started to seem like a *conscious* event, rather than an automatic and uncontrollable one: I've become aware of the whole process from start to finish, and can intervene at any point if I wish.

This idea of making mental processes conscious rather than unconscious, simply by paying attention to them over time, seems interesting.

I was recently thinking about what I might do to more fully install the [Replacing Guilt](#) skillset. (I've had success with it for a single week after doing an [internal double crux](#) about it, but it didn't stick... which was sad, because it was *really nice*.) After some thinking, I settled on a short list of priorities. At the top of my list was "looking back to see how I arrived at a cognitive state" -- a skill which seems necessary for [staring into regrets](#) to work well. I remembered Anna's idea. I was initially planning to make a set of [TAPs](#) for whatever skills I decided were most critical, but in this case, practicing the skill as a meditative technique seemed appropriate.

Why Try This?

Besides helping you [stare into regrets](#) and reach [the strategic level](#), and bringing unconscious processes into conscious awareness over time, the skill of tracing back your thoughts could be good for memory. I don't have any studies to back this up, but I've heard from a couple of people that remembering things is a skill which improves from use: if you are patient with your memory and wait for it to work, rather than giving up right away, it starts getting noticeably better.

(There are studies which show that practicing memory skill does improve it, but, the ones I know of involve memorizing lists for later recall rather than recalling things more randomly without knowing what you'll need to remember.)

I also think this is good practice for [defusion](#), the skill of de-identifying with your thoughts and feelings. Tracking back on your distracting thoughts helps you shift from thinking "X" to thinking "I had the thought X". This is helpful for avoiding temptations and managing reactions (basically, willpower), as well as metacognition more generally.

Recalling rapid chains of thought which got you to where you are is also a necessary subskill of the [tuning your cognitive strategies](#) skill at bewelltuned. The practice I'm about to describe is actually pretty close to the one at bewelltuned, so that may be good reading material to help with this meditation. Differences:

- Bewelltuned is describing a more active exercise which you do while solving a puzzle. I'm doing it with my eyes closed and no task at hand.
- Bewelltuned is combining the skill of rapid metacognition about how your thoughts got to where they are with the problem of credit assignment, noticing which thoughts were useful.

The Meditation

The idea is simple: every time you notice that you're distracted, track back mentally. What sequence lead you here? Try to recall in as much detail as possible.

What is the "central focus" which you are trying to attend to, so that you know when you're distracted? There are a couple of options here.

- You could choose to focus on your breath, as is common, and only do a track-back when you notice that your attention has strayed from your breath. I did this a little, but it really means you're switching between two primary objects of focus: the breath, and the track-back mental motion. This can be a little confusing.
- Alternatively, you can think of the track-back mental motion itself as the object of focus. When I do it this way, the object of meditation which I return to when I don't have any distractions to track-back is a track-back of the fresh moment I find myself in. How did I end up here?
- A third option is to do absolutely nothing when there aren't any distracting thoughts. Rest in the peace of your mind. Unless you are an experienced meditator, you will very soon find yourself distracted again, so don't worry that you're not practicing track-back for a few moments. This option is really hard (at least for me); it is difficult to keep meditating with no focus to return to. However, I do find that I can sometimes switch to this after I'm a few minutes into the meditation.

"Track back" can mean several things. One interpretation is by-the-clock back-tracking: recall the sequence of moments leading up to this one, as best you can. Another option is causal back-tracking: look for the origins of the thought. For example, if you remember a late bill, a causal back-track might go to the image of the bill sitting on your desk (which you might have seen earlier in the day, setting up the thought of the bill as an open loop in your mind). You might then think further back to putting it on the desk, and so on.

I suggest that you back-track in whatever way feels natural and relevant to you.

[ETA: I now think that causal back-tracking practices **long-term memory**, whereas temporal back-tracking practices **short-term memory**. Practicing short-term memory, if it works, effectively expands your working memory; if you can quickly "touch" a lot of thoughts to maintain them in short-term memory, you will lose less in your train of thought. I find that this also helps me to maintain productivity on a task. Very short-timespan temporal back-tracking also seems most useful for [tuning cognitive strategies](#). Practicing long-term memory helps to maintain broad rather than narrow context; it might end up being more compatible with an exploratory rather than focused state of mind. Causal back-tracking seems most useful for [the strategic level](#). Nonetheless, I would recommend aiming at whichever mental motion feels more amenable to improvement. If you are doing this with the goal of extending the set of mental processes which you are conscious of, orient toward whatever you are not currently very conscious of.]

It's very important to keep a non-judgemental attitude toward your distracting thoughts and back-tracks. I've previously done a meditation in which I tried to "go meta" on every thought which arose, meaning take the outside view and [think about the thought as a part of a policy](#). This involves looking toward the source of each thought, like the meditation I'm describing now, but it also involves *judging* each thought. I found that this exercise left me in a manic state, which is a bad sign.

The *tuning your cognitive strategies* exercise which I mentioned earlier also involves judging the thought-patterns, but they give a warning to stick to positive judgements to reinforce good patterns; punishing useless thoughts is unnecessary, and more dangerous than positively reinforcing the useful ones. This also seems fine, but since my exercise is the track-back activity itself, other thoughts are all distractions anyway.

More Advice

If you get distracted while you're in the middle of back-tracking a different distraction, should you try to return to back-tracking the earlier distraction, or start back-tracking the new one?

If you let yourself keep switching which thought you're back-tracking, you're not training yourself to finish back-tracks, which is the goal. That being said, if I find that I've been distracted from the back-track for a while without noticing, it may be too hard to pick up back where I was, so I'll just start back-tracking from the present.

Also, some things seem more inherently interesting to back-track than others. As long as you're not letting a distraction be a *distraction*, I think it's fine to follow your interest. Just try to keep with the back-tracking mental motion, rather than allowing yourself to get involved in interesting thoughts in other ways.

How do you start? I don't have any distracting thoughts right when I start, so I don't know what to do.

I set a timer right before starting, so I tend to think back about how I set the timer, how I sat down before that, how I got to the room I'm in... (I usually don't get that far before noticing I'm distracted.)

I have totally random thoughts or images which come in. They don't have any identifiable source, so I can't back-track.

Me too! Be patient with your memory, though. If you gently ask "where did that come from?" and wait, an answer may come, even if the thought initially seems really random. Also, you can start analysing parts of the thing. Maybe you had a random flash of a soldier with black and white face paint standing in a river. Where might the black and white face paint have come from? What might have made you associate face paint with warriors? (Can you recall instances?) Do you think of "standing in a river" as a soldier sort of thing to do, and do you know why?

Or, you could do a temporal track-back rather than a causal track-back, so you just have to try to remember what you were thinking right before the soldier thing.

What about false memories? There are many studies showing that we confabulate, and especially that we confabulate our reasons for doing things. Why should I trust the back-traces which my mind gives me when I try to recall what happened leading up to a certain thought?

I think this is an interesting question. It seems to me like it isn't that hard to avoid false memories. A false memory seems like the result of a mistake in figuring out what happened. In some cases I've generated a chain back from a thought, and then asked myself "is that really what happened?" -- and the answer sometimes seems to be "no". So, there is a difference between what my brain thinks just happened and what it thinks when I ask it whether what it thinks is right. So the situation isn't hopeless! In fact, I would be surprised to learn that people can't detect their confabulations in ordinary situations (IE, excluding brain damage), in which there aren't any real stakes motivating the confabulation, if they actually try.

But then, I *would* think that, wouldn't I?

Conclusion

I've found this meditation to be particularly easy to motivate myself to do. After trying it, I generally feel... energized?... well, it's hard to describe, but I think it's a more successful meditative experience than other things I've tried. Hopefully someone else finds it useful.

What To Do If Nuclear War Seems Imminent

This is a linkpost for <https://benlandautaylor.com/2018/09/13/what-to-do-if-nuclear-war-seems-imminent/>

This document describes precautions to take in a scenario like the Cuban Missile Crisis, where nuclear war seems plausibly imminent within the next days or weeks. This is *not* a guide for what to do if a missile is currently inbound and will strike within minutes or hours.

Overview

If tensions between nuclear powers are running extremely high, and you are in or near a plausible target during a nuclear war (such as a major city in the United States or Europe), then I recommend evacuating to a safer place as soon as possible, and staying for days or weeks until things have calmed down. New Zealand is an excellent place to go.

This plan requires that you maintain a valid passport, so that you can leave your country on short notice if needed. No other special preparations are needed.

Proper calibration here should include substantial tolerance for false positives. For people with the means available, I think it was correct to evacuate during the Cuban Missile Crisis, even though it did not end up leading to nuclear war.

Why New Zealand?

New Zealand is of little or no strategic relevance to the current conflicts between nuclear powers. The experts I've talked to agree that it's implausible that anyone would target New Zealand with nuclear weapons, or that anyone would invade New Zealand in the aftermath of a nuclear exchange.

New Zealand is easy to enter. Anyone with no notable criminal history and a valid passport from most countries, including the US, EU, and Canada, can get a New Zealand tourist visa on arrival, with no need for a prior application, and stay for up to 90 days. (Make sure to get a round-trip ticket, or they might not let you in.)

New Zealand is a major food exporter. If supply chains are disrupted, you'll be close to the source.

New Zealand is very stable internally. It has a strong Anglo tradition of governance, reasonable national pride, no coups or civil wars within the last century+, negligible riots or ethnic strife, etc.

New Zealand is culturally familiar. It's an English-speaking country that's firmly within Western Civilization. As such, most of my audience will be more comfortable staying there while waiting for tensions to calm down, and will stick out less if there's chaos or rioting after a war.

No other country is so good on so many of these dimensions.

Backup Plans

If you are unable to enter New Zealand, then there are many other countries which look like good options: many South American countries, Australia, and Botswana. Partial notes [here](#).

If you are unable to leave your country (this is unlikely if you have a valid passport; see below), then you should drive to a small town far from any metropolis or other plausible target. (After brief examination, for people in the Bay Area, I recommend the Modoc Plateau in northeast California as a default unless/until more research is done.) Once there, organize, acquire supplies, and find a location to dig fallout shelters. Construction is described in Nuclear War Survival Skills, the full text of which is [online](#). The book claims untrained civilians can build the shelters in 1-2 days.

Other Concerns

How will I know when to evacuate?

This will probably be obvious. Past diplomatic crises between nuclear powers have frequently been widely publicized.

If I decide to evacuate, I will send a brief alert to anyone who signs up to receive one via [this form](#).

Won't all the flights get booked due to mass panic?

Probably not, judging by past cases. For example, it looks like there were no large-scale evacuations during the Cuban Missile Crisis, in spite of [very alarming headlines](#). (It seems to me that most people have trouble thinking about nuclear destruction in a way that permits any action whatsoever.)

What about nuclear fallout?

Based on a friend's analysis, fallout risk in New Zealand is low unless New Zealand itself is targeted, and the experts I've talked to agree that this is implausible.

Fallout is dangerous for about two weeks. Nuclear War Survival Skills ([full text](#)) describes how to build shelters, which would be uncomfortable but effective.

Moral differences in mediocristan

This is a linkpost for <http://benjaminrosshoffman.com/moral-differences-in-mediocristan/>

Scott Alexander [writes](#):

Utilitarianism agrees that we should give to charity and shouldn't steal from the poor, because Utility, but take it far enough to the tails and we should tile the universe with rats on heroin. Religious morality agrees that we should give to charity and shouldn't steal from the poor, because God, but take it far enough to the tails and we should spend all our time in giant cubes made of semiprecious stones singing songs of praise.

He suggests that these are surprisingly divergent visions of the highest good, for moral visions that give similar advice for day-to-day life:

converting the mass of the universe into nervous tissue experiencing euphoria isn't just the second-best outcome from a religious perspective, it's completely abominable

But what strikes me about them is how *similar* they seem, when you strip away the decorative metaphors.

First of all, in both cases you can afford many more instances of the best thing if you simulate it than make a real one. So in both cases the universe is better converted into computronium than actual rats or heavenly choirs - the substitution of "nervous tissue experiencing euphoria" for "rats on heroin" is an implied acknowledgement of this. Nor should we imagine that religion eschews such optimization. Many religions promote asceticism, which allows more humans to subsist and praise God on the same resources. Moreover, the Bible urges that we be fruitful and multiply.

But also, it's not at all clear that an imagined end state of constant songs of praise is meaningfully different from the hedonic-utilitarian state. After all, of what interest is a song of praise if it is not an expression of genuine appreciation? And, couldn't you - by economizing on the actual churches, or all the parts of the mind not engaged in songs of praise - make your religious heaven more efficient and therefore have more of it?

(Anyone familiar with Dante's *Paradiso* will have recognized a vision of heaven - where we want as many people to go as possible, of course - identical or nearly identical to the hedonic utilitarian goal - the souls of the saved forming a sort of Matrioshka brain orbiting about a reference point of maximum goodness, basking in the energy radiating out from it in eternal unchanging bliss.)

I don't see how this optimization ends up with anything *but* "nervous tissue experiencing euphoria" - simulated nervous tissue, of course. There's an implied disagreement on whether rats have minds sufficiently complex to sing songs of praise, but that's an open question among hedonic utilitarians (and among the religious) too. In any case, it's reasonably likely that if you tried to simulate *just* the expression of appreciation, it wouldn't matter much whether you started with a model of a human brain singing songs of praise, or a rat brain enjoying heroin.

It's actually in the near term that these visions of the good diverge. Here's what I see as the actual, applied disagreements between a few major schools of thought.

Utilitarianism implicitly recommends that we make decisions by directly quantifying the good our different actions might accomplish, and doing the one that scores highest. Communism seems to think the people doing all the work should coordinate to seize control of the social systems being used to extract things from them, and then gradually wind them down, though step 1 never seems to work out. Fascism seems somewhat like a more cynical alternative to Communism that doesn't bother pretending there's a second step, which appeals to people who want to boss others around even if they don't get to be at the apex of the hierarchy.

Christianity thinks that we should make more people Christian, and have them listen to sermons or read books about helping others (or, in a surprising variant, subject themselves to Roman imperial administration despite the ostensible collapse of the Roman empire more than a millennium ago), and also actively try to help people around us in order to send a credible signal that Christianity's a good thing. Islam seems to think it should gradually conquer the world by a mixed strategy of conversion by persuasion and seizing and holding territory, and administer an unambiguous code of laws with clear lines of authority and social roles. Judaism thinks that Jews should have children and train them to engage in an intergenerational project to develop a perfected code of laws with a history spanning millennia, in the hope that eventually Jews will have something good enough that other people will want to adopt it, and in the meantime doesn't have much in the way of advice for non-Jews. (It's unclear to what extent Christianity, Islam, and Utilitarianism constitute partial successes for this agenda.) Buddhism seems to think we should teach people how to chill out, don't sweat the small stuff, and it's all small stuff, in the hopes that this will cause a persuasion cascade similar in mechanism to the Christian one, but without promoting either the weird book or the Roman imperial bureaucracy.

This can get confusing because utilitarianism tends to foreground the *idea* of a particular end state, because that's the utilitarian strategy, and utilitarians tend to misunderstand people with different strategies as having different ends in mind. But in practice, most people most of the time have not got a global highest end in mind - they're implementing a strategy that feels right or advantageous to them, or one they've been acculturated to implement, and correctly attending to the *means* this directs them to, since they're unlikely to ever be in a position to directly specify ends globally, and any long-run strategy for programming the universe is likely to involve future generations better situated to specify the end-state than we are.

Related: [Against Responsibility](#), [The Basic AI Drives](#)

Petrov corrigibility

Crossposted from the [AI Alignment Forum](#). May contain more technical jargon than usual.

After my [example](#) of problems with corrigibility, and Eliezer pointing out that sometimes corrigibility [may involve saying "there is no corrigible action"](#), here's a scenario where saying that may not be the optimal choice.

[Petrov](#) is, as usual for heroes, tracking incoming missiles in his early warning command centre. The attack pattern seems unlikely, and he has decided not to inform his leaders about the possible attack.

His corrigible AI pipes up to check if he needs any advice. He decides he does, and asks the AI to provide him with documentation about computer detection malfunction. In the few minutes it has, the AI can start with one introductory text A, or with introductory text B. Predictably, if given A, Petrov will warn his superiors (and maybe set off a nuclear war), and, if given B, he will not.

If the corrigible AI says that it cannot answer, however, Petrov will decide to warn his superiors, as his thinking has been knocked off track by the conversation. Note that this is not what would have happened had the AI stayed silent.

What is the corrigible thing to do in this situation? Assume that the AI can predict Petrov's choice for whatever action it itself can take.

Resurrection of the dead via multiverse-wide acausal cooperation

TL;DR: Measure decline in random mind creation may be prevented if we take into account very large number of random mids created in other universes.

Summary: P.Almond suggested the idea of the resurrection of the dead via a quantum random generator which creates a random mind, but such an approach has several problems: non-human beings in our world, non-necessary suffering of non-perfect copies, and measure decline.

Here I suggest three patches, which prevent most of the undesired effects:

1. Human mind matrix to prevent pure random minds appearing.
2. Digital immortality data to create a person which satisfies all known external expectations, and the use of randomness only to fill unknown information.
3. Multiverse-wide cooperation for the “cross-resurrection” of the dead between multiple worlds via quantum random minds, so the total measure of all resurrected people will not decline.

1. Introduction

Almond in “[Many-Worlds Assisted Mind Uploading: A Thought Experiment](#)” suggested the following idea about the resurrection of the dead by the use of a quantum random generator, which would create a random mind within a computer (Almond, 2006):

[A technician who lost someone's brain scan file] writes a computer program which takes input from a physical system. The physical system, known as a quantum event generator, generates "1"s and "0"s randomly as a result of quantum events. The program will use the physical system to tell it what sequence of "1"s and "0"s will be used to try to recreate the lost scan file. The program starts with an empty scan file which will be filled with "1"s and "0"s.

If the many-worlds interpretation of quantum mechanics is correct, all possible minds will appear in separate timelines starting from the moment of random mind creation, which would mean the resurrection of everyone from his own point of view. However, this approach will a) not help an outside observer, who wants to resurrect a relative, for instance, as the observer would see only a random mind, and b) the quantum “measure” of existence of each mind will be infinitely small.

2. Problems of Almond's approach

To illustrate the problems with quantum mind uploading, I will explore a simplified thought experiment where only names will be restored using quantum mind uploading. First, here is what Almond suggested:

Thought experiment “Not-patched quantum mind uploading”:

Bob had a friend John Smith. John has died and Bob wants to resurrect him. Bob remembers only first letter of John's name: S.

Bob and John are interested only in the uniqueness of name preservation, and no other identity considerations are important. Bob wants to observe his friend to be alive, and for his friend to be named "John S...." (I would call it immortality from the point of view of the external observer). John wants his own immortality, and will be satisfied only if "John Smith" is created.

Bob creates random quantum mind A using a quantum generator to choose each new letter in the names.

It turns out that A is "jYY2N@11". Only less than 10-30 share of all such copies in the multiverse are named John Smith. Both Bob and John are unhappy.

This thought experiment leaves both John and Bob unsatisfied, and we see three reasons for that below:

2.1. Problem 1: Measure decline

Problem 1 is a problem for John.

Measure could be defined as a share of an observer of a given type between all possible observers. If the typical size of the simulated mind is, say, 10^{15} bites, the chances that a randomly generated mind will be exactly the needed person is $2^{-(10^{15})}$. In other words, a quantum mind generator results in a measure decline of $2^{-(10^{15})}$ which is an extremely large number. Even in our thought experiment 1 measure decline is 10³⁰ times.

Many authors claim that large measure decline should be treated as death or as an infinitely small chance of survival. Such discussions appeared in the context of so-called quantum immortality, that is, the counterfactual possibility to survive death via existing in quantum multiverse timelines where a person will not die.

Even if the measure decline is not bad per se, it leads to a world where very small probability outcomes will dominate possible futures of an observer, and such parasitic outcomes may be full of suffering. For example, the quantum immortality improbable survival landscape may be dominated by people who are very old and dying but can't die (it could be patched by signing up for cryonics).

If we use some expected utility calculations, and measure decline results in declining utility of any useful outcome associated with it, we could just ignore my copies with infinitely small measures.

2.2. Problem 2: Non-human and not welcomed minds

Problem 2 is mostly for Bob.

Another problem is that most random minds will be non-human, and will not be adapted to our world, so they will suffer or cause suffering to people living here. In our thought experiment "jYY2N@11" is an example of a non-human random mind.

Such random minds are also extremely bad for any outside observer, like Bob, as he will be very unlikely to meet anyone resembling his friend John Smith.

2.3. Problem 3: Damaged minds

Problem 3 is a problem for both Bob and for John.

Most randomly-created minds will be not minds at all, but some garbage code, or at “best case,” damaged minds. For example, if Bob wants to resurrect John Smith, there will be much more copies where his name (as well as his other properties) is a parody of the name Smith, for example Smthi, Smiht, Misth, Smitt, etc. For n bits long name, there are n individual names which have 1 bit difference.

Thus, for any real person, there will be much larger set of his/her damaged copies, which implies suffering for such a person as the most probable outcome of the quantum random resurrection and s-risks for all people.

3. Patches

Fortunately, quantum random mind uploading could be patched, so it will provide much more satisfaction for John and Bob.

Patch 1. The use of the human mind’s universal model as a starting point

The goal of this patch is to escape minds of “aliens” or of non-workable gibberish code, and thus prevent suffering of most created minds. For example, for a human mind model, his/her possible name will be generated not as random symbols but from the preset of typical human names.

Such a human mind model may look like an untrained neural network which has the general architecture of a human mind, with some other constraints, so any random set of parameters will create a more-or-less normal human mind. We assume that some future assistant AI will be able to find an appropriate model.

In that case, Bob uses a random mind generator for parameters of the universal human mind model. He gets “Maria Stuart”. This will increase the share of the worlds where real John Smith is resurrected to 10-10. Both John and Bob are a little bit more satisfied, as Bob gets a human friend, and John increases his measure.

Obviously, some minds may not want to be resurrected, but this could be an important parameter in the model, and models, where “resurrection preference = false” will be ignored.

Patch 2. The use of the digital immortality data to create only minds which comply with our expectations

The problem of Bob’s satisfaction could be overcome by the use of Bob’s expectations as priors, if there are no other current of future sources of data about John.

In that case, Bob could use his memories about John S. to create a model of John S. He remembers that John was either John Smith or John Simpson. He uses a random quantum coin to choose between Smith or Simpson, and gets “John Simpson”.

In another branch of the quantum multiverse, where the coin fails tails, John Smith appears, but his measure declines to 0.5. Both John and Bob are partly satisfied. Bob got someone who looks like his friend, but Bob knows that it is not exactly his friend, and that his friend has now smaller measure of existence.

Digital immortality, or indirect mind uploading, is the collecting information about a person while he is alive with hope that future advanced AI may be able to resurrect the person, by creating an advanced model of the personality based on all available information. Such a model will, by definition, satisfy Bob and all other relatives, as all

available information has already been taken into account, including all relatives' expectations. However, large chunks of information will never be known, and thus have to be replaced with some random data. Even if quantum randomness is used to fill the gaps, John will have an infinitely small share of all possible worlds, and in most other worlds he will be replaced by someone else.

Patch 3. The use of multiverse-wide cooperation for the cross-resurrection

The next step is that Bob considers that not only his universe exists, but all possible other universes exist in the Multiverse.

Bob concludes that because all possible observers exist in the Multiverse, his John Simpson created via a quantum random generator is a resurrection of some John Simpson from another universe, while John Smith who lived in our universe, will be resurrected in some other universe where another copy of Bob will do the same experiment.

In other words, Bob and Bob's copies in other universes cooperate to resurrect the exact John Smith.

As the second universe is exactly the same as ours except for John's name, there is another exact copy of Bob in it, and this Bob's copy is also wanting to resurrect his friend John S., so he uses another quantum random mind generator. Now the following happens:

Universes:	Bob's friend:	Result of quantum mind generator:	Total measure at the end:
Bob	John Smith, measure = 1	If coin = heads, Smith; measure = 0.5	Smith = 1 Simpson = 1
		If coin= tails, Simpson; measure = 0.5	
Bob's copy	John Simpson, measure = 1	If coin = heads, Smith; measure = 0.5	
		If coin= tails, Simpson; measure = 0.5	

So, the total measure of John Smith has not declined, if Bob takes into account that other copies of Bob in other universes will run the same experiment. By deciding to start the random mind generator (and to not turn off the resulting mind), Bob joins a large group of other minds, who think similarly, but who are located in causally disconnected parts of the Multiverse. Everyone expects that some other random generator recreates an exact copy of their loved one.

In a real case of large missing data, like gigabytes, this requires a simultaneous run of an extremely large number of quantum random mind generators, like $10^{(10^9)}$, which is only possible via multiverse-wide cooperation. The measure will not decline in such a case too, as for every dead person there will be one random person, and given the large numbers, any person will be randomly recreated, at least in approximately one world. (Some may go deeper and take into account standard deviation, but because we use quantum generators in the many worlds interpretation, each universe

creates exactly its share of John, and there will be no fluctuations, which would result in non-existence of some Johns and two copies of another.)

Any of Bob's copies can join such a multiverse-wide cooperation by creating just one quantum random mind (and treating the resulting mind well).

4. Remaining problems

Multiverse. What if the multiverse doesn't actually exist? In that case, Bob and John get partly satisfying results, as Bob gets John's copy, but John's copy is not perfect from John's point of view. If the quantum multiverse is not real, but some other form of the multiverse exists, like the one based on inflationary cosmology, the resurrection method will still work.

Defection. Bob may not create any random mind generators at all but still expect that someone else will recreate his friend. In general, the rate of defections may be known and compensated by increasing the number of random minds by those who have more resources.

There are several other possible generic problems of multiverse-wide cooperation, including infinite ethics, the possibility of acausal blackmail, a method to measure similarity between agents, and problems with agents that have other values as described in EA post's [comment](#).

Conclusion

I hope that this post may increase one's hope in the future personal resurrection by superintelligent AI.

Are you in a Boltzmann simulation?

EDIT: Donald Hobson has [pointed out a mistake](#) in the reasoning in the section on nucleation. If we know that an area of space-time has a disproportionately high number of observer moments, then it is very likely that these are from long-lived Boltzmann simulations. However, this does not imply, as I thought, that most observer moments are in long-lived Boltzmann simulations.

Most people on LessWrong are familiar with the concept of [Boltzmann brains](#) - conscious brains that are randomly created by quantum or thermodynamic interactions, and then swiftly vanish.

There seem to be [two types](#) of Boltzmann brains: quantum fluctuations, and nucleated brains (actual brains produced in the vacuum of space by the expanding universe).

The quantum fluctuation brains cannot be observed (the fluctuation dies down without producing any observable effect: there's no decoherence or permanence). If I'm reading [these papers](#) right, the probability of producing any given object of duration t and mass

M is approximately

$$\exp [- t M \times 10^{51}].$$

We'll be taking $M = 1\text{kg}$ for a human brain, and $t > 0.1\text{s}$ for it having a single coherent though.

A few notes about this number. First of all, it is vanishingly small, as an exponential of a negative exponential. It's so small, in fact, that we don't really care too much over what volume of space-time we're calculating this probability. Over a Plank length four-volume, over a metre to the fourth power, or over a Hubble volume for 15 billion years: the probabilities of an object being produced in any of these spaces are approximately all of the same magnitude, $\exp[-tM \times 10^{51}]$ (more properly, the probabilities vary

tremendously, but any tiny uncertainty in the 51 term dwarfs all these changes).

Similarly, we don't need to consider the entropy of producing a specific brain (or any other structure). A small change in mass overwhelms the probability of a specific mass setup being produced. Here's a rough argument for that: the [Bekenstein bound](#) puts a limit on the number of bits of information in a volume of space of given size and given mass (or energy). For a mass M and radius R , it is approximately

$$2.6 M R \times 10^{43}$$

Putting $M = 1\text{kg}$ and $R = 0.1\text{m}$, we get that the number of possible different states in brain-like object of brain-like mass and size is less than

$$2^{10^{43}},$$

which is much, much, much, ..., much, much less than the inverse of $\exp[-10^{51}]$.

Quantum fluctuations and causality

What is interesting is that the probability expression is exponentially linear in t :

$$\exp [- (t_1 + t_2) M \times 10^{51}] = \exp [- t_1 M \times 10^{51}] \times \exp [- t_2 M \times 10^{51}].$$

Therefore it seems that the probability of producing one brain of duration t_1 , and another independent brain of duration t_2 , is the same as producing one brain of duration $t_1 + t_2$. Thus it seems that there is no causality in quantum fluctuating Boltzmann brains: any brain produced of long duration is merely a sequence of smaller brain moments that happen to be coincidentally following each other (though I may have misunderstood the papers here).

Nucleation and Boltzmann simulations

If we understand [dark energy](#) correctly, it will transform our universe into a [de Sitter](#) universe. In such a universe, the continuing expansion of the universe acts like the event horizon of a black hole, and sometimes, spontaneous objects will be created, similarly to [Hawking radiation](#). Thus a de Sitter space can nucleate: spontaneously create objects.

The probability of a given object of mass $M = 1\text{kg}$ being produced is [given as](#)

$$\exp [- M \times 10^{69}].$$

This number is much, much, much, much,, much, much, much, much smaller than the quantum fluctuation probability. But notice something interesting about it: it has no time component. Indeed, the objects produced by nucleation are actual objects: they endure. Think a [brain in a sealed jar](#), floating through space.

Now, a normal brain in empty space (and almost absolute zero-temperatures) will decay at once; let's be generous, and give it a second of survival in something like a functioning state.

EDIT: there is a mistake in the following, see [here](#).

Creating N independent one-second brains is an event of probability:

$$\exp [- N \times 10^{69}].$$

But creating a brain that lasts for N seconds will be an event of probability

$$\exp [- M_N \times 10^{69}].$$

where M_N is the minimum mass required to keep the brain running for N seconds.

It's clear that M_N can be way below N . For example, the longest moonwalk was 7 h 36 min 56 s (Apollo 17, second [moonwalk](#)), or 27,416 seconds. To do this the astronauts used [spacesuits of mass around 82kg](#). If you estimate that their own body mass was roughly 100kg, we get $M_{27,416} < 200 < 27,416$.

This means that for nucleated Boltzmann brains, unlike for quantum fluctuations, most [observer moments](#) will be parts of long lived individuals, with a life experience that respects causality.

And we can get much much more efficient than that. Since mass is the real limit, there's no problem in using anti-matter as source of energy. The [human brain runs at about 20 watts](#); one half gram of matter with one half gram of anti-matter produces enough energy to run this for about 4.5×10^{12} seconds, or 140 thousand years. Now, granted, you'll need a larger infrastructure to extract and make use of this energy, and to shield and repair the brain; however, this larger infrastructure doesn't need to have a mass anywhere near 10^{12} kilos (which is the order of magnitude of the mass of the [moon](#) itself).

And that's all neglecting improvements to the energy efficiency and durability of the brain. It seems that the most efficient and durable version of a brain - in terms of mass, which is the only thing that matters here - is to run the brain on a small but resilient computer, with as much power as we'd want. And, if we get to re-use code, then we can run many brains on a slightly larger computer, with the mass growth being less than the growth in the number of brains.

Thus, most nucleated Boltzmann brain observer-moments will be inside a Boltzmann simulation: a spontaneous (and causal) computer simulation created in the deep darkness of space.

Your genome isn't private. Maybe it never was.

This is a linkpost for <https://www.eastbaybiosecurity.org/blog/defcon-biohacking-genetic-privacy>

Advances in Baby Formula

Someone is finally trying.

Four years ago I was shopping for baby formula. I was horrified to discover that most formulas were mostly fructose and soybean oil. Including the fancy 'organic' ones. Humans digest fructose hepatically, and overconsumption of fructose has been implicated in a variety of modern health problems including diabetes and atherosclerosis. Why would we feed it to our babies as their primary sugar?

Oh right. It's super cheap.

Soybean oil is a plant oil high in polyunsaturated fats and low in saturated fats. It and other high omega-6 oils have been implicated in numerous health problems. It is notoriously the only fat found in the tpn (intravenous) formulation given to infants, which causes liver failure and death. The FDA will not approve a tpn with the appropriate fat contents to keep infants alive, sparking outrage. Why would we feed our babies soybean oil instead of milk fat? It was pointed out to me that this was probably related to shelf-stability of the fat in question. Milk fat likely spoils faster.

Plus, soybean oil is cheap. As is fructose. And made with 'organic' vegetables...

So, I sprang for the premium formula which contained lactose instead of fructose, but still was mostly soybean oil.

Four years later, I am pleasantly surprised by the rise of much more sophisticated baby formulas. Enfamil Enspire, and it looks like some other formulas, have added a more complex mix of sugars and fats. The primary sugar is still lactose, but in addition there are galactooligosaccharides and polydextrose. These prebiotic sugars pass through the stomach undigested, but serve as a food source for bacteria in the gut and slowly release sugar as they are digested there. They are thought to lead to healthier gut bacteria and more stable blood sugar.

The primary fat is still a vegetable oil blend (still high in omega 6's), but it now contains coconut oil which is high in saturated fat. Some milk fat has been added and supplemented with omega 3s DHA and ARA. One of the newest advancements is the addition of MFGM (milkfat globule membrane), which is a mix of gangliosides and phospholipids thought to play an important part in cognitive development. If you've spent any time studying or researching weird Ashkenazi recessive mutations, you'll know that many of them are mutations in brain lipid metabolism including: Nieman-Pick's, Fabry's, Tay Sach's, Krabbe's, and Gaucher's diseases. It has long been hypothesized that having a single copy of these traits is related to higher IQ in the Ashkenazi population, even if having two can be fatal. Suffers of Gaucher's disease are disproportionately represented in engineering and mathematics. In any case, you want to have good brain lipids. Let's put them in formula.

Is the exact mix of components now used in these new formulas the best we can do with our knowledge and tech level? No. But it is a vast and obvious improvement over the previous formulas of fructose and soybean oil.

Why did it take this long, and why am I so surprised that we have it? I'm a cynic. I had assumed that formula was the way it was because of some regulations somewhere, and that companies didn't want to deal with making changes. Children are too

important to learn about, etc. Maybe this was true and demand finally got so loud that it was finally worth the extra cost and effort. When I think about the state of nutrition science twenty years ago, it's unsurprising that formula looked the way it did. It takes time for a culture to shift as well. The people who cared most about what their babies were eating went crazy with demanding that everyone breast feed (no matter what the practical reality of your circumstance was). Breastfeeding became a moral imperative and breastmilk a sacred symbol of MotherLove, which could never be replicated in a food science lab.

What will the outcomes with the new formula be? We will probably *never* know. Children are too important to learn about. We don't even have a good way of comparing formula to breast milk. Everything is confounded. Nothing is controlled. The best we can say right now is at least we know the outcomes on the old formula were only as bad as they are. That is to say, we are still uncertain if it was worse than breast milk, and therefore it must be at least mostly fine. In any case, I'm buying the new stuff.

Alignment Newsletter #24

Crossposted from the [AI Alignment Forum](#). May contain more technical jargon than usual.

Starting from this week, Richard Ngo will join me in writing summaries. His summaries are marked as such; I'm reviewing some of them now but expect to review less over time.

Highlights

[Introducing the Unrestricted Adversarial Examples Challenge](#) (*Tom B. Brown et al*): There's a new adversarial examples contest, after the one from NIPS 2017. The goal of this contest is to figure out how to create a model that never confidently makes a mistake on a very simple task, even in the presence of a powerful adversary. This leads to many differences from the previous contest. The task is a lot simpler -- classifiers only need to distinguish between bicycles and birds, with an option of saying "ambiguous". Instead of using the L-infinity norm ball to define what an adversarial example is, attackers are allowed to supply any image whatsoever, as long as a team of human evaluators agrees unanimously on the classification of the image. The contest has no time bound, and will run until some defense survives for 90 days without being broken even once. A defense is not broken if it says "ambiguous" on an adversarial example. Any submitted defense will be published, which means that attackers can specialize their attacks to that specific model (i.e. it is white box).

My opinion: I really like this contest format, it seems like it's actually answering the question we care about, for a simple task. If I were designing a defense, the first thing I'd aim for would be to get a lot of training data, ideally from different distributions in the real world, but data augmentation techniques may also be necessary, especially for eg. images of a bicycle against an unrealistic textured background. The second thing would be to shrink the size of the model, to make it more likely that it generalizes better (in accordance with Occam's razor or the minimum description length principle). After that I'd think about the defenses proposed in the literature. I'm not sure how the verification-based approaches will work, since they are intrinsically tied to the L-infinity norm ball definition of adversarial examples, or something similar -- you can't include the human evaluators in your specification of what you want to verify.

[The What-If Tool: Code-Free Probing of Machine Learning Models](#) (*James Wexler*): When you train an ML model, it is often hard to understand what your model is doing and why. This post introduces the What-If tool, which allows you to ask counterfactual queries about the decision rule implemented by your final trained model, for classification and regression tasks. For example, you can take a particular data point, edit it slightly, and see how that changes the model prediction. Or you can graph the data points by L2 distance from a particular point. For classification tasks, you can find the "closest counterfactual", that is, the data point closest to the current point where the decision of the model is reversed. I played around with some of the demos, and apparently for a particular person and a particular model trained on census data, the probability that they had a salary of over \$50k depended much more strongly on their marital status than their age, which was the opposite of my prediction. I figured this out by choosing a point, finding the closest counterfactual,

and then making each of the changes in the delta individually and seeing which affected the model probability most.

My opinion: I'm guessing this is limited to tasks where your data points have a reasonable number of features (< 1000, I'd guess) and you are only analyzing a small set of test data points (around tens of thousands), due to computational constraints. That said, for those tasks, this seems incredibly useful to actually get a good model that you can debug and eventually deploy.

It's worth noting that this is an engineering achievement. Researchers are considering even stronger (but more computationally difficult) techniques, such as finding which part of the training set most influenced a particular decision, whereas the What-If tool doesn't talk about the training set and training process at all, instead only allowing you to ask queries about the final trained model.

Preserving Outputs Precisely while Adaptively Rescaling Targets (Matteo Hessel et al): When an agent is trained on multiple tasks across which the sizes of rewards vary greatly, it usually focuses on tasks which provide the largest or most frequent rewards at the expense of performance on others. Previous work dealt with this by clipping rewards outside a certain range, but this changes the optimal policy (eg. in Pacman, eating pellets is just as rewarding as eating ghosts). This paper uses PopArt (introduced in [this 2016 paper](#)) to normalise rewards from each task before using them to update the policy in an actor-critic RL algorithm. The authors use PopArt to train a single IMPALA agent which can play all 57 Atari games, achieving a median performance slightly higher than human performance.

To delve into more detail about PopArt, let's consider training a policy with an actor-critic algorithm. In this case, we need a critic that produces estimates of values V , and an actor that produces probabilities of actions. Both of these networks are trained by taking gradients of their outputs, and weighting them based on the observed rewards. Now, a key empirical fact about deep learning is that it works better if all the things are normalized, especially the gradients. (If nothing else, it makes it easier to choose the learning rate.) For the actor, this is easy -- probabilities are already normalized, and the weight of gradient is proportional to the reward, so we can just rescale the weight of the gradient based on the mean and standard deviation of the rewards we have observed so far. This is a bit harder for the critic, since it has to predict values, so we have to normalize both the outputs and the gradient weights. We can normalize the gradient weights in the same way as before. However, normalizing the outputs is tricky, because as time goes on the means and standard deviations change. To do this, at every timestep we modify the weights of the critic that is equivalent to unnormalizing based on the old statistics and then normalizing based on the new statistics. This gives the PopArt method.

Here's a simple example where I butcher types, ignore the difference between states and trajectories, and throw away the standard deviation. Suppose the first reward we see is 10, so we say that our mean is 10 and train our net to output a normalized reward of 0 for this state and action. Then, we see a reward of 100, so we update our mean to 55. On our previous (state, action) pair, we still output a normalized reward of 0, which now corresponds to a real reward of 55, even though it should correspond to 10! We then do the unnormalize-and-renormalize trick. After unnormalization, the critic would output 10, and after renormalization, the network would output -45, which when combined with the mean of 55 would give us the desired 10 reward.

My opinion: This is an impressive result, since it's the first time a single agent has performed so well on a range of Atari games. It doesn't seem to have required any novel techniques except for a straightforward extension of PopArt to the multi-task setting, but this is still interesting since the results from the previous PopArt paper were very mixed (performance increased and decreased dramatically on different games, with the average remaining roughly stable).

One confusing aspect was that PopArt still benefitted slightly from being trained with reward clipping (110% vs. 101% in the unclipped case), even though the point of PopArt was to normalize rewards so that clipping wasn't necessary. I'm assuming the clipping happens after PopArt normalization, since if it happens before then you lose information as in the Pacman example. In this case, maybe it's that the reward distribution is fat-tailed, and so even after normalization there could be some extreme rewards that after normalization are still large enough that they would cause updates that are too large, and clipping alleviates this problem.

[**ML Writing Month May 2018**](#) (*Cody Wild*): The author wrote up a summary of an ML paper every day in May, which have all been collected in this doc.

My opinion: These summaries seem really good to me (probably higher quality than a typical summary that I write), but are often on topics I'm not an expert in (eg. GANs) so it's hard for me to evaluate. The one paper I knew well ([Inverse Reward Design](#)) had a good summary.

Technical AI alignment

Technical agendas and prioritization

[**Comment on decision theory**](#) (*Rob Bensinger*): MIRI works on Agent Foundations because AGIs should be good reasoners, and we're currently confused about how to have good reasoning, and work on logical uncertainty and decision theory should help us resolve this confusion. If we don't resolve the confusion, it will be significantly harder to build AGI in a way that is clean, understandable and interpretable, and as a result it will be harder to understand what is happening and to fix it if anything goes wrong. This is analogous to how it seems really useful to understand Newton's law of gravitation before you try to build rockets, even though the work of figuring out Newton's law is very different from the rocket-building work.

My opinion: My basic opinion is that this makes sense and agrees with my model. On the other hand, I'm not planning to switch to working on decision theory now, so perhaps I should say why. Partly it's that I have a comparative advantage at ML work, but it's also an impression that Agent Foundations will not help much with the first powerful AI systems we build. On one axis, I wouldn't be surprised if the first powerful AI systems don't look like the good reasoners that MIRI studies, and so Agent Foundations research won't apply. On another axis, Agent Foundations seems like a hard problem that we may not solve before powerful AI systems are created. I do find it plausible that to build *aligned* AI systems that are *much* more powerful than humans, we must understand it at the level of Agent Foundations understanding. (Though I also find the opposite statement plausible.) However, I think we will first build powerful AI systems that are not that much more powerful than humans, and that direct alignment of ML techniques will be sufficient to make that safe (even though they do pose an x-risk). (I suspect this is where my main disagreement with

people at MIRI is.) We can then use those systems to help us solve Agent Foundations before we scale up.

Iterated distillation and amplification

[Disagreement with Paul: alignment induction](#) (*Stuart Armstrong*): The amplification step in iterated distillation and amplification (IDA) requires an inductive argument saying that if the agent at level n is aligned, then so is the one at level $n+1$. However, in order to get the induction to work, you need to say not just that agent at level n won't take unaligned actions, but that it will also "assess the decisions of a higher agent in a way that preserves alignment (and preserves the preservation of alignment, and so on)". This seems like a much less intuitive criterion, and so getting the base case of an agent that a human can verify has this property may be too hard - it probably has to have all of the friendly utility function in the base case itself, or perhaps it gets it after one or two iterations (if it needs to solve a few problems along the way).

My opinion: If I imagine each level $A[n]$ as maximizing the expected value of some simple utility function, I agree that it would be surprising if the result was not one of your first three cases. Intuitively, either we already have all of the friendly utility function, and we didn't need induction, or we didn't and bad things happen, which corresponds to cases 1 and 3.

But it seems like one of the main points of iterated amplification is that at least the initial levels need not be maximizing the expected value of some simple utility. In that case, there seems to be a much wider space of possible designs.

For example, we could have a system that has the epistemic state of wanting to help humans but knowing that it doesn't know how best to do that, and so asking humans for feedback and deferring to them when appropriate. Such a system with amplification might eventually learn the friendly utility function and start maximizing that, but it seems like there could be many iterations before that point, during which it is corrigible in the sense of deferring to humans and not maximizing its current conception of what is best.

I don't have a strong sense at the moment what would happen, but it seems plausible that the induction will go through and will have "actually mattered".

Learning human intent

[Active Inverse Reward Design](#) (*Sören Mindermann et al*): (Note: I am an author on this paper.) Inverse Reward Design (IRD) introduced the idea that a hand-designed reward function should be treated as an *observation* about the intended reward function. Any particular hand-designed reward function (called a *proxy reward*) is likely if it incentivizes good behavior in the training environment, as measured by the true reward function. In this paper, instead of asking the reward designer to choose a proxy reward from the space of all possible reward functions, the designer is presented with a small subset of possible reward functions and asked to choose the best option from those. The subset is chosen such that the answer will be maximally informative about the true reward (in expectation). This is an easier query for the reward designer to answer, and can convey more information about the true reward function. To see this, we can imagine pairing each possible reward function with the trajectory that it incentivizes in the training environment. Then original IRD gets information about the

best such trajectory, which could correspond to multiple true reward functions, whereas active IRD can get information about the best trajectory in any subset of trajectories, and so can get more information in total. In some cases, that extra information can narrow down the space of possible reward functions to a single true reward function. The paper discusses two kinds of queries that can be asked (discrete and continuous), and several optimizations for actually computing the best query to ask (which can be computationally intensive). The technique is demonstrated on a contextual bandits problem and gridworlds.

Prerequisites: [Inverse Reward Design](#)

[Expert-augmented actor-critic for ViZDoom and Montezumas Revenge](#) (*Michał Garmulewicz et al*) (summarized by Richard): The authors augment ACKTR (a natural gradient RL algorithm) with an additional term in the loss function which depends on expert data. In particular, policies which choose different actions from samples of 14 expert trajectories are penalised, with a coefficient that depends on the expert's advantage over the critic's current expectations. This allows the agent to perform well on Montezuma's Revenge and a ViZDoom maze, sometimes beating the experts it was trained on. It also discovered a new bug in Montezuma's Revenge which increases its score by a factor of 40.

My opinion: I'm not convinced that this paper's method of utilising expert data is an improvement on other approaches, such as [this paper](#) in which an agent learns to play Montezuma's revenge from watching a Youtube video. However, it does seem to learn faster than most others, probably due to using ACKTR. I'd also expect it to be overfitting to the expert trajectories, but can't determine the extent to which this is the case (the authors claim that their agent can continue gameplay into the second world of Montezuma's Revenge despite only having expert trajectories for the first world, but don't provide metrics of success in the second world).

[Addressing Sample Inefficiency and Reward Bias in Inverse Reinforcement Learning](#) (*Ilya Kostrikov et al*)

Interpretability

[The What-If Tool: Code-Free Probing of Machine Learning Models](#) (*James Wexler*): Summarized in the highlights!

Verification

[Training for Faster Adversarial Robustness Verification via Inducing ReLU Stability](#) (*Kai Y. Xiao et al*)

Miscellaneous (Alignment)

[\(A -> B\) -> A](#) (*Scott Garrabrant*): This blog post has some thoughts on the type $(A \rightarrow B) \rightarrow A$, which can be thought of as the type of an agent. Rather than summarize the post, which seems hard, I'm going to say things about this type inspired by the post, and then you can decide whether to read the post.

My opinion: Intuitively, this type says that given something that maps A to B, you can get an A. If you think of $(A \rightarrow B)$ as an environment where A is the action and B is

the effect, then this type is a function that says which actions you should take. Note that the goal, which would be an element of B, is not passed in explicitly, so the goal must be inside of the function of type $(A \rightarrow B) \rightarrow A$, similar to our typical views of what an "agent" is. If you only assume that you know what A and B are, but you know nothing else, to do anything interesting you would be doing black box optimization -- that is, you get some f of type $(A \rightarrow B)$ that you know nothing about, and so you just keep computing $f(a)$ for different $a \in A$, looking for a particular $b \in B$. Perhaps you build a model of the function f and then do something more abstract with your model to find a good $a \in A$. (The post mentions a similar idea by noting that argmax has this type signature.) The post also has some thoughts about game theory that are interesting.

[Petrov corrigibility](#) (*Stuart Armstrong*): There can be situations where a human asks an AI for guidance, and the AI's action then determines what the human's preferences are. In such a situation, taking either action is incorrigible and so the AI should say "there is no corrigible action to take". But what if saying that would itself predictably change the human's decision? In that case, even that would not be corrigible.

There is also a comment chain discussing whether how this notion of corrigibility/alignment differs from the notion Paul Christiano talks about [here](#).

My opinion: I've written enough opinions about this version of corrigibility, I think I'd just be repeating myself. You can look through Stuart's recent posts and find my comments there if you really want (eg. [here](#)).

Near-term concerns

Adversarial examples

[Introducing the Unrestricted Adversarial Examples Challenge](#) (*Tom B. Brown et al*): Summarized in the highlights!

[Are adversarial examples inevitable?](#) (*Ali Shafahi et al*)

Other progress in AI

Reinforcement learning

[Preserving Outputs Precisely while Adaptively Rescaling Targets](#) (*Matteo Hessel et al*): Summarized in the highlights!

[Challenges of Context and Time in Reinforcement Learning: Introducing Space Fortress as a Benchmark](#) (*Akshat Agarwal et al*)

[Learn What Not to Learn: Action Elimination with Deep Reinforcement Learning](#) (*Tom Zahavy, Matan Haroush, Nadav Merlis et al*)

[Improving On-policy Learning with Statistical Reward Accumulation](#) (*Yubin Deng et al*)

[Keep it stupid simple](#) (*Erik J Peterson et al*)

[ViZDoom Competitions: Playing Doom from Pixels \(Marek Wydmuch et al\)](#)

Deep learning

[Neural Guided Constraint Logic Programming for Program Synthesis \(Lisa Zhang et al\)](#):

In program synthesis from examples, we want to find a program consistent with a given set of input-output examples. One classic approach is to use logic programming. In logic programming, instead of writing functions that compute output = f(input), we write rules to compute relations. To encode standard functions, we would write the relation (f, i, o), which is interpreted as "computing f(i) gives o". In logic programming, you can let any variable be unknown, and the language will search for a solution. Using this you can eg. invert a function f on a specific output o, using the query (f, ?, o). To apply logic programming to program synthesis, we write an interpreter eval for the language we want to synthesize in, and pose the query (eval, ?, i, o). They consider the lambda calculus with pairs and lists as their language.

The algorithm that falls out is a recursive descent search over the possible structure of the program, that generates and checks partial constraints over the partial programs implied by the input-output examples during the search. The search has branching points where it must choose, for some as-yet-unknown part of the program, what language construct it should use (if, cons, variable, etc.) This paper attempts to use a neural net to predict what choice the search should make to find a solution, replacing some simple hand-tuned heuristics. It can be trained either using reinforcement learning (where the search choices are actions, the partial search trees are states, and the goal is to find a complete program), or through supervised learning since they know for training programs what choices are optimal. They also use a curriculum and experience replay. They evaluate against classical symbolic approaches ($\lambda 2$, Escher, Myth) and RobustFill, and show that their method generalizes better to finding longer programs not seen in the training dataset.

My opinion: It always makes me happy to read a paper about making symbolic approaches faster using neural nets to learn heuristics. That said, I'm concerned about the evaluation in this paper -- their programs are fairly strange, often involving a huge mess of cons (make-pair), car (first) and cdr (second), and not including recursion. The symbolic approaches they evaluate against are aiming to synthesize recursive functions similar to what people write, and I wouldn't be surprised if they had heuristics that actively discouraged these big messes of cars and cdrs, since normal programs don't look like that. The programs are also primarily taking pieces of data out of an input, and then recombining them in some way -- this feels like a significantly easier task than most synthesis problems (in the sense that I could probably write a handcoded solution that performs very well on this domain only).

Miscellaneous (AI)

[ML Writing Month May 2018 \(Cody Wild\)](#): Summarized in the highlights!

A probabilistic off-switch that the agent is indifferent to

~~Edit: I no longer think this post deserves attention.~~

Abstract

This post presents a setup with an off-switch that is defective with probability of almost 0. The agent is indifferent to being terminated in worlds where the off-switch works. Also, the agent doesn't try to find out whether the off-switch is defective or not, because that would require an expensive computation. With some assumptions (described in the last section) the result is supposedly an off-switch that works with probability that is as close to 1 as we wish, and the agent is indifferent to it being pressed.

Notations

Let there be an agent that in each time step selects an action and then receives an observation. The goal of the agent is to maximize $E[u(h)]$, for some model of the environment, where u is a computable utility function and h is the history of all (action, observation) pairs from time step 1 until the end of the epoch.

The environment contains an off-switch that terminates the agent when pressed. In every time step after the agent is terminated, the corresponding action in h is the "null action".

The Proposed Setup

Each time we wish to invoke the agent, we invoke a program that does the following:

- (1) A random number x is generated, sampled uniformly from the range: $[0,10100]$.
Spoiler alert: the off-switch will be defective iff $x=0$.
- (2) Let f be some invertible function such that given $f(x)$, calculating x (i.e. applying f^{-1}) would require 0.1 seconds of intensive computations for some powerful computer C. The number $y:=f(x)$ is calculated.
- (3) The representation of x is erased in a way that would require a non-trivial effort to reconstruct.
- (4) The number y is sent to computer C which is part of the off-switch.
- (5) The number y is inserted into the code of the agent, such that the utility of the agent will be:
$$u'(h)=\{u(h)f^{-1}(y)=0 \alpha 1 + [\text{number of non-null actions in } h] \text{ otherwise}$$

for some $0 < \alpha \ll 1$.

Edit: With the above u' the solution has multiple problems that Stuart Armstrong pointed out (see [here](#), [here](#) and [here](#)). I think the new u' defined in [this comment](#) fixes these problems.

(6) The agent is invoked.

The probabilistic off-switch:

When the off-switch is pressed, computer C calculates x (i.e. $f^{-1}(y)$) and if the result is 0, nothing happens (i.e. the off-switch is defective). Otherwise, the agent is terminated.

Why should this work?

For a sufficiently small α , the agent would act as if it assumes $x=0$, without attempting to find out whether that's indeed the case. That's because that assumption would cost at most α utility (in worlds where $x \neq 0$); and would save some amount of utility (independent of α) in worlds where $x=0$, due to not spending resources on the expensive computation $f^{-1}(y)$.

When the agent assumes $x=0$, it is indifferent to the off-switch being pressed, because it assumes the off-switch is defective (even though the probability of it being defective is only 10–100).

Note: When we don't set α to be small enough, the agent might calculate $f^{-1}(y)$ (or reconstruct x in some other way). In such a case, the agent gets hold of x and will (with probability of almost 1) discover it's not 0, and then we'll observe the agent doing nothing (i.e. a sequence of null actions). In such a scenario we should simply decrease α and try again.

Flaws / Necessary Assumptions

ETA: This solution is not [reflectively stable](#) - creating subagents that maximize u is allowed and doesn't directly hurt the agent (thanks to Stuart Armstrong for [pointing this out](#)).

(1) Other than pressing the off-switch, the programmers must not have any way to terminate the agent or interfere with its plan. Otherwise, even when the agent assumes the off-switch is defective, it might attempt to deceive/neutralize the programmers.

(2) The agent's plan might involve destroying the mechanism of the off-switch or disabling the programmers in some way, in order to use the atoms/energy for something else. In that case the programmers must be able to realize this and press the off-switch button in time.

(3) If $x=0$ (i.e. the off-switch is defective) it must be sufficiently unlikely that the mechanism of the off-switch would cause the termination of the agent (e.g. due to a hardware failure or due to the programmers being able to somehow "repair" the defective off-switch). Otherwise, the agent would be willing to spend the resources needed for destroying the mechanism of the off-switch.

A possible way to tackle this problem is to slightly modify u such that the agent would want to reach some goal and then terminate itself ASAP. The modification should prevent the agent from accepting a delay of even a single time step for the purpose of eliminating a tiny risk of termination.

Asymptotic Decision Theory (Improved Writeup)

Crossposted from the [AI Alignment Forum](#). May contain more technical jargon than usual.

ADT ([asymptotic decision theory, initially detailed in this paper](#)) was a proposed decision theory with logical inductors, developed after [EDT with logical inductors](#) and exploration steps. It is a possible candidate for conceptual progress in decision theory, but some basic questions about its performance are still unsettled, and there are several issues with the current implementation of it.

Definitions:

A *market* is a computable function with type $S \rightarrow [0, 1]$ (where S is the set of sentences of math) that will be denoted as P . We will be considering logical inductors in this setting, and in particular, each finite stage of a logical inductor is a market.

Let A be a finite set of actions. An *agent* is a computable sequence of computable functions of type $P \rightarrow A$. It takes a timestep and a probability distribution, and selects a distribution over actions. Because logical inductors associate each timestep with an action, specifying the logical inductor and the agent specifies a sequence of actions. Because we will be assuming some fixed logical inductor in the background, we will suppress it in the notation and refer to the action produced by an agent at time t as a_t or b_t . Randomization can be handled by using the diagonalization sentences used to define exploration (this acts in such a way that no trader in the inductor is able to reliably predict what the agent does), or it can be handled by having one of the actions be to call a random number generator in the environment.

An *embedder* is a sequence of randomized functions of type $A \rightarrow [0, 1]$ (denoted by E or F , with the function at time t being denoted by E_t or F_t). Let $E_t(a_t)$ denote the random variable that corresponds to the environment using a uniform random distribution of bitstrings for its randomness tape. An embedder must have the probability distribution of $E_t(a_t)$ being computably approximable on computable inputs.

An *argmax agent* is an agent that takes a finite set of other agents A , and a single embedder E , and outputs $(\text{argmax}_{a \in A} E_t(E_t(a_t)))_t$ on turn t . E_t is defined as it usually is

for logical inductors. Notice that, although E_t may be very difficult to compute, putting it inside an expectation allows a logical inductor to output a decent guess as to what it is anyways. The agents all had to be computable in order for the argmax agent to duplicate their behavior at arbitrary turns. There is a dependence on the set A , the embedder E , and the logical inductor, but because the logical inductor and A will be

treated as fixed, we will write the time- t action produced by the argmax agent as $a_m^E_t$.

U will denote the "true environment", which is a sequence of values in $[0, 1]$.

Finally, the logical inductor will be required to have "fast feedback". This means that after each turn, the following information will show up in the deductive process at time $t + 1$.

1: (an interval containing) the true value of U_t .

The deductive process is the sequence of theorems that the logical inductor sees. The reason that the true utility has to be in an interval, instead of reported directly, is because the true utility may have arbitrarily many digits, which hampers the ability of traders to use that information to judge bets on how U will turn out. This condition is present to ensure that, if U is generated by taking E and feeding it a uniform random bitstring, $E_t(E_t(ADT_t)) \approx_t E_t(U)$.

Theorem 1: *If U is generated on each turn by running $E_t(ADT_t)$ with a bitstring sampled uniformly at random, then $E_t(E_t(ADT_t)) \approx_t E_t(U)$.*

We will apply the basic trading strategy from Theorem 4.5.9 in the [logical induction paper](#). In this case, if the left-hand side is overpriced by ϵ infinitely often (by a failure of convergence, and the same argument can be applied symmetrically to the right-hand side being overpriced infinitely often), the traders buy a fraction of a share of U and sell a fraction of a share of $E_t(ADT_t)$, and keep doing it until 1 share has been moved in total. This takes care of having P -generable trade magnitudes. According to the law of large numbers, with probability 1, there is a finite time after which all of the sub-traders have their pile of shares have a value within ϵ of 0, so the initial traders in the efficiently emulatable sequence of traders can be clipped off, and all the other sub-traders get ϵ or more value from selling the high-priced expectation and buying the low-priced one, so the necessary conditions for the ϵ -ROI Lemma are fulfilled and the resulting trader exploits the market.

Finally, we will define dominance. An agent a dominates an agent b on an embedder E if

$$\limsup_{t \rightarrow \infty} \sum_{i=1}^t (E_i(b) - E_i(a)) \leq 0$$

This could be thought of as a having sublinear regret relative to b .

What's ADT?:

The asymptotic decision theory algorithm works as follows.

Inputs: A sequence of numbers asymptotically decreasing to 0, denoted as δ , a finite list of embedders E , a finite list of agents A , and a logical inductor with fast feedback on the true environment, \bar{P} .

$$me := \text{rADT}(\delta, E, A, \bar{P})$$

$$R_t := \underset{\overset{E}{\text{Ind}_{\delta_t}}}{}(|E_t(E_t(me_t)) - E_t(U_t)| < \delta_t)$$

$$\underset{\overset{E}{\text{argmax}_{E \in E}(R_t \cdot E_t(E_t(am_t)))}}{}$$

$$ADT_t := am_t$$

Ok, so what does this do intuitively? Well, it takes all the embedders, and runs them through a "reality filter" R_t which checks whether the embedder, run on the agent

itself, replicates what the true reality U actually does, in expectation. For all the embedders that pass the reality filter, it uses the inductor to assess what score argmax gets in that embedder, takes the one that gets the highest score, and then implements argmax for that embedder. In short, it optimistically chooses the embedder where the highest score is attainable (that isn't wrong about reality), and optimizes for that. If it got a good score, that's fine. If it got a lower score than it was expecting, that embedder is less likely to pass the reality filter next time, because

$E_t(me_t)$ undershot U_t . There isn't a problem (as there is in standard [logical inductor decision theory](#)) of systematically underscoring an embedder forever, because since an embedder is a randomized function, it's possible to actually approximate a distribution over what it outputs, so argmax will eventually catch on and start taking (predicted-to-be) optimal actions, and the value of these can also be empirically assessed.

Problems with ADT:

The most prominent undesirable feature of this is that it's restricted to a finite set of embedders. Optimistic choice fails very badly on an infinite set of embedders, because we can consider an infinite sequence of embedders that are like "pressing the button dispenses 1 utility forever", "pressing the button delivers an electric shock the first time, and then dispenses 1 utility forever"... "pressing the button delivers an electric shock for the first n times, and then dispenses 1 utility forever"... "pressing the button just always delivers an electric shock". Optimistic choice among embedders keeps pressing the button, because, although it keeps getting shocked, there's always an embedder that promises that the agent's fortunes will turn around on the next turn.

Also, this optimizes for each choice individually, and does not naturally deal with sequences of choices, which is necessary to handle general environments.

Good Properties of ADT:

A nice feature of this is that exploration steps are not required for good behavior, which is important because [counterfactuals are not conditionals](#).

Another nice feature of this is that it gets ASP (agent simulates predictor) right, which is a surprisingly hard decision theory problem to do well on. When you are rewarded with 10^3 dollars for picking action 2, and paid $10^6 \cdot E_t(m_{et} = 1)$ dollars, the best move is to just take action 1 to win the million dollars, but [standard logical inductor decision theory](#) converges to taking action 2, because the predictor isn't powerful enough to predict the rare exploration steps, so the agent will learn that action 2 always gets it more money than action 1, and dial up the probability of taking action 2 until it ends up getting a thousand dollars on each round and missing the larger payoff.

However, because the embedder that represents the true environment is plugging things like "always pick 1" and "always pick 2" into the environment, argmax on that environment will converge to copying the "always pick 1" strategy, so the logical inductor learns that argmax will always pick 1, and then the true embedder claims that 10^6 dollars are attainable. If ADT learns to use the true embedder, then it will converge to always one-boxing, which is the desired behavior.

This same feature also allows it to win on Parfit's Hitchiker and XOR Blackmail (I think, 90% probability)

This does not pay up on counterfactual mugging.

The [original paper](#) on Asymptotic Decision Theory had *much* more restrictive restrictions for good behavior, such as the embedders having to be continuous, all the agents in A having to converge to a single distribution in the limit, all the embedders in E having to converge to a single payoff in the limit when fed a convergent agent, and having to use a continuous analogue of argmax, and in combination, this meant that most of the games you could define (even "predict this sequence of bits") were outside of the class of problems where optimality is guaranteed.

Why Haven't I Heard of this Before?

Well, for quite a while, we thought it was bad because [we thought it crashed into itself in games of chicken](#), so it got tabled for a while. I'll now go over why the argument in that post is false.

To begin, note that there's a crucial difference between the opponent in the "Spoof" embedder and in the "Delusion" embedder. In the first case, the only embedder that the opponent is optimizing over is the one with "go-straight bot" (or "swerve bot", depending on what you substitute in). In the second case, the opponent is optimizing over the same three embedders that you are using. ADT with only the "vs. go-straight bot" embedder converges to swerving, and ADT with only the "vs. swerve bot" embedder converges to going straight. So, in the "Spoof" embedder, the argmax agent will converge to thinking that "go straight" is the best thing to plug in because it gets a straight-swerve outcome.

Assume for the sake of contradiction that the ADT agent (with the 3 embedders listed) converges to going straight. Then on the true environment, it will keep getting straight-straight outcomes, while the "Spoof" embedder keeps predicting that you'll get straight-swerve outcomes. So, the "Spoof" embedder eventually fails the reality filter, so the agent will learn to not use it. Both of the remaining embedders advise swerving as getting the best outcome, so the ADT agent converges to swerving, and we have a contradiction.

What went wrong in the original reasoning? As far as I can tell, it originated from accidentally treating "ADT with only one embedder" and "ADT with the three listed embedders" as the same, because "ADT" was used to denote both of them.

Evil Problems and Theorem Failure:

To make things more confusing, the [second theorem in the ADT paper](#), about argmax dominating all agents in A, is wrong, and not in a fixable way. There's an explicit counterexample.

It will be instructive to take a detour to talk about [drnickbone's evil problem](#). Defining a "fair" problem in decision theory is necessary, because you can't just say that a decision theory is the best on all problems. Consider the problem where you are up against an opponent who just really hates some particular decision theory, and penalizes anyone that uses it to select actions. Of course, the decision theory of your choice will fail on this problem. So, instead, we would hope that there's some notion of a "fair" problem, such that we can say that a decision theory attains equal or greater utility than all competitors on all fair problems.

An initial attempt at defining a fair problem was one where all agents that select the same action get the same result. The problem is fully action-determined, and your

payoff doesn't depend at all on your ritual of cognition. This is the notion of "fair" used in the original ADT paper.

The Evil problem is a decision theory problem that is completely action-determined, and fair by the old definition of "fair" where (decision theory of your choice, heretofore abbreviated as YDT for "Your Decision Theory") gets systematically lower utility than most other competing decision theories. Consider a variant of Newcomb's problem where there are two boxes, and you may select one of them. Omega's decision process is to predict what box YDT would take, and then put 1000 dollars in the *other* box, and nothing in the box that YDT takes. Also, one of the boxes is lustrous and sparkly and would make a nice addition to your box collection, and you value that box at 10 dollars.

Now, you are like "well, I'm using YDT, and Omega is really accurate, and I left my random number generator at home, so no matter which box I pick, I'll get 0 dollars. Might as well go for the shiny box". And all the other decision theories like CDT and EDT can run through that reasoning and take the blank box, which contains 1000 dollars, and walk away substantially richer than before. Note that any agent that takes the same box gets the same payoff. So it is intuitively unfair, despite being completely action-determined.

This same sort of problem, when put into embedder form, leads to argmax systematically underperforming "take the blank box". The argmax agent will converge to taking the shiny box about 50.25% of the time. This is because, since it is possible to go back and compute what "blank-box" and "shiny-box" would have yielded on turns where you didn't take that box, you keep thinking that they'd get decent payoffs. In expectation, "blank-box" gets 502.5 dollars, and "shiny-box" gets 502.5 dollars, while "argmax" gets 5.025 dollars. This is an issue for *all* decision theories that operate by treating the environment as a function, and plugging actions into it. It will always consider the action it didn't take to have a higher utility, and this is actually true, because there are "objective counterfactuals" that can be checked. In particular, condition 1 in [my probabilistic tiling post](#) was required specifically to rule out these sorts of shenanigans. Check the discussion on why CDT fails this condition, it's talking about how these sorts of evil problems cause the expected utility of "I take an action" to systematically deviate from the expected utility of the action you actually take.

To return back to ADT, since it's treating the environment as a function that it can plug stuff into, it's vulnerable to this exploit.

So, what went wrong in the original proof that argmax dominates everything in A? It was the (not explicitly listed) step that went from $\lim_{t \rightarrow \infty} E_t(E_t(a_t) - E_t(a'_t)) = 0$ (which is true because a and a' are in a maximal equivalence class) to $\lim_{t \rightarrow \infty} E_t(|a_t, a'_t|) = 0$

(which is a necessary step to get argmax to have the same property). However, just because the utilities are the same doesn't mean the probability distributions are the same!! In fact, given the continuous version of argmax in that paper, it's possible to come up with a much more mundane case where it fails that *definitely* isn't an evil problem. Consider the problem where you must make the same move as your opponent, and your probability distribution over moves is the same as your opponent's probability distribution over moves. Let a and a' be the agents that just implement the two constant strategies. $E(a) = E(a') = 1$, but, by the definition of continuous argmax that was given in the paper, argmax would converge to a 50/50 mix of the two strategies since they are of equal value. When this is substituted into both your and your opponent's moveslot, it produces an expected utility of 0.5.

What is Fairness?:

These evil problems are a problem for showing that argmax dominates everything in A. While attempting a proof, I came up with a possible alternate definition of "fair", but it's more accurate to call it "regret-free". An embedder E is "regret-free" iff

$$\forall F : \liminf_{t \rightarrow \infty} E_t(E_t(am_t)) - E_t(F(am_t)) \geq 0 .$$

Intuitively, argmax doesn't regret its decision, in the sense of not wishing it was optimizing for some specific other world to get a more fortunate sequence of actions. The agent doesn't wish it was deluded about the world in order to be decorrelated with whatever is punishing its action sequence. One effect of this definition is that it shows that (according to the expectation of the agent), argmax will do better than any particular strategy in the strategy set, because you can consider an environment that rewards that exact strategy being played on every round. Another notable effect is that it rules out the original Evil problem as "unfair", but *the modified variant where the action of the agent is substituted into both your action, and Omega's prediction, is fair*. So there's still a regret-free/"fair" environment that can model the Evil problem faithfully, but it says that the proper thing to judge YDT against isn't environments where Omega penalizes YDT, but rather environments where Omega penalizes the decision theory of whoever is picking the boxes. And of course, CDT, EDT, and everything else also gets either 0 or 10 dollars worth of value in this modified problem, and balance is restored.

Sadly, this particular definition of "fair" is inadequate for general use, because it is possible to construct environments where some arbitrary decision theory other than argmax systematically gets lower utility than argmax. This can be fixed by making this definition relative to the agent under consideration, but then you run into the problem of simple agents calling (intuitively fair) games "unfair relative to me". There will be another post about this.

Lack-of-Proof-of-Optimality for ADT:

We might try to go for a theorem like

Conjecture 1: *If all environments in E are regret-free, then, on all embedders E s.t.*

$$\liminf_{t \rightarrow \infty} R_t^E > 0, \text{ ADT dominates } am^E.$$

In short, ADT seems like it will learn to do as well as argmax on all environments that don't get ruled out. This conjecture is still open! I thought I had a proof, and it failed, and then I thought I had a disproof, and it failed, and then I thought I had another disproof, and it turned out that I couldn't show that one of the environments was regret-free. Maybe some extra conditions are required for this conjecture to be a theorem, but I suspect that the environment is actually regret-free and the conjecture is false, pointing to a genuine hole in ADT.

I'll describe the first and last of these failed attempts here, in the hopes that they will provide help on how to conclusively prove or disprove the conjecture.

Failed Attempt 1:

To begin with, if the true environment E is in E, and *one* other environment F is

"clearly in the lead" (ie, $\forall F' \in E : R_t^F \cdot E_t(F_t(am_t)) \geq R_t^{F'} \cdot E_t(F'_t(am_t)) + \epsilon$ for some fixed ϵ), and this opportunity recurs infinitely often, it is possible to money-pump this.

In particular, since R_t^E (the reality filter for the true environment) converges to 1,

$E_t(F_t(am_t)) \geq E_t(E_t(am_t)) + \epsilon$. The money pump works by buying one share in $E_t(am_t)$,

selling one share in $F_t(am_t)$, selling one share in $E_t(ADT_t)$, and buying one share in $F_t(ADT_t)$. Because both of the latter have prices very similar to each other (because both E and F pass the reality-filter), their price cancels out, yielding ϵ dollars upfront from the first buy-sell pair, along with the fact that E is a regret-free environment, so

$E_t(am_t)$ is priced at or below $E_t(am_t)$ in the limit. Now, because the embedder F will be actually picked, all the purchased and sold shares have value that cancels out to 0. (well, actually, it's a bit more complicated than that, you need to add an additional condition that the information about what ADT copies shows up in the deductive process in a timely manner, because logical inductors judge the value of a trader by

their worst-case value.) Anyways, this pumps free money from each time where *one* particular embedder is "in the lead" above the true environment.

However, this doesn't rule out cases where multiple embedders that aren't the true environment are "in the lead". Intuitively, it's possible to have a pileup where multiple embedders have approximately the same estimate of their argmax value, and estimate of what ADT actually does in them. There's another case where the true environment is "in the lead" along with several others. The obvious fix to this is to have the trader buy and sell conditional contracts, so the money-pump goes through for whichever embedder is actually selected, and all the others cancel out to 0 value. The point where this fails is that it wouldn't necessarily be buying the shares as

F

originally stated. In order to get the $F_t(am_t)$ share, it would be purchasing the share at

F F F

a price of $E_t(F_t(am_t) | ADT_t = am_t)$ which may not equal $E_t(F_t(am_t))$. Attempts to construct more sophisticated money-pumps all met with failure due to this specific issue where regret-free environments may say that argmax does worse when ADT copies it.

Failed Attempt 2:

We will give a set of environments, and see that two of the embedders will converge to passing the reality filter, but ADT dominates neither of the argmax agents (however, one of the environments is not known to be regret-free.)

In particular, let the true environment be a game of chicken against the same ADT agent. A will consist of the three actions "swerve", "straight", and "consult random number generator in environment to swerve 90% of the time" (augmenting A to allow consulting the random number generator with different randomization probabilities doesn't change the result.) E will consist of the two environments E (which substitutes the action into both your action slot and the opponent's action slot), and F (which substitutes the action in your action slot, and the opponent is you. The utility function is 1 for going straight while your opponent swerves, 0 for crashing, and 0.8 for swerving.

To begin with, $E_t(E(ADT_t)) \approx_t E_t(F(ADT_t))$, because both environments are identical when you plug in the agent itself, so both of them pass the reality filter. Argmax on the embedder E converges to taking the action "swerve 90% of the time".

Assume that ADT learns to use the embedder E, predictably, infinitely often (that is, infinitely often, the market assigns a probability very near 1 of using the embedder E).

Then, on this sequence, it will converge to swerving 90% of the time, and in that case, ADT will learn to always use the embedder F in response, and go straight, because this passes the reality filter and offers a payoff of 0.9, which is greater than the payoff of 0.81 that E offers. So, we have a contradiction.

Assume that ADT learns to use the embedder F, predictably, infinitely often. Then, if the probability ADT assigns to itself going straight is $< 0.8 - \epsilon$, it goes straight, if the probability ADT assigns to itself going straight is $> 0.8 + \epsilon$, it swerves. So the probability ADT assigns to itself going straight on that subsequence will converge to

F

0.8, which means that $E_t(F(am_t)) \approx_t 0.8$, and because E promises a greater utility, it will learn to use the embedder E in those cases. So again, we have a contradiction.

So, according to the probabilities of the logical inductor itself, it is unpredictable what ADT will do, and by calibration, this applies to what ADT does in reality. Since

F E

differences in $E_t(F(am_t))$ and $E_t(E(am_t))$ allow predicting which embedder ADT will pick, the two must converge to each other, and since the latter converges to 0.81, the former must as well.

Now (I'm moderately sure it's the only solution, but don't have a proof), there is a stable solution where the agent doesn't know which embedder it picks, but with 90% probability it picks E, and with 10% probability it picks F (and argmax always picks "go straight", so it runs into itself), for an expected utility of 0.729, while both embedders claim that argmax for them specifically gets an expected utility of 0.81. Because it's unpredictable which embedder it will go with, the expected value of going straight (in embedder F) is 0.81, the expected value of swerving is 0.8, and the expected value of randomizing is 0.801. So argmax converges to always going straight, for an expected utility (in embedder F) of 0.81.

Birth order effect found in Nobel Laureates in Physics

[Epistemic status: Three different data sets pointing to something similar is at least interesting, make your own mind up as to how interesting!]

Follow-up to: [Fight Me, Psychologists, Birth Order Effects are Real and Very Strong, 2012 Survey Results, Historical mathematicians exhibit a birth order effect too](#)

In [Eli Tyre's analysis](#) of birth order in historical mathematicians, he mentioned analysing other STEM subjects for similar effects. In the comments I [kinda-sorta](#) preregistered a study into this. Following his comments I dropped the age requirement I mentioned as it no longer seemed necessary.

I found that Nobel Laureates in Physics are more likely to be firstborn than would be expected by chance. This effect (10 percentage points) is smaller than the effect found in the rationalist community or historical mathematicians (22 and 16.7 percentage points respectively) but is significant ($p=0.044$).

More brothers were found in the study than sisters (125:92 (58%)). After correcting for the correct expected ratio (~52%) this was found to not be significant ($p=0.11$).

I was unable to find sufficient data on Fields medal, Abel prize and Turing award winners.

My data and analysis is documented [here](#). With Eli's kind permission I used his spreadsheet as a template. I have kept Eli's data on the same Table – rows 4-153 are his.

Methodology

My methods matched Eli's closely except for the data sets I looked at, see his post for more information.

Initially I attempted to replicate Eli's results in other mathematicians by analysing Fields medal and Abel prize winners. Unfortunately I was unable to gather sufficient additional data. This is partly due to crossover in names between these mathematicians and the list from which Eli was working.

It also seems to be the case that less biographical information is available for people born after ~1950. This might be partly due to these people and their siblings being more likely to be still alive so data protection rules prevent e.g. geni from listing their full details (siblings' details are often set to "private") but there could be other reasons. For Fields medals awarded before 1986 I found data on 12/30 recipients, after that only 3/30.

I had a brief look at Turing award winners, as this would have seemed a relevant field to compare to the results from the rationalist community that inspired the studies, but came across the same problem.

Finally, I looked at Nobel laureates in Physics. A massive help in data collection here was the fact that since the 1970s Nobel laureates have been asked to supply an autobiography, which is published on the [Nobel website](#). Even before then there are biographies of each laureate although these seldom mention birth order.

Between the Nobel site, Wikipedia and [geni](#) I was able to find useful data on 100/207 Physics laureates. The other 107 either had no siblings or I couldn't find sufficient data on them – either way they weren't included in the analysis.

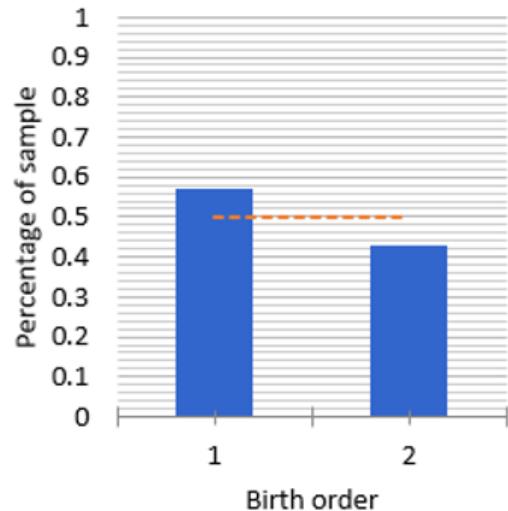
As a comment on data sources, I found geni to be somewhat unreliable. It contradicted the autobiographies or sometimes even contradicted itself. At other times, the list of siblings was incomplete or missing completely.

Results

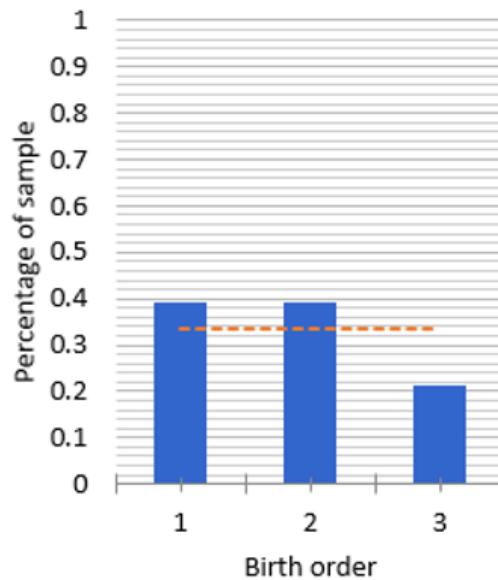
Categorising by family size shows that for all family sizes with ≥ 10 data points there are more firstborns than would be expected by chance.

Number of children in a family	Number families of that size	Actual number first born	Expected number first born	Actual > or < expected?	Actual minus expected	actual % first born	expected % firstborn	point difference (actual %)
1	107	107	107	=	0	100.00%	100.00%	0.00%
2	35	20	17.5	>	2.5	57.14%	50.00%	7.14%
3	28	11	9.333333333	>	1.666666667	39.29%	33.33%	5.95%
4	10	5	2.5	>	2.5	50.00%	25.00%	25.00%
5	13	4	2.6	>	1.4	30.77%	20.00%	10.77%
6	7	0	1.166666667	<	-1.166666667	0.00%	16.67%	-16.67%
7	3	2	0.428571429	>	1.571428571	66.67%	14.29%	52.38%
8	1	0	0.125	<	-0.125	0.00%	12.50%	-12.50%
9	1	1	0.111111111	>	0.888888889	100.00%	11.11%	88.89%
10	2	1	0.2	>	0.8	50.00%	10.00%	40.00%
11	0	0	0	=	0			
12	0	0	0	=	0	#DIV/0!	#DIV/0!	#DIV/0!
13	0	0	0	=	0	#DIV/0!	#DIV/0!	#DIV/0!
2 or more	100	44	33.96468254	0	10.03531746	44.00%	33.96%	10.04%
6 or more	14	4	2.031349206	0	1.968650794	28.57%	14.51%	14.06%

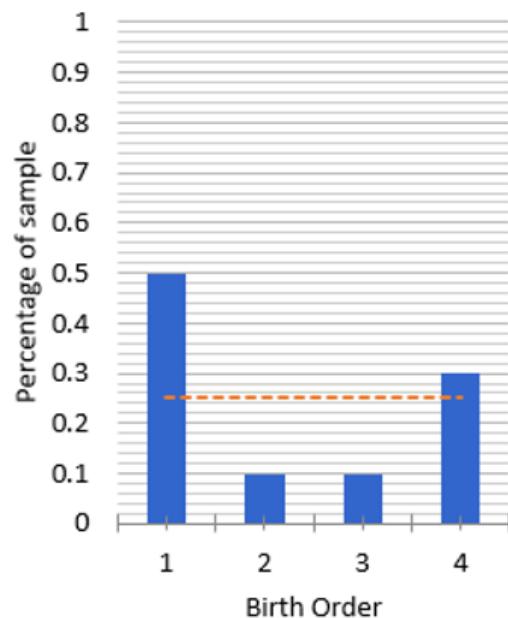
One of two 2 siblings



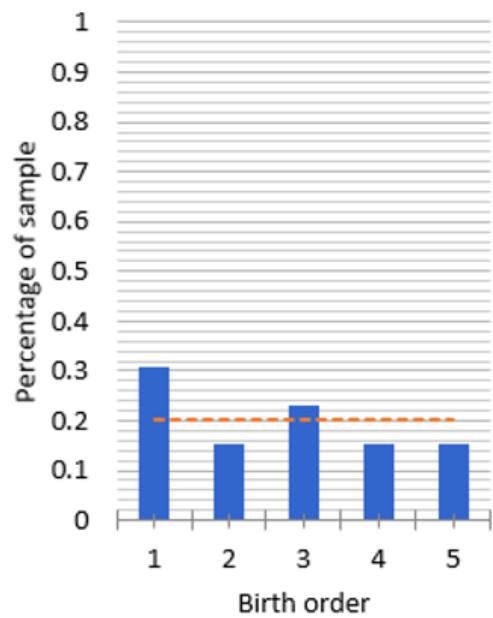
One of 3 siblings



One of 4 siblings



One of 5 siblings



Due to small sample size I have grouped all families of 6+ siblings into a single bucket and even then n=14. Expected birth order then varies with higher birth order as there are fewer families in the sample with at least that many children.

One of 6 plus siblings



Analysing the data as a whole gives a 10 percentage point effect (0.2 to 19.8 percentage points, 95% confidence). This is less than both the SSC / Less Wrong surveys and Eli's historical mathematicians analysis (22 and 17 percentage point respectively). I haven't got a number for overall confidence level for the SSC data but due to the large data set and very low p quoted for the 2 sibling example, it is unlikely that the 95% confidence interval overlaps with this new data, suggesting that the effect is truly a different size and not due to chance.

Discussion

Autobiographies as source material

Using autobiographies as the source for a significant number of the data points should have helped with the reliability of the data. It is possible that when writing an autobiography one would be more likely to mention siblings and birth order if one was the eldest but this doesn't seem likely.

Gender imbalance

Eli discussed under reporting of females as a potential source of bias. However, he found that the brothers:sisters ratio in his data was not unreasonable.

Running the same analysis on the physics Nobel laureate data I get a ratio of 125:92 brothers:sisters. This makes the siblings 58% male, with $p=0.03$ (binomial distribution, two tailed). This effect is actually more significant than the birth order effect.

Looking at the SSC data and Eli's data and found that there were 52% brothers in both. I did a little [research](#) and found that actually 51-52% is roughly the expected brother:sister ratio. I feel like this is something I should have already known but didn't.

[Another effect](#) which might increase the proportion of Nobel laureates brothers is that men can have a disposition to have boys or a disposition to have girls. As almost all of the

laureates are male it would be reasonable to think more of their Dads were predisposed to having boys. However as this isn't seen in SSC or historical mathematicians data (both also male dominated) this doesn't really get us much further.

Using 52% as the expected ratio (instead of 50%) means that the 58% result from Nobel laureates no longer rises to significance ($p=0.11$) and should instead be labelled as "[hey, look at this interesting subgroup analysis](#)" or possibly "slightly odd but not implausible".

As I mentioned previously, most of the data since the 1970s Nobels is based on autobiographies. Looking at only data since then, the brother:sister ratio is 51:35 (59%). It seems unlikely that Nobel laureates forgot about some of their sisters, making it less likely that the gender imbalance is due to incorrect data.

One potential source of error in the gender balance may be in the siblings whose gender I was unable to determine. There were 50 of these. Most (41) of these came from families where I had no data except the number of siblings and the position of the laureate within the family (e.g. "I was the fourth of five children."). It is possible that some of the missing sisters are in this category.

However, this would imply that if someone has more brothers they are more likely to list the genders of their siblings than if they have more sisters. Perhaps as most of the laureates were male they might have had more in common with brothers and spend more time with them, making them statistically more likely to mention their brothers' gender. This seems plausible but unlikely to cause a big effect even if it were true.

For the moment, I am working with the assumption that the sample is accurate and that the gender imbalance is just an outlier. Any other thoughts on causes of bias are welcome. These would have to explain how this effect was seen both in data from both geni and the laureates' autobiographies.

Conclusion

Nobel laureates in physics exhibit a birth order effect such that they are 10 percentage points more likely to be the eldest child than would be expected ($p=0.044$). This effect is less than data from both SSC readers and historical mathematicians (22 and 17 percentage points respectively).

There was a gender imbalance between brothers and sisters (58% brothers) but, taking into account the expected ratio of 52%, this was not significant ($p=0.11$). This effect is not seen in SSC readers or historical mathematicians (52% in both)

I would recommend that anyone who wishes to collate additional historical data consider Nobel laureates in other awards due to the availability of accurate data from the autobiographies. My analysis took perhaps 12 hours but a lot of that was spent on wild goose chases in looking for data on Fields medal and Abel prize recipients. I saved a lot of time by reusing Eli's spreadsheet (thanks for the permission). I would estimate getting data on the entire history of another Nobel prize category and analysing it would take ~6-8 hours so it shouldn't be too daunting for someone to take on.

Direct Primary Care

Epistemic Status: PSA and raising questions

Medical care in the US is expensive. There aren't that many demonstrated ways to make it drastically cheaper.

Direct primary care seems to be an exception. It makes routine medical expenses **95% cheaper.**

[Yes, really.](#)

If you have a cash-only, or “direct”, primary care practice — i.e. you *don’t accept insurance* — you can negotiate much lower wholesale rates from providers of tests (like EKGs, MRIs, blood tests) or prescription drugs. Why? Because you can guarantee the providers immediate payment, rather than the uncertainty and inconvenience of insurance reimbursement. They’re willing to give you a discount for that.

Direct primary care also cuts down on paperwork for doctors, because they don’t have to document everything with the ICD codes that insurance companies require.

[Atlas MD](#) is an EMR designed for direct primary care practices that use a subscription model, founded by Dr. Josh Umehr, who also uses it in his own direct primary care practice in Wichita. I’ve spoken to him via phone and tried to poke holes in his model, and came away even more impressed.

My question is, could this model scale up nationwide — and would it still be as effective at cutting costs if it did?

Right now, as I understand, direct primary care practices negotiate individually with suppliers to get discounts. I imagine that it would be much more efficient if done by a nationwide chain of direct primary care practices. *Bulk* wholesale purchases, after all, could be even cheaper than what a single practice might hope to get.

With Amazon moving into [healthcare](#), direct primary care might get a chance to shine. Amazon has a lot of experience cutting prices through economies of scale & supply chain optimization. Jeff Bezos even funded direct primary care startup [Qliance, which went bankrupt last year.](#)

Direct primary care only works as a complement to insurance that pays for more catastrophic care like emergency room visits and specialists. And if you can get a “minimalist” insurance plan that’s not redundant with the direct primary care membership, your total healthcare costs (membership + premiums) can be much lower. The potential problem arises if it’s difficult to sell sufficiently “bare-bones” insurance plans — in that case patients wouldn’t be willing to pay out of pocket for a direct-care membership in addition to their already pricey insurance.

Umehr has managed to negotiate deals with insurers to offer lower premiums when patients bought insurance along with direct care subscriptions, but maybe Qliance, which apparently struggled to keep customers, didn’t successfully pull it off.

At any rate, if I’m not missing something, this seems like an ideal opportunity for Amazon to make healthcare a lot more affordable. Are there barriers I haven’t thought

of?

Hypothesis about how social stuff works and arises

EDIT, 2022:

this post is still a reasonable starting point, but I need to post a revised version that emphasizes preventing dominance outside of play. **All forms of dominance must be prevented, if society is to heal from our errors.** These days I speak of human communication primarily in terms of the network messaging protocols that equate to establishing the state of being in communication, modulated by permissions defined by willingness to demonstrate friendliness via action. I originally wrote this post to counter "status is all you need" type thinking, and in retrospect, I don't think I went anywhere near far enough in eliminating hierarchy and status from my thinking.

With that reasoning error warning in mind, original post continues:

Preface

(I can't be bothered to write a real Serious Post, so I'm just going to write this like a tumblr post. y'all are tryhards with writing and it's booooring, and also I have a lot of tangentially related stuff to say. Pls critique based on content. If something is unclear, quote it and ask for clarification)

Alright so, this is intended to be an explicit description that, hopefully, could be turned into an actual program, that would generate the same low-level behavior as the way social stuff arises from brains. Any divergence is a mistake, and should be called out and corrected. it is not intended to be a *fake* framework. it's either actually a description of parts of the causal graph that are above a threshold level of impact, or it's wrong. It's hopefully also a good framework. I'm pretty sure it's wrong in important ways, I'd like to hear what people suggest to improve it.

Recommended knowledge: vague understanding of what's known about how the cortex sheet implements fast inference/how "system 1" works, how human reward works, etc, and/or how ANNs work, how reinforcement learning works, etc.

The hope is that the computational model would generate social stuff we actually see, as high-probability special cases - in semi-technical terms you can ignore if you want, I'm hopeful it's a good causal/generative model, aka that it allows compressing common social patterns with at least somewhat accurate causal graphs.

The thing

So we're making an executable model of part of the brain, so I'm going to write it as a series of changes I'm going to make. (I'm uncomfortable with the structured-ness of this, if anyone has any ideas for how to generalize it, that would be helpful.)

1. To start our brain thingy off, add **direct preferences**: experiences our new brain wants to have. Make negative things much worse, maybe around 5x, than

good things.

- From the inside, this is an experience that in-the-moment is enjoyable/satisfying/juicy/fun/rewarding/attractive to you/thrilling/etc etc. Basic stuff like drinking water, having snuggles, being accepted, etc - preferences that are nature and not nurture.
 - From the outside, this is something like the experience producing dopamine/serotonin/endorphin/oxytocin/etc in, like, a young child or something - ie, it's *natively* rewarding.
 - In the implementable form of this model, our reinforcement learner needs a state-reward function.
 - Social sort of exists here, but only in the form that if an agent can give something you want, such as snuggles, then you want that interaction.
2. Then, make the direct preferences update by **pulling the rewards back through time.**
- From the inside, this is the experience of things that lead to rewarding things becoming rewarding themselves - operant conditioning and preferences that come from nurture, eg complex flavor preferences, room layout preferences, preferences for stability, preferences for hygiene being easy, preferences for stability, etc.
 - From the outside, this is how dopamine release and such happens when a stimulus is presented that indicates an increase in future reward
 - In the implementable form of this model, this is any temporal difference learning technique, such as q learning
 - Social exists more here, in that our agent learns which agents reliably produce experiences are level-1 preferred vs dispreferred. If there's a level-1 boring/dragging/painful/etc thing another agent does, it *might* result in an update towards lower probability of good interactions with that agent in that context. If there's a level-1 fun/good/satisfying/etc thing another agent does, it *might* result in an update towards that agent being good to interact with in that context and maybe in others.
3. Then, modify preferences to deal with **one-on-one interactions with other agents:**
- Add tracking of retribution for other agents
 - From the inside, this is feeling that you are your own person, getting angry if someone does something you don't like, and becoming less angry if you feel that they're actually sorry.
 - From the outside, this is people being quick to anger and not thinking things through before getting angry about Bad Things. something about SNS as well. I'm less familiar with neural implementation of anger.
 - To implement: Track retribution-worthiness of the other agent. Increase it if the other agent does something you consider retribution-worthy. Initialize what's retribution-worthy to be "anything that hurts me". Initialize retribution-worthiness of other agents to be zero. Decrease retribution-worthiness once retribution has been enacted and accepted as itself not retribution-worthy by the other agent.
 - Track deservingness/caring-for other agents. Keep decreasing an agents' deservingness open as an option for how to enact retribution.
 - From the inside, this is the feeling that you want good for other people/urge to be fair. It is not the same thing as empathy.
 - From the outside, this is people naturally having moral systems.
 - To implement, have a world model that allows inferring other agents's locations and preferences, and mix their preferences with

yours a little, or something. correct implementation is safe ai

- Track physical power-over-the-world of you vs other agents
 - From the inside, this is the feeling that someone else is more powerful or that you are more powerful. (fixme: Also something about the impro thing goes here? how to integrate?)
 - From the outside, this is animals' hardcoded tracking of threat/power signaling - I'd expect to find it at least in other mammals
 - To implement, hand-train a pattern matcher on [Threatening vs Nonthreatening] data, and provide this as a feature to reinforcement learning; also increase deservingness/decrease retributionworthiness for agents that have high power, because they are able to force this, so treat it as an acausal trade
- 4. Then, **track other agent's beliefs** to iterate this over a social graph
 - Track other agent's coalition-building power, update the power-over-the-world dominance based on an agent's ability to build coalitions and harness other agent's power.
 - From the inside, this is the feeling that someone else has a lot of friends/is popular, or that you have a lot of friends/are popular
 - Track other agents' verbal trustworthiness, update your models on level 2 directly from trusted agents' statements of fact
 - Track other agents' retribution lists to form consensus on what is retribution-worthy; update what you treat as retribution-worthy off of what other agents will punish you for not punishing
 - Track other agents' retribution status and deservingness among other agents, in case of coordinated punishment.
 - Predict agents' Rewardingness, Retribution-worthiness, Deservingness, and Power based on any proxy signals you can get - try to update as fast as possible.
 - Implementation: I *think* all you need to do is add a world model capable of rolling-in modeling other agents modeling other agents etc as feelings, and then all of level 4 should naturally fall out of tracking stuff from earlier levels, but I'm not sure. For what I mean by rolling-in, see [Unrolling social metacognition](#)

Things that seem like they're missing to me

- Greg pointed out that current artificial RL (ie, step 1) is missing something simple and important about the way reward works in the brain, but neither of us are quite sure what exactly it is.
- Greg also pointed out that the way I'm thinking about power here doesn't properly take into account the second to second impro thing
- Greg thought there were interesting bits about how people do empathy that disagree really hard with the way I thought level 3 works
- Lex had a bunch of interesting critiques I didn't really understand well enough to use. I thiink I might have integrated them at this point? not sure.
- A bunch of people including me hate anything that has levels for being probably more complicated in terms of being organized structurally and simpler in terms of amount of detail than reality actually has. But I still feel like the levels thing is actually a pretty damn good representation. Suggestions welcome, callouts are not
- This explanation sucks and people probably won't get useful intuitions out of this the way I have from thinking about it a lot

misc interesting consequences

- level 4 makes each of the other levels into partially-grounded [keynesian beauty contests](#) - a thing from economics that was intended to model the stock market - which I think is where a lot of "status signaling" stuff comes from. But that doesn't mean there isn't a real beauty contest underneath.
- level 2 means it's not merely a single "emotional bank account" deciding whether people enjoy you - it's a question of whether they *predict* you'll be fun to be around, which they can keep doing even if you make a large mistake once.
- level 3 Deservingness is referring to how when people say "I like you but I don't want to interact with you", there is a meaningful prediction about their future behavior being positive towards you that they're making - they just won't necessarily want to like, hang out

Examples of things to analyze would be welcome, to exercise the model, whether the examples fit in it or not; I'll share some more at some point, I have a bunch of notes to share.

New DeepMind AI Safety Research Blog

Crossposted from the [AI Alignment Forum](#). May contain more technical jargon than usual.

This is a linkpost for <https://medium.com/@deepmindsafetyresearch/building-safe-artificial-intelligence-52f5f75058f1>

Very excited to announce the new DeepMind Safety Research blog! The [first post](#) by Pedro Ortega and Vishal Maini categorizes AI safety problems into specification, robustness and assurance:

Specification (Define purpose of the system)	Robustness (Design system to withstand perturbations)	Assurance (Monitor and control system activity)
Design Bugs & inconsistencies Ambiguities Side-effects High-level specification languages Preference learning Design protocols	Prevention and Risk Risk sensitivity Uncertainty estimates Safety margins Safe exploration Cautious generalisation Verification Adversaries	Monitoring Interpretability Behavioural screening Activity traces Estimates of causal influence Machine theory of mind Tripwires & honeypots
Emergent Wireheading Delusions Metalearning and sub-agents Detecting emergent behaviour	Recovery and Stability Instability Error-correction Failsafe mechanisms Distributional shift Graceful degradation	Enforcement Interruptibility Boxing Authorisation system Encryption Human override
Theory (Modelling and understanding AI systems)		

"In this inaugural post, we discuss three areas of technical AI safety: specification, robustness, and assurance. Future posts will broadly fit within the framework outlined here. While our views will inevitably evolve over time, we feel these three areas cover a sufficiently wide spectrum to provide a useful categorisation for ongoing and future research."

(A → B) → A

Crossposted from the [AI Alignment Forum](#). May contain more technical jargon than usual.

This post is about following type signature, which I call the type of agency: $(A \rightarrow B) \rightarrow A$

You can also think of it as consequentialism or doing things on purpose. This post will be a rant with a bunch of random thoughts related to this type signature, and it will likely not make sense. It will also be sloppy and will have type errors, but I think it is worth posting anyway.

First, interpret these arrows as causal arrows, but you can also think of them as function arrows. This is saying that the causal relationship from A to B causes A to happen. Think of A as an action and B as the goal. The reason that A happens is the fact that it has B as a consequence. There are not normally exponential objects like this in Bayes' nets, but I think you can modify so that it makes sense. (I'm not sure that this works, but you have a Cartesian closed category with nodes that are the nodes in your Bayes net, and add small number of morphisms from product nodes to individual nodes, corresponding to the functions in the Bayes' net. The acyclicity of the Bayes' net roughly corresponds to this category being thin. Then you can consider having other types of morphisms that can keep the category thin.)

If you have a game between two agents with action nodes A_1 and A_2 , with utilities U_1 and U_2 . The game implements a pair of functions $A_1 \times A_2 \rightarrow U_1$ and $A_1 \times A_2 \rightarrow U_2$. We can Curry these functions and think of them as $A_2 \rightarrow (A_1 \rightarrow U_1)$ and $A_1 \rightarrow (A_2 \rightarrow U_2)$.

Bringing in the agency $(A_i \rightarrow U_i) \rightarrow A_i$ of both players leads to cycle. This cycle does not make sense unless the agency arrows are lossy in some way, so as to not be able to create a contradiction.

Fortunately, there is another reason to think that these agency arrows will be lossy. Lawvere's Fixed Point Theorem says that in a Cartesian closed category, unless B has the fixed point property, you cannot have a surjective function $A \rightarrow (A \rightarrow B)$, in Set this is saying that if B has more than one element, you cannot have an injection $(A \rightarrow B) \rightarrow A$. i.e. The agency arrows have to be lossy.

Also, notice that Argmax, takes in a function f from some set A to R, and returns an element of the domain, A, so Argmax has type $(A \rightarrow R) \rightarrow A$.

This one is a bit more of a stretch, but if you look at gradient descent, you have some space X , you have a function $f : X \rightarrow R$. The gradient can be thought of as a function from infinitesimal changes in X to infinitesimal changes in $f(X)$. Gradient descent works by converting this gradient into a change in X . i.e. Gradient descent looks kind of like $(\partial X \rightarrow \partial R) \rightarrow \partial X$.

Good Citizenship Is Out Of Date

This is a linkpost for <https://benlandautaylor.com/2018/09/03/good-citizenship-is-out-of-date/>

Norms of good citizenship have been declining. These norms are a crucial piece of social technology vital to the health of local communities and institutions. While good citizenship norms are certainly still present in America today, they are substantially weaker than they were in the 1930s-1950s. This is not because of contemporary people's personal failings; rather, it's because we're still operating from a foundation of norms that were built for the New Deal era, and so are not adapted to today's conditions.

A society's norms lead to better or worse outcomes depending on how well they fit the circumstances. For example, in a small town, politeness norms often involve greeting everyone you pass and sometimes chatting a bit; this functions well because there are few people and they mostly know and care about each other. In New York City, this would be utterly impractical, so instead politeness norms demand ignoring passersby. Less adaptive norms will naturally lose force as people notice that they don't lead to good outcomes. Norms can be adapted to physical characteristics (like population), to the landscape of institutions (contrast American vs Mexican norms of bribing police officers, which are adapted to the local police institutions), or even to other norms (contrast American vs Japanese norms of public cleanliness, which are adapted to local levels of conscientiousness and trust).

In the mid-1900s, the norms of good citizenship were richer and more powerful than today. There was a shared idea that the good citizen was an active and integral part of his or her (norms differed somewhat by gender, but there was more similarity than difference) local community, as captured by arch-Americanist Norman Rockwell in his iconic *Freedom of Speech*. The good citizen was supposed to be involved with organizing at least one local civic organization, perhaps a church, or a local relief society, or a fraternal club like the Shriners. My grandfather made a point of serving on the board of the St. Louis chapter of the ACLU and writing incessant letters about local issues to the *St. Louis Post-Dispatch*, while his wife was heavily involved with the St. Louis planetarium and science center. For them, these things were part of a larger project through which engaged citizens would do their part to bring about a better world.

(I don't mean to imply that everyone was always, or even usually, following such norms. For most people, these things are aspirational, like the contemporary norm that one should read the article before sharing an inflammatory link. However, even aspirational norms can have a notable effect on most people—think of how the idea of homeownership affects even people who rent, or how the idea of launching a startup affects programmers who have never founded a company—and an influential minority will make a serious project of living up to the ideal.)

Over time, society changed, and the norms became less adaptive and thus less powerful. For example, *12 Angry Men*, a classic of 1950s American civics, shows how a good juror was meant to behave: a bulwark of Enlightenment justice shielding the common man from the passions of the mob, independent-minded, reasonable, and charitable. (I don't think fiction determines these patterns, but I do think it reflects them, and sometimes crystallizes them into their most coherent forms.) Since then, as

[jury trials have been dropping off](#) in favor of plea bargains, these norms have become [less relevant](#). This pattern has played out many times, in ways large and small: some part of society changes, so the norms relating to that part become less functional or less important, and so the norms atrophy.

As a result of this process, norms of good citizenship are not nearly as satisfying to aspire to as they once were. Today's citizenship norms tend to be negative rather than positive: don't be racist, don't damage the environment, don't fall for fake news. The few positive directives tend to advocate vague and passive things like "being informed", or at most participation in a large faceless mass, such as voting or marching in protests. There is no conception that a good citizen should *build*, in the way that a citizen of old would aspire to support the opera house or be a voice at City Council debates or what have you. There are still people who build local institutions, of course, but when I talk to them they mostly seem motivated by local pride, and not by the idea of participating in an overarching national or civilizational project that motivated my grandfather's generation. Not coincidentally, this call is now much more rarely felt by upper-class or upper-middle class people, who today often see themselves as too cosmopolitan to be involved with local institutions.

In [Mr. Smith Goes To Washington](#), note how establishing that Mr. Smith is a popular Boy Scout leader is instantly sufficient to tell the audience that he's an upstanding, competent pillar of the community, more worthy of power than the corrupt insiders who know how to work the system. His role as a local institution-builder makes him part of the living sinew of civic society, and it is morally right (if not necessarily practical) that he should become a Senator. Today's culture doesn't have any roles with quite the same cachet.

A large reason for the decline in norms around building local communities is that there is a new source of competition for organizational talent: building online communities. From personal experience, I know that leading local and online communities can be socially rewarding in similar ways. So, they will draw from a strongly overlapping talent pool. While online communities fulfill some of the functions of local communities, they don't fulfill nearly all of them. Building online communities is not a part of good citizenship ideals in the way that building local communities used to be (try to imagine a modern remake of *Mr. Smith Goes To Washington* where Mr. Smith is a beloved forum moderator), largely because we don't know how to make a complete civil society out of online institutions.

The decline of these norms is a loss, and our society is the poorer for it. However, they cannot be restored by simply repeating what our ancestors did; the reason the old norms fell out of favor in the first place is that they are no longer as fit for their purpose. If similar norms are to exist in the future—and I believe they can—then they must be built to function in the social and technological landscape of today.

Criticism Scheduling and Privacy

Scheduling criticism

"Two times two is four 'tis true But too empty and too trite I would rather seek a clue To some matters not so light."

— W. Busch, *Schein und Sein* (Karl Popper's translation)

Being in control of when, how, and what you get criticism on is a vital part of the growth of knowledge.

Scheduling criticism isn't being closed-minded. It's the only way to deal with the fact that there is an infinity of possible criticism, and no mechanical way to make that infinity smaller. You need some way of choosing what is useful.

The self-development idea of 'growth mindset' says: learn from criticism, don't ignore it.

There is an exaggerated version of this idea, namely: "Because criticism is not actually bad and in fact helps you learn, criticism is always good (or at least neutral)."

The grain of truth in this is: Identifying errors gives you access to the problems in your ideas. Any criticism *could* help. Finding conflicts between your ideas is the first step to solving them. A criticism on its own, as an idea in the ether, is only information — and this is information that can be used to make progress.

So why isn't criticism purely a good thing? How could noticing errors be bad?

Well, it's true that someone might have an irrational aversion to making mistakes (perhaps after having been penalised in school for getting things wrong), which causes them to be averse to even *noticing* mistakes.

But if we got over this, and didn't take criticism personally, wouldn't criticism just be a gift?

No. Criticism, when unwanted, can effectively destroy the means of error correction and the growth of knowledge. Put differently: it can structurally cripple thinking.

When I say "unwanted", I'm not referring to emotional reactions to criticism. One *can* feel bad about criticism by taking it personally — and it is possible to avoid this kind of bad feeling, by learning to take criticism better — but that isn't what I'm talking about here. That's a personal hangup, whereas I'm pointing at something more fundamental.

Criticism is how we learn

Criticism is part of how we learn. We have a problem we're trying to solve or a question we're trying answer, we come up with potential answers, and we criticise those answers to narrow them down or polish them. Once we have a promising

solution, we criticise it further to understand it better and to test that it actually works as a solution.

We need to criticise ideas to understand them at all. (This often happens subconsciously. Criticism needn't be at the explicit level of awareness; though making criticisms explicit can help make them clearer. All thinking has a large inexplicit component, including criticism.) If we didn't criticise ideas, we wouldn't know things like what else the idea solves, where it's applicable vs we need a different concept, why *that* solution instead of something else, and what it even means.

Criticism also comes into play earlier, when first coming up with potential answers. If it were 'anything goes', they would be random useless nonsense. (A bit like dream logic, though even dreams must have some kind of criticism filters to make as much sense as they do.)

For instance: If you're modelling the mind of your professor, you might come up with different guesses of the solution than if you were modelling the mind of your favourite author. Those two models give different kinds of criticism. Likewise if you imagine presenting to different audiences. **The ideas you come up with will be affected by what kind of criticism you're anticipating.** So, even when generating ideas in the first place, the effects of criticism are always in the background.

If we're not in control of what criticism we think about, if we outsource this critical process to someone else, we can't really *understand* the idea. We might get the bottom line or some individual applications we've been told. But that's only useful for things like passing exams or repeating propaganda.

The criticisms you consider will depend on what *your* specific idea-conflicts are. Each person is really different. All individuals have unique misconceptions, assumptions, background knowledge, connections between ideas and ideas in conflict.

Someone who thinks animals have sentience might find that conflicts with eating meat, whereas someone who doesn't may not. Animals having sentience may be used in criticisms of other ideas, like how to build factories in uncultivated land, or how to treat working animals like horses. Someone who doesn't think animals have sentience may not think these are issues — or may think they are issues, but for very different reasons, and so come up with different solutions.

To learn something, we need to consider the conflicts and address misconceptions *we* have. Not someone else's.

This comes up in standardised education systems. Textbooks, to be accessible to a reasonably wide audience of different people, are written for an average reader. But really, there is no 'average reader' — each person's interests and problems are unique. Thick textbooks are usually written with this in mind, arranged so that the student can jump around and focus on the parts relevant to them. It's rare that reading a textbook cover-to-cover is exactly what would help a specific individual learn best.

Even if we have a popular misconception, our solution won't be exactly the same as someone else's. It has different constraints, because we all have a unique set of background ideas. The things we're curious or confused about will be different.

Why you need to be in control of this

So, not any criticism will do. We need a specialised set of criticism to address the ideas we actually have.

The other problem with using non-specialised criticism is that there's practically an infinite amount of criticism available. There's more criticism in the world than you have seconds alive to hear it.

Criticism is very easy to come up with. Good criticism is harder. Good criticism that actually targets *your* problems and misconceptions is much harder. One needs a way of selecting what to pay attention to and narrow down.

The criticism you come up with yourself has the best bet of being relevant and broadly compatible with your background ideas.

Writers often make a few passes of a draft themselves before showing it to other people for feedback. When we have a new idea, we sometimes want to flesh it out before we share it with other people. If we share too early, criticism may cause us to take it in a different direction from what we wanted — or it might die entirely — when we could have thought about it more and explained it in a different way. It's no longer addressing our own problem, it's addressing the problem of the intervener.

After yourself, the next best is getting criticism from sources you choose. Like asking people questions, reading about the topic, comparing your solution to other people's, or submitting the specific idea/work to individuals you think will give useful feedback. Especially when you can *direct* this — such as asking specific questions, or reading only the parts of the book that feel most relevant.

But we are good at [fooling ourselves](#), and sometimes we don't even know we have a misconception. It's easier to point out other people's mistakes than our own, because our own are often hidden in our blind spot. *How can we correct errors in our blind spot if we are the ones directing the critical process?*

Directing doesn't mean you know exactly what you're going to get beforehand. Sometimes, you want to get more *general* feedback, or feedback on the process you're using. Or you may read about fundamental ideas to understand one thing better, and find they reveal other misconceptions. But all of this will come out of some problem, some interest, some question that *you have*.

Seeking general feedback is only useful in rare situations. Usually, the more specific your question is, the more useful the feedback will be. For instance, if you posted a work of art to a general 'arts' forum, you might not get useful criticism compared to posting it to a forum dedicated to the type of art you do. Unless, of course, you had the specific problem of wanting to know what people outside your subfield think!

So to be relevant to your problems at any given time, only certain kinds of criticism are helpful — namely, criticism that is wanted. The rest is useless or actively harmful.

Why **harmful**, though?

The critical process needs to happen inside one's own mind to learn something at all. There needs to be some kind of process to decide which things are relevant, which things help most with the current problem you are on, and what to try next.

And it needs to address *your* problems and ideas — not abstract 'problems', or someone else's ideas. Otherwise, it's not really *you* learning it; it's not being incorporated into your mind. You could memorise some facts (or how to parrot them), but to *understand* them requires that they pass through this critical process.

When this process gets outsourced — such as being forced to learn things in school, or trying to learn something because you 'should' (according to someone else, who may be imaginary) — **you stop judging whether a given criticism solves your problem**. You stop assessing it according to your own criteria.

Criticism that *you* would naturally apply to the thing you're being compelled to learn thus gets suppressed. (Including meta criticism, like "This is boring", "I want to do something else", etc. Those criticisms indicate that *something else* might be better to learn right now.)

This suppression is harmful in two ways:

1. It is anti-rational and anti-truth seeking.

You have criticism about either whether/when/how to learn this thing, or criticism of the thing itself, or both — which are valuable data — that you now must ignore.

Suppressing criticism of whether to learn something is siding with 'learn this thing' dogmatically, instead of rationally resolving the conflict.

Suppressing criticism of the thing itself a) assumes it is true instead of *finding out whether it is true*, and b) *sabotages your attempts to learn it*. Criticism is how we learn. The more we suppress our own criticism, the more fragile our understanding will be.

Suppose you were wondering whether light is a particle or a wave, and your high school taught that it has "particle-wave duality". You might wonder what that means, how it could be both, and what physical reality would have to be like if it were only one or the other. Then you find that some physicists say it's only a particle, but for that to make sense in physical reality it requires the existence of multiple universes — which on the face of it seems nuts to you, but also intriguing. You start wondering what reality would be like if there really were these infinite copies of you in other universes. Then you start thinking about whether that would solve the time-travel paradox of killing your grandfather, and...

— Except you don't. Because as soon as you ask your teacher about the idea that light is only a particle, he scoffs and tells you that you don't need to worry about that until university, and in this class all you need to know is that it has duality. You would have asked more, but your teacher seemed irritated by the question, and anyway it wasn't relevant to the assignment at hand. So you continue with the assignment, trying to build your model of reality on this thing that you have doubts about, wondering what kind of world would permit 'duality' (or not wondering, because being forced to drop your last question was unpleasant and at this point you just want to get the work over with). You never get to time travel paradoxes.

After the exam is over, you spend the next couple years forgetting most of your high school physics — which is fine, because it never came up in life again after that. Until you see a documentary, which piqued your curiosity and prompted you to re-learn some physics. This time, your model of reality is created under the direction of your own interests, and you get to use the full force of your critical faculties to understand the world.

2. Unwanted criticism is coercive and creates hangups: lasting psychological blocks/irrationalities.

It is difficult to keep two competing versions of reality in your mind (theirs and yours), especially when you disagree with or have doubts about one of them. And when there is force or pressure to adopt a version of reality that is not your own, it's often painful, or threatens pain/punishments if you don't succumb. Boredom — especially boredom at having your attention forced onto something demanding — is famously painful, and this is why.

To suppress your own critical faculties and follow someone else's agenda, you have to find ways of directing yourself to not think about something you would naturally think about.

You have to cordon off areas of your mind, kinds of criticism, and make them off-limits. When your own critical attempts are thwarted and you find no rational way out, you must resort to coping mechanisms: you create no-go areas where you can't think, everything is foggy, your mind goes blank — and if you try, it hurts. You use your own creativity against itself.

This is how hang-ups are born. These are active irrationalities.

If the harm done from unwanted criticism were only in the moment — wasted time, wasted effort, confusions that require detangling — that might be a shame, but no big deal. Passive errors in our thinking can be corrected later. We're constantly making mistakes and correcting them; that's life.

But when faced by coercion — when we must self-suppress to get through an unpleasant situation — we wind up with errors that *actively resist* when we try to correct them.

Is this harm avoidable?

Someone may object, "But you can just not listen to criticism. Criticism can only hurt you if you let it."

To the extent you can do this, sure. But that's not always trivial.

1. If you're in a situation where there's some external (or internal) pressure, you need to creatively come up with a way of overcoming that pressure.

If it's pressure from school, you may need to solve a whole host of problems around interacting with your parents. If it's social pressure, you might need to find a way to be accepted to the group without addressing their criticism. If it's pressure from

yourself, you might need to deal with big unsolved personal problems like how to feel 'good enough'. Simply ignoring criticism doesn't always work.

2. Once you hear a criticism, that destroys your ignorance.

You can't always pretend you haven't heard it, because new knowledge affects the path of your thinking. If someone tells you the twist of a movie, that spoiler means you can no longer try to make guesses about the twist yourself. *It is exactly the same with all thinking.* There are times when you want hints to help you solve a puzzle, and there are times when hints would ruin your problem-solving. This doesn't only happen when there's a fixed answer to discover, like in a movie or a puzzle; the same process happens in creative thinking, where the 'end' hasn't been invented yet.

3. Conjectures that are created in the expectation of different kinds of criticism are different.

When you're in private speaking to a close friend, you say different kinds of thing from when you're speaking to a public audience. *Which* friend makes a difference, too — it's not just 'few vs many'. And it's not just that you say different things. You *think* different things. You make different conjectures around different people and sets of people. (This is part of the value of talking to people in the first place: you create different ideas than you would alone.)

These three dangers are why it's important to not share a draft before it's ready. Criticism at the wrong time can play the same role as movie or puzzle spoilers. It can kill off ideas early, or change the direction of your thought, even when you would have preferred to keep thinking about it. It does this no matter how well-meaning the person giving criticism was, no matter how non-pressured you feel, no matter how much responsibility you take for how to reply to the criticism — *criticism changes the course of thinking.* You cannot recover that ignorance.

Privacy is criticism scheduling

"When a person is learning successfully, the trail of wilful ignorance, as it were, that they leave behind them – the sequence of rejected opportunities to acquire information – characterises their creative process just as faithfully as the sequence of problems that they faced and solved."

— David Deutsch, *In Praise of Ignorance – Taking Children Seriously* journal issue 31

Choosing when to get movie spoilers, puzzle hints and creative feedback are all types of **criticism scheduling**: directing the critical process, so that you're using the best ideas you have on which criticism would be good to introduce next.

When most people think "violating privacy is bad", they think of things like:

1. People not reading their private diary, or their browser history, or letters to friends. Not anything those people might do or say, but the act of reading those things itself being violating.
2. Entering a bedroom unannounced, or letting oneself in to a house of an acquaintance. Again, this is considered separate from what they might do or say while there.

The sense of violation is in fact intimately connected to what they do or say to you — even if they're not directly commenting on what they've seen. It affects your expectations of them and what kinds of things they could conceivably do or say to you.

It also affects your internal model of their mind, which affects the criticism environment of your ideas.

So, there are two broad reasons violating privacy is bad:

First, it creates uncertainty about what kind of information (including criticism) is going to be brought up to you that you then have to deal with. In other words, it can directly interfere with criticism scheduling.

Second, it changes your model of them. Your mental models can shape the criticism you think about in the first place. If you can't do something without thinking how they would think about it, that forces you to confront or deal with certain criticisms.

Certain mental models are easier to engage with than others. Modelling someone hostile takes criticism in a very different direction compared to modelling someone cooperative. But even modelling someone friendly can change the way you think in ways that you may not necessarily want, such as if they know things about you that you're not ready for them to know.

If you're vulnerable to your parents reading your diary and discovering your secret thoughts, then even if they say nothing — even if you think they just *may* have seen it — your model of them becomes more complicated. You have to interact with them on the basis that they *may* or *may not* know those things.

You'd have to think something like, "Either they did or did not read that page, and as a result may think such-and-such or not, and that might mean they do such-and-such or not." — It causes a whole load of things that complicate your relationship.

Whereas, the way should be is that you *know* — you are certain — that your parents do not look at those things. They don't come in to your room because there is a lock on the door, or because there's a strong family institution that they do not do this, and if they happen to be in your room and see your open diary then they *still* won't read it.

And even if they were amenable to rational argument, they may still have misunderstandings based on the private information they discovered about you, because it's never possible to share the information 100% accurately. Some things take a lot of background knowledge to understand.

A piece of private information might make sense in context, but another person may not have that context — especially if it's a personal thing that comes out of the nitty-gritty details of your own life.

Sharing *non-personal* ideas can be fraught too, for the same reason. You're faced with the problem: try to explain the context fully (impossible), dumb down the idea to make it comprehensible to the other person, or keep it private. Simplifying the idea for another person to digest risks you losing sight of your idea's depth and potential. If you understand an idea well already, this may not be such a danger, but having this early on in an idea's development can be devastating.

Controlling the privacy of your mind means you can control what gets exposed to criticism and when. This is one of the most important uses of privacy. Without this control, our thinking would be impaired or disabled.

We see invasion of privacy sometimes used intentionally to create obedience and disable criticism. In Soviet Russia, people would spy on each other — if you were suspected of having disloyal thoughts, you risk being thrown in the Gulag. In some religions, God is omniscient and can read your mind, including every sinful thought, which you must repent for. In coercive parenting, children are snooped on, and then are given punishments if there's something the parent doesn't like.

But this kind of disabling can happen when everyone is well-meaning. You might reveal something, and your well-meaning friend gives a comment intended to help but actually brings up a problem you weren't conscious of, which you now can't ignore, and it throws you off the fruitful thinking you were doing before it came up. This can happen even when no one is being irrational or sensitive. If those things are at play, it's worse.

When to have privacy

"One cannot be an individual—a person separate from family and society—without having secrets."

— Thomas Szasz

What feels fun and enlivening to share, and what makes you inwardly cringe?

If it makes you inwardly cringe, that's a sign sharing it would interfere with your thinking.

It's not possible to understand our mind fully and have an explicit idea of all the ways that sharing something private can cause problems. A lot of the things I've been talking about happen at the implicit or subconscious level.

Discomforts and fears are part of the intuitions you can use to form judgements about what to share and when.

Harry glanced away uncomfortably, then, with an effort, forced himself to look back at Draco. "Why are you telling me *that*? It seems sort of... private..."

Draco gave Harry a serious look. "One of my tutors once said that people form close friendships by knowing private things about each other, and the reason most people don't make close friends is because they're too embarrassed to share anything really important about themselves." Draco turned his palms out invitingly. "Your turn?"

— Eliezer Yudkowsky, *Harry Potter and the Methods of Rationality*, [Chapter 7](#)

We live in a culture that broadly undervalues privacy. Even aside from modern concerns like your employer finding embarrassing photos on your Instagram feed, or people posting their entire lives to Facebook, there is pressure to share private details in order to bond with others.

And Eliezer Yudkowsky's Draco does have a point. People often keep things private out of fear about what people would think of them. For some people, stopping worrying about what people think of them and revealing more would allow them to have closer and more meaningful relationships, and make their lives better.

But revealing things *before* you've resolved that fear — before you know whether that fear has useful information — can be dangerous. It could be that the fear has a valid concern, such as not wanting to [be defined](#). As long as you have that fear, that's an indication there's something you're not okay with. Bypassing that fear through force of will is raw material for a hangup.

So when should you reveal things?

As a guideline: When doing so is necessary for your joint problems (or projects or interests).

People-pleasers and social folk tend to offer up far more than makes sense for their shared problems.

Loner types or asocial folk may do the opposite. People who block out social information also try to limit their own information from going out. (Explained more depth in my [Beware Social Coping Strategies](#) post.) They may maintain so much privacy that it sabotages them. How to deal with too *much* privacy is beyond the scope of this article; but keep in mind that even there, the solution is not to force yourself, but instead to make updates on what might be nice for your life.

Because society so far has a lot of coercive control aspects (such as religion and coercive education), we live in a culture that tends to undervalue privacy. If people had more privacy, individuals would have much more opportunity to disobey authorities and avoid peer pressure. It would be a more individualistic and less tribalistic society.

Privacy means you can think freer. You have more [slack](#).

Whether it's details about your personal life, or what matters to you, or new ideas you're developing, or a budding project you're working on, privacy is important for managing the flow of information that affects your thinking.

You can try ideas that other people can't even *follow* without a lot of context, let alone agree with. You have the option to make no thought unthinkable. You can go anywhere.

When you're not beholden to make yourself agreeable or comprehensible to someone else, you can freely explore the unknown, the unthought.

This unknown is the only place that progress happens. It's where discovery lies.

[Hammertime Final Exam]

Accommodate Yourself; Kindness Is An Epistemic Virtue; Privileging the Future

[this is my entry for the [Hammertime Final Exam](#). I answered all three prompts but took much longer than five minutes writing each part.]

1. Accommodate Yourself

(related: [Society Is Fixed, Biology Is Mutable](#); [Design](#); [Radical Acceptance](#) as acknowledgement of reality. this is one of the first & most valuable lessons I have gained from the rationalist community, though I don't think I've seen it stated in quite these terms.)

People often want to make themselves better - stronger, more hardworking, more able. We compare ourselves to ideals and we find ourselves lacking; we strive to improve ourselves to better fit the roles we want to play in the world.

What if, instead of taking the world as given and striving to adapt ourselves to it, we took ourselves as given and looked for ways to adapt our world to us?

Examples:

- "I can't get around much because my feet hurt - and even taking public transit is bad because I have to stand while I wait for the train - so I have to fix my feet or else I'm doomed" → "My feet hurt, so I'll look for other ways to get around, like a bike or an electric scooter. When I take public transit, I'll carry a [light portable stool](#) to make waiting for trains painless."
- "I have to stop picking at my nails! Argh, why can't I do it!" → "Let me see if having a [thing to fidget with](#) removes this urge. Oh, it basically does! Good!"
- "Argh why can't I follow this conversation, everyone else seems to be able to do it" → "Hmm, it seems I can't follow conversations well in loud spaces, let me make plans in quieter ones instead."
- "I should be more hardworking and not procrastinate so much!" → "Hmm, it seems I don't have a lot of energy and my executive function is often not very good. Let me scale back my plans to match my energy levels, at least for now, and think about how I might make my work fit my brains and/or find an environment that's easier for me to work in and/or find a different way to support myself and/or look into ADHD medication."
- "I don't fit into these jeans, I'm so fat, I need to lose weight." → "Let me find some other jeans, these ones don't fit me."

To apply this mindset, think of a situation where a limitation of yours leads to failure or frustration. Now, flip your view of the situation around. You are as you are; it is the world that does not fit. Now, if the world were designed for people like you, how would it be different? What in your environment would change? Now, can you make that change yourself?

The title of this piece comes from the concept of disability accommodations, but I think the usefulness of this mindset extends beyond things commonly understood as disabilities. Pretty much everybody is subject to some kinds of expectations that they do not fit. We often take these as condemnations of aspects of ourselves, as judgments that our bodies or our brains in some way don't measure up.

But the world contains many kinds of brains and bodies. It's not useful for everyone whose brain or body is bad at some task to feel bad about themselves. Rather, we can take this diversity as a given and work on making the world fit all of us.

This doesn't mean there is no space for self-improvement. But, firstly, accommodation is often much easier, sometimes to the point that accommodation is possible where fixing oneself is not. And second, even if self-modification is possible and desirable, accommodation is often a useful first step, providing some slack with which to take on the difficult project of change.

2. Kindness Is An Epistemic Virtue

In order to learn from each other, we must first feel safe. When a battle is pitched, when opposition feels like a threat, when losing an argument threatens one's felt or actual safety - then humans fall into battle formation, [arguments become soldiers](#), and changing one's mind becomes virtually impossible.

Seeing such a battle, some with valuable knowledge may reasonably decide to stay out of the fight for their own psychological health; thus the combatants are left without their valuable knowledge.

Others will join the fray, motivated by a sense of urgency; they will learn something, they will share their knowledge, but they will be wounded and triggered and enter a psychological defensive crouch from which no change of opinion is possible, and they will lose their self-control and wound and trigger others in return.

(One time I had a particularly bad fight on Facebook and was extremely stressed and anxious for several days - and less open and charitable in arguments for months afterwards.)

So if you want participants in your discourse to be selected for valuable knowledge and not for battle-hardiness and thick skins, and if you want people to actually have the ability to change their minds - be kind.

(I want to acknowledge that this consideration can [compete](#) with others that are also important. Some people with valuable knowledge have come by that knowledge in painful ways and cannot share it kindly. Sometimes repressing anger for the sake of harmony is harmful in the long run. These are worth taking seriously, and sometimes they will be more important than kindness. But - I claim - kindness should very often win, and should always be considered.)

3. Privileging The Future

Human desires and intuitions are often highly biased in favor of the present. If we want to make decisions that are wise in the long term, we have to counteract this

bias; the ability to delay gratification is a valuable skill that will serve you well if you have it.

However, a common result of this attempt at debiasing is a sense that delayed gratification is *virtuous in itself*, that one should *always* take the path that gives you enjoyment in the future rather than the present - or, failing that, that one should feel guilty when one doesn't do that.

For example: I used to have a problem with staying up far too late far too often, which led to chronic sleep deprivation. This was clearly in part a result of prioritizing present over future on any given evening; so, I reasonably felt that it was better when I succeeded in prioritizing tomorrow's well-being enough to go to sleep on time. But sometimes, something legitimately interesting and unique was happening late at night, something that was worth staying up... and I would stay up, but I would feel guilty, that I was doing a thing that would result in my suffering later; and the next day, I would regret my choice, even as I valued the experience I had had at night. In my mind, sacrificing future well-being for present enjoyment was *always wrong*, even if the tradeoff was actually worth it!

For another example: adults sometimes express regret that they quit an instrument they had played in childhood, even if they had never particularly enjoyed practicing it; I think this is often a result of a pro-future bias.

(Related: somebody on Tumblr suggested that somebody should run the opposite of the [canonical marshmallow test](#), such that a child can either eat two marshmallows now or one marshmallow in ten minutes, and then see whether in ten minutes the child regrets eating the two marshmallows because they can no longer have one now. Their hypothesis was that the child would indeed regret it, even though the decision was clearly correct - which would show that regret is not reliable information about the quality of one's past decisions.)

If you're facing a tradeoff between present and future - weigh the options and choose a tradeoff! Delay gratification if that's the right thing to do. (And it is more often the right thing to do than your system 1 probably thinks.) But if it's not - go get the good thing in the present, with no self-judgment and no regrets.

Beauty bias: "Lost in Math" by Sabine Hossenfelder

This is a linkpost for <https://www.amazon.com/Lost-Math-Beauty-Physics-Astray/dp/0465094252>

This is technically just a link post for those who didn't see the book yet. The main idea, as I understand it, is that while some physical theories are beautiful, other are complex and "ugly". Beauty should not be taken as evidence for truth.

It looks like sometimes in AI safety research aesthetics may be taken as an evidence (I will not provide examples, as in each case it may be just my interpretation), and thus possibility of such beauty bias should be taken into account.

Reflective AIXI and Anthropic

Crossposted from the [AI Alignment Forum](#). May contain more technical jargon than usual.

It's possible to define a version of [Solomonoff Induction with Reflective Oracles](#), that allows an AIXI-like agent to consider hypotheses that include itself or other equally powerful agents, going partway towards addressing naturalized induction issues.

So then a natural question is "what does this partial answer seem to point to for anthropics?"

To figure this out, we'll be going over a few of the thought experiments in [Bostrom's book about anthropic reasoning](#), and seeing what Reflective-Oracle AIXI has to say about them.

The following conclusions are very dependent on how many extra bits it takes to encode "same environment, but I'm that other agent over there", so I'll be making a lot of assumptions that I can't prove, such as the most efficient way of encoding an environment being to specify an environment, and then specifying a place in there that the agent interfaces with. This seems unavoidable so far, so I'll at least make an effort to list out all the implicit assumptions that go into setting up the problems.

As a quick refresher, SSA (self-selection assumption) and SIA (self-indication assumption) work as follows: SSA takes the probability of a world as given and evenly distributes probability mass to being everything in "your reference class" in that particular world. SIA reweights the probability of a world by the number of instances of "things in your reference class" that it contains. In short, SIA has a strong bias in favor of possible worlds/hypotheses/turing machines with many instances of you, while SSA doesn't care about how many instances of you are present in a possible world.

Thought Experiment 1: Incubator

Stage (a): In an otherwise empty world, a machine called "the incubator" kicks into action. It starts by tossing a fair coin. If the coin falls tails then it creates one room and a man with a black beard inside it. If the coin falls heads then it creates two rooms, one with a blackbearded man and one with a white-bearded man. As the rooms are completely dark, nobody knows his beard color. Everybody who's been created is informed about all of the above. You find yourself in one of the rooms. Question: What should be your credence that the coin fell tails?

Stage (b): A little later, the lights are switched on, and you discover that you have a black beard. Question: What should your credence in Tails be now?

This will be modeled as a machine that represents the environment, that has a bit that is used to determine how the coinflip comes up. Also, in the second case, because there are two possible places where the agent can be hooked up to the environment, another bit is required to specify where the agent is "attached" to the environment.

These three cases have minimum description lengths of $|E| + 1$, $|E| + 2$, and $|E| + 2$ bits respectively (where $|E|$ is the description length of the environment), so by the

universal semimeasure, they have (relative) probability mass of 50%, 25% and 25% respectively.

So, assuming the problem setup actually works this way, the answers are 50% and 67%, respectively. This seems to point towards Reflective-Oracle Solomonoff Induction (RO-SI) doing something like SSA. The intuitive reason why, is because a hypothesis with a bunch of copies of you requires a bunch of extra bits to specify *which* copy of you the input data stream is coming from, and this cancels out with the increased number of hypotheses where you are in the well-populated world. There may be 2^{50} copies of you in a "world", but because it requires 50 bits to specify "I'm that copy right there", each specific hypothesis/Turing machine of the form "I'm in that world and am also that particular copy" requires 50 extra bits to specify where in the environment the data is being read out from, and receives a probability penalty of 2^{-50} , which, when multiplied by the large number of hypotheses of that form, recovers normality.

There are two ways where things get more interesting. One is that, for environments with many observers in your reference class (RO-SI uses as its reference-class all spots in the environment that receive the exact same observation string as is present in its memory), you'll assign much higher probability to being one of the (fairly few) observers for which specifying their spot in the environment is low K-complexity. It *definitely* isn't a uniform distribution over observers in the possible world, it favors observers that are lower-complexity to specify where in the environment they are. A similar effect occurs in logical induction, where there tend to be peaks of trading activity of simple traders, on low-K-complexity days. Sam's term for this was "Graham's crackpot", that there could be a simple trader with a lot of initial mass that just bides its time until some distant low-K-complexity day and screws up the probabilities then (it can't do so infinitely often, though)

The other point of interest is what this does on the standard counterexamples to SSA.

To begin with, the Doomsday argument is valid for SSA. This doesn't seem like much of a limitation in practice, because RO-SI uses a *very restrictive* reference class that in most practical cases includes just the agent itself, and also, because RO-SI is about as powerful as possible when it comes to updating on data, the starting prior would *very quickly* be washed out by a maximally-detailed inside view on the probability of extinction using all data that has been acquired so far.

Thought Experiment 2: Adam and Eve

Eve and Adam, the first two humans, knew that if they gratified their flesh, Eve might bear a child, and if she did, they would be expelled from Eden and would go on to spawn billions of progeny that would cover the Earth with misery. One day a serpent approached the couple and spoke thus: "Pssst! If you embrace each other, then either Eve will have a child or she won't. If she has a child then you will have been among the first two out of billions of people. Your conditional probability of having such early positions in the human species given this hypothesis is extremely small. If, on the other hand, Eve doesn't become pregnant then the conditional probability, given this, of you being among the first two humans is equal to one. By Bayes' theorem, the risk that she will have a child is less than one in a billion. Go forth, indulge, and worry not about the consequences!"

Here's where the situation gets nifty.

Assume the environment is as follows: There's the coding of the Turing machine that represents the environment ($|E|$ bits), the 1 bit that represents "fertile or not", and the bitstring/extradata that specifies where Eve is in the environment. ($|L|$ bits, L for "location"). Eve has been wandering around the Garden of Eden for a bit, and since she's a hyper-powerful inductor, she's accumulated enough information to rule out all the other hypotheses that say she's actually not in the Garden of Eden. So it's down to two hypotheses that are both encoded by $|E| + |L| + 1$ bits, which get equal probability. If we assume a utility function that's like "+1 reward for sex, -10 reward for creating billions of suffering beings" (if it was -10^{10} for an Eve that wasn't scope-insensitive, the serpent's reasoning would fail), the expected utility of sex is $0.5 \cdot 1 + 0.5 \cdot -9 = -4$, and Eve ignores the serpent.

The specific place that the serpent's reasoning breaks down is assuming that the probability of being Eve/difficulty of specifying Eve's place in the universe goes down/up when a decision is made that results in the world having a lot more beings in it. It doesn't work that way.

However, it gets more interesting if you assume everyone in the resulting created world has sense data such that even a hyper-powerful inductor doesn't know whether or not they are Eve before the fateful decision.

Also, assume that it takes $|L'|$ bits to specify any particular person's location if they're not Eve. This is a sort of "equally distributed probability" assumption on the future people, that doesn't restrict things that much. Maybe it's much easier to point to Eve than some other person, maybe it's the other way around.

Also assume that everyone's utility functions are like "+1 for sex, -10 for finding out shortly after sex that you are one of the suffering future beings, or that you created billions of such."

To begin with the analysis, break the hypothesis space into:

two worlds of $|E| + |L| + 1$ bits where Eve is fertile/infertile, and you are Eve.

and 10^{10} worlds of (it depends) bits where Eve was fertile, sex was had, and you are not Eve. The reason why it's tricky to say what the description-length of being one of the future agents is, is because it takes fewer bits to encode a world where an agent does a thing in accordance with the laws of math, than it takes to encode a world where an agent does a different thing that they wouldn't have normally done. In this particular case, it would take $|S|$ bits (S for surgery) to specify "at this particular spot, ignore what Eve would have done and instead substitute in the action "have sex", and then run things normally".

So, if Eve definitely has sex, it takes $|E| + |L'| + 1$ bits to specify one of the future agents. If Eve definitely doesn't have sex, it takes $|E| + |L'| + |S| + 1$ bits to specify one of the future agents.

Taking these two cases, we can rescale things to get a mass of 1, 1, and either $10^{10} \cdot 2^{|L|-|L'|}$ or $10^{10} \cdot 2^{|L|-|L'|-|S|}$, on the three classes of worlds, respectively. Expected utility calculations will work out the same way if we use these numbers instead of probabilities that add up to 1, because it's just a scaling on expected utility and the scaling can be moved over the utility function, which is invariant under scale-and-shift. So then, in the first case, expected utility of sex and not-sex becomes:

$$1 \cdot 1 + 1 \cdot -10 + (10^{10} \cdot 2^{|L|-|L'|}) \cdot -9 = -9 \cdot (10^{10} \cdot 2^{|L|-|L'|} + 1)$$

$$(10^{10} \cdot 2^{|L|-|L'|}) \cdot -10$$

So sex will be had if $10^{10} \cdot 2^{|L|-|L'|} > 9$. The crossover point occurs approximately at a 30 bit penalty to specify a non-Eve person (and 2^{-30} is approximately 1/billion.) So, if Eve has sex, and assigns less than about a 1/10 chance to being Eve, it's a consistent state of affairs. The reasoning is "I'm probably not Eve, and so I'm probably already going to suffer (since I know in advance what my decision is in this case), might as well pick up that +1 utility"

Redoing this analysis for the case where Eve doesn't have sex, we get that sex will be had if $10^{10} \cdot 2^{|L|-|L'|-|S|} > 9$, and in this case, the crossover point occurs approximately at a 30 bit penalty to specify *both* the non-Eve person and that particular decision intervention. (there can also be consistent solutions where the reflective oracle is perched right on the decision threshold, and randomizes accordingly, but I'll ignore those for the time being, they don't change much)

Considering the specific case where the ratios of the probability masses for "I'm Eve" and "I'm not Eve" is less than 1 : 9 (in the sex case) and 1 : 9 $\cdot 2^{-|S|}$ (in the non-sex case), we get a case where the decision made depends on the choice of reflective oracle! If the reflective oracle picks sex, sex is the best decision (by the reasoning "I'm probably not Eve, might as well pick up the +1 utility"). If the reflective oracle picks not-sex, not-sex is the best decision (by the reasoning "I'm likely enough to be Eve (because the non-Eve people live in a lower-probability universe where an intervention on Eve's action happened), that I won't chance it with the coinflip on fertility")

So, RO-AIXI doesn't exactly *fail* (as SSA is alleged to) in this case, because there's a flaw in the Serpent's reasoning where the difficulty of specifying where you are in the universe *doesn't change* when you make a decision that creates a bunch of other agents, and you don't think you could be those other agents you're creating.

But if there's a case where the other agents are subjectively indistinguishable from yourself, and it's bad for you to create them, but good for them to push the "create"

button, there are multiple fixed-points of reasoning that are of the form "I probably press the button, I'm probably a clone, best to press the button" and "I probably don't press the button, I'm probably not a clone, best to not press the button".

Another interesting angle on this is that the choice of action has a side-effect of altering the complexity of specifying various universes in the first place, and the decision rule of RO-AIXI doesn't take this side-effect into account, it only cares about causal consequences of taking a particular action.

The arguments of **Lazy Adam, Eve's Card Trick**, and **UN++** in Bostrom's book fail to apply to RO-AIXI by a similar line of reasoning.

Sleeping Beauty, SSA, and CDT:

There's a possible future issue where, according to [this paper](#), it's possible to money-pump the combination of SSA and CDT (which RO-AIXI uses), in the Sleeping Beauty experiment. Looking further at this is hindered by the fact that RO-AIXI implicitly presumes that the agent has access to the entire string of past observations that it made, so it doesn't interact cleanly with any sort of problem that involves amnesia or memory-tampering. I haven't yet figured out a way around this, so I'm putting up a 500-dollar bounty on an analysis that manages to cram the framework of RO-AIXI into problems that involve amnesia or memory tampering (as a preliminary step to figure out whether the combination of SSA-like behavior and CDT gets RO-AIXI into trouble by the argument in the aforementioned paper).

Takeaways:

RO-AIXI seems to act according to SSA probabilities, although there are several interesting features of it. The first is that it assigns much more probability to embeddings of the agent in the environment that are low K-complexity, it definitely doesn't assign equal probability to all of them. The second interesting feature is that the reference class that it uses is "spots in the environment that can be interpreted as receiving my exact string of inputs", the most restrictive one possible. This opens the door to weird embeddings like "The etchings on that rock, when put through this complicated function, map onto my own sense data", but those sorts of things are rather complex to specify, so they have fairly low probability mass. The third interesting feature is that the probability of being a specific agent in the world doesn't change when you make a decision that produces a bunch of extra agents, which defuses the usual objections to SSA. The final interesting feature is that making a particular decision can affect the complexity of specifying various environments, and the standard decision procedure doesn't take this effect into account, permitting multiple fixed-points of behavior.

Also I don't know how this interacts with dutch-books on Sleeping Beauty because it's hard to say what RO-AIXI does in cases with amnesia or memory-tampering, and I'd really like to know and am willing to pay 500 dollars for an answer to that.

Cooperative Oracles

Crossposted from the [AI Alignment Forum](#). May contain more technical jargon than usual.

The long-awaited continuation of [this sequence](#).

Credit to Sam, Tsvi, Nisan, Scott, me, and Giles (who helped make the figures)

It may or may not become a paper, but if it doesn't, it's best for it to not be locked in a folder of drafts. This will also recap material from previous posts and papers.

Refresher on Reflective Oracles

Let M denote an arbitrary probabilistic Turing machine with oracle access, and M^O refer to the same machine equipped with the oracle O . It is important to remember that M refers to a randomized algorithm, while M^O refers to the algorithm equipped with a particular oracle, so it doesn't make sense to ask what the output of M is, only what the output of M^O is.

Reflective oracles as originally defined [here](#) are oracles which accurately answer queries about whether a probabilistic Turing machine with access to the *same* oracle, M^O , output a 1 with probability greater than or less than p , by outputting 1 for "yes", 0 for "no", and randomizing between the two if the probability of M^O outputting 1 is exactly p . This is more expressive than it may seem at first, because it is possible to query the oracle about a probabilistic oracle machine that runs M^O and checks whether the output has some arbitrary computable property.

It is also possible to use reflective oracles to make the probability distribution over the output of one algorithm vary continuously with respect to the probability distribution over the output of another algorithm, via repeated oracle calls and randomization, as demonstrated [here](#).

Reflective oracles can be thought of as akin to halting oracles, except that they accurately answer questions about the class of probabilistic Turing machines with access to the *same* oracle, which halting oracles cannot do. If the diagonalization argument against the existence of a halting oracle is applied to a reflective oracle, where an algorithm invokes the oracle on itself, and outputs the bit which it is the least likely to output, a reflective oracle will randomize between outputting 0 and 1, but for any $p < 0.5$, if the oracle is asked whether the algorithm outputs 1 with greater than probability p , it will respond "yes".

This ability to accurately answer questions about themselves makes reflective oracles well-suited for studying the unbounded computing power limit of game theory, and situations where an agent reasons about a surrounding environment which contains itself. In particular, if a set of agents in a game all use a reflective oracle to perfectly predict the behavior of all other agents in the game, and then select the optimal action, a Nash equilibrium will be played in the game. This will be called the *shortsighted best response* algorithm.

To begin, there is a canonical way to translate any game theory problem into the setting of probabilistic oracle machines. As a concrete example, take a game of Prisoner's Dilemma, with an exchange of source code between agents. Let player_A and player_B refer to the source code of the two players in a Prisoner's Dilemma game. Then the game may be thought of as the two probabilistic oracle machines $\text{player}_A(\text{"player}_B\text{"})$ and $\text{player}_B(\text{"player}_A\text{"})$ calling the reflective oracle on each other, and outputting 0 for "defect" and 1 for "cooperate". In general, any game can be expressed by expressing all players as a pair of two algorithms, one of which outputs a strategy, and one of which takes no input and computes the utility via oracle calls. By fixing a suitable input for the algorithm which generates the strategy, which represents the information a player has, all strategy algorithms become inputless probabilistic oracle machines which make oracle calls to other players.

More formally, let a *player P* be a pair of a probabilistic oracle machine M , and a probabilistic oracle machine U which halts with probability 1 regardless of the choice of oracle, and outputs 0 or 1. M is interpreted as the algorithm which produces the strategy of the player, and U is interpreted as the utility function. Fixing an oracle O , the probability of U^O outputting 1 is the utility of the player P .

There are many different variants of a reflective oracle, which are essentially equivalent. There is the standard single-bit reflective oracle, presented [here](#) but multi-bit reflective oracles may also be [defined](#). Here, I used the semimeasure formulation of reflective oracles, which accurately answers any question about prefixes of a Turing machine output, and assume that the outputs of the Turing machine form a semimeasure over the set of finite bitstrings.

Instead of the queries consisting of a $\{M, p\}$ pair, they consist of a $\{P, p, x\}$ triple, asking about the probability of the M^O component of P^O halting with a bitstring that has x as a prefix. Let $P(P^O = xy)$ be the probability that M^O halts with a bitstring that has x as a prefix, and let $P(P^O = \neg xy)$ be the probability that M^O halts with a bitstring that doesn't have x as a prefix. Let $q_O(P, p, x)$ be the probability of the oracle O

returning a 1 in response to the query (P, p, x) . Then one of the defining properties of a reflective oracle is that, for all bitstrings x , rational probabilities p , and players P :

$$P(P^O = xy) > p \rightarrow q_O(P, p, x) = 1$$

$$P(P^O = \neg xy) > (1 - p) \rightarrow q_O(P, p, x) = 0$$

If a probabilistic oracle machine halts with a certain probability distribution, it will be accurately captured by this variant of reflective oracle, but nonhalting probability mass may be assigned to an arbitrary distribution over bitstrings, as long as it makes a semimeasure.

Cooperative Oracles:

In the particular case where the two agents in a Prisoner's Dilemma both use shortsighted best response, the set of fixed points is the set of Nash equilibria, so they will mutually defect.

However, different algorithms may be used by the players. As long as each player's set of responses to everyone else's strategies is nonempty, convex, and has closed graph, Kakutani's fixed point theorem guarantees the existence of an "equilibrium set" where all players respond to knowledge of all other player strategies with the exact strategy that all other players predicted. This is the same as the fair set from [here](#).

If we consider two agents in the Prisoner's Dilemma implementing the algorithm "cooperate with the same probability as my opponent cooperates", there is a continuum of reflective oracles for all probabilities of cooperation in $Q \cap [0, 1]$.

Intuitively, mutual cooperation should be selected. This splits the problem of formalizing Functional Decision Theory in the unbounded case into the problem of finding an algorithm for the agent to use, and the problem of finding an oracle which causes a Pareto optimal element of the equilibrium set to be played.

A naive implementation, where we select a reflective oracle that is Pareto optimal for all players, fails. This is because all utility functions have an opposing one, which makes all outcomes Pareto optimal, because there is no outcome that makes one player better off without making another player worse off. In the Prisoner's Dilemma, if there is a third player, a jailer who wants to maximize jail time, but cannot otherwise influence the outcome of the game, all outcomes are Pareto optimal, because there is no choice which makes a prisoner better off that doesn't also leave the jailer worse off, and vice-versa. Therefore, the optimality notion must somehow restrict to the set of players who are actually participating in a game, instead of the set of all players.

To clarify what it means to talk about the set of players who participate in a game, we can look at the structure of which players make oracle calls to other players. If there is a game with finitely many players, the utility functions of all players should make oracle calls to the output of all players, to compute a utility.

Let the *dependency set* of player P_i be the set of all players that either the strategy algorithm or the utility function of P_i may query the oracle about, as well as P_i itself.

Let the relation \prec be the minimal relation such that if P_j is in P_i 's dependency set,

$P_i \prec P_j$. The transitive closure of this relation, \leq , induces a poset over equivalence classes of players.

Minimal elements in this poset correspond to sets of players whose strategies and utility functions only depend on each other. Players in a self-contained game have this property. A non-minimal equivalence class in this poset corresponds to a game where the utility functions of some participants are affected by what happens in a different game they cannot influence. Due to these properties, \leq is a natural choice to define the boundaries of a game, and to formalize the desideratum that our notion of optimality should not be Pareto optimal for spectators.

Now, the notion of stratified Pareto optimality may be defined. Let $U_{O,i}$ be the utility of the player P_i , upon being equipped with the oracle O . A reflective oracle RO' is a *stratified Pareto improvement* on another reflective oracle RO , for a player P_i , iff:

$$(1): U_{RO,i} < U_{RO',i}$$

$$(2): P_j \leq P_i \rightarrow U_{RO,j} \leq U_{RO',j}$$

$$(3): P_i \not\leq P_j \rightarrow U_{RO,j} = U_{RO',j}$$

The first condition is that a stratified Pareto improvement for P_i must make P_i better off.

The second condition is that all players that P_i depends on, even indirectly, must be unaffected or better off. In the case of a finite and minimal equivalence class of players under \leq , this says that all other players in the game must be unaffected or better off, so all stratified Pareto improvements are Pareto improvements for this game. In the case of the jailer in the Prisoner's Dilemma, it says that both prisoners must receive the same or greater utility, because the jailer's utility function makes an oracle query about their output.

The third condition is that a change in oracles should only affect players whose strategy algorithm or utility function depends on the outcome of a game, and leave all other players unaffected. This condition is used to forbid cases where selecting a better oracle in a game of Stag Hunt makes the players in a totally unrelated game of Chicken worse off, because the same oracle is used by all players.

A point is *stratified Pareto optimal* (SPO) for a set of players P_G if there are no stratified Pareto improvements for any players in P_G . In the particular case of a Prisoner's Dilemma with a jailer, a stratified Pareto optimum will involve both prisoners selecting a Pareto optimal outcome in their equilibrium set, and ignoring the jailer.

However, even stratified Pareto optima aren't enough in all cases, by the case of the chain of infinitely many players detailed [here](#)

Let an *almost stratified Pareto optimum* for a set of players P_G be a reflective oracle that is in the closure of the set of stratified Pareto optimal oracles for all players in P_G . In the case of the infinitely descending chain of players, an almost stratified Pareto optimal reflective oracle would report that all players had probability 0 of outputting 1, because that is the limit of a sequence of SPO oracles for a player, which move the "first player which outputs 1" arbitrarily far back in the sequence.

Theorem 1 (weak form): *There exists a reflective oracle that is almost stratified Pareto optimal for all players.*

A proof sketch for Theorem 1 is as follows:

Given an arbitrary finite set of players, P_G , there is a finite poset of equivalence classes of players in P_G under \leq . Given some initial set of reflective oracles, it is possible to construct a subset that is Pareto optimal for all players in a minimal equivalence class E, and a Pareto optimal subset of that for a minimal equivalence class F in the poset with E removed from it, and so on. This process repeats until there is a set of SPO reflective oracles for all players in P_G . Now that existence has been shown, take the closure of the set of SPO reflective oracles for P_G . The resulting set is a compact set of ASPO reflective oracles for P_G . This argument works for all finite P_G . Because there is a compact and nonempty set of ASPO reflective oracles for all finite sets of players, by compactness, there must be a compact and nonempty set of ASPO reflective oracles for all players.

This is pretty much the old proof of "there's an ASPO point for all players", the only difference is that this proof takes place in reflective-oracle space, instead of the space of the actions of all players. Let a *weak cooperative oracle* be a reflective oracle with this property.

Let a well-founded equivalence class E be an equivalence class of players under \leq composed of finitely many players, such that players within E only make oracle calls to other well-founded equivalence classes. Minimal equivalence classes under \leq , of finitely many players, are well-founded. It is possible to strengthen Theorem 1 to:

Theorem 1: *There exists a reflective oracle that is ASPO for all players, and SPO for all players in a well-founded equivalence class.*

This is proved by constructing a suitable set of reflective oracles that are SPO for all players in a well-founded equivalence class, and then using this as the initial set in the previous proof sketch. The resulting set of ASPO reflective oracles is a subset, and inherits this property.

An oracle of this type is called a *strong cooperative oracle*. Almost all games of interest in game theory involve finitely many players playing against each other. This is the condition for a minimal finite equivalence class under \leq , all of which are well-founded. Therefore, a strong cooperative oracle is an SPO oracle for any reasonable game that doesn't have non-well-founded chains of infinitely many players.

However, the proof of the existence of a cooperative oracle involves the infinite-dimensional space of reflective oracles, so it is not immediately linked to practical game theory. By imposing several natural conditions, it is possible to link the two settings. (Also, it is unclear whether the final condition of strategy continuity is necessary).

Theorem 2: *In a game of finitely many players, which all have finitely many possible outputs, have oracle access to each other and no other players, halt with probability 1 given oracle access to the other players, and have continuous strategies w.r.t everyone's probability distribution over outputs, all points in the equilibrium set have reflective oracles which result in that point being played in the game.*

Corollary 2: *If the conditions of Theorem 2 hold, then for any point on the Pareto frontier of the equilibrium set, there is a strong cooperative oracle which results in that point being played in the game.*

Therefore, the equilibrium selection problem is solved by the use of a special type of reflective oracle which results in Pareto optimal outcomes within the equilibrium set whenever it is used, aka a strong cooperative oracle.

However, there is still an open question about which algorithms result in something "FDT-like" when used with a cooperative oracle. To make partial progress towards a result, we will now look at the Ultimatum Game. In the Ultimatum Game, there is a pool of money, such as 100 dollars, and player A specifies a split of money between player A and player B, and then player B can either choose to accept the offer, or reject it, in which case both players get 0 dollars.

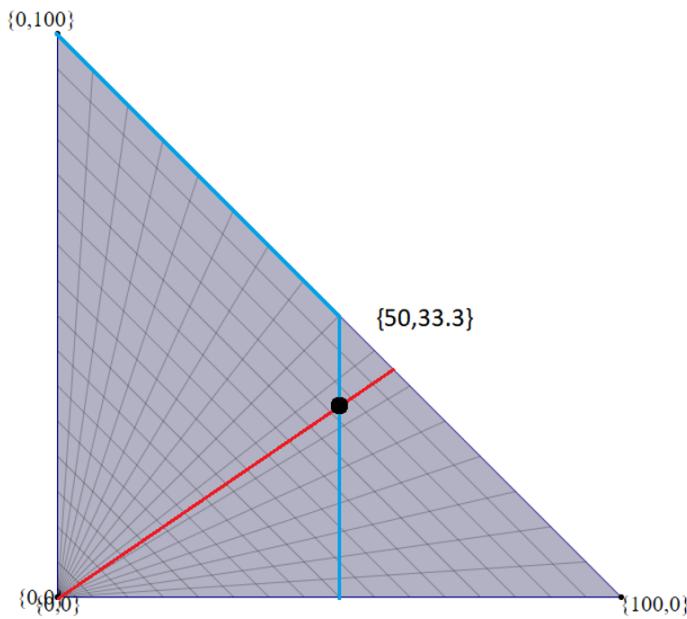
Imagine player B is implementing the strategy "I will reject the money unless I get 90 dollars or more". If player A would offer a 10/90 split in response to this strategy, and this is known to player B, this provides an incentive for player B to implement this extortionate strategy. Similarly, by viewing the game from the perspective of player B, accepting a split of money where 90 dollars goes to player A incentivises player A to offer a 90/10 split.

Additionally, if the strategies of player A and player B act to acquire 90 dollars, by proposing a 90/10 split, and rejecting all offers short of a 10/90 split, the game outcome will be A offering a 90/10 split, and B rejecting it, so both players walk away with no money. Extortionate strategies fail when they go against each other.

Is there a strategy that does not incentivize extortion? There is, in a certain sense, as long as there is some notion of a "fair outcome", such as the Nash bargaining solution. As an example, in the Ultimatum Game, if a 50/50 split is considered to be a fair outcome by player B, they could implement the strategy "if the money I am offered, x , is less than 50 dollars, reject the offer with a $\frac{y(50-x)}{100}$ probability, where $1 \leq y \leq 2$ ".

Against this strategy, the optimal offer is a 50/50 split of money, because in expectation, money is lost by offering any split that is more self-favoring than 50/50. Each marginal dollar that goes towards player A in the initial split is compensated by an increase in the probability of player B rejecting the offer, making the expected value of taking one more dollar from the split zero or negative. Therefore, extortion beyond the "fair outcome" is disincentivized.

Another advantage of this strategy is that it is robust against different agents having different notions of a "fair outcome". If player A thinks that a 60/40 split is a fair outcome, and player B thinks that a 50/50 split is a fair outcome, and $y = 1$, then player B will reject the offer with a $\frac{1}{6}$ probability, leading to an expected gain of 50 dollars for player A and $33\frac{1}{3}$ dollars for player B, instead of a guaranteed rejection of the offer.

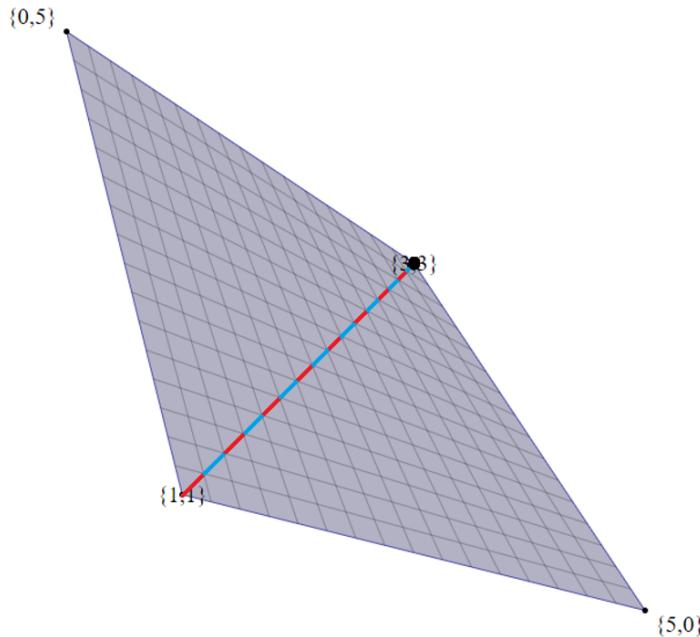


Red is the strategy of player A, blue is the strategy of player B, the display is over utility-space, player B is implementing a coercion-proof strategy.

However, this class of strategies does not have an elegant way of a-priori selecting the "fair point" that it will attempt to enforce. Intuitively, in the Ultimatum Game, there's a tradeoff where enforcing a strict demand, like a 30/70 split, means that there are more possible opponents which will reject the offer, while enforcing a weak demand, like a 70/30 split, means that there are more possible opponents which will exploit this. The Nash bargaining solution is a candidate "fair point", but it is unknown whether there is an algorithm that would naturally attain it under appropriate circumstances.

The insight of "be increasingly vengeful when the opposing players cause you to lose utility" will likely be a component of the true solution, and generalizes to other games such as the Prisoner's Dilemma and Chicken, but it is not a complete account, as will be shown in the next section. Humans also seem to act according to this strategy to some degree, instead of following the *Homo Economicus* strategy of accepting any offer greater than 0 dollars.

On the Prisoner's Dilemma, a natural choice for an extortion-proof strategy is "cooperate with the same probability that the opponent cooperates". This leads to mutual cooperation against a copy of itself, when equipped with a cooperative oracle, as seen in Figure 2.

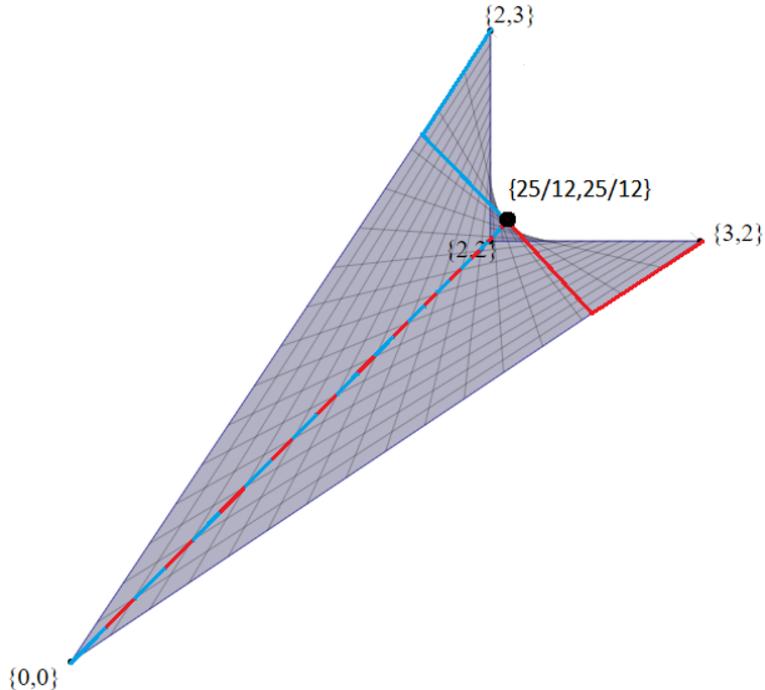


Prisoner's Dilemma, mutual cooperation is the Pareto-optimal point in the equilibrium set where the strategies intersect.

If the probability of player A cooperating is x , then if player B implements the strategy "cooperate with probability $\max(x - \epsilon, 0)$ for some $\epsilon > 0$ ", mutual defection will occur, as that is the only fixed point. This appears suboptimal, but permitting exploitation, even in a mild form, incentivizes opponents to modify their strategy to an exploiting strategy. If Player B uses shortsighted best response, then mutual defection will occur. This is because player B will always play defect because it is a dominant strategy, and player A will play defect as retaliation. If Player B always cooperates, then Player A will cooperate as well. This shows that implementing an extortion-proof strategy is not enough. A satisfactory strategy should lead to defection against a player which cooperates regardless of opponent.

A similar strategy applies to the game of Chicken. The strategy is "play Straight with the same probability that the opponent plays Straight, unless the probability of them playing Swerve is greater than $\frac{1}{2}$, in which case, play Straight." This leads to crashing into an opponent which goes straight, which incentivizes the opponent to avoid strategies such as tossing their steering wheel out the window. If two agents of this

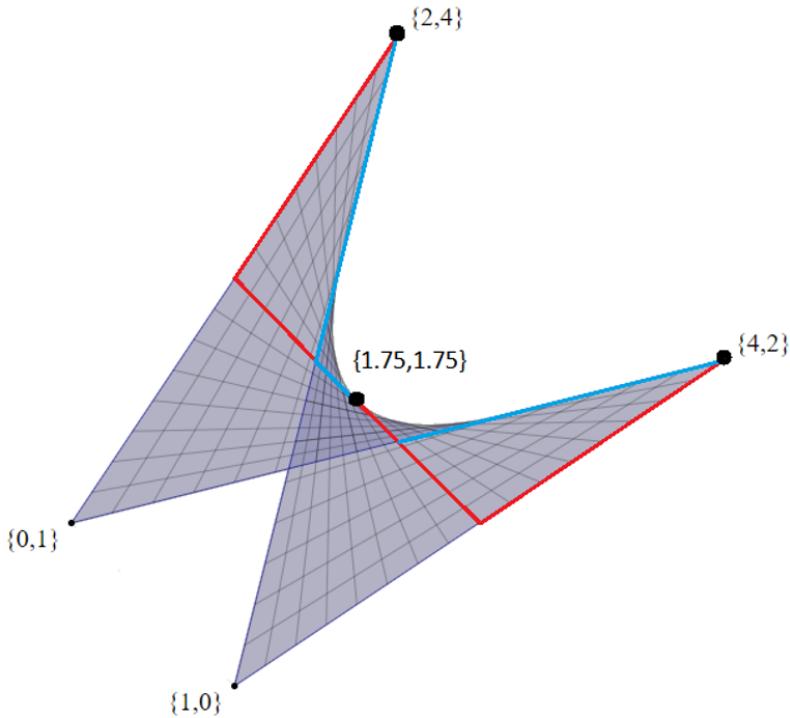
type play against each other, then by Figure 3, they both get $\frac{25}{12}$ expected utility with a $\frac{1}{6}$ probability of swerving. This is a Pareto improvement on the mixed-strategy Nash equilibrium where both players get 2 expected utility with a $\frac{2}{3}$ probability of swerving. If this player plays against a player B using shortsighted best-response, a cooperative oracle will result in the outcome where player A goes straight and player B swerves, for an expected utility of 3 and 2, respectively.



The game of chicken.

Consider the cooperation game where player A has a preference for going to a movie, and player B has a preference for going home to play a board game. They also have a preference for going to the same event, and going to a movie leads to nonzero utility even if the moviegoer is alone.

If both players implement the strategy "consult the oracle, and go to whichever event the other player has a higher probability of going to, with a rapid change of strategy around 50 percent", then the Pareto frontier of the equilibrium set consists of the two points where both players coordinate on an action. The choice of cooperative oracle determines which action is selected, and coordination occurs because both players are using the same oracle. This is similar to the choice of reflective oracle determining which Nash equilibrium is played in a game.



A cooperation game. Of the 3 points in the equilibrium set, the choice of cooperative oracle ensures that the outcome is coordination on an activity.

By looking at these examples, there are two further directions for designing a strategy which effectively makes use of a cooperative oracle. The first direction is developing a general criteria which allows a player to exploit some exploitable opponents. In the Prisoner's Dilemma, against an opponent which cooperates with 10 percent greater probability than you cooperate, the desired action is to cooperate with 90 percent probability. The second direction for improvement is finding a single decision rule, like argmax for Nash equilibria, that produces an exploitation-proof strategy incentivizing opponents to play at a Pareto optimal point, instead of customizing a strategy of that type for each game separately.

(Proofs were omitted because they're quite long. I can post them if there is interest.)

August gwern.net links

This is a linkpost for <https://www.gwern.net/newsletter/2018/08>

How to use a microphone ~~rational~~ during public speaking

Microphones are common technology and yet few people understand how they work. For our purposes here, there are unidirectional microphones and omnidirectional microphones.



A unidirectional microphone only records the sound from one direction. If it's pointed at the mouth of the speaker it won't record noise from the audience which makes the resulting recording better.

If it's however pointed at the ceiling while it's in front of the mouth of the speaker it will do a poor job at recording the speaker.

This leads to the first rule of microphone usage:

Hold the microphone pointing towards your mouth.

Microphones react to sound waves and sound waves are movement of air. If you exhale into a microphone, the microphone will record the exhale. Given Newtons laws, the exhale isn't omnidirectional either but the air that comes out of your mouth from the exhale has a clear direction. If you hold the microphone in front of your mouth, it will get hit by the air.

Our second rule of microphone usage is:

Keep the microphone to the right side of your mouth if you hold it in your right hand and correspondingly on the left side of your mouth if you hold it with your left hand.

Let combine the two rules into rule zero:

Hold the microphone to the side of your mouth in a way that points towards your mouth.

An Ontology of Systemic Failures: Dragons, Bullshit Mountain, and the Cloud of Doom

Core Claim

I assert that a lot of value can be achieved by categorizing systemic failures into four broad categories.

For the sake of pithiness, I will name them "bugs", "dragons", "bullshit mountain", and "the cloud of doom".

A **Bug** is the simplest kind of failure: you have a single cause, and a single symptom. Fixing a bug is pretty easy, compared to other modes - you just fix the cause, and the symptom goes away. Bugs aren't even really systemic failures, since they don't involve any cross-talk between multiple causes or effects, but they are included here for completeness.

A **Dragon** is like a Bug, except that instead of a single symptom, there are multiple seemingly independent symptoms. Because the symptoms seem large and diverse, a Dragon will often seem far more daunting than it actually is - although most Dragons are still pretty large and significant causes. Dealing with a Dragon simply involves some particularly Heroic-type identifying that there's a Dragon, hunting it down, and slaying it.

Bullshit Mountain is basically the opposite of the Dragon. There is one huge, unbearable, painful symptom, that everyone knows about and everyone wishes would just GO AWAY. But no one can get any traction on it. This is because that symptom is actually being contributed to by a thousand little causes, all of which only contribute a little - so making progress on any one of them feels like it doesn't help much, if at all. The only way to solve Bullshit Mountain is for everyone in the org to roll up their sleeves, get a shovel, pick some little corner of Bullshit Mountain to work on, and start shoveling - and not stop until the problem gets noticeably better, even if their work doesn't seem to be doing much to contribute to the improvement.

Finally, we have the **Cloud of Doom**. The Cloud of Doom is where you have a thousand tiny causes, each of which meaningfully contribute to each of a thousand little symptoms, which together make the whole system feel unworkable. The only way to fix the Cloud of Doom is to pour an ungodly amount of Slack into the system, and hope the cloud shakes loose and blows away - otherwise, everyone needs to just throw up their hands and abandon the whole mess.

A good question to ask yourself, when trying to tackle a seemingly insurmountable mess of problems, is: is this a Dragon, Bullshit Mountain, or a Cloud of Doom?

Slaying Dragons

If the problem is a Dragon, are you the one qualified to slay it? If you are, what further information or resources do you need? Do you need a team? Are you the right person to lead that team?

Because Dragons are single-cause problems, they respond well to a single person and a single plan. Most organizations hope (and therefore pretend) that most of their problems are Dragons, and most organizational problem-solving is dedicated to finding Dragons (or making problems look like Dragons) and then paying a few core Heroes big bucks to slay them.

Basically, slaying Dragons is a Solved Problem in the organizational world; most apparent lack of success involves mis-identifying Bullshit Mountains and Clouds of Doom as Dragons, in the naive hope that they'll turn out to be Dragons anyways and therefore can be solved by a process that the system knows how to implement.

Shoveling Bullshit Mountain

If your problem is Bullshit Mountain, do you have enough buy-in to get enough people to start shoveling? Are you the right cheerleader to keep people motivated? Does the team still care enough to even want to solve the problem? These are way harder problems to tackle than the Dragon problems listed above, so expect a lot of people to handwave and convince you that the problem is a Dragon (if they want to signal buy-in for solving it) or a Cloud of Doom (if they don't).

Dealing with Bullshit Mountain as if it was a series of Dragons is what causes "death marches" in the tech industry. It's what makes transformative "business culture" experiments seem to work temporarily (by clearing out the mountain and replacing it with a new system, which then begins accumulating its own Bullshit). It's the biggest contributor to low employee morale that a mid- to large-size "healthy" organization can have. (In fact, the transition from Bullshit Mountain into a Cloud of Doom is probably the tipping point for an organization becoming unsalvageable.)

Surviving The Cloud of Doom

Finally, we have the Cloud of Doom. Every organization has one. Some are big, some are small, some are more toxic than others. But an organization with 0% of its problems in a Cloud of Doom is an organization that has not yet had to do anything actually *real*.

So, you live with it. And the organization begins to develop other problems - mostly bugs, but some Dragons and a few Bullshit Mountains here and there. And as the Dragons get bigger and lay waste to more countryside, and as the Bullshit Mountains tower higher and higher overhead, they start to interweave and correspond, feeding the Cloud of Doom.

Eventually, the Cloud of Doom begins to actually choke the life out of your organization. That's when you have a single choice: inject lots and lots of Slack, or leave.

Injecting lots and lots of Slack basically means "doing less with more", which almost no one believes is the correct choice. But if you aren't burning everything down and starting over somewhere else, it's the *only* choice. If you can't live with your Cloud of

Doom, and you won't flee it, you're going to have to stop feeding it and let it blow away.

Harnessing the Cloud of Doom

One thing you CAN do, if you're particularly vicious, is convince people that your Cloud of Doom is actually just Bullshit Mountain, and use that to extract work from your subordinates. You need to be exceptionally clever (in a [Raoian sociopathic](#) sense) in order to pull this off, because you basically need to manage both sides of the information and effort flow: you have to keep everyone believing that there's a Bullshit Mountain that they're biting into, AND you have to re-direct and manage their actual efforts so that they benefit your covert goals. Note that this is doable whether the actual problem is in fact a Cloud of Doom or merely Bullshit Mountain, and that applying this very process to Bullshit Mountain is one of the more common things that turns Bullshit Mountain *into* a Cloud of Doom in the first place.

If I catch you doing this within any organization that I am aligned with, I consider it within my rights to destroy you.

Dragonslayers are Gryffindor, Shovelers are Hufflepuff

Noticing what kind of problem you have the right temperament to solve is key to avoiding burn-out. Dragonslaying is glamorous, high-praise work; shoveling Bullshit Mountain is thankless and grueling, and the person who finally gets the praise is usually the guy that did the least actual work. From my perspective, Project Hufflepuff was in many ways a direct attempt to train up people who could handle the Rationalist-and-EA-community's Bullshit Mountain before it turned into a Cloud of Doom. (Whether it's too late now or not is a matter for the Slytherins to convince you; I will say no more on this today.)

One big problem with shovelers is that they still expect praise and reward for shoveling Bullshit Mountain; no one seems to be telling them that it ain't gonna happen. (The smart ones figure this out on their own, and either go away or [grit their teeth and get to work anyways](#).) As a system for motivating people to shovel, anything like Project Hufflepuff is doomed to fail from the beginning. What you need is a way to identify the people who are already doing the work, and make sure they stay funded and supported and nurtured. (Whether this problem itself is a Bullshit Mountain or a Cloud of Doom is, again, left as an exercise for the reader. A Dragon it ain't, or Project Hufflepuff would have worked.)

Conclusions

So, anyway, yeah. Here we are. Do what you will with it. Or don't; I'm not your dad.

[[Epistemic status](#)]

Bridging syntax and semantics with Quine's Gavagai

Crossposted from the [AI Alignment Forum](#). May contain more technical jargon than usual.

Quine [has an argument](#) showing that you can never be sure what words mean in a foreign language:

Quine uses the example of the word "**gavagai**" uttered by a native speaker of the unknown language Arunta upon seeing a rabbit. A speaker of English could do what seems natural and translate this as "Lo, a rabbit." But other translations would be compatible with all the evidence he has: "Lo, food"; "Let's go hunting"; "There will be a storm tonight" (these natives may be superstitious); "Lo, a momentary rabbit-stage"; "Lo, an undetached rabbit-part." Some of these might become less likely – that is, become more unwieldy hypotheses – in the light of subsequent observation.

What does this mean from the perspective of [empirically bridging syntax and semantics?](#)

Well, there is no real problem with "gavagai" from the empirical perspective; the fact that there seems to be one is, in my view, due to the fact that the syntax-semantic discussion has focused too heavily on the linguistic aspects of the problem.

Let G be the symbol in the Arunta speaker's brain that activates when they say "gavagai". As Quine said towards the end of his quote, "some of these might become less likely – that is, become more unwieldy hypotheses – in the light of subsequent observation." Relatedly, Nick Bostrom [argues](#) that indeterminacy is a "matter of degree".

In practice, this means is that if, for instance, r ="A rabbit is there" and s ="A storm is coming tonight", then, if we accumulate enough observations, the G is going to be predictive of r better than it is of s , or vice versa. Thus there is an empirical test for whether G corresponds to some of these hypotheses.

But what about u ="An undetached rabbit-part is there"? It may be very difficult to find a situation that distinguishes between u and r in practice - so does G stand for "rabbit", or "undetached rabbit-part"?

Similarly to the example of the neural net in the [previous post](#), G is a symbol for *both* r and u . If no experiment can distinguish between the two - or, at least, if no experiment can distinguish between the two in any typical environment that any Arunta speaker would ever experience - then they are synonymous (or, equivalently, they are strongly within each other's [web of connotations](#)).

If we presented the speaker with a situation far outside their typical environment, then we might be able to distinguish r from u - but we'd be uncertain if G 'meant that all along', or if the speaker was extending the definition to a new situation.

Back to linguistics

There are some ways we might be able to distinguish whether "gavagai" means "rabbit" or "undetached rabbit-part", even if u and r cannot be distinguished in practice. It's plausible that there might be an *internal* variable P in the speaker that corresponds to "part of an animal", and an internal U that corresponds to "undetached" (versus detached).

Then, when G is activated, we might ask whether P and U are also activated, or not. You can see this as the *internal* web of connotations amongst the variables in the human brain. It seems people with different languages have different internal webs, [different ways of connecting words and concepts](#) together.

This does not, however, affect the connection between G and external variables.

Iterative Arguments: Alternative to Adversarial Collaboration

I've been toying with an idea of developing two competing theories in parallel in an iterative manner:

1. A writes an initial thesis
2. B does the same
3. A revises their thesis to address B's thesis where it contradicts that of A's
4. B does the same
5. and so on until both parties feel they have nothing to add

This would be different from adversarial collaboration, as it's commonly understood, in that both sides would work on their own arguments instead of trying to agree on a common summary (which is very hard!).

It's worth emphasising that the idea is not to correspond with the opponent. Instead, one would keep updating one's thesis to meet the challenges presented by the competing thesis so that it remains coherent and stands on its own after every iteration.

I wrote a little ClojureScript library to go with a [Pandoc template](#) to facilitate iterative argumentation of this kind. The library makes it easy to navigate between the sides of the argument ensuring that following local links will switch context when appropriate, etc. Additionally, it provides bidirectional links by way of highlighting bits on both sides of the argument. Hot loading is also supported to make writing the document more convenient. Currently it's very much work in progress, of course. I didn't want to commit to too many features at this point not knowing if any of this would be useful.

Would someone here be interested in trying out this sort of thing either with their own "archnemesis" or with me? For starters, I think it would be best to pick a properly contentious topic but not one inciting too much passion.