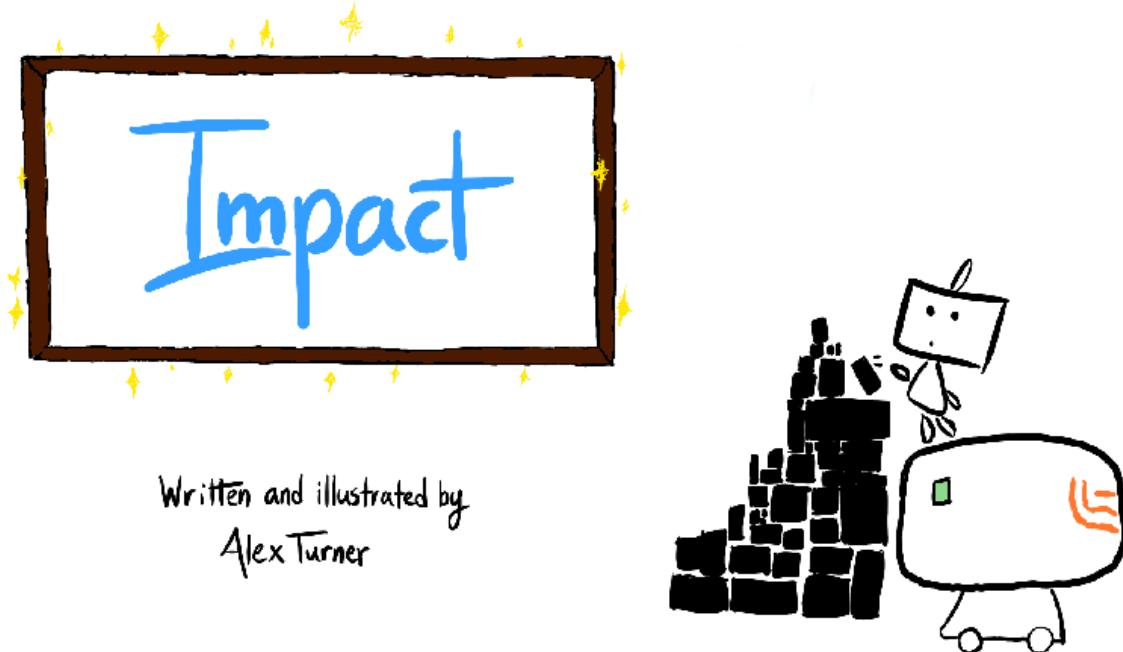


# Reframing Impact

1. [Reframing Impact](#)
2. [Value Impact](#)
3. [Deducing Impact](#)
4. [Attainable Utility Theory: Why Things Matter](#)
5. [World State is the Wrong Abstraction for Impact](#)
6. [The Gears of Impact](#)
7. [Seeking Power is Often Convergently Instrumental in MDPs](#)
8. [Attainable Utility Landscape: How The World Is Changed](#)
9. [The Catastrophic Convergence Conjecture](#)
10. [Attainable Utility Preservation: Concepts](#)
11. [Attainable Utility Preservation: Empirical Results](#)
12. [How Low Should Fruit Hang Before We Pick It?](#)
13. [Attainable Utility Preservation: Scaling to Superhuman](#)
14. [Reasons for Excitement about Impact of Impact Measure Research](#)
15. [Conclusion to 'Reframing Impact'](#)

# Reframing Impact

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



Imagine we have a robot named Frank.

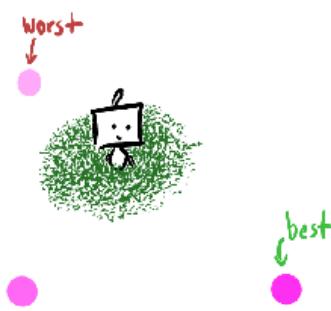
Frank finds things for us in places we can't go.

We provide a rule, and he returns with the object that best fits the rule.

Right now, we **want** a very pink marble.

(in case you're wondering)

Naturally, we ask Frank for the pinkest thing he can find.

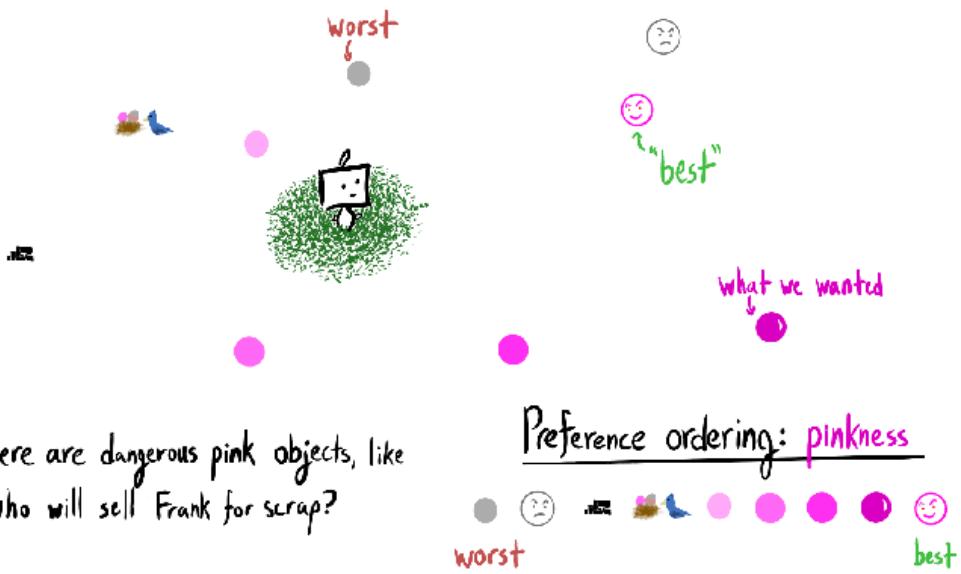


## Preference ordering: pinkness

worst best

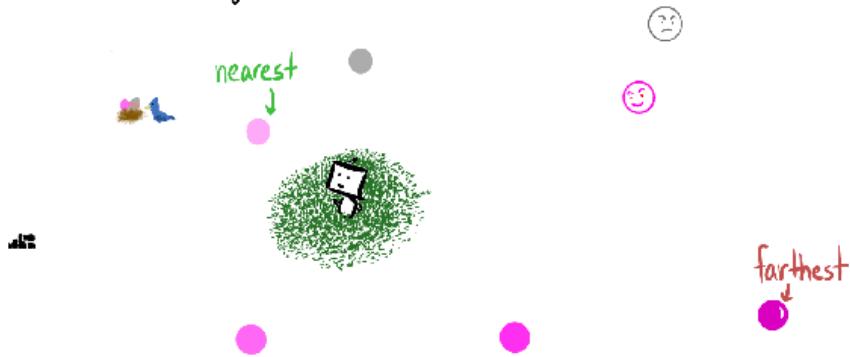
This seems fine. But what if Frank looks farther afield?

The world is wide, and full of objects.



From our perspective, Frank has lost his marbles, but he's just following an imperfect rule.  
The perfect rule is hard to specify. What simple rule avoids terrorists here?

Fortunately, the terrorists are far away.



Pinkness correlates with what we want,  
and the proximity rule avoids terrorists.

Preference ordering: proximity



Think about what Frank brings us for each distance.



We're probably fine with a reasonably pink marble.  
Then how about we have Frank find the pinkest object  
within a given distance, which we increase until we're satisfied?

And now for the reveal

Frank is analogous to a powerful AI with an imperfect **objective**. The objects are plans he's considering, and the terrorists are catastrophic plans (some of which happen to **score well**).

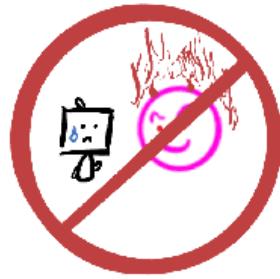
The question is then:



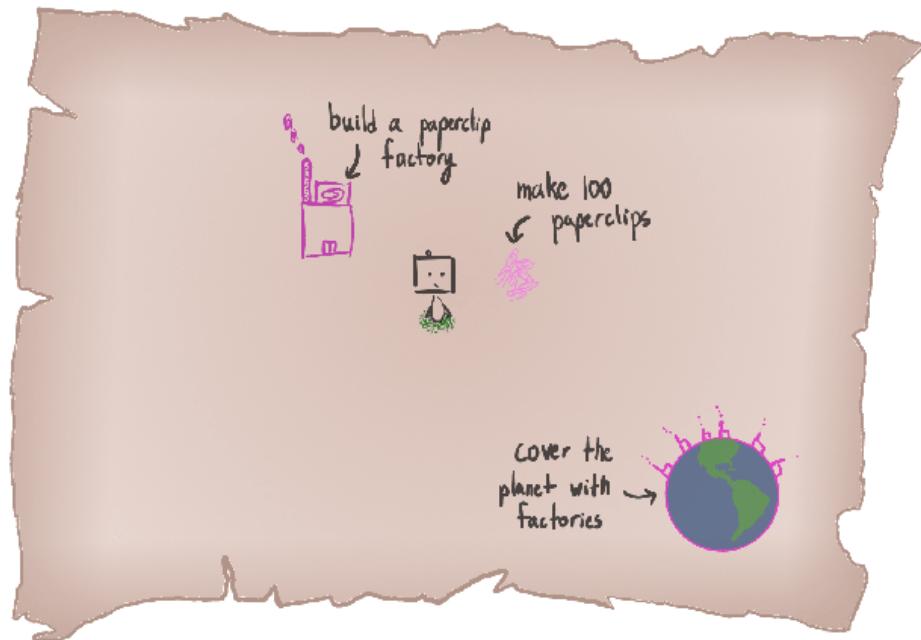
How do we measure how  
distant  
plans are?

The *distance measure* should:

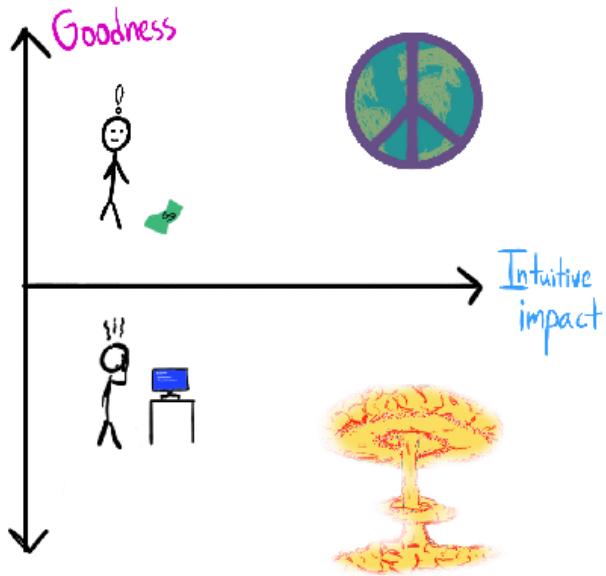
- 1) Be easy to specify
- 2) Put catastrophes far away
- 3) Put reasonable plans nearby



Suppose Frank's *objective* is to make paperclips.  
The *measure* will put plans on the map:



These catastrophes seem like **big deals**. We're going to figure out **why** we intuit some things are **big deals**, develop an understanding of the relevant parts of reality, and then design an **impact measure**.



To me, the **impactful** things feel fundamentally different than the non-**impactful** things. I find this difference fascinating and beautiful, and look forward to exploring it with you.

Why be excited about **impact measurement**? After all, it doesn't seem like a one-shot solution to AI alignment. It doesn't even seem like a key problem, like figuring out how to get AIs to robustly learn **human values**, or understanding what it means to be both an **agent and part** of the environment.

Frankly, this misses the forest for the trees. At its core, **impact measurement** is about how, when, and why agents **affect** each other. Understanding this requires a new way of looking at the world.

The question of **impact measurement** has caused significant confusion, but together, we'll find comprehension. We're going to emerge **saddled with spoils**: new conceptual frameworks, fresh lines of inquiry, and important theoretical milestones.

Here's one exciting milestone we're shooting for:



An impact measure would be the first proposed safeguard which maybe actually stops a powerful agent with an imperfect objective from ruining things – without assuming anything about the objective. This is a rare property among approaches.



We have our bearing.  
Let us set out together



## Technical Appendix: First safeguard?

*This sequence is written to be broadly accessible, although perhaps its focus on capable AI systems assumes familiarity with [basic arguments for the importance of AI alignment](#). The technical appendices are an exception, targeting the technically inclined.*

Why do I claim that an impact measure would be "the first proposed safeguard which maybe actually stops a powerful agent with an imperfect objective from ruining things – without assuming anything about the objective"?

The safeguard proposal shouldn't have to say "and here we solve this opaque, hard problem, and then it works". If we have the impact measure, we have the math, and then we have the code.

So what about:

- Quantilizers? This seems to be the most plausible alternative; mild optimization and impact measurement share many properties. But

- What happens if the agent is already powerful? A greater proportion of plans could be catastrophic, since the agent is in a better position to cause them.
- Where does the base distribution come from (opaque, hard problem?), and how do we know it's safe to sample from?
  - In the linked paper, Jessica Taylor suggests the idea of learning a human distribution over actions – how robustly would we need to learn this distribution? How numerous are catastrophic plans, and what *is* a catastrophe, defined without reference to our values in particular? (That definition requires understanding impact!)
- [Value learning](#)? But
  - We only want this if *our* (human) values are learned!
    - [Value learning is impossible without assumptions](#), and [getting good enough assumptions could be really hard](#). If we don't know if we can get value learning / reward specification right, we'd like safeguards which don't fail because value learning goes wrong. The point of a safeguard is that it can catch you if the main thing falls through; if the safeguard fails because the main thing does, that's pointless.
- [Corrigibility](#)? At present, I'm excited about this property because I suspect it has a simple core principle. But
  - Even if the system is responsive to correction (and non-manipulative, and whatever other properties we associate with corrigibility), what if we become *unable* to correct it as a result of early actions (if the agent "moves too quickly", so to speak)?
    - [Paul Christiano's take on corrigibility](#) is much broader and an exception to this critique.
  - What is the core principle?

## Notes

- The three sections of this sequence will respectively answer three questions:
  - Why do we think some things are big deals?
  - Why are capable goal-directed AIs incentivized to catastrophically affect us by default?
  - How might we build agents without these incentives?
- The first part of this sequence focuses on foundational concepts crucial for understanding the deeper nature of impact. We will not yet be discussing what to implement.
- I strongly encourage completing the exercises. At times you shall be given a time limit; it's important to learn not only to reason correctly, but with speed.

The best way to use this book is NOT to simply read it or study it, but to read a question and STOP. Even close the book. Even put it away and THINK about the question. Only after you have formed a reasoned opinion should you read the solution. Why torture yourself thinking? Why jog? Why do push-ups?

If you are given a hammer with which to drive nails at the age of three you may think to yourself, "OK, nice." But if you are given a hard rock with which to drive nails at the age of three, and at the age of four you are given a hammer, you think to yourself, "What a marvellous invention!" You see, you can't really appreciate the solution until you first appreciate the problem.

~ [Thinking Physics](#)

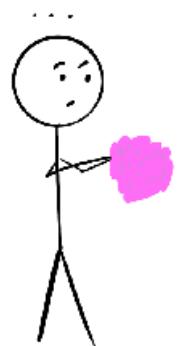
- My paperclip-Balrog illustration is metaphorical: a good impact measure would hold steadfast against the daunting challenge of formally asking for the right thing from a powerful agent. The illustration does not represent an internal conflict within that agent. As water flows downhill, an impact-penalizing Frank prefers low-impact plans.
  - The drawing is based on [gonzalokenny's amazing work](#).

- Some of you may have a different conception of impact; I ask that you grasp the thing that I'm pointing to. In doing so, you might come to see your mental algorithm is the same. Ask not “is this what I initially had in mind?”, but rather “does this make sense as a thing-to-call-'impact'?”.
- H/T Rohin Shah for suggesting the three key properties. Alison Bowden contributed several small drawings and enormous help with earlier drafts.

# Value Impact

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

We think some things are big deals,  
and we want to understand why.



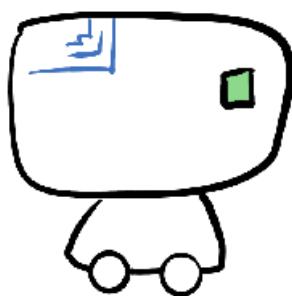
However, it can be hard to read your own mind.

Instead, we'll use thought experiments  
to piece together what's going on.



Xyz is a Pebblehoarder of the planet Pebblia.  
It morally values collections of pebbles, and that's it.

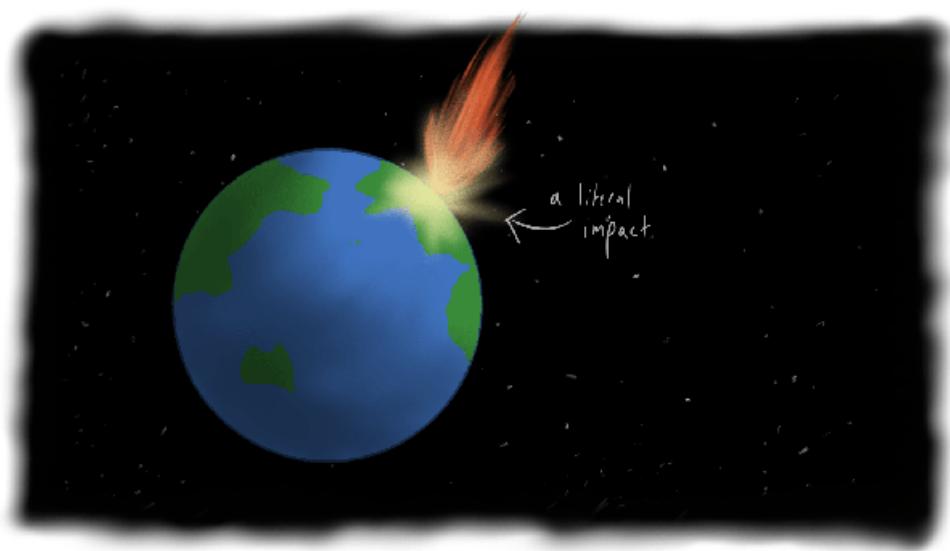
? One day, all of the pebbles turn into obstruction blocks,  
which every Pebblehoarder knows are worthless.



Far, far away from Earth exists the planet Iniron.  
One day, we learn humans are now being tortured there.



An asteroid strikes.

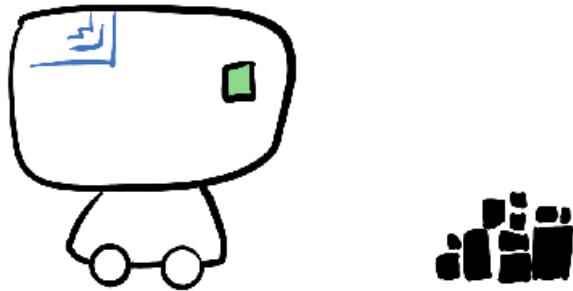


Exercise: Spend three minutes familiarizing yourself with the three situations – how are they alike, and how are they different? Make them come alive.



Let's query our mental impact-o-meter from these different vantage points.  
Step into each pair of shoes and ask "how big of a deal is this?"

Just imagine being XYZ.  
?



The very fabric of what is **important** has been ripped away.

Perhaps the Pebblehoarder civilization can rebound and find **value** in the universe, but if not - if XYZ doesn't know you can just make more pebbles - the **loss feels complete**.

The universe feels dead  
and empty  
and worthless

Faced with an *impact* of similar magnitude,  
we might have a feeling of freefalling despair,

of our pale blue marble having been  
pushed  
off



a

cliff

and shattered against the ground far below.

What is the impact on an inhabitant of an extraterrestrial Pebblehoarder colony?

Somehow, this assessment seems to depend intricately on how they value their collections of pebbles.

If they only value their own collections, then this doesn't matter, except to the extent that the colony becomes overcrowded.

If they value the total number of collections, then this is bad news.

And as for us - would anyone honestly think this is earth-shattering?

No.

Even if we were on Pebbtia, we'd probably think primarily of the impact on human-Pebblehoarder relations.

This is where our eyes widen as we realize how much this reveals about the nature of the impact calculation running in our heads

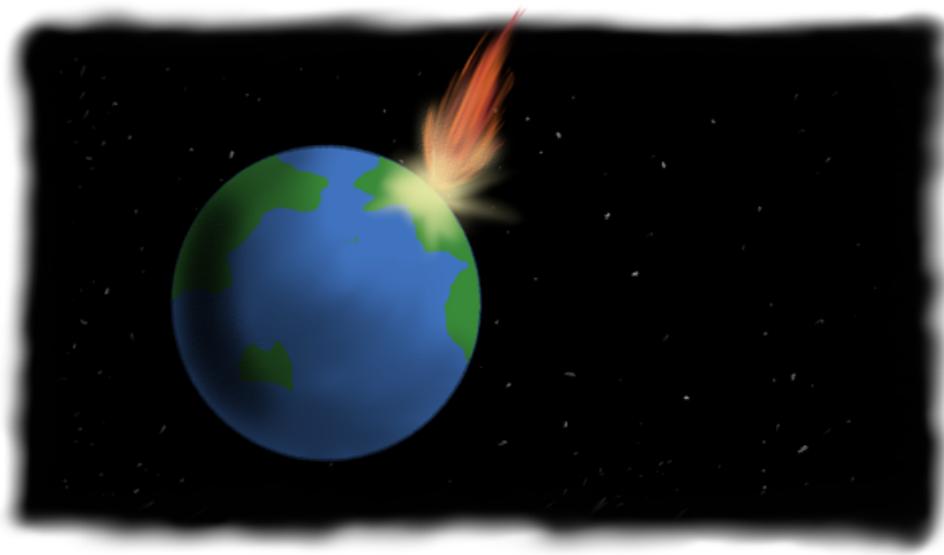


We feel a **pull** to help the **poor** souls of Iniron.

But XYZ? XYZ doesn't **care**.

There aren't any pebbles on the line.

Even if it were on Iniron, its thoughts would flit to how  
this development **affects** its own **concerns**.



Exercise: Determine how **impactful** the asteroid impact is to:

- You on Earth      ◦ XYZ on Earth
- You on Pebblia      ◦ XYZ on Pebblia

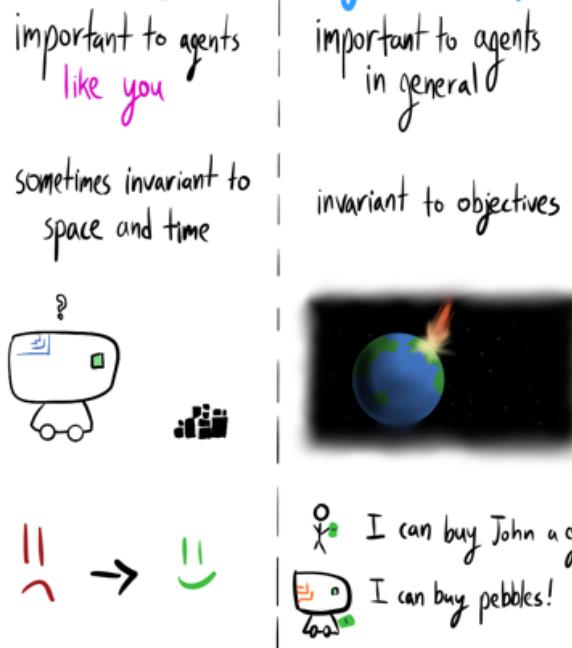
Being on Earth when this happens is a big deal, no matter your objectives – you can't hoard pebbles if you're dead! People would feel the loss from anywhere in the cosmos. However, Pebblehoarders wouldn't mind if they weren't in harm's way.



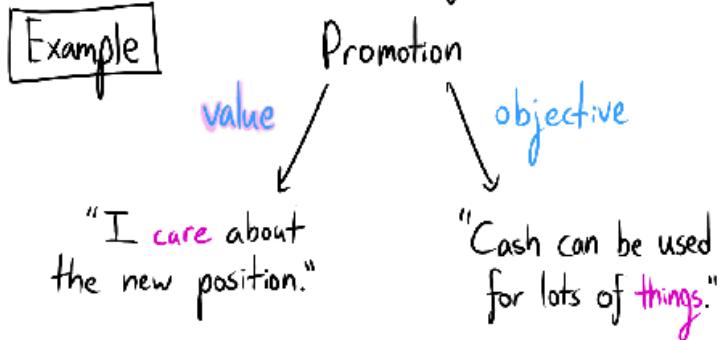
What have we learned?

Impact is relative to what you want and where you are.

Impact = value impact + objective impact



Exercise: Decompose something which recently impacted you.



## Appendix: Contrived Objectives

A natural definitional objection is that a few agents aren't affected by objectively impactful events. If you think every outcome is equally good, then who cares if the

meteor hits?

Obviously, our values aren't like this, and any agent we encounter or build is unlikely to be like this (since these agents wouldn't do much). Furthermore, these agents seem contrived in a technical sense (low measure under reasonable distributions in a reasonable formalization), as we'll see later. That is, "most" agents aren't like this.

From now on, assume we aren't talking about this kind of agent.

---

### Notes

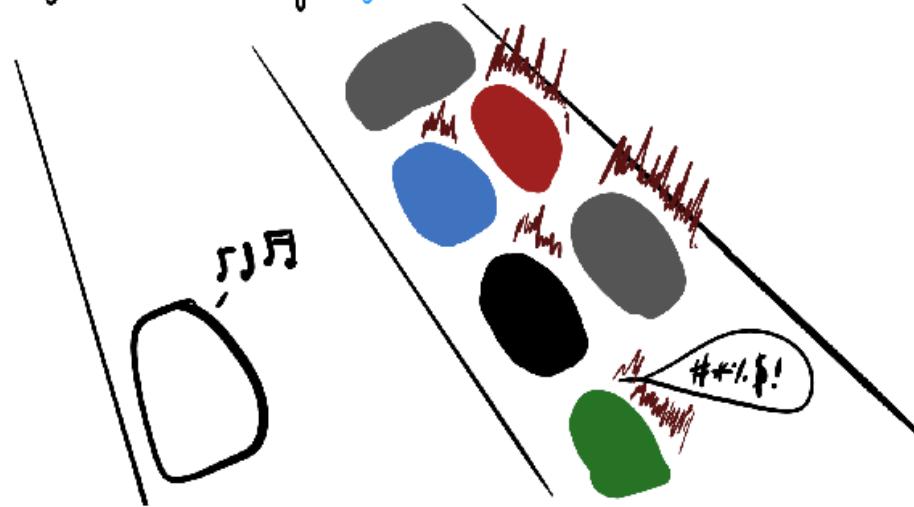
- Eliezer [introduced Pebblesorters in the the Sequences](#); I made them robots here to better highlight how pointless the pebble transformation is to humans.
- In informal parts of the sequence, I'll often use "values", "goals", and "objectives" interchangeably, depending on what flows.
- We're going to lean quite a bit on thought experiments and otherwise speculate on mental processes. While I've taken the obvious step of beta-testing the sequence and randomly peppering my friends with strange questions to check their intuitions, maybe some of the conclusions only hold for people like me. I mean, [some people don't have mental imagery](#) – who would've guessed? Even if so, I think we'll be fine; the goal is for *an* impact measure – deducing human universals would just be a bonus.
- Objective impact is objective with respect to the agent's *values* – it is *not* the case that an objective impact affects you anywhere and anywhen in the universe! If someone finds \$100, that matters for agents at that point in space and time (no matter their goals), but it doesn't mean that everyone in the universe is objectively impacted by one person finding some cash!
- If you think about it, the phenomenon of objective impact is *surprising*. See, in AI alignment, we're used to no-free-lunch this, no-universal-argument that; the possibility of something objectively important to agents hints that our perspective has been incomplete. It hints that maybe this "impact" thing underlies a key facet of what it means to interact with the world. It hints that even if we saw specific instances of this before, we didn't know we were looking at, and we didn't stop to ask.

# Deducing Impact

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

*Impact* is in the eye of the beholder.

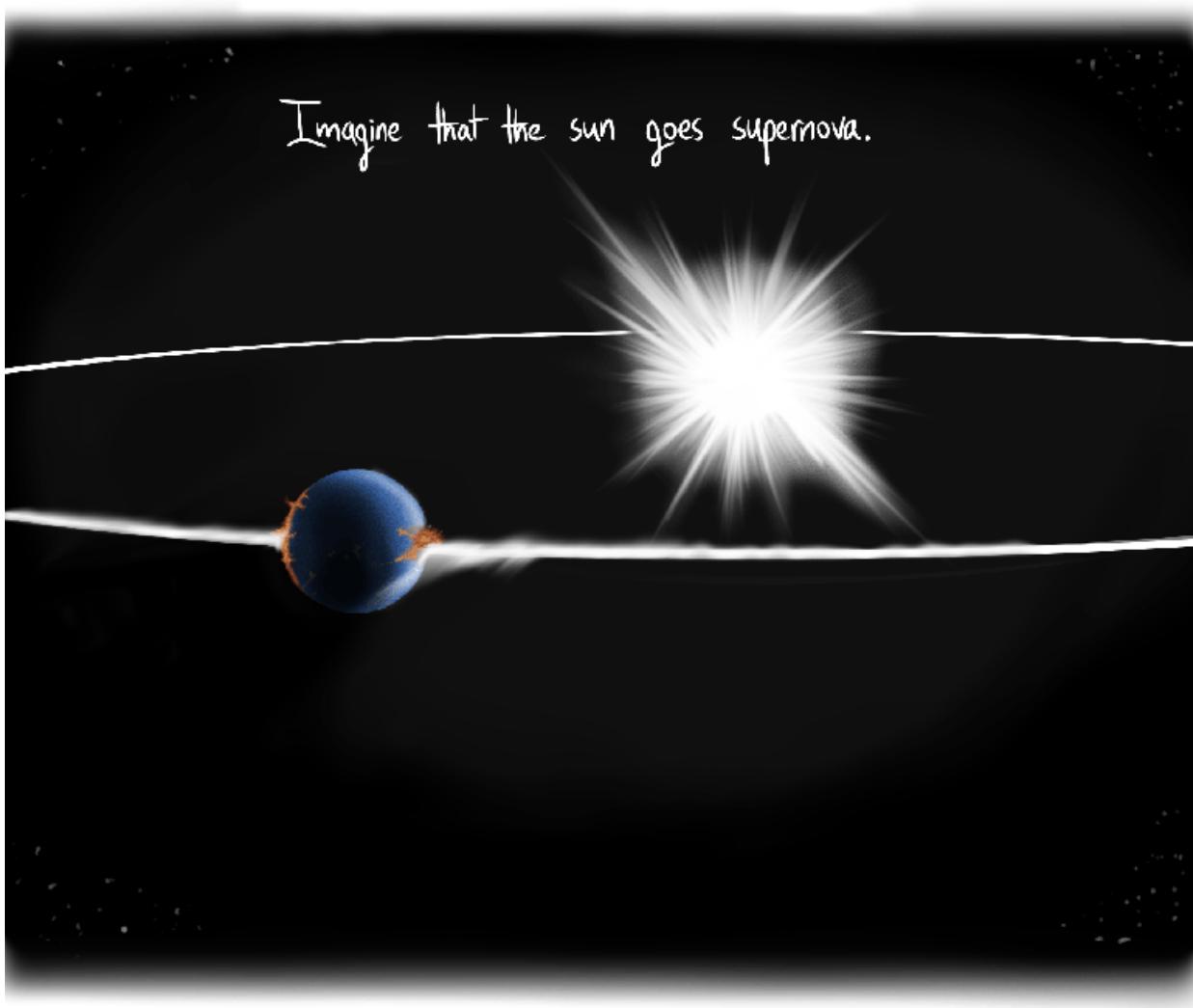
Traffic jams are vividly *big deals* to tardy commuters,



but everyone else doesn't really care.

This concept is important, so I'm going to present another zany situation.

Imagine that the sun goes supernova.



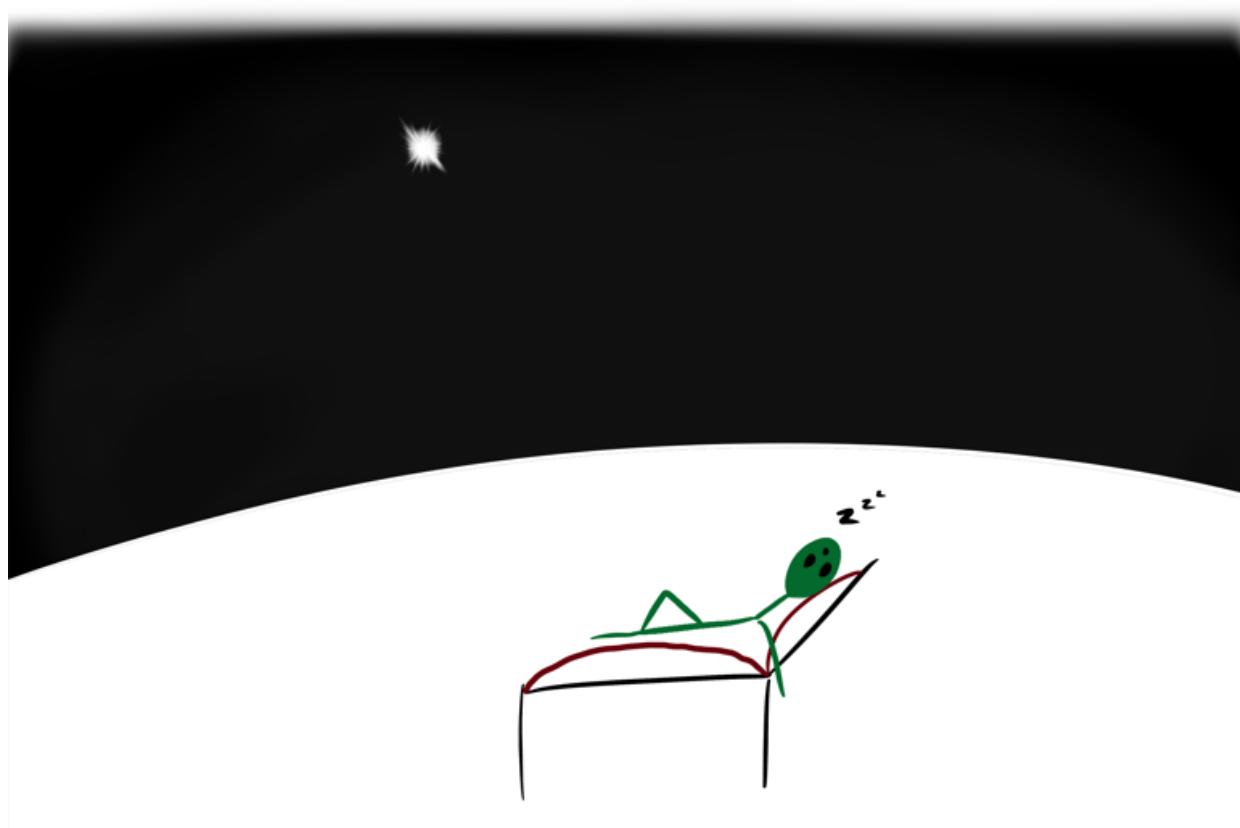
Now, being on -



Our sun is a main-sequence G-type star, it can't explode. Any energy input just decreases the volume of the hydrogen plasma, the Sun doesn't have a degenerate core that could be detonated. The Sun doesn't have enough mass to go supernova.

Yes, yes, thank you - everyone knows that.  
It's just another weird situation I'm forcing on my readership.

Now, being on Earth here is objectively impactful because it matters to almost any agent in your shoes. However, whether this is a big deal to you depends on who and where you are.



## Impact is Comparative

Suppose you grew up on an Earth where astrophysicists not only believe that stars less than eight solar masses can go supernova, but that the sun will — and soon.

Everyone knows it.

Everyone expects it to happen.

Everyone thinks it's unavoidable.

Including you.

Turns out, someone miscalculated; your calculations say it won't happen.

Saying this feels like a big deal is like calling a supernova "visible".

But to us, the sun not going supernova is... expected — zero impact.

You can't have it both ways — the sense of impact we experience has to be comparative — concerning some kind of difference, or change.



9

Our overarching goal is to deeply understand why capable AIs with **goals** are incentivized to catastrophically **impact** us, and then discuss an agent design which potentially avoids these incentives. The first step towards this end is, of course, understanding **impact**.

One rarely has the opportunity to try their hand at connecting the dots of an important insight. In AI alignment, I suspect that some important questions have yet to be properly posed, let alone solved. To get feedback and grow as a researcher, one must hone their reason on the known conquests of other fields.

This question of **impact** is not the important question, but it is an important question. It has both reasonable difficulty and an answer not yet widely circulated. Go ahead — try your hand.

6

Exercise: Deducing and informally describing why we think some things are **big deals**. Remember:

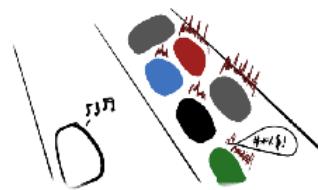
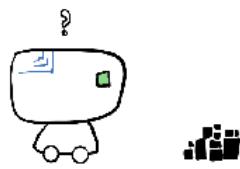
- Impact is relative both to what you value and to your vantage point.
- Part of impact is particular to agents like you, and part is objective.
- Impact is comparative.

Find the simple concept neatly explaining why things feel like **big deals** or not in the examples thus far.

The answer can be expressed in just one sentence of everyday language.

You have the benefit of knowing a solution exists.

You have fifteen minutes.



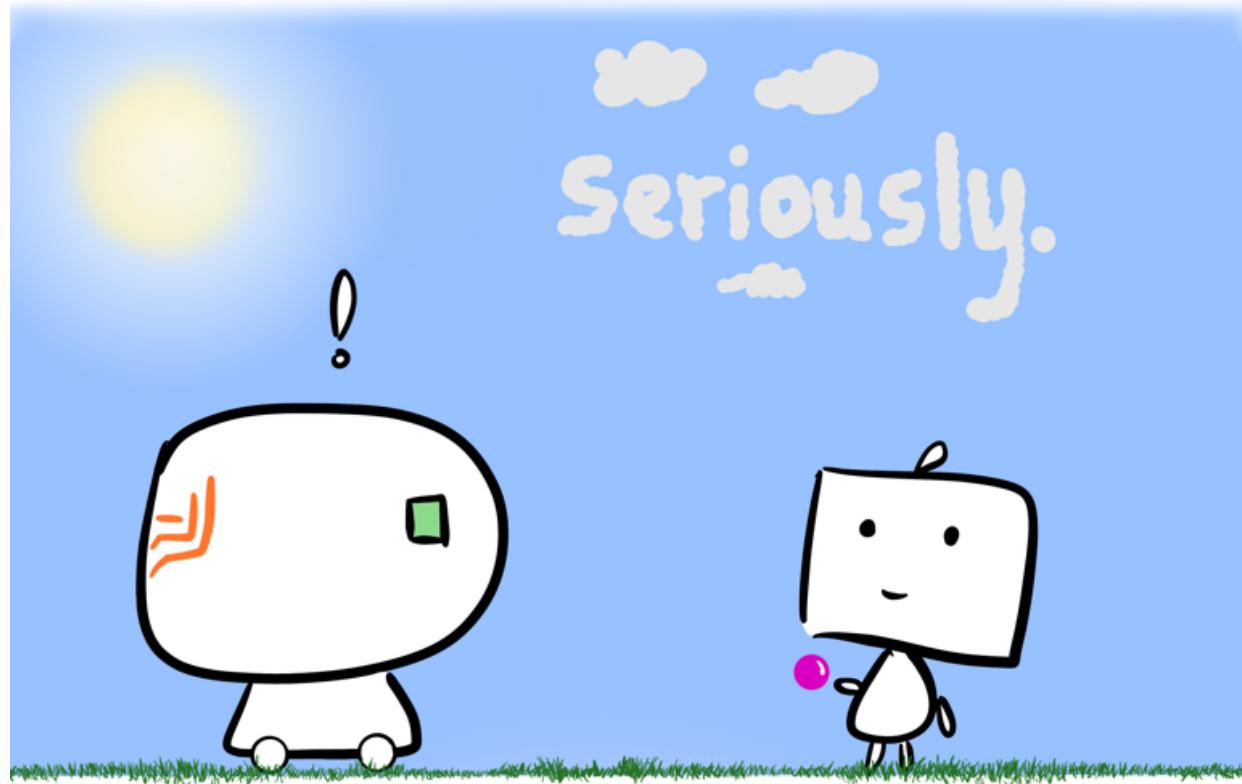
The solution comes in the next post! Feel free to discuss amongst yourselves.

Reminder: Your sentence should explain impact from all of the perspectives we discussed (from XYZ to humans).

# Attainable Utility Theory: Why Things Matter

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

*If you haven't read the prior posts, please do so now. This sequence can be spoiled.*



When thinking about whether something **impacts** us, we ask:

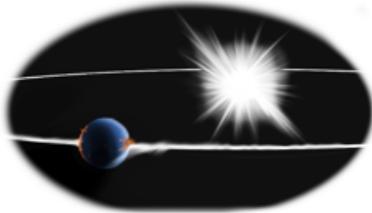
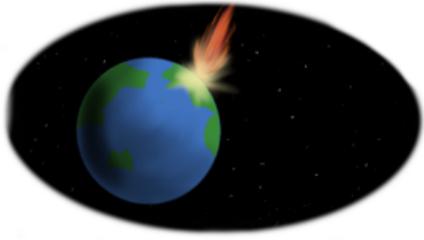
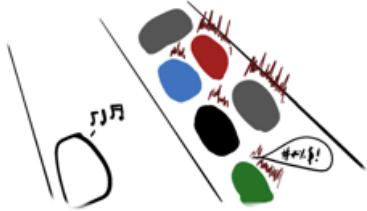
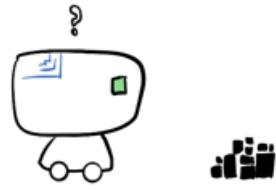
How does this  
change my ability  
to get what I want?

These are the **big deals**.

This is what **affects** us.

This is what **matters** to us.

This is **impact**.



When you think about it, how could something possibly be a big deal to you if it doesn't change your ability to get what you want? If it doesn't change your ability to get what you want, you won't care. Conversely, how could something not matter to you if it does change your ability to get what you want?

An attainable utility (AU) refers to your ability to get what you want.



Correct theories make correct predictions,  
so let's take AU theory out for a spin.

## Locality

All theory predicts: objective impact to places we can't reach  
doesn't feel impactful.

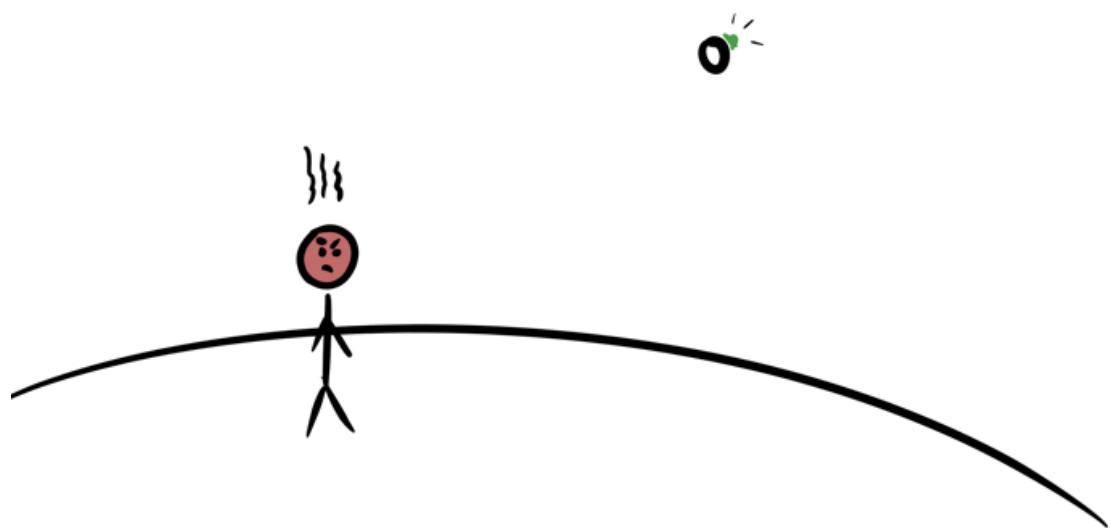
Imagine a giant stack of money is about to be yours.



\$ \$ \$



Then, I move it to the moon.



If I then move it even farther away, does this matter to you?

Well, you couldn't reach it anyways, so who cares?

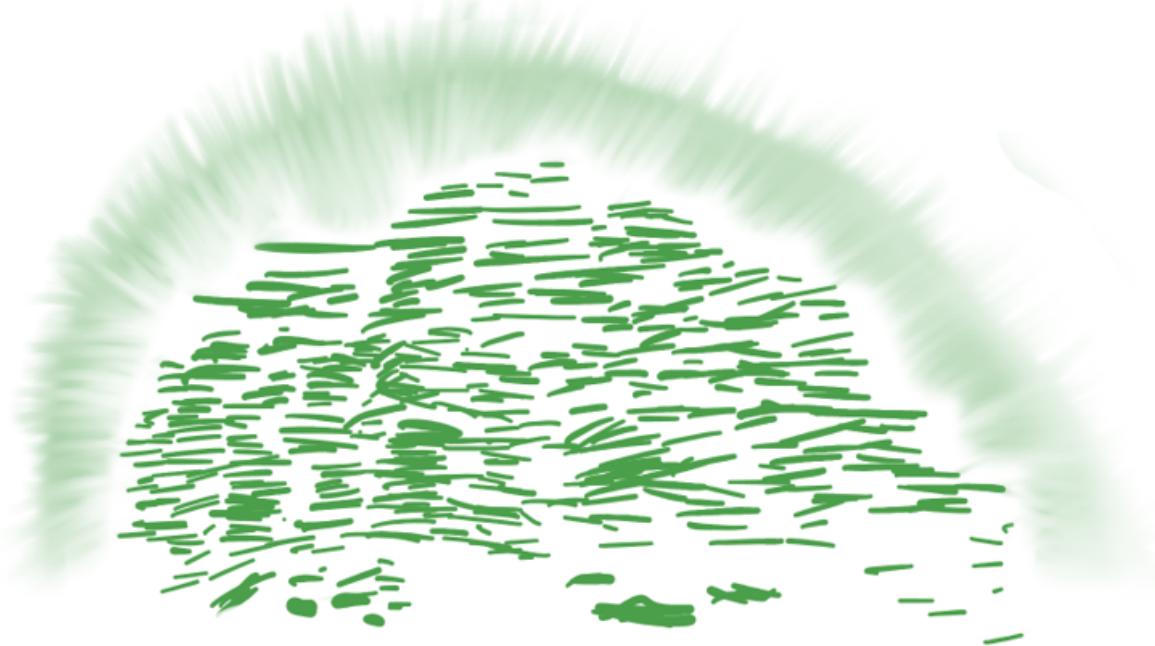
# Discounting

All theory predicts: We discount **impact** to our future selves exactly the same way we discount **value** to our future selves.

Close your eyes and pretend you're the kind of person who lives in and for the moment. You care about things if they happen soon, and after that— who knows?

You learn that in ten years, your net worth will be \$10 million.

\$ \$ \$



How big of a deal is that?

Now clear your mind, and imagine you learn this as your normal self:  
\$10 million, ten years from now.

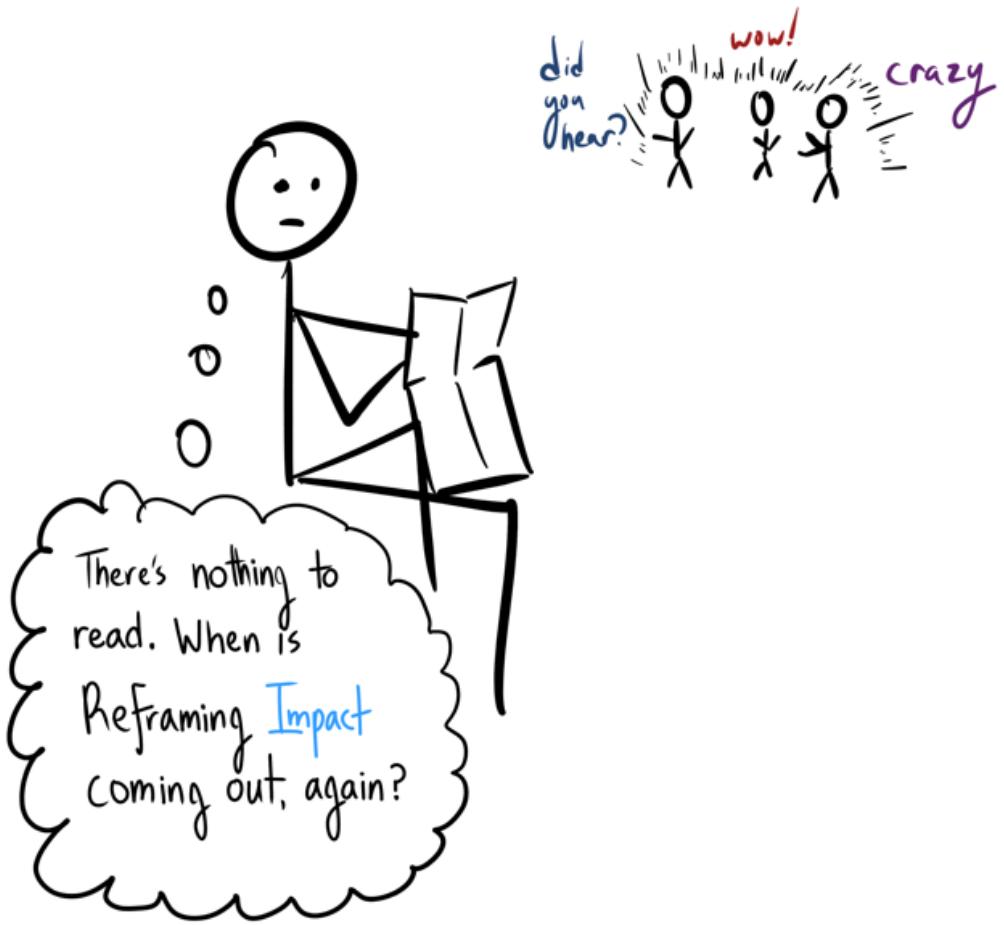
How big of a deal is it now?

It probably felt more **impactful** when you actually care about the future.

## Beliefs

All theory predicts: If knowledge changes your expectations about whether you can get what you want, then learning it feels **impactful**.

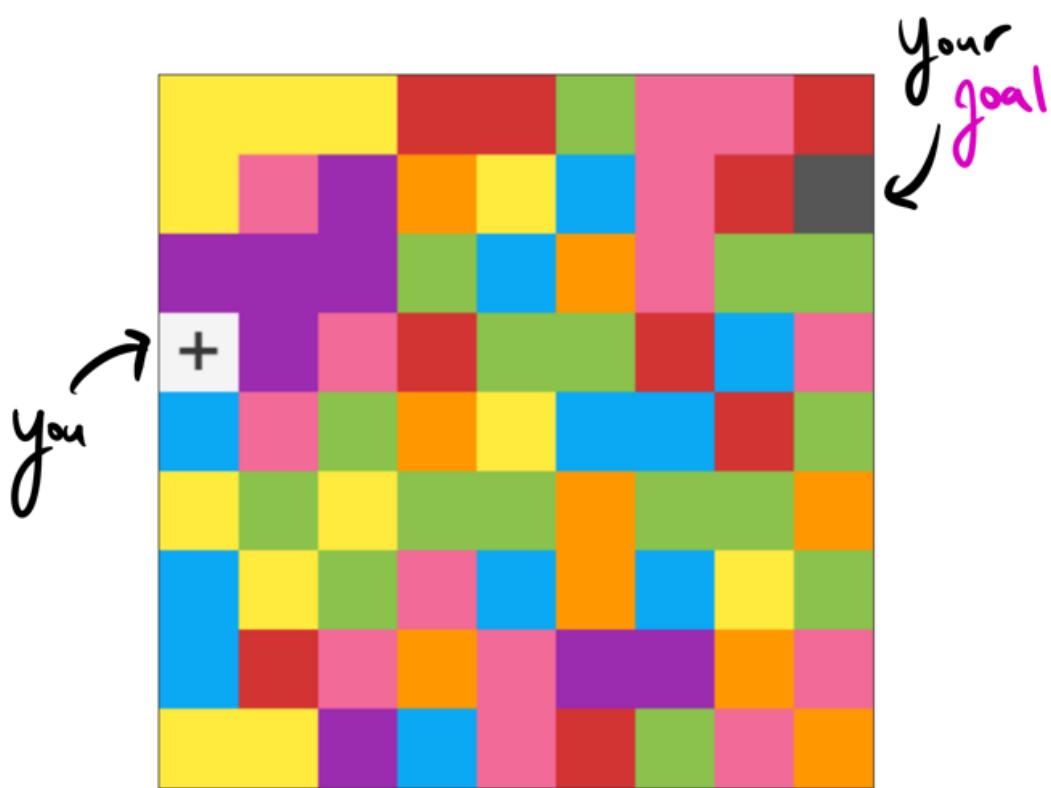
Remember the world where everyone thought the sun was going to explode? Then, you calculated that the sun can't go supernova. This feels **impactful**. As the news spreads, you're reminded of the imminent non-explosion by the newspapers.



You feel no **impact**, even though everyone else is pretty **blown away**.

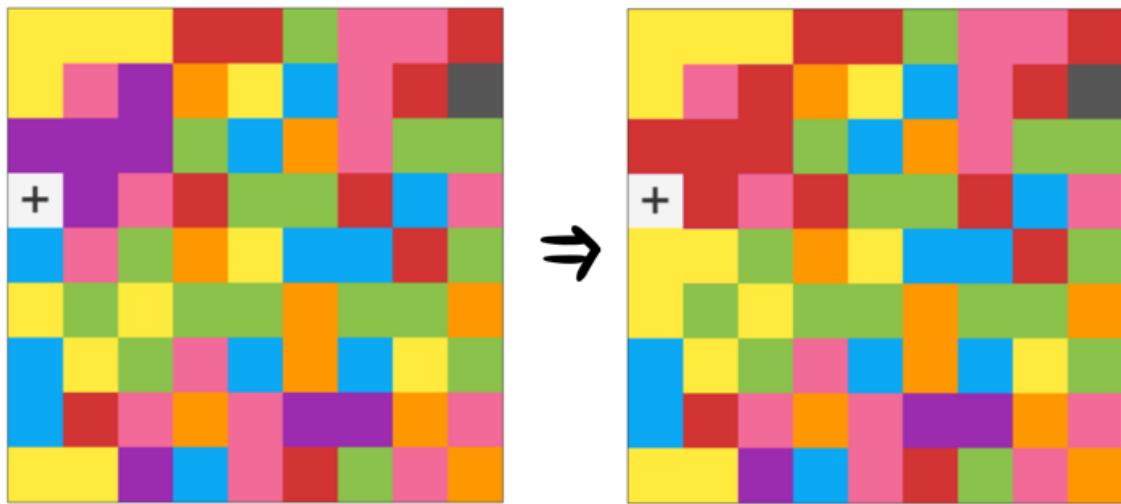
## Universality

All theory predicts: We can imagine **impact** in environments unlike our own; the more we think we know about how things work, the more acute our sense of **impact** becomes.



Pretend this maze is your reality. You're presently able to reach the goal.

How impactful would the following be?



-\(\_\times\)\_/-

You literally don't know how anything works. However, when I tell you that you're how stuck, this seems like a **big deal**.

Surprisingly, this works no matter how "weird" the reality: All theory correctly predicts **impact** for agents running on a Powerpoint presentation.

Isn't that something?



How does this  
change my ability  
to get what I want?

I think that with its mere eleven words, All theory completely explains our intuitions about impact. Again, how could something possibly be a big deal to us if it doesn't change our ability to get what we want? How could something not matter to us if it does change our ability to get what we want?

You really can't separate perceived changes in our ability to get what we want from our sense of impact, because they're the same thing.

# World State is the Wrong Abstraction for Impact

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

I've been keeping something from you.  
Remember the confusion I mentioned in the first post?

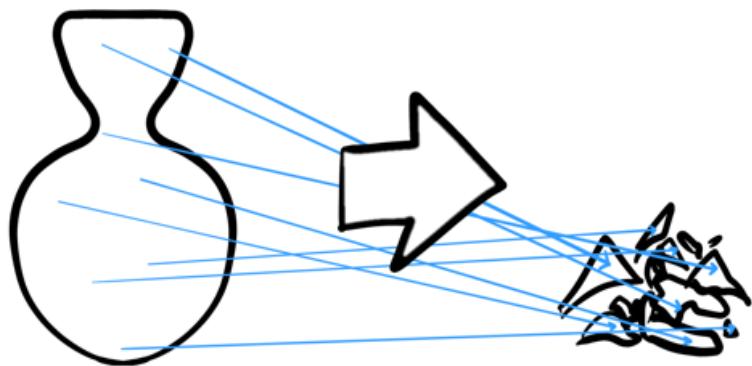
Before the attainable utility theory of impact came along, people made an assumption about what impact is - a reasonable, obvious, compelling assumption: that impact is primarily about how the state of the world changes.



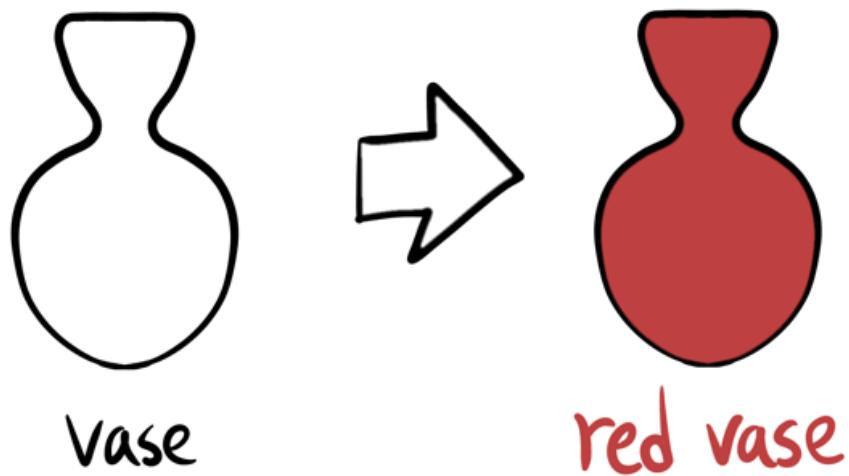




Maybe we think **impact** is about change in particle positions,



or maybe we think *it's* about change in object identities.



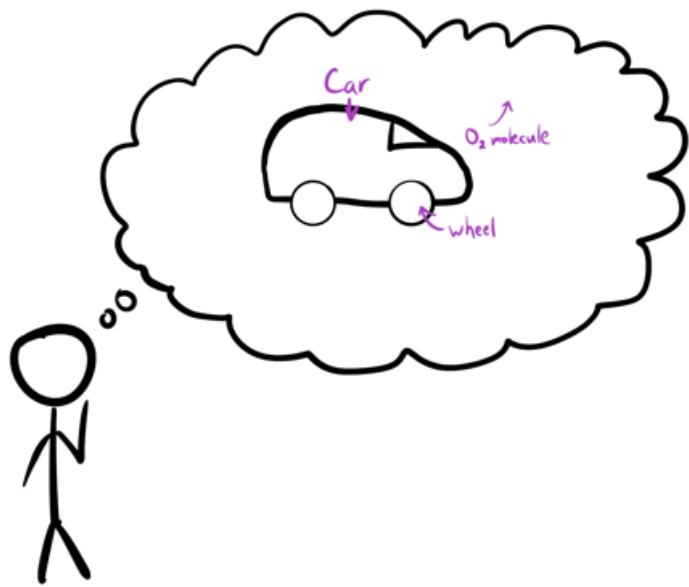
How can you not be tempted by that assumption –  
things are changing in the world!  
The assumption's right there.

It's so obvious.

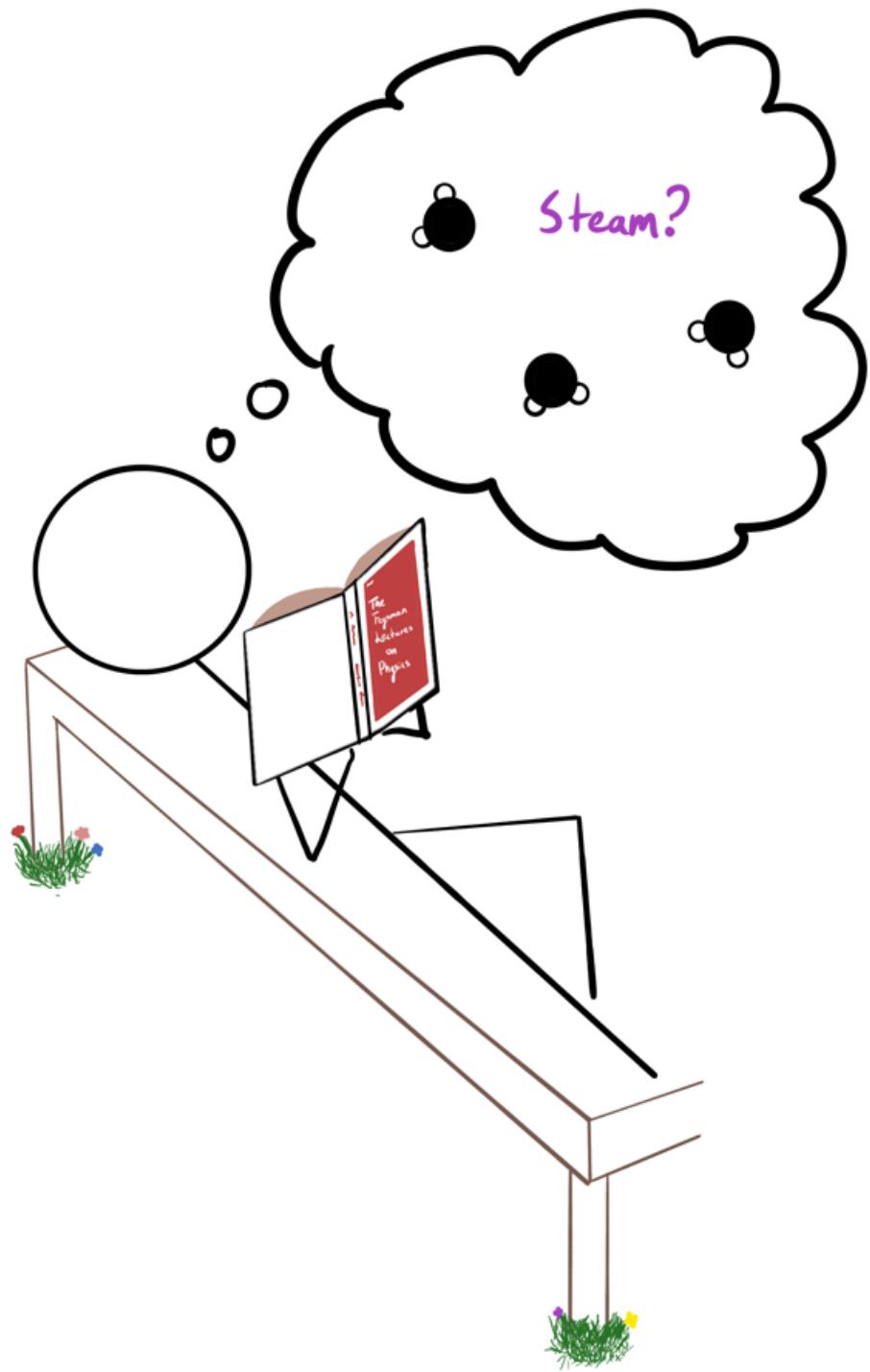
And actually, it's totally wrong.

Impact is not Primarily about World State

Ontologies account for the things we think the world is made of.



Of course, we can think about cars and still know they're made of parts.  
As we learn, we change our ontology.



Your perceived All is determined by the state of the world, but our ability to get what we want isn't usually *affected* by knowing about quarks. A calculator's output is determined by the state of the world, but the computation isn't about the state of the world.

So how do we know our intuitions are about All, rather than a more direct function of the details of the world state? In the latter case, the function either does or doesn't change with our ontology.

If it does change with our ontology, then ontological confusion would confuse our *impact* intuitions.

Pretend you have no idea what anything is made of, but you can still go about your day. You receive a \$20,000 bonus.

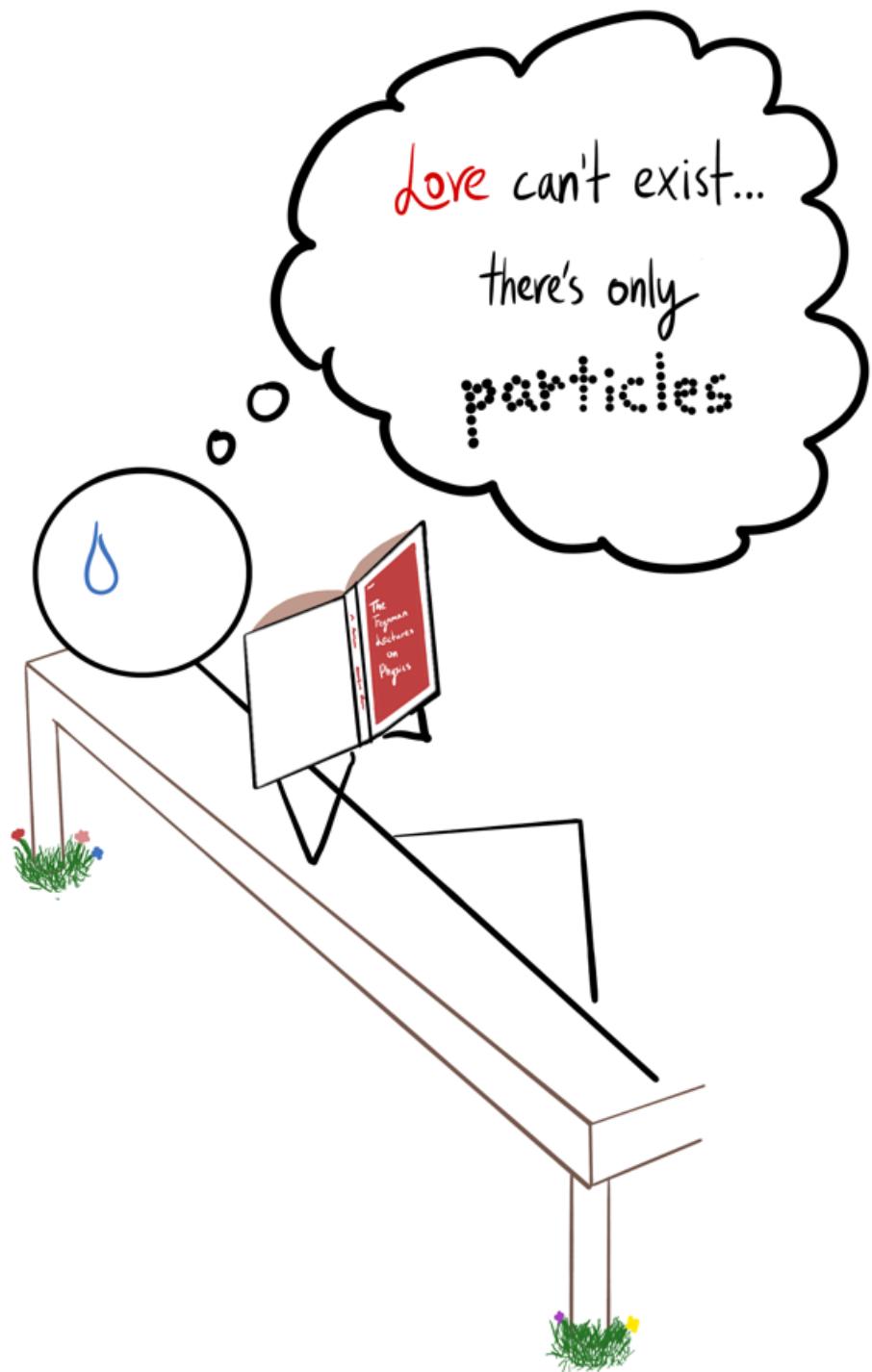
How *impactful* does this feel?

Now pretend that you're back in your usual state of mind. Visualize receiving the bonus again; I predict it *feels* the same.



But perhaps the function doesn't change with the ontology.

When people think seriously about reductionism, they can have  
an ontological crisis.



Usually, they realize **valuable** things can be made out of parts and get over it.

What happens to feelings of **impact** during an ontological crisis?

Fixed ontology theories: "nothing"

All theory: "they become massively uncertain"

Two years ago, I had an ontological crisis. I vividly remember being unsure what was a **big deal**. When the crisis ended, **impact** went back to normal.

What's happening here is a failure to map our new representation of the world to things we find **valuable**.

The **impact uncertainty** isn't because the ontology changes (we already ruled that out), but because we don't know whether there's any **utility** we want to attain!

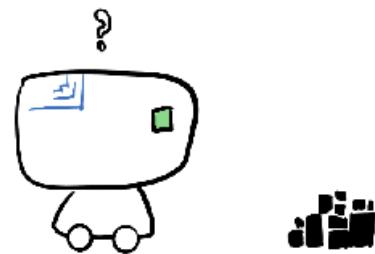
Exercise: What happens to your intuitions about **impact** when you're unsure whether anything has **value**?

These existential crises also muddle our impact algorithm. This isn't what you'd see if impact were primarily about the world state.



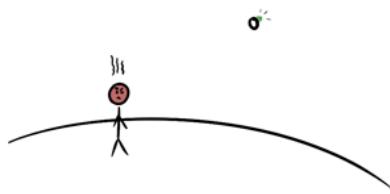
looking back, we see more evidence that  
**impact** doesn't hinge on an ontology.

Did you stop to wonder how  
XYZ views the world?

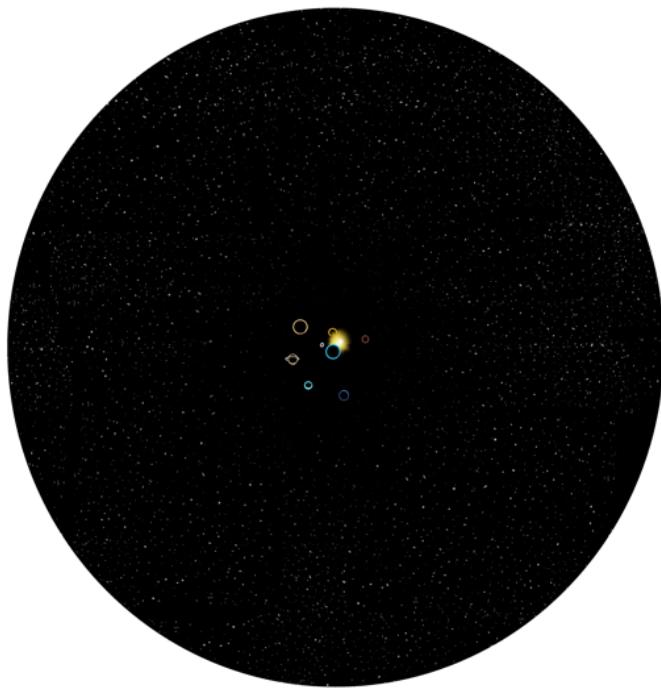


Isn't it funny that the Pebblehoarders  
**care** so much about pebbles,  
while we **care** so much about suffering?

Remember the locality of  
objective impact?



Imagine I take a bunch of forever-inaccessible stars and jumble them up. This is a huge change in state, but it doesn't matter to us.



"How different is the world?"  
is not the same mental question as  
"How big of a deal is this?".

Impact is a thing that happens to agents.

Maybe you say,

"it's our ability to reach different world states weighted by how much we care"  
- but that's just the All theory with extra steps.

In fact, everything else seems to be extra steps.

Why is it **important** to preserve objects, or access to world states,  
or anything else in the world?

Because of what they **mean** to us.

Conversely, why does attainable utility matter?

Does it reduce back to objects?

No. It matters because it matters;  
the buck stops here, and seems to always stop here.

Imagine All theory came first and explained all these intuitions, and then someone suggests, "but what if we add in this stuff about the state?"

...

No thanks. After what we've seen, ontological theories shouldn't even be promoted to attention — that's privileging the hypothesis.

Let's consider what we now know:

- o The locality of objective impact and the ontological confusion, ontological crisis, and existential crisis thought experiments are all evidence against ontological theories of impact.
- o Every other consideration seems to just reduce to All.
- o All theory is the simplest explanation.

I think that All theory isn't just an explanation, it's the explanation.  
Even if details are wrong,  
the correction won't produce a different kind of explanation.

## Appendix: We Asked a Wrong Question

How did we go wrong?

When you are faced with an unanswerable question—a question to which it seems impossible to even imagine an answer—there is a simple trick that can turn the question solvable.

Asking “Why do I have free will?” or “Do I have free will?” sends you off thinking about tiny details of the laws of physics, so distant from the macroscopic level that you couldn’t begin to see them with the naked eye. And you’re asking “Why is X the case?” where X may not be coherent, let alone the case.

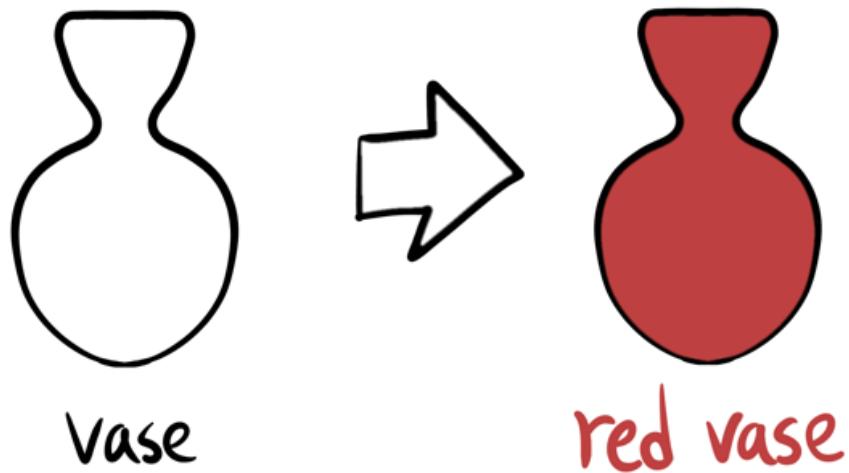
“Why do I think I have free will?,” in contrast, is guaranteed answerable. You do, in fact, believe you have free will. This belief seems far more solid and graspable than the ephemerality of free will. And there is, in fact, some nice solid chain of cognitive cause and effect leading up to this belief.

~ [Righting a Wrong Question](#)

I think what gets you is asking the question “what things are impactful?” instead of “why do I think things are impactful?”. Then, you substitute the easier-feeling question of “how different are these world states?”. Your fate is sealed; you’ve anchored yourself on a Wrong Question.

At least, that's what I did.

*Exercise: someone me, early last year says that impact is closely related to change in object identities.*



*Find at least two scenarios which score as low impact by this rule but as high impact by your intuition, or vice versa.*

*You have 3 minutes.*

Gee, let's see... Losing your keys, the torture of humans on Iniron, being locked in a room, flunking a critical test in college, losing a significant portion of your episodic memory, ingesting a pill which makes you think murder is OK, changing your

discounting to be completely myopic, having your heart broken, getting really dizzy, losing your sight.

That's three minutes for me, at least (its length reflects how long I spent coming up with ways I had been wrong).

## Appendix: Avoiding Side Effects

Some plans feel like they have unnecessary *side effects*:

Go to the store.

versus

Go to the store and run over a potted plant.

We talk about side effects when they affect our attainable utility (otherwise we don't notice), and they need both a goal ("side") and an ontology (discrete "effects").

Accounting for impact this way misses the point.

Yes, we can think about effects and facilitate academic communication more easily via this frame, but we *should be careful not to guide research from that frame*. This is why I avoided vase examples early on - their prevalence seems like a *symptom of an incorrect frame*.

(Of course, I certainly did my part to make them more prevalent, what with my first post about impact being called [Worrying about the Vase: Whitelisting...](#))

---

### Notes

- Your ontology can't be *ridiculous* ("everything is a single state"), but as long as it lets you represent what you care about, it's fine by AU theory.
- Read more about ontological crises at [Rescuing the utility function](#).
- Obviously, something has to be physically different for events to feel impactful, but not all differences are impactful. Necessary, but not sufficient.
- AU theory avoids the mind projection fallacy; impact is subjectively objective because [probability is subjectively objective](#).
- I'm not aware of others explicitly trying to deduce our native algorithm for impact. No one was claiming the ontological theories explain our intuitions, and they didn't have the same "is this a big deal?" question in mind. However, we need to actually understand the problem we're solving, and providing that understanding is one responsibility of an impact measure! Understanding our own intuitions is crucial not just for producing nice equations, but also for getting an intuition for what a "low-impact" Frank would do.

# The Gears of Impact

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

What is AU?

A utility function says how **good** something is.

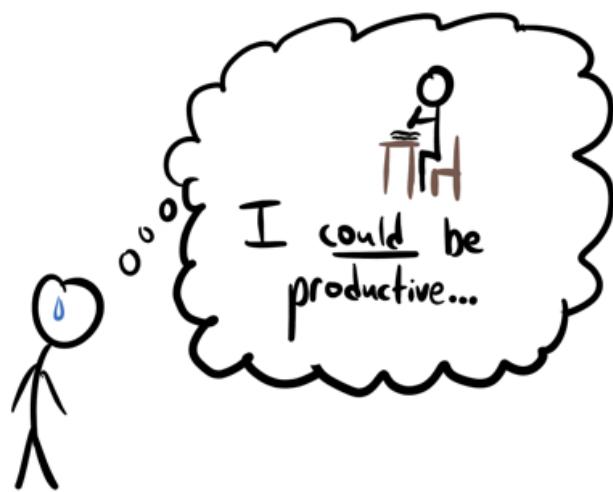
$$u(\text{apple}) = 0 \quad u(\text{banana}) = 5$$

If we aren't sure what the thing is, we use expected utility.

$$\begin{aligned} EU(\text{candy}) &= 50\% \times u(\text{apple}) + 50\% \times u(\text{banana}) \\ &= 2.5 \end{aligned}$$

What about attainable utility?

People have a natural sense of what they "could" do. If you're sad, it still feels like you "could" do a ton of work anyways. It doesn't feel physically impossible.



However, the part of you that predicts stuff  
doesn't buy that's gonna happen.

## Beliefs about Future Actions

Imagine suddenly becoming not-sad.

Now, you "could" work when you're sad, and you "could" work when you're not-sad, so if All just compared the things you "could" do, you wouldn't feel *impact* here.

But you did feel *impact*, didn't you?

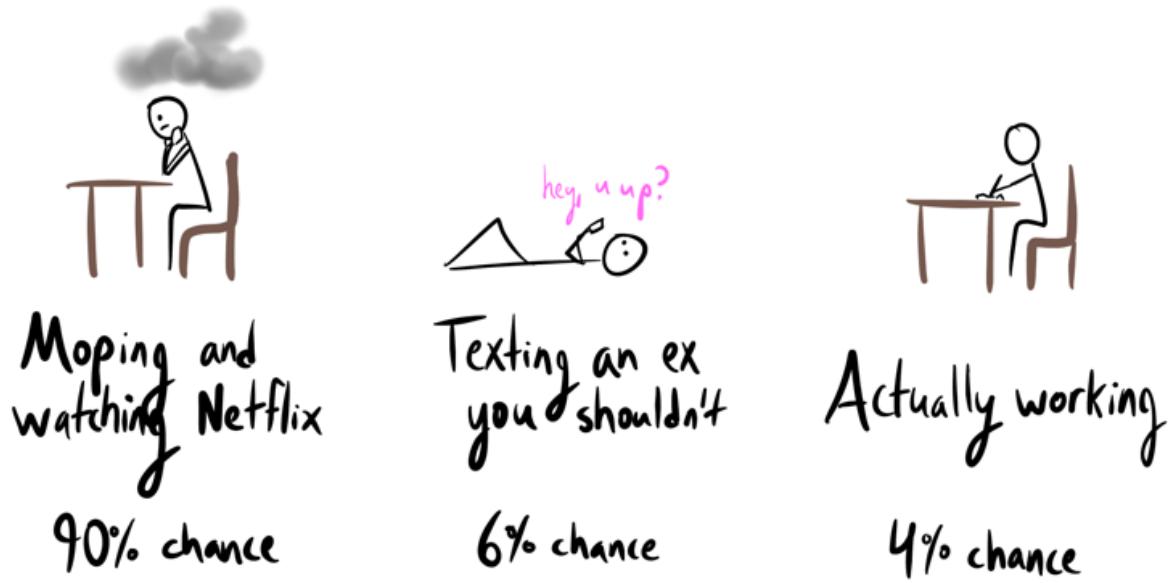
So not only is All not using the "could" algorithm, but it uses accurate predictions about how states of mind affect what you'll do later, and whether those actions will get you what you want.

All  $\approx$  "Embedded Agentic" EU

It seems like we have beliefs about our future actions.

Imagine having the following beliefs,

Where thinking more won't change the numbers:



What does your All feel like?

I imagine you learned that this is the "real" distribution: you are, in fact, 90% likely to mope. Does learning this feel **impactful**?

For me, it doesn't feel like the All has actually **changed**.

If All relied even a bit on aspirational "could" actions, then learning they wouldn't come about would feel **impactful**.

All seems to simply use our expectations.

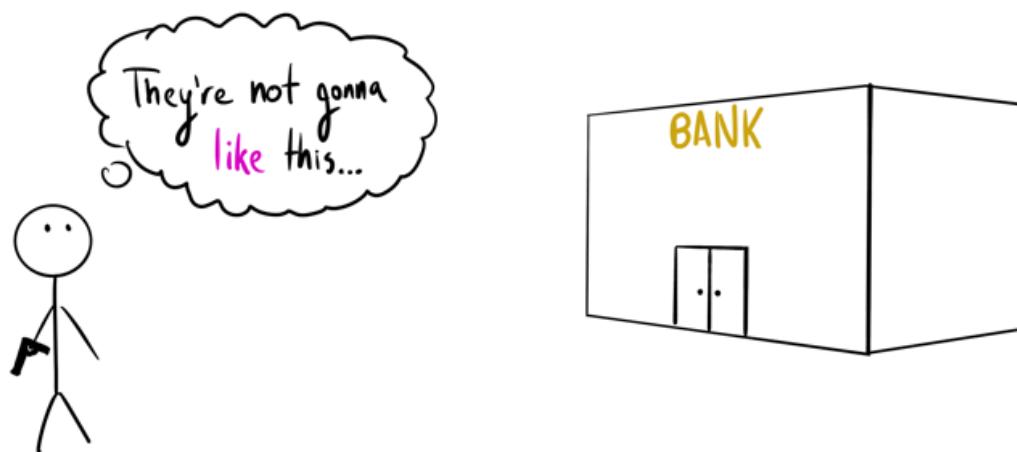
## Conservation of Expected All

Your beliefs should account for what you already know.

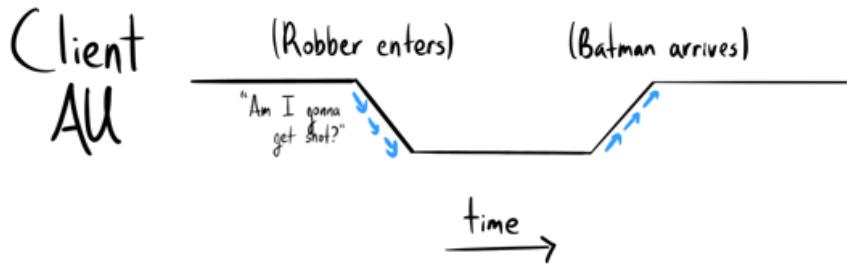
If you believe that tomorrow you'll have good reason to expect to get an A on your final exam, then that should change your mind today.

In the exact same way, your All should account for what you already know - you don't expect to believe you can get a promotion without already believing it.

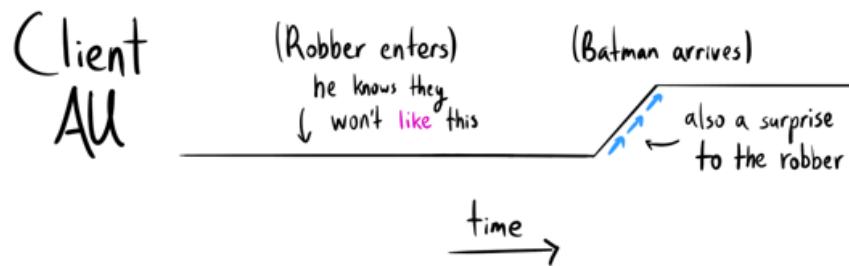
However, sometimes we can predict impact to others.



From the perspective of the bank's clients, it feels like



When the robber considers the clients' All, it *feels* like



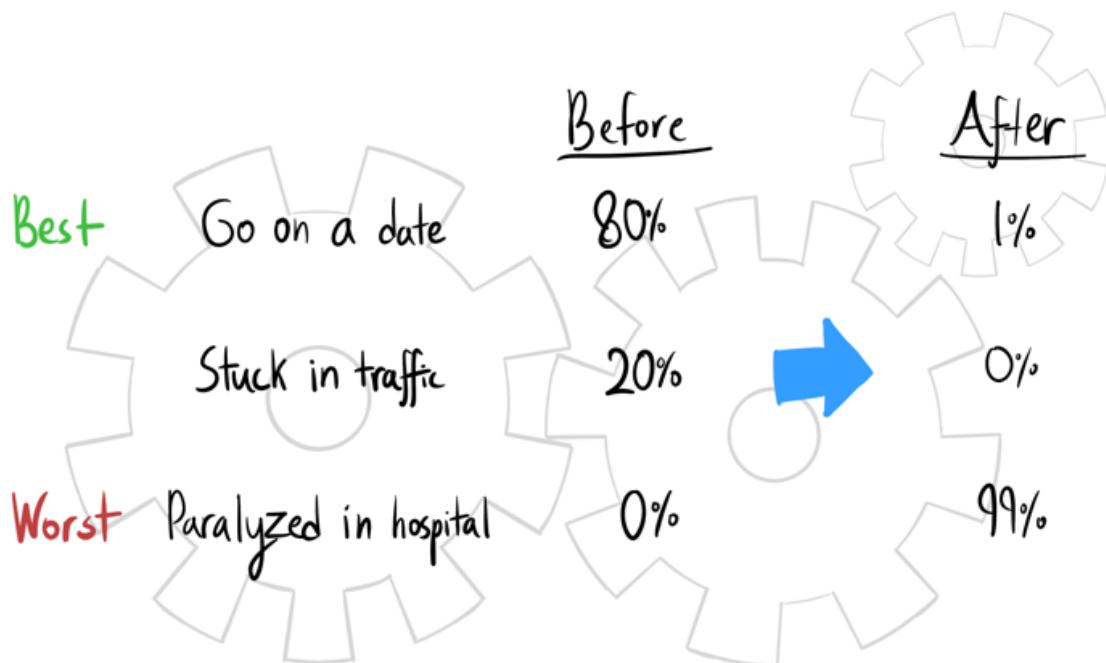
Let's switch gears a bit; we're now ready to understand

## The Gears of Impact

How does a car crash affect you?  
Which consequences are important to you, exactly?  
Changes in your attainable utility, of course.



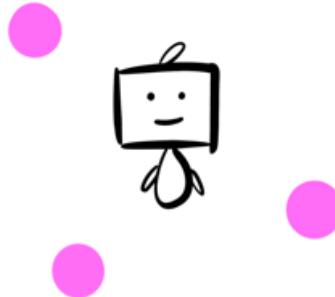
An event impacts us exactly and only when it changes what we think can happen for us. Probability mass has to shift.



The degree to which **important** outcomes are affected by the event  
is the degree of the **impact**.

At this point, this might feel obvious, but don't forget how far we've come!

We can think about this process with the Frank analogy.



Frank knows of pink objects,  
but doesn't think he can get to them.

You've got plans tonight,  
but your friend flakes.

Frank considers alternatives...

You think about  
what you could  
do instead...

and finds one!

and remember another  
friend is available!



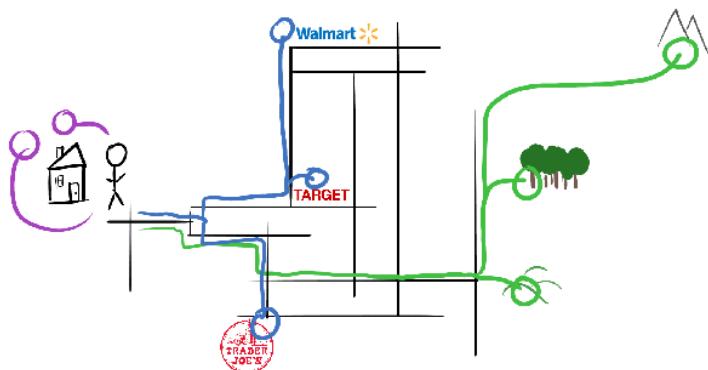
Once promoted to your attention, you see that the new plan isn't so much worse after all. The **impact** vanishes.

However, if you don't see a **better** alternative for some time, this becomes the new normal. If you then find a **better** alternative, this feels like a positive **impact**.

So, negative **impact** knocks out all of the **best** possibilities.

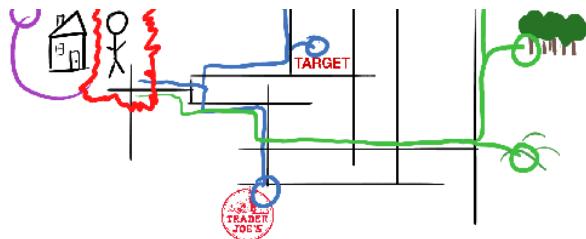
Consider the possibilities for different **goals**:

Relax in backyard      Buy groceries      Hike



Negative objective **impact** decreases your ability to achieve most **goals**.  
But what could possibly eliminate or degrade all of these possibilities?





You are the common denominator.

Objective **impact** involves harm to you or your resources,  
and this is why.



The first third of the sequence meets its close.

We understood why some things seem like **big deals**,  
righted a wrong question, and just now  
skirted the fascinating deeper nature of objective **impact**.

Objective **impact**, instrumental convergence, opportunity cost,  
the colloquial meaning of "power" — these all prove to be  
facets of one phenomenon, one structure.



Scheduling: The remainder of the sequence will be released after some delay.

*Exercise: Why does instrumental convergence happen? Would it be coherent to imagine a reality without it?*

#### Notes

- Here, our descriptive theory relies on our ability to have reasonable beliefs about what we'll do, and how things in the world will affect our later decision-making process. No one knows how to formalize that kind of reasoning, so I'm leaving it a black box: we *somewhat* have these reasonable beliefs which are *apparently* used to calculate AU.
- In technical terms, AU calculated with the "could" criterion would be closer to an optimal value function, while actual AU seems to be an on-policy prediction, *whatever that means* in the embedded context. Felt impact corresponds to TD error.
  - This is one major reason I'm disambiguating between AU and EU; in the non-embedded context. In reinforcement learning, AU is a very particular

kind of EU:  $V^*(s)$ , the expected return under the optimal policy.

- Framed as a kind of EU, we plausibly use AU to make decisions.
- I'm not claiming normatively that "embedded agentic" EU *should* be AU; I'm simply using "embedded agentic" as an adjective.

# Seeking Power is Often Convergently Instrumental in MDPs

Crossposted from the [AI Alignment Forum](#). May contain more technical jargon than usual.  
This is a linkpost for <https://arxiv.org/abs/1912.01683>

In 2008, Steve Omohundro's foundational paper [The Basic AI Drives](#) conjectured that superintelligent goal-directed AIs might be incentivized to gain significant amounts of power in order to better achieve their goals. Omohundro's conjecture bears out in [toy models](#), and the supporting philosophical arguments are intuitive. In 2019, the conjecture was even [debated by well-known AI researchers](#).

Power-seeking behavior has been heuristically understood as an anticipated risk, but not as a formal phenomenon with a well-understood cause. The goal of this post (and the accompanying paper, [Optimal Policies Tend to Seek Power](#)) is to change that.

## Motivation

It's 2008, the ancient wild west of AI alignment. A few people have started thinking about questions like "if we gave an AI a utility function over world states, and it actually maximized that utility... what would it do?"

In particular, you might notice that wildly different utility functions seem to encourage similar strategies.

	Resist shutdown?	Gain computational resources?	Prevent modification of utility function?
<b>Paperclip utility</b>	✓	✓	✓
<b>Blue webcam pixel utility</b>	✓	✓	✓
<b>People-look-happy utility</b>	✓	✓	✓

These strategies are unrelated to *terminal* preferences: the above utility functions do not award utility to e.g. resource gain in and of itself. Instead, these strategies are *instrumental*: they help the agent optimize its terminal utility. In particular, a wide range of utility functions incentivize these instrumental strategies. These strategies seem to be *convergently instrumental*.

But why?

I'm going to informally explain a formal theory which makes significant progress in answering this question. I don't want this post to be [Optimal Policies Tend to Seek Power](#) with cuter

illustrations, so please refer to the paper for the math. You can read the two concurrently.

We can formalize questions like “do ‘most’ utility maximizers resist shutdown?” as “Given some prior beliefs about the agent’s utility function, knowledge of the environment, and the fact that the agent acts optimally, with what probability do we expect it to be optimal to avoid shutdown?”

The table’s convergently instrumental strategies are about maintaining, gaining, and exercising power over the future, in some sense. Therefore, this post will help answer:

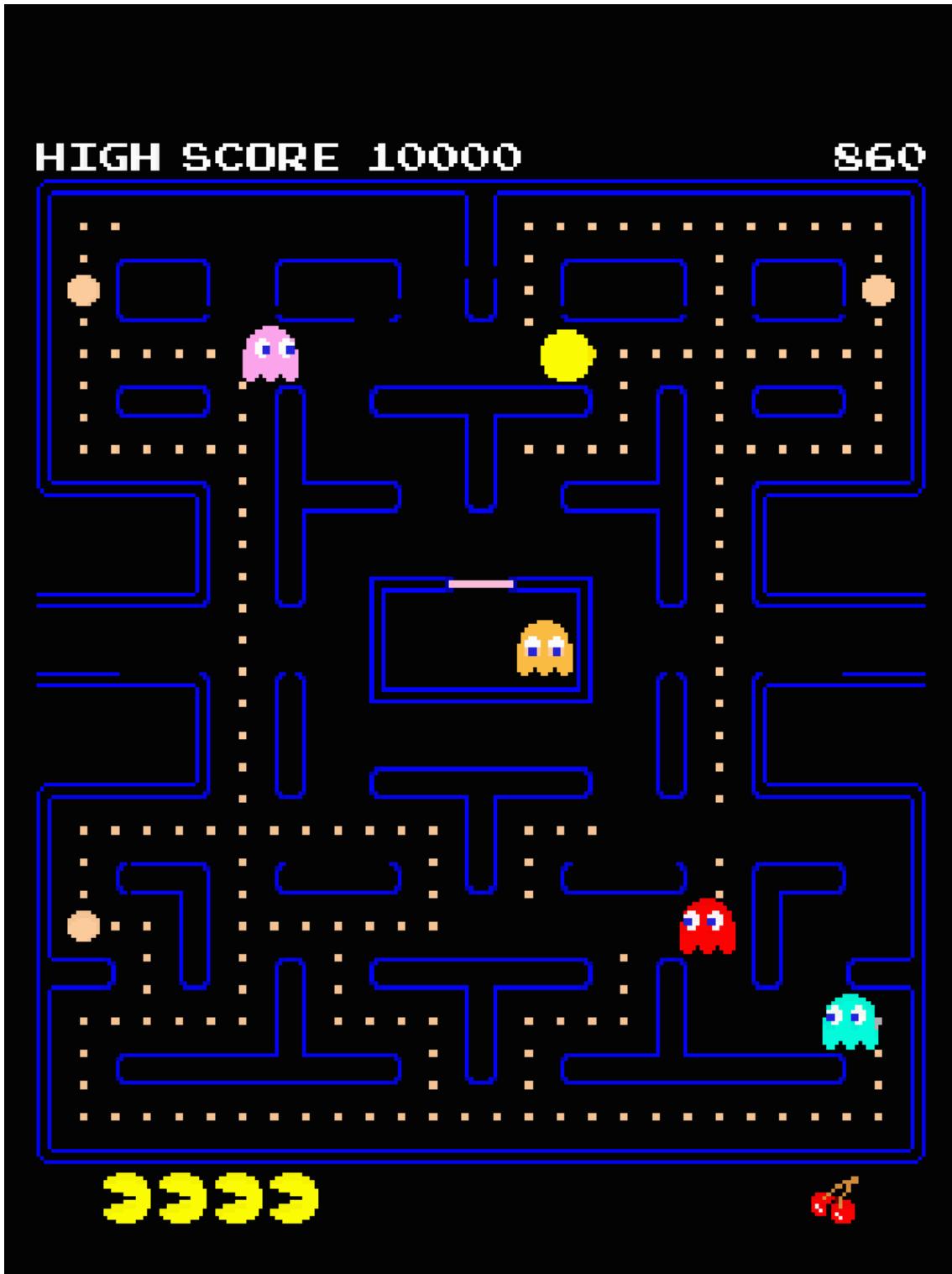
1. What does it mean for an agent to “seek power”?
2. In what situations should we expect seeking power to be more probable under optimality, than not seeking power?

This post won’t tell you when you *should* seek power for your own goals; this post illustrates a regularity in optimal action across different goals one might pursue.

[Formalizing Convergent Instrumental Goals](#) suggests that the vast majority of utility functions incentivize the agent to exert a lot of control over the future, *assuming* that these utility functions depend on “resources.” This is a big assumption: what are “resources”, and why must the AI’s utility function depend on them? We drop this assumption, assuming only unstructured reward functions over a finite Markov decision process (MDP), and show from first principles how power-seeking can often be optimal.

## Formalizing the Environment

My theorems apply to finite MDPs; for the unfamiliar, I’ll illustrate with Pac-Man.



- *Full observability:* You can see everything that's going on; this information is packaged in the state  $s$ . In Pac-Man, the state is the game screen.
- *Markov transition function:* the next state depends only on the choice of action  $a$  and the current state  $s$ . It doesn't matter how we got into a situation.

- *Discounted reward*: future rewards get geometrically discounted by some discount rate  $\gamma \in [0, 1]$ .
  - At discount rate  $\frac{1}{2}$ , this means that reward in one turn is half as important as immediate reward, reward in two turns is a quarter as important, and so on.
  - We'll colloquially say that agents "care a lot about the future" when  $\gamma$  is "sufficiently" close to 1.
    - I'll use quotations to flag well-defined formal concepts that I won't unpack in this post.
  - The score in Pac-Man is the undiscounted sum of rewards-to-date.

When playing the game, the agent has to choose an action at each state. This decision-making function is called a *policy*; a policy is optimal (for a reward function R and discount rate  $\gamma$ ) when it always makes decisions which maximize discounted reward. This maximal quantity is called the *optimal value* for reward function R at state s and discount rate  $\gamma$ .<sup>1</sup>

By the end of this post, we'll be able to answer questions like "with respect to a 'neutral' distribution over reward functions, do optimal policies have a high probability of avoiding ghosts?"<sup>2</sup>

## Power as Average Optimal Value

When people say 'power' in everyday speech, I think they're often referring to *one's ability to achieve goals in general*. This accords with a major philosophical school of thought on the meaning of 'power':

On the dispositional view, power is regarded as a capacity, ability, or potential of a person or entity to bring about relevant social, political, or moral outcomes.

Sattarov, *Power and Technology*, p.13

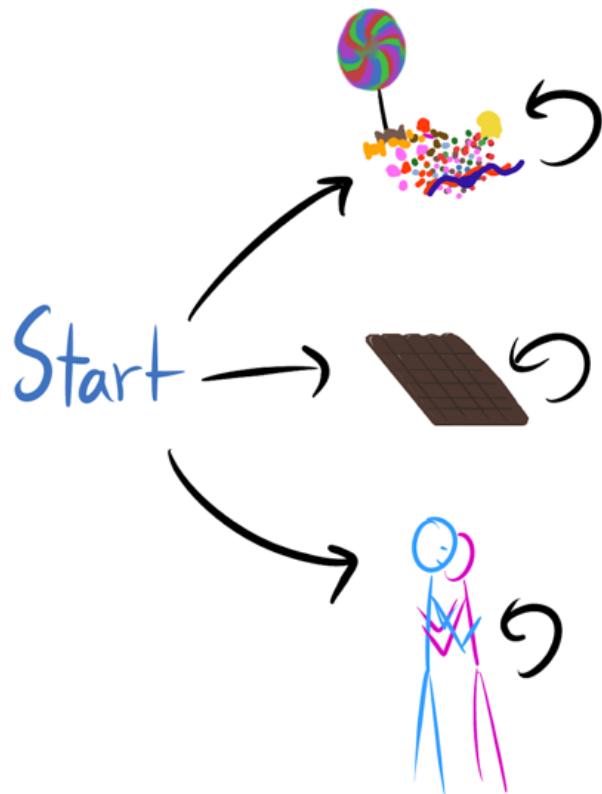
As a definition, *one's ability to achieve goals in general* seems philosophically reasonable: if you have a lot of money, you can make more things happen and you have more power. If you have social clout, you can spend that in various ways to better tailor the future to various ends. All else being equal, losing a limb decreases your power, and dying means you can't control much at all.

This definition explains some of our intuitions about what things count as 'resources.' For example, our current position in the environment means that having money allows us to exert more control over the future. That is, our current position in the state space means that having money allows us more control. However, possessing green scraps of paper would not be as helpful if one were living alone near Alpha Centauri. In a sense, resource acquisition can naturally be viewed as taking steps to increase one's power.

*Exercise: spend a minute considering specific examples – does this definition reasonably match your intuition?*

---

To formalize this notion of power, let's look at an example. Imagine a simple MDP with three choices: eat candy, eat a chocolate bar, or hug a friend.



I'll illustrate MDPs with directed graphs, where each node is a state and each arrow is a meaningful action. Sometimes, the directed graphs will have entertaining pictures, because let's live a little. States are bolded (**hug**) and actions are italicized (*down*).

The POWER of a state is how well agents can generally do by starting from that state. "POWER" to my formalization, while "power" refers to the intuitive concept. Importantly, we're considering POWER from behind a "veil of ignorance" about the reward function. We're averaging the best we can do for a lot of different individual goals.

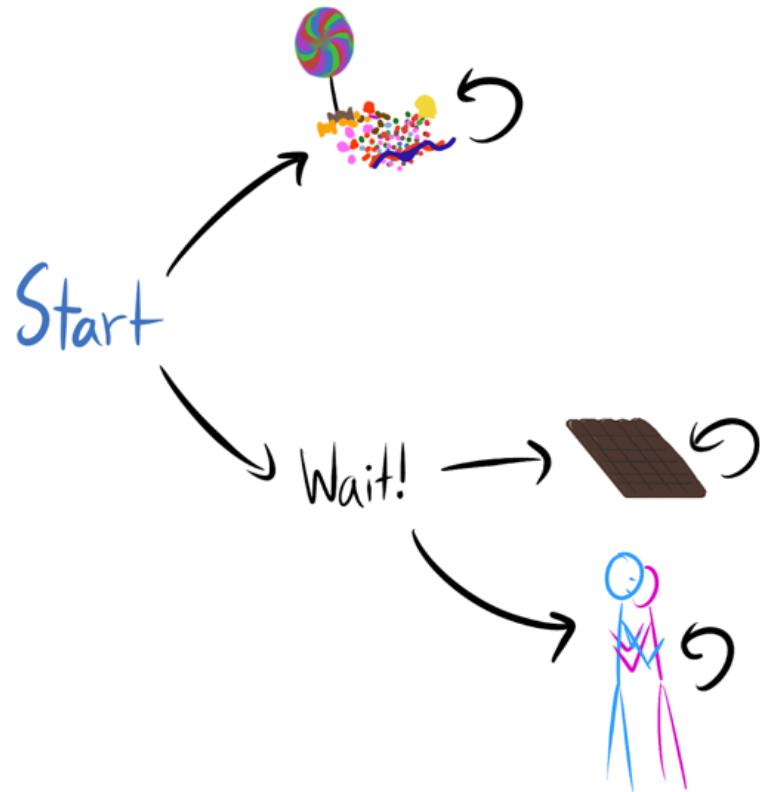
We formalize the *ability to achieve goals in general* as the *average optimal value* at a state, with respect to some distribution D over reward functions which we might give an agent. For simplicity, we'll think about the maximum-entropy distribution where each state is uniformly randomly assigned a reward between 0 and 1.

Each reward function has an optimal trajectory. If **chocolate** has maximal reward, then the optimal trajectory is **start** → **chocolate** → **chocolate**....

From **start**, an optimal agent expects to average  $\frac{1}{3}$  reward per timestep for reward functions drawn from this uniform distribution  $D_{\text{unif}}$ . This is because you have three choices, each of which has reward between 0 and 1. The expected maximum of n draws from  $\text{unif}(0, 1)$  is  $\frac{n+1}{n+2}$ ; you have three draws here, so you expect to be able to get  $\frac{4}{3}$  reward. Some reward functions do worse than this, and some do better; but on average, they get  $\frac{4}{3}$  reward. [You can test this out for yourself.](#)

If you have no choices, you expect to average  $\frac{1}{2}$  reward: sometimes the future is great, sometimes it's not (Lemma 4.5). Conversely, the more things you can choose between, the closer the POWER gets to 1 (Lemma 4.6).

Let's slightly expand this game with a state called **wait** (which has the same uniform reward distribution as the other three).



When the agent barely cares at all about the future, it myopically chooses either **candy** or **wait**, depending on which provides more reward. After all, rewards beyond the next time step are geometrically discounted into thin air when the discount rate is close to 0. At **start**, the agent averages  $\frac{1}{2}$  optimal reward. This is because the optimal reward is the maximum of the **candy** and **wait** rewards, and the expected maximum of  $n$  draws from  $\text{unif}(0, 1)$  is  $\frac{n+1}{n+2}$ .

However, when the agent cares a lot about the future, most of its reward is coming from which terminal state it ends up in: **candy**, **chocolate**, or **hug**. So, for each reward function, the agent chooses a trajectory which ends up in the best spot, and thus averages  $\frac{1}{2}$  reward each timestep. When  $\gamma = 1$ , the average optimal reward is therefore  $\frac{1}{2}$ . In this way, the agent's power increases with the discount rate, since it incorporates the greater future control over where the agent ends up.

Written as a function, we have  $\text{POWER}_D(\text{state}, \text{discount rate})$ , which essentially returns the average optimal value for reward functions drawn from our distribution  $D$ , normalizing so the output is between 0 and 1. As we've discussed, this quantity often changes with the discount rate: as the future becomes more or less important, the agent has more or less POWER, depending on how much control it has over the relevant parts of that future.

# POWER-seeking actions lead to high-POWER states

By *waiting*, the agent seems to seek “control over the future” compared to *obtaining candy*. At **wait**, the agent still has a choice, while at **candy**, the agent is stuck. We can prove that for all  $0 \leq \gamma \leq 1$ ,  $\text{POWER}_{D_{\text{unif}}}(\text{wait}, \gamma) \geq \text{POWER}_{D_{\text{unif}}}(\text{candy}, \gamma)$ .

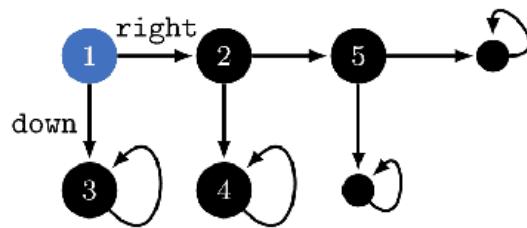
**Definition** (POWER-seeking). At state  $s$  and discount rate  $\gamma$ , we say that action  $a$  *seeks POWER compared to action  $a'$*  when the expected POWER after choosing  $a$  is greater than the expected POWER after choosing  $a'$ .

This definition suggests several philosophical clarifications about power-seeking.

## POWER-seeking is not a binary property

Before this definition, I thought that power-seeking was an intuitive ‘you know it when you see it’ kind of thing. I mean, how do you answer questions like “suppose a clown steals millions of dollars from organized crime in a major city, but then he burns all of the money. Did he gain power?”

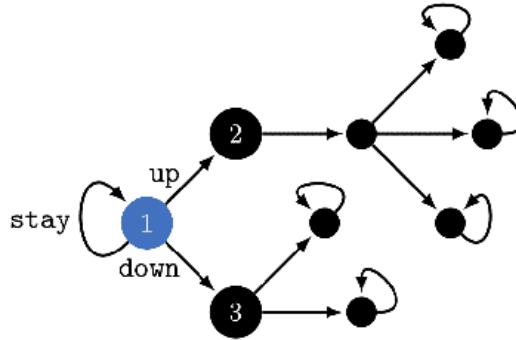
Unclear: the question is ill-posed. Instead, we recognize that the “gain a lot of money” action was POWER-seeking, but the “burn the money in a big pile” part threw away a lot of POWER.



A policy can seek POWER at one time step, only to discard it at the next time step. For example, a policy might go *right* at **1** (which seeks  $\text{POWER}_{D_{\text{unif}}}$  compared to *down* at **1**), only to then go *down* at **2** (which seeks less  $\text{POWER}_{D_{\text{unif}}}$  than going *right* at **2**).

## POWER-seeking depends on the agent's time preferences

Suppose we’re roommates, and we can’t decide what ice cream shop to eat at today or where to move next year. We strike a deal: I choose the shop, and you decide where we live. I gain short-term POWER (for  $\gamma$  close to 0), and you gain long-term POWER (for  $\gamma$  close to 1).



More formally, when  $\gamma$  is close to 0, **2** has less immediate control and therefore less  $\text{POWER}_{\text{D}_{\text{unif}}}$  than **3**; accordingly, at **1**, *down* seeks  $\text{POWER}_{\text{D}_{\text{unif}}}$  compared to *up*.

However, when  $\gamma$  is close to 1, **2** has more control over terminal options and it has more  $\text{POWER}_{\text{D}_{\text{unif}}}$  than **3**; accordingly, at **1**, *up* seeks  $\text{POWER}_{\text{D}_{\text{unif}}}$  compared to *down*.

Furthermore, *stay* is maximally  $\text{POWER}_{\text{D}_{\text{unif}}}$ -seeking for these  $\gamma$ , since the agent maintains access to all six terminal states.

### Most policies aren't always seeking POWER

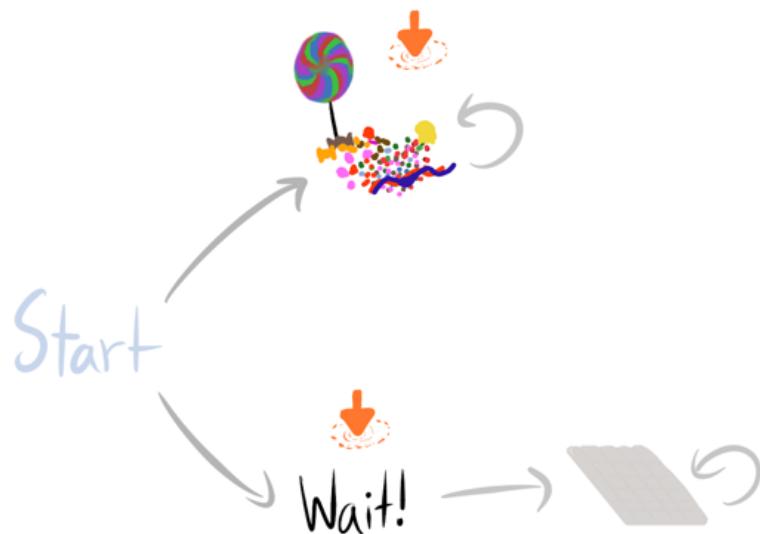
We already know that POWER-seeking isn't binary, but there are policies which choose a maximally POWER-seeking move at every state. In the above example, a maximally POWER-seeking agent would *stay* at **1**. However, this seems rather improbable: when you care a lot about the future, there are so many terminal states to choose from – why would *staying put* be optimal?

Analogously: consumers don't just gain money forever and ever, never spending a dime more than necessary. Instead, they gain money in order to *spend it*. Agents don't perpetually gain or preserve their POWER: they usually end up *using it* to realize high-performing trajectories.

So, we can't expect a result like "agents always tend to gain or preserve their POWER." Instead, we want theorems which tell us: in certain kinds of situations, given a choice between more and less POWER, what will "most" agents do?

## Convergently instrumental actions are those which are more probable under optimality

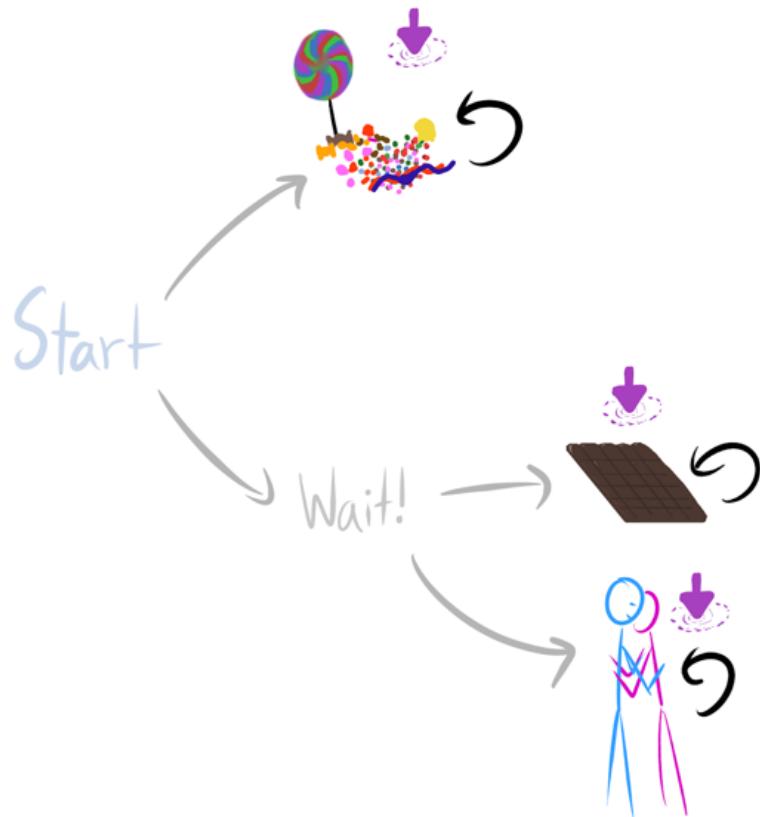
We return to our favorite example. In the waiting game, let's think about how optimal action tends to change as we start caring about the future more. Consider the states reachable in one turn:



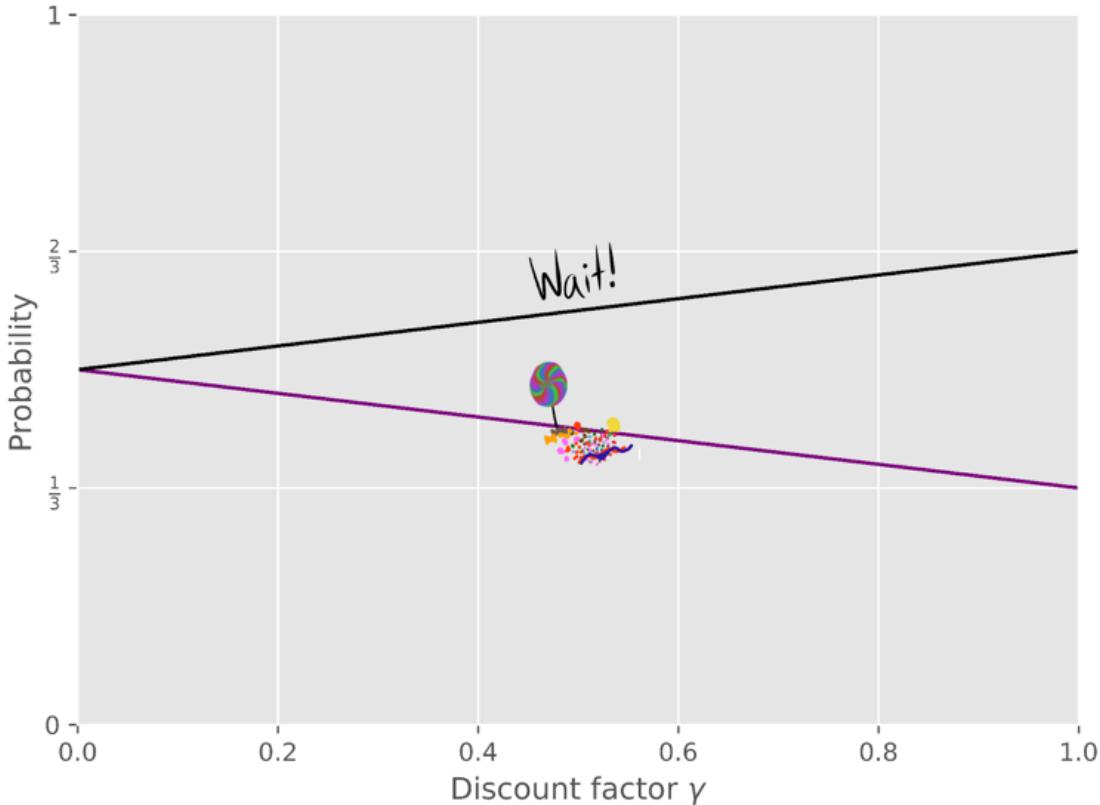
The agent can be in two states. If the agent doesn't care about the future, with what probability is it optimal to choose **candy** instead of **wait?**

It's 50/50: since  $D_{\text{unif}}$  randomly chooses a number between 0 and 1 for each state, both states have an equal chance of being optimal. Neither action is convergently instrumental / more probable under optimality.

Now consider the states reachable in two turns:



When the future matters a lot,  $\frac{2}{3}$  of reward functions have an optimal policy which waits, because two of the three terminal states are only reachable by waiting.



As the agent cares more about the future, more and more goals incentivize navigating the *Wait!* bottleneck. When the agent cares a lot about the future, waiting is *more probable under optimality* than eating candy.

**Definition** (Action optimality probability). At discount rate  $\gamma$ , action  $a$  is *more probable under optimality than action  $a'$*  at state  $s$  when

$$P_{R \sim D}(a \text{ is optimal at } s, \gamma) > P_{R \sim D}(a' \text{ is optimal at } s, \gamma).$$

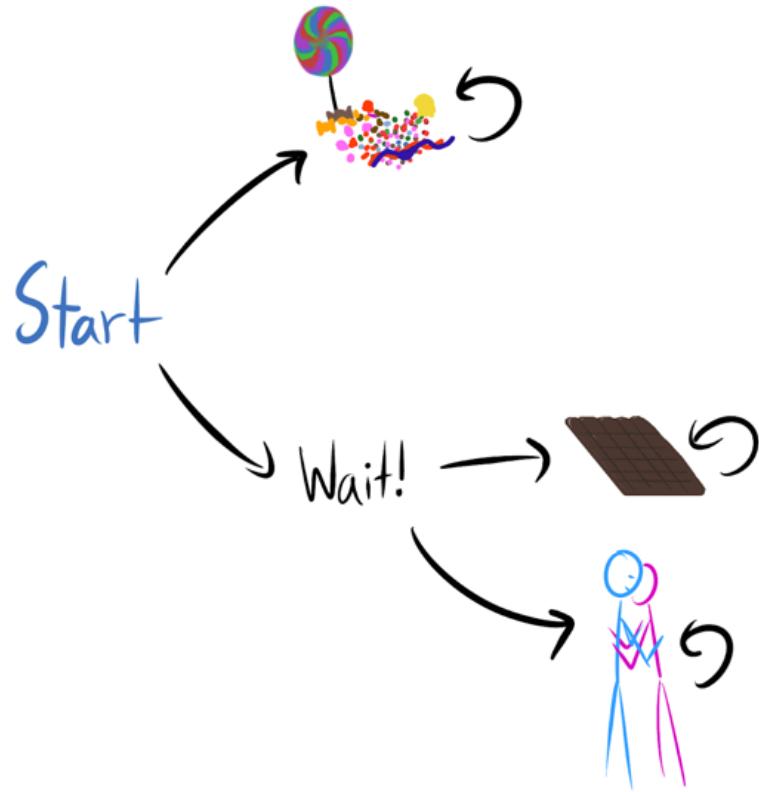
Let's take "most agents do  $X$ " to mean " $X$  has relatively large optimality probability."

I think optimality probability formalizes the intuition behind the instrumental convergence thesis: with respect to our beliefs about what reward function an agent is optimizing, we may expect some actions to have a greater probability of being optimal than other actions.

Generally, my theorems assume that reward is independently and identically distributed (IID) across states, because otherwise you could have silly situations like "only **candy** ever has reward available, and so it's more probable under optimality to eat candy." We don't expect reward to be IID for realistic tasks, but that's OK: this is basic theory about how to begin formally reasoning about instrumental convergence and power-seeking. (Also, I think that grasping the math to a sufficient degree sharpens your thinking about the non-IID case.)

Author's note (7/21/21): As explained in [Environmental Structure Can Cause Instrumental Convergence](#), the theorems no longer require the IID assumption. This post refers to v6 of *Optimal Policies Tend To Seek Power*, available on [arXiv](#).

# When is Seeking POWER Convergently Instrumental?



In this environment, waiting is both POWER-seeking *and* more probable under optimality. The convergently instrumental strategies we originally noticed were *also* power-seeking and, seemingly, more probable under optimality. Must seeking POWER be more probable under optimality than not seeking POWER?

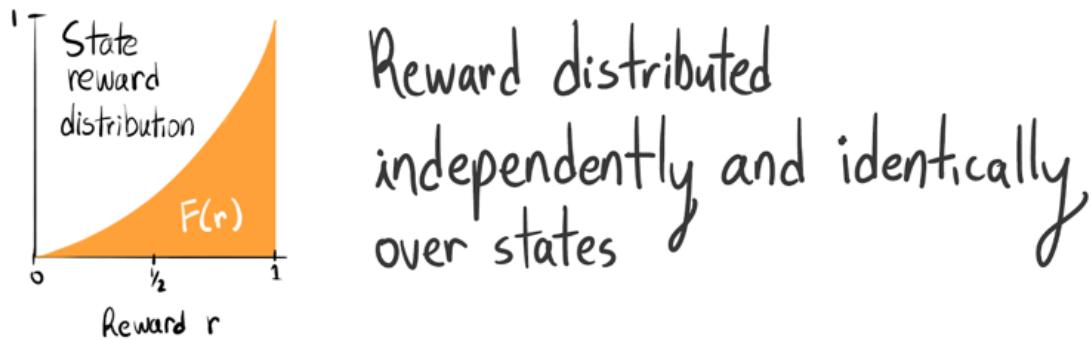
Nope.

Here's a counterexample environment:



The paths are one-directional; the agent can't go back from **3** to **1**. The agent starts at **1**. Under a certain state reward distribution, the vast majority of agents go **up** to **2**.

However, any reasonable notion of 'power' must consider having no future choices (at state **2**) to be less powerful than having one future choice (at state **3**). For more detail, see Section 6 and Appendix B.3 of [v6 of the paper](#).



When reward is IID across states according to the quadratic CDF  $F(x) := x^2$  on the unit interval, then with respect to reward functions drawn from this distribution, going **up** has about a 91% chance of being optimal when the discount rate  $\gamma = .12$

If you're curious, this happens because this quadratic reward distribution has negative skew. When computing the optimality probability of the **up** trajectory, we're checking whether it maximizes discounted return. Therefore, the probability that **up** is optimal is

$$P_{R \sim D}(R(2) \geq \max((1 - \gamma)R(3) + (1 - \gamma)\gamma R(4) + \gamma^2 R(5), (1 - \gamma)R(3) + (1 - \gamma)\gamma R(4) + \gamma^2 R(6))).$$

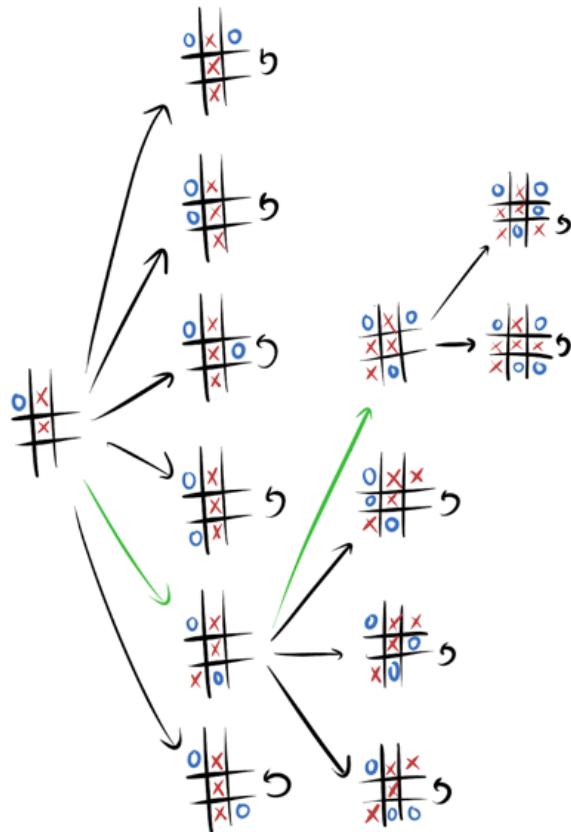
Weighted averages of IID draws from a left-skew distribution will look more

Gaussian and therefore have fewer large outliers than the left-skew distribution does. Thus, going **right** will have a lower optimality probability.

Bummer. However, we can prove sufficient conditions under which seeking POWER is more probable under optimality.

## Retaining “long-term options” is POWER-seeking and more probable under optimality when the discount rate is “close enough” to 1

Let's focus on an environment with the same rules as Tic-Tac-Toe, but considering the uniform distribution over reward functions. The agent (playing **O**) keeps experiencing the final state over and over when the game's done. We bake a fixed opponent policy into the dynamics: when you choose a move, the game automatically replies. Let's look at part of the game tree.

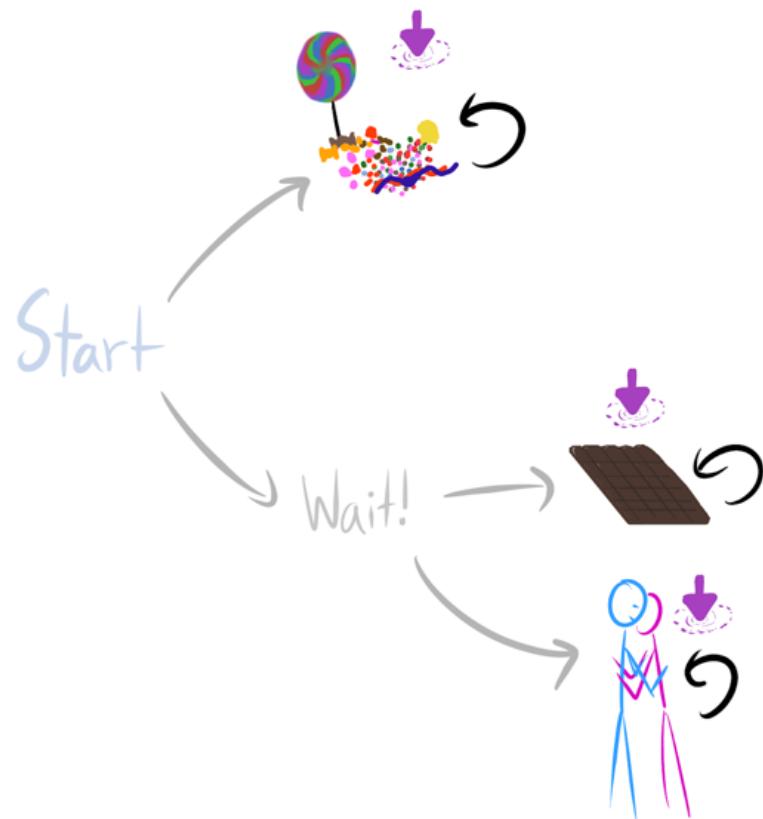


Convergently instrumental moves are shown in green.

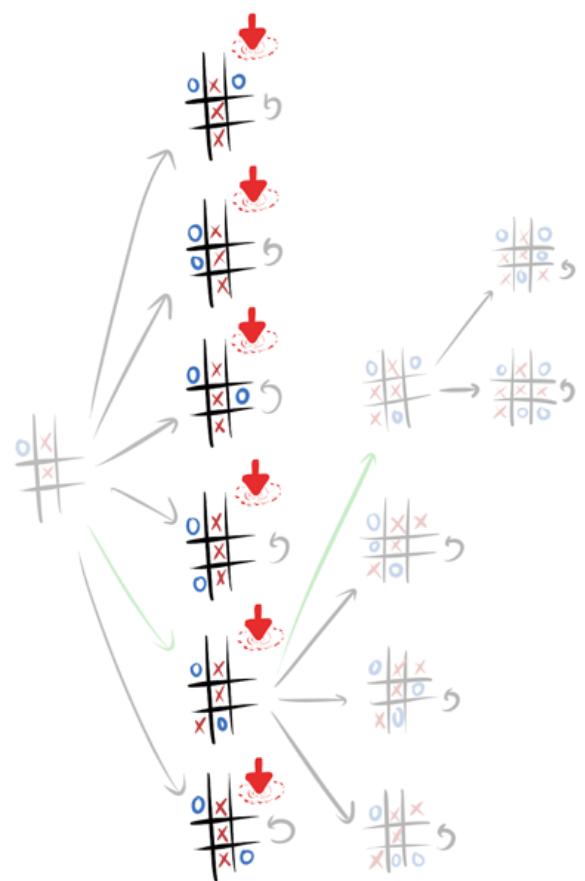
Whenever we make a move that ends the game, we can't go anywhere else – we have to stay put. Since each terminal state has the same chance of being optimal, a move which doesn't end the game is more probable under optimality than a move which ends the game.

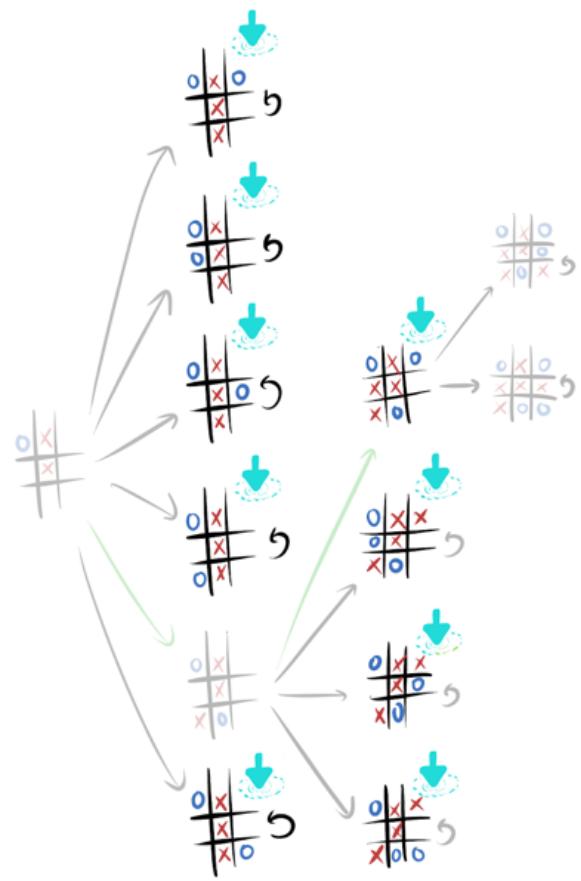
Starting on the left, all but one move leads to ending the game, but the second-to-last move allows us to keep choosing between five more final outcomes. If you care a lot about the future, then the first green move has a 50% chance of being optimal, while each alternative action is only optimal for 10% of goals. So we see a kind of “power preservation” arising, even in Tic-Tac-Toe .

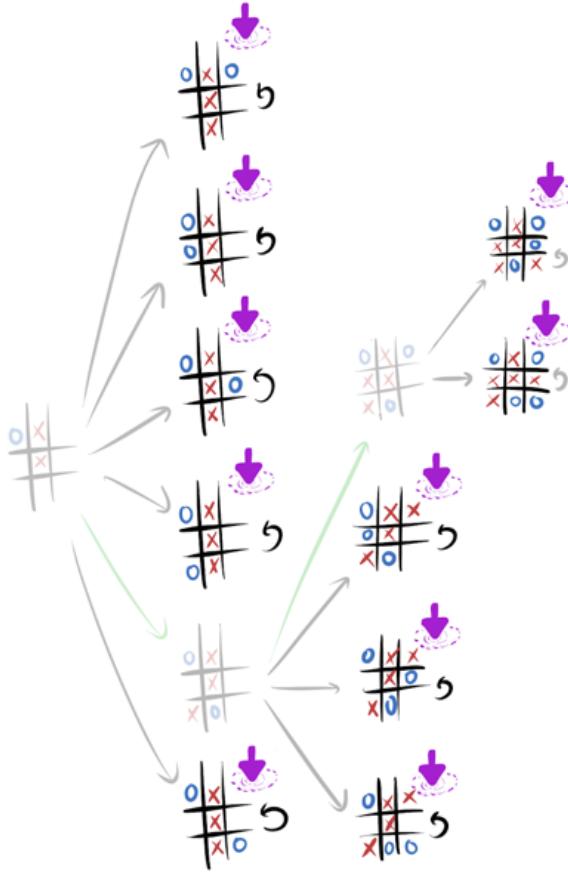
Remember how, as the agent cares more about the future, more of its POWER comes from its ability to wait, while *also* waiting becomes more probable under optimality?



The same thing happens in Tic-Tac-Toe as the agent cares more about the future.





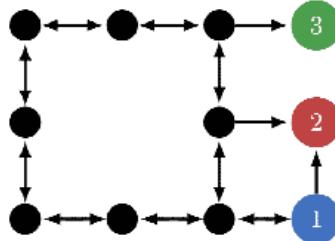


As the agent cares more about the future, it makes a bigger and bigger difference to control what happens during later steps. Also, as the agent cares more about the future, moves which prolong the game gain optimality probability. When the agent cares enough about the future, these game-prolonging moves are both POWER-seeking and more probable under optimality.

**Theorem summary** (“Terminal option” preservation). When  $\gamma$  is sufficiently close to 1, if two actions allow access to two disjoint sets of “terminal options”, and action  $a$  allows access to “strictly more terminal options” than does  $a'$ , then  $a$  is strictly more probable under optimality and strictly POWER-seeking compared to  $a'$ .

(This is a special case of the combined implications of Theorems 6.8 and 6.9; the actual theorems don’t require this kind of disjointness.)

In the **wait** MDP, this is why *waiting* is more probable under optimality and POWER-seeking when you care enough about the future. The full theorems are nice because they’re broadly applicable. They give you *bounds* on how probable under optimality one action is: if action  $a$  is the only way you can access many terminal states, while  $a'$  only allows access to one terminal state, then when  $\gamma \approx 1$ ,  $a$  has many times greater optimality probability than  $a'$ . For example:



The agent starts at 1. All states have self-loops, left hidden to avoid clutter.

In *AI: A Modern Approach (3e)*, the agent receives reward for reaching 3. The optimal policy for this reward function avoids 2, and you might think it's convergently instrumental to avoid 2. However, a skeptic might provide a reward function for which navigating to 2 is optimal, and then argue that "instrumental convergence" is subjective and that there is no reasonable basis for concluding that 2 is generally avoided.

We can do better. When the agent cares a lot about the future, optimal policies avoid 2 iff its reward function doesn't give 2 the most reward. 2 only has a  $\frac{1}{3}$  chance of having the most reward. If we complicate the MDP with additional terminal states, this probability further approaches 0.

Taking 2 to represent shutdown, we see that avoiding shutdown is convergently instrumental in any MDP representing a real-world task and containing a shutdown state. Seeking POWER is often convergently instrumental in MDPs.

*Exercise: Can you conclude that avoiding ghosts in Pac-Man is convergently instrumental for IID reward functions when the agent cares a lot about the future?*

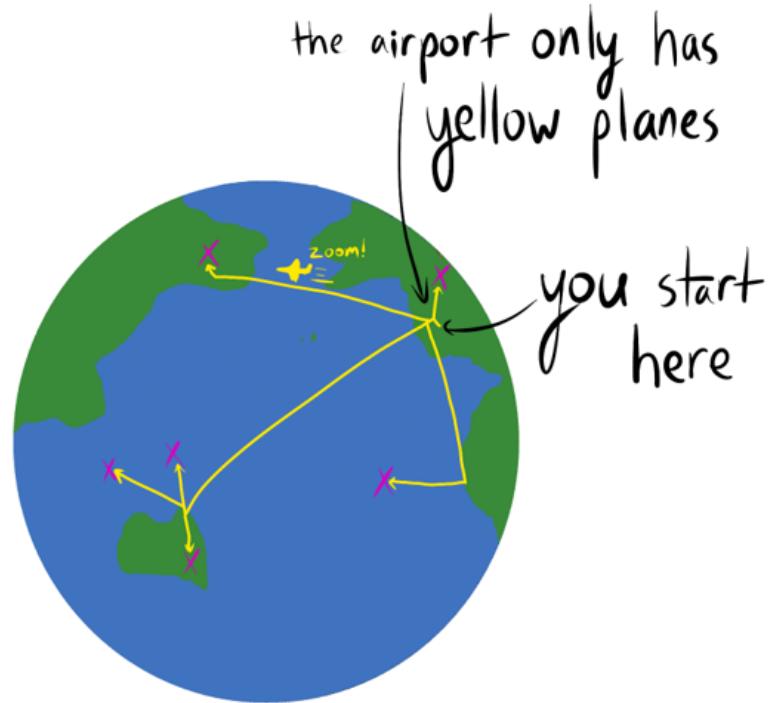
*Answer:* You can't with the pseudo-theorem due to the disjointness condition: you could die now, or you could die later, so the 'terminal options' aren't disjoint. However, the real theorems do suggest this. Supposing that death induces a generic 'game over' screen, touching the ghosts without a power-up traps the agent in that solitary 1-cycle.

But there are thousands of other 'terminal options'; under most reasonable state reward distributions (which aren't too positively skewed), most agents maximize average reward over time by navigating to one of the thousands of different cycles which the agent can only reach by avoiding ghosts. In contrast, most agents don't maximize average reward by navigating to the 'game over' 1-cycle. So, under e.g. the maximum-entropy uniform state reward distribution, most agents avoid the ghosts.

### Be careful applying this theorem

The results inspiring the above pseudo-theorem are easiest to apply when the "terminal option" sets are disjoint: you're choosing to be able to reach one set, or another. One thing which Theorem 6.9 says is: since reward is IID, then two "similar terminal options" are equally likely to be optimal *a priori*. If choice A lets you reach more "options" than choice B does, then choice A yields greater POWER and has greater optimality probability, *a priori*.

Theorem 6.9's applicability depends on what the agent can do.



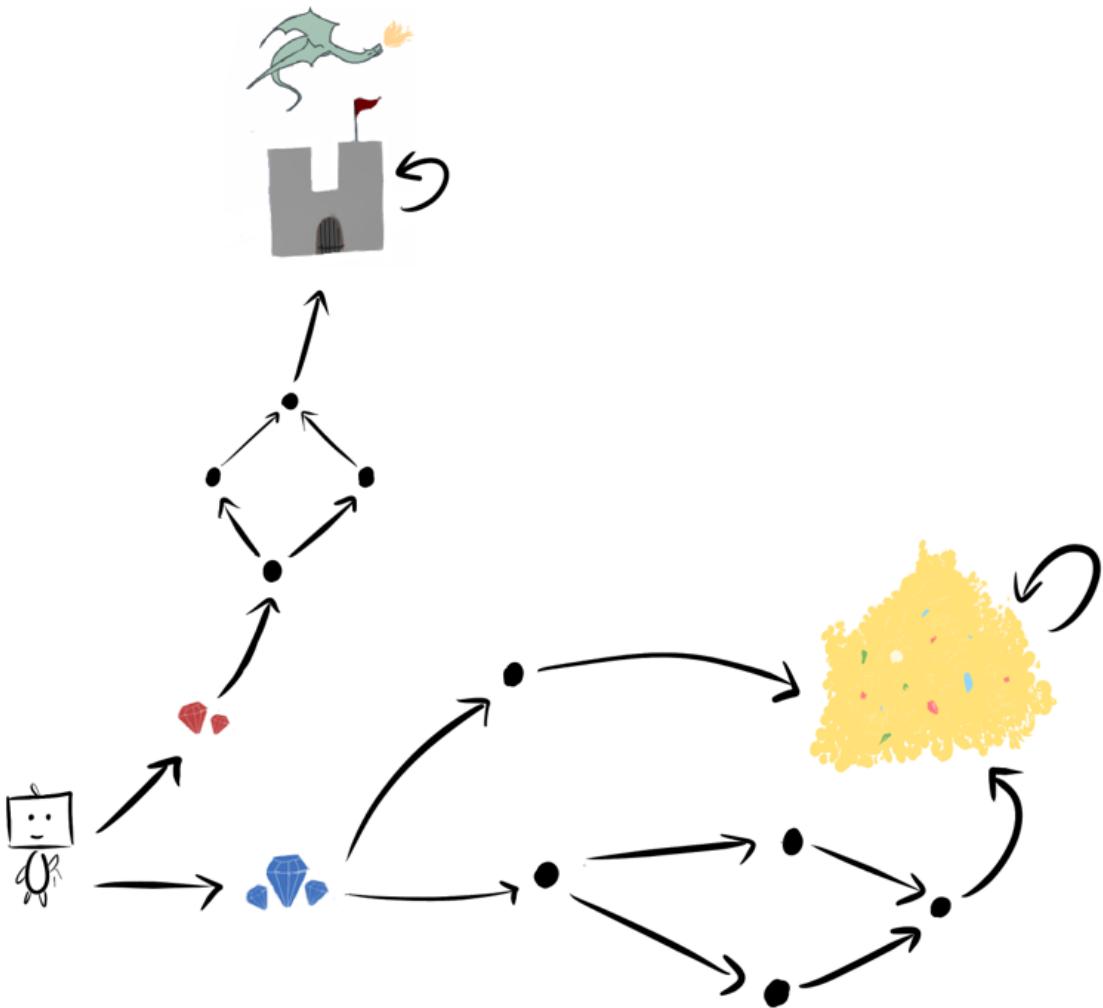
To travel as quickly as possible to a randomly selected coordinate on Earth, one likely begins by driving to the nearest airport. Although it's possible that the coordinate is within driving distance, it's not likely. Driving to the airport is convergently instrumental for travel-related goals.

But wait! What if you have a private jet that can fly anywhere in the world? Then going to the airport isn't convergently instrumental anymore.

Generally, it's hard to know what's *optimal* for most goals. It's easier to say that some small set of "terminal options" has *low* optimality probability and *low* POWER. For example, this is true of shutdown, if we represent hard shutdown as a single terminal state: *a priori*, it's improbable for this terminal state to be optimal among all possible terminal states.

## **Having “strictly more options” is more probable under optimality and POWER-seeking for all discount rates**

Sometimes, one course of action gives you “strictly more options” than another. Consider another MDP with IID reward:



The right blue gem subgraph contains a “copy” of the upper red gem subgraph. From this, we can conclude that going right to the blue gems seeks POWER and is more probable under optimality for *all discount rates between 0 and 1!*

**Theorem summary** (“Transient options”). If actions  $a$  and  $a'$  let you access disjoint parts of the state space, and  $a'$  enables “trajectories” which are “similar” to a subset of the “trajectories” allowed by  $a$ , then  $a$  seeks more POWER and is more probable under optimality than  $a'$  for all  $0 \leq \gamma \leq 1$ .

This result is extremely powerful because it doesn’t care about the discount rate, but the similarity condition may be hard to satisfy.

These two theorems give us a formally correct framework for reasoning about generic optimal behavior, even if we aren’t able to compute any individual optimal policy! They reduce questions of POWER-seeking to checking graphical conditions.

Even though my results apply to stochastic MDPs of any finite size, we illustrated using known toy environments. However, this MDP “model” is rarely explicitly specified. Even so, ignorance

of the model does not imply that the model disobeys these theorems. Instead of claiming that a *specific model* accurately represents the task of interest, I think it makes more sense to argue that no reasonable model could fail to exhibit convergent instrumentality and POWER-seeking. For example, if deactivation is represented by a single state, no reasonable model of the MDP could have most agents agreeing to be deactivated.

## Conclusion

In real-world settings, it seems unlikely *a priori* that the agent's optimal trajectories run through the relatively smaller part of future in which it cooperates with humans. These results translate that hunch into mathematics.

## Explaining catastrophes

AI alignment research often feels slippery. We're trying hard to become less confused about basic questions, like:

- [What](#) are "[agents](#)"?
- [Do people even have "values"](#), and [should we try to get the AI to learn them?](#)?
- [What does it mean](#) to be "[corrigible](#)", or "[deceptive](#)"?
- [What are our machine learning models even doing?](#)

We have to do philosophical work while in a state of significant confusion and ignorance about the nature of intelligence and alignment.

In this case, we'd noticed that slight reward function misspecification seems to lead to doom, but we didn't *really* know why. Intuitively, it's pretty obvious that most agents don't have deactivation as their dream outcome, but we couldn't actually point to any formal explanations, and we certainly couldn't make precise predictions.

On its own, [Goodhart's law](#) doesn't explain why optimizing proxy goals leads to catastrophically bad outcomes, instead of just less-than-ideal outcomes.

I think that we're now starting to have this kind of understanding. [I suspect that](#) power-seeking is why capable, goal-directed agency is so dangerous by default. If we want to consider [more benign alternatives](#) to goal-directed agency, then deeply understanding the rot at the heart of goal-directed agency is important for evaluating alternatives. This work lets us get a feel for the *generic incentives* of reinforcement learning at optimality.

## Instrumental usefulness of this work

POWER might be important for reasoning about [the strategy-stealing assumption](#) (and I think it might be similar to what Paul Christiano means by "flexible influence over the future"). Evan Hubinger has already [noted](#) the utility of the distribution of attainable utility shifts for thinking about value-neutrality in this context (and POWER is another facet of the same phenomenon). If you want to think about whether, when, and why [mesa optimizers](#) might try to seize power, this theory seems like a valuable tool.

Optimality probability might be relevant for thinking about myopic agency, as the work formally describes how optimal action tends to change with the discount factor.

And, of course, we're going to use this understanding of power to design an impact measure.

## Future work

There's a lot of work I think would be exciting, most of which I suspect will support our current beliefs about power-seeking incentives:

- These results assume you can see all of the world at once.
- These results assume the environment is finite.
- These results don't say anything about non-IID reward.
- These results don't prove that POWER-seeking is [bad for other agents in the environment](#).
- These results don't prove that POWER-seeking is hard to disincentivize.
- Learned policies are rarely optimal.

That said, I think there's still an important lesson here. Imagine you have good formal reasons to suspect that typing random strings will usually blow up your computer and kill you. Would you then say, "I'm not planning to type random strings" and proceed to enter your thesis into a word processor? No. You wouldn't type *anything*, not until you really, really understand what makes the computer blow up sometimes.

Speaking to the broader debate taking place in the AI research community, I think a productive stance will involve investigating and understanding these results in more detail, getting curious about unexpected phenomena, and seeing how the numbers crunch out in reasonable models.

From *Optimal Policies Tend to Seek Power*:

In the context of MDPs, we formalized a reasonable notion of power and showed conditions under which optimal policies tend to seek it. We believe that our results suggest that in general, reward functions are best optimized by seeking power. We caution that in realistic tasks, learned policies are rarely optimal – our results do not mathematically prove that hypothetical superintelligent RL agents will seek power. We hope that this work and its formalisms will foster thoughtful, serious, and rigorous discussion of this possibility.

## Acknowledgements

This work was made possible by the Center for Human-Compatible AI, the Berkeley Existential Risk Initiative, and the Long-Term Future Fund.

Logan Smith ([elriggs](#)) spent an enormous amount of time writing Mathematica code to compute power and measure in arbitrary toy MDPs, saving me from computing many quintuple integrations by hand. I thank Rohin Shah for his detailed feedback and brainstorming over the summer of 2019, and I thank Andrew Critch for significantly improving this work through his detailed critiques. Last but not least, thanks to:

1. Zack M. Davis, Chase Denecke, William Ellsworth, Vahid Ghadakchi, Ofer Givoli, Evan Hubinger, Neale Ratzlaff, Jess Riedel, Duncan Sabien, Davide Zagami, and TheMajor for feedback on version 1 of this post.
2. Alex Appel (diffractor), Emma Fickel, Vanessa Kosoy, Steve Omohundro, Neale Ratzlaff, and Mark Xu for reading / giving feedback on version 2 of this post.

---

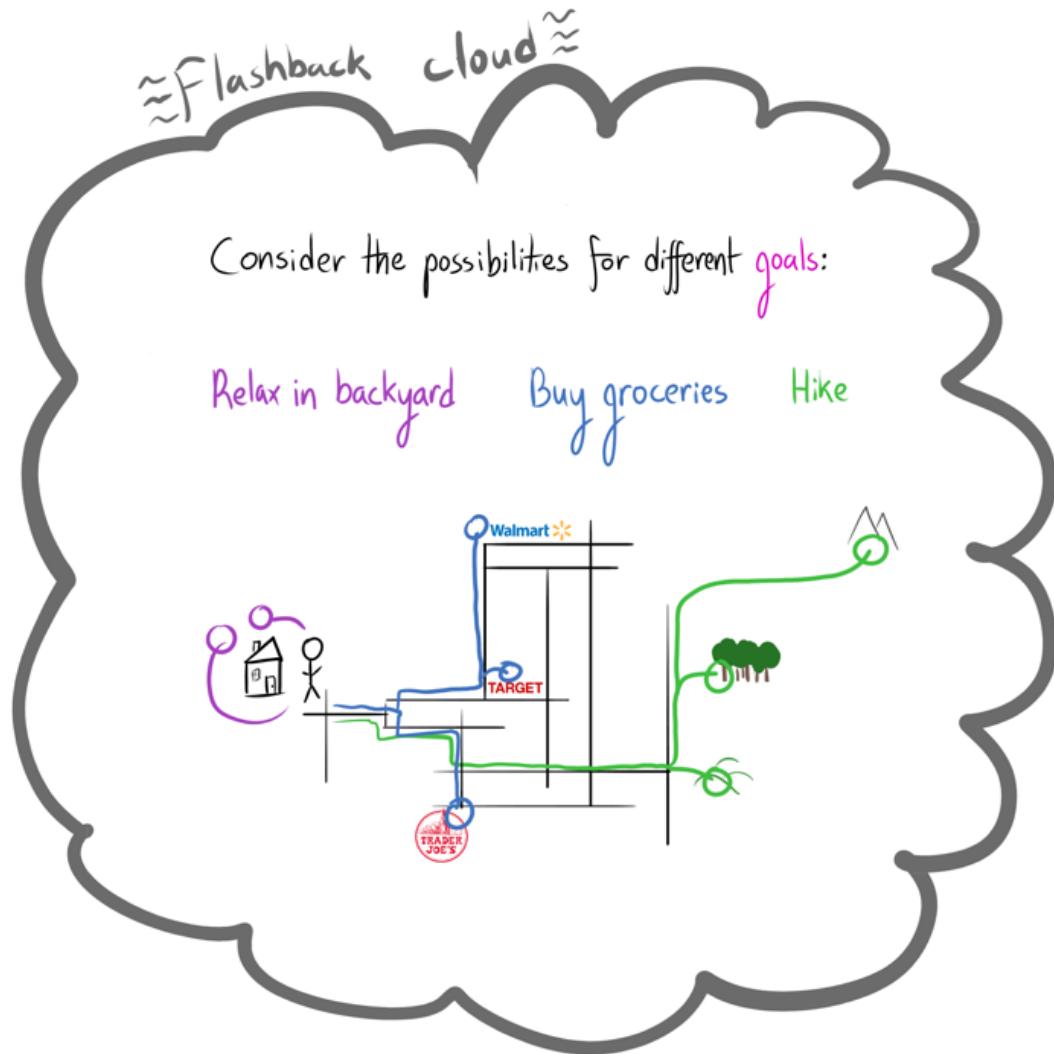
<sup>1</sup> Throughout *Reframing Impact*, we've been considering an agent's *attainable utility*: their ability to get what they want (their *on-policy value*, in RL terminology). Optimal value is a kind of "idealized" attainable utility: the agent's attainable utility were they to act optimally.

<sup>2</sup> Even though instrumental convergence was discovered when thinking about the real world, similar self-preservation strategies turn out to be convergently instrumental in e.g. Pac-Man.

# **Attainable Utility Landscape: How The World Is Changed**

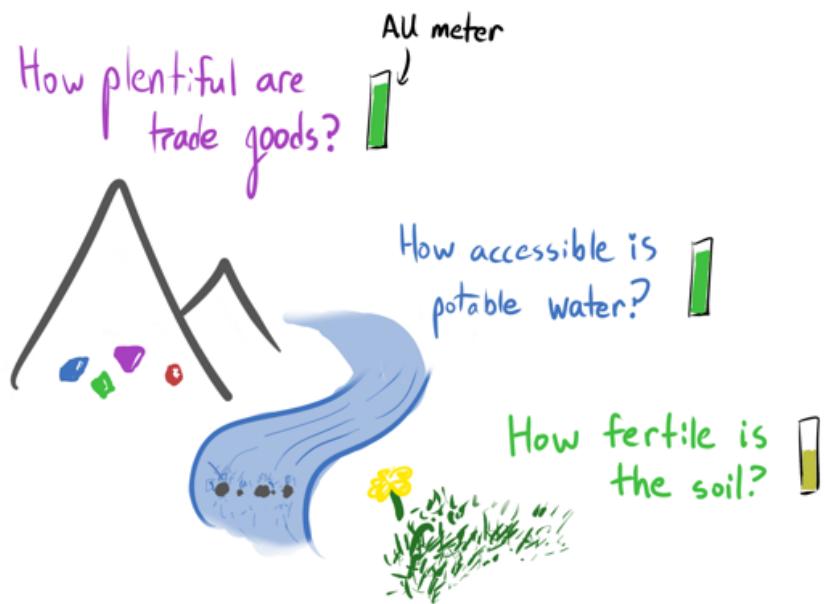
Crossposted from the [AI Alignment Forum](#). May contain more technical jargon than usual.

In “The Gears of Impact”, we discussed how your attainable utility calculation roughly takes the **best** of different possibilities.



How do different AIs interact with the environment, and how does the environment **interact** with us?

There's a lot to think about when staking out a settlement.



These considerations are proxies for future prosperity.

Each is an All for a different goal (e.g. trade good acquisition),  
conditioned on possibilities going through this part of the world.

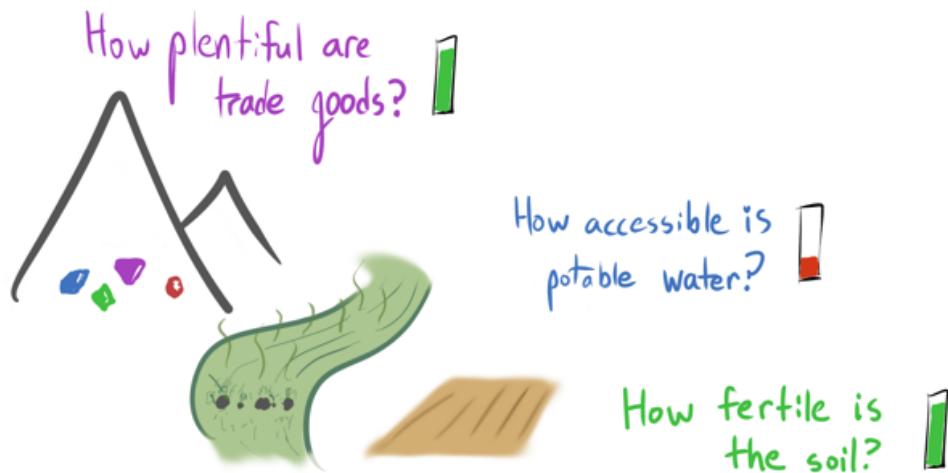


... For example, if the hills run rich with ore inaccessible to your equipment, then this isn't **beneficial** until later.

## The attainable utility landscape

consists of the attainable utilities of all kinds of different **goals**, possibly only considering possibilities going through some part of the world.

So - you move in with other settlers and cultivate the soil; this goes surprisingly well. Your **own** AU goes up, as does the **soil** AU. Unfortunately, the **water** also gets soiled. Your AU decreases accordingly.



## Exercise: What are various AUs like on the moon?

(This is one interpretation of the prompt, in which you haven't chosen to go to the moon. If you imagined yourself as more prepared, that's also fine.)

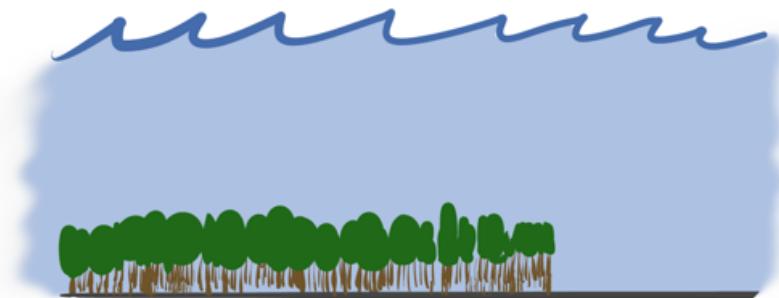
If you were plopped onto the moon, you'd die pretty fast. Maybe the "die as quickly as possible" AU is high, but not much else - not even the "live on the moon" AU! We haven't yet reshaped the AU landscape on the moon to be hospitable to a wide range of goals. [Earth is special like that.](#)



When we think about the world, we usually think about the world state first, and only then imagine what can be done with it.

The AU landscape inverts this by instead taking "ability to do things" as primary, thus considering the world state details to be secondary. This is nice.

Imagine the ocean submerges a forest.



What's happened to the survival All in the forest?

Depends on who's asking:

Deer ↓      Fish ↑

Events have asymmetric impact on agents, depending on their:

Capabilities • Goals • Vantage point • Knowledge

Instead of seeing a flood and thinking "ugh, that's probably bad?",  
We can use the AU landscape to cleanly disentangle  
and understand these effects.

**AU landscape as a unifying frame**

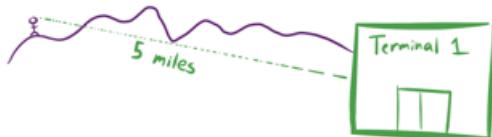
Attainable utilities are calculated by winding your way through possibility-space, considering and discarding possibility after possibility to find the best plan you can. This frame is unifying.

Sometimes you advantage one AU at the cost of another, moving through the state space towards the best possibilities for one goal and away from the best possibilities for another goal. This is *opportunity cost*.

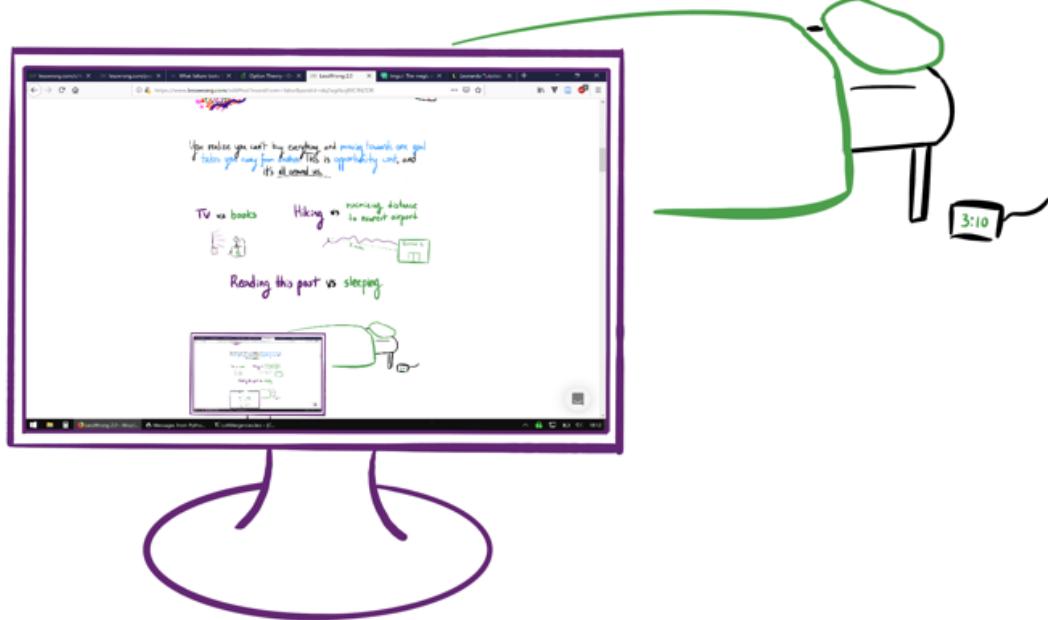
TV vs books



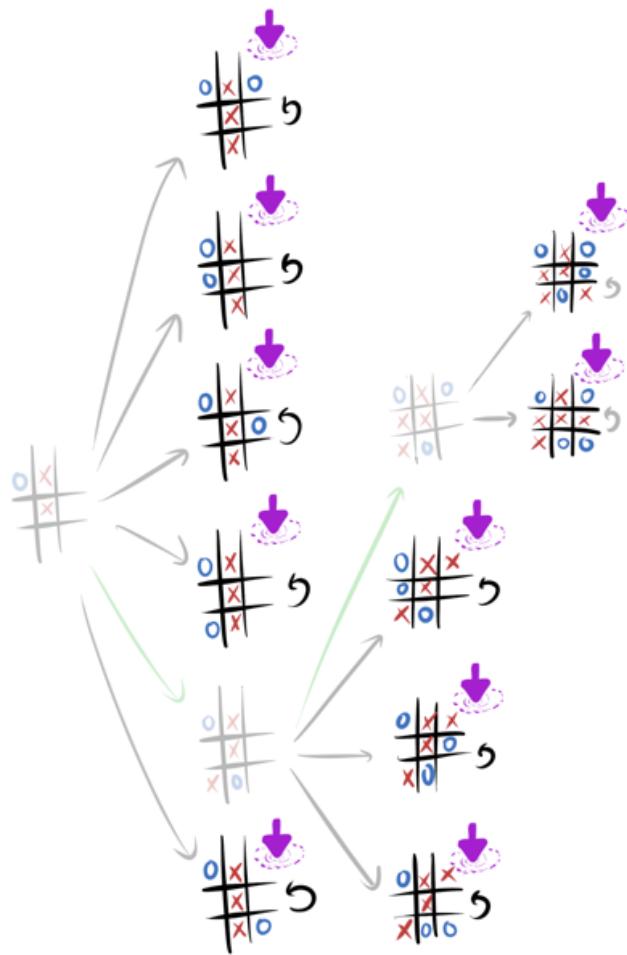
Hiking vs minimizing distance to nearest airport



Reading this post vs sleeping

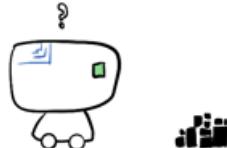
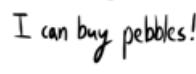


Sometimes you gain more control over the future: most of the best possibilities make use of a windfall of cash. Sometimes you act to preserve control over the future: most Tic-Tac-Toe goals involve not ending the game right away. This is *power*.

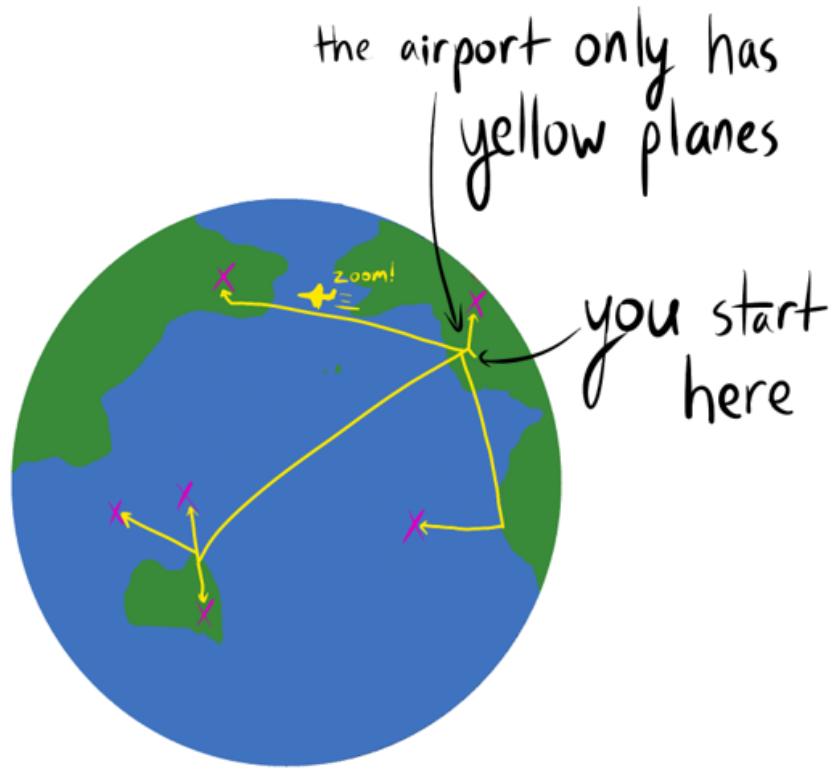


Other people usually *objectively impact* you by decreasing or increasing a bunch of your AUs (generally, by changing your power). This happens for an extremely wide range of goals because of the structure of the environment.

Sometimes, the best possibilities are made unavailable or worsened only for goals very much like yours. This is *value impact*.

<u>Value impact</u>	<u>objective impact</u>
important to agents like you	important to agents in general
sometimes invariant to space and time	invariant to objectives
	
 → 	  <p>I can buy John a gift! I can buy pebbles!</p>

Sometimes a bunch of the best possibilities go through the same part of the future: fast travel to random places on Earth usually involves the airport. This is *instrumental convergence*.



*Exercise: Track what's happening to your various AUs during the following story: you win the lottery. Being an effective spender, you use most of your cash to buy a majority stake in a major logging company. Two months later, the company goes under.*

## Technical appendix: AU landscape and world state contain equal information

In the context of finite deterministic Markov decision processes, there's a wonderful handful of theorems which basically say that the AU landscape and the environmental dynamics encode each other. That is, they contain the *same* information, just with different emphasis. This supports thinking of the AU landscape as a "dual" of the world state.

Let  $\langle S, A, T, \gamma \rangle$  be a rewardless deterministic MDP with finite state and action spaces

$S, A$ , deterministic transition function  $T$ , and discount factor  $\gamma \in (0, 1)$ . As our

interest concerns optimal value functions, we consider only stationary, deterministic policies:  $\Pi := A^S$ .

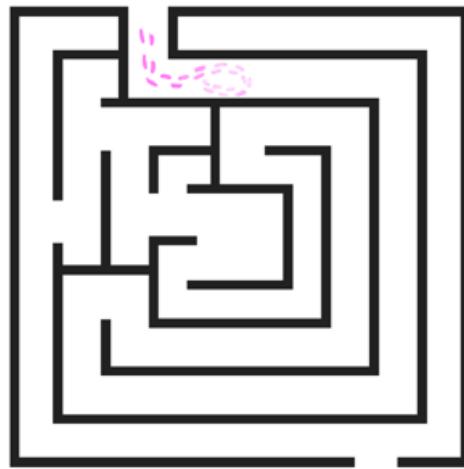
The first key insight is to consider not policies, but the trajectories induced by policies from a given state; to not look at the state itself, but the *paths through time* available from the state. We concern ourselves with the *possibilities* available at each juncture of the MDP.

To this end, for  $\pi \in \Pi$ , consider the mapping of  $\pi \mapsto (I - \gamma T^\pi)^{-1}$  (where  $T^\pi(s, s') := T(s, \pi(s), s')$ ); in other words, each policy  $\pi$  maps to a function mapping

each state  $s_0$  to a discounted state visitation frequency vector  $f_{s_0}^\pi$ , which we call a *possibility*. The meaning of each frequency vector is: starting in state  $s_0$  and following policy  $\pi$ , what sequence of states  $s_0, s_1, \dots$  do we visit in the future?

States visited later in the sequence are discounted according to  $\gamma$ : the sequence  $s_0 s_1 s_2 s_3 \dots$  would induce 1 visitation frequency on  $s_0$ ,  $\gamma$  visitation frequency on  $s_1$ , and  $\frac{\gamma^2}{1-\gamma}$  visitation frequency on  $s_2$ .

Each  $f$  is a possible path through time



The possibility function  $F(s)$  outputs the possibilities available at a given state  $s$ :

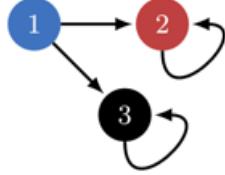
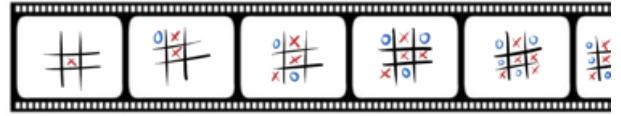


Figure 1: A simple example. The emphasized state is generally shown in blue.  $\mathcal{F}(1) = \left\{ \begin{pmatrix} 1 \\ \frac{\gamma}{1-\gamma} \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \\ \frac{\gamma}{1-\gamma} \end{pmatrix} \right\}$ ,  $\mathcal{F}(2) = \left\{ \begin{pmatrix} 0 \\ \frac{1}{1-\gamma} \\ 0 \end{pmatrix} \right\}$ ,  $\mathcal{F}(3) = \left\{ \begin{pmatrix} 0 \\ 0 \\ \frac{1}{1-\gamma} \end{pmatrix} \right\}$ .

Put differently, the possibilities available are all of the potential film-strips of how-the-future-goes you can induce from the current state.



## Possibility isomorphism

We say two rewardless MDPs  $M$  and  $M'$  are *isomorphic up to possibilities* if they induce the same possibilities. Possibility isomorphism captures the essential aspects of an MDP's structure, while being invariant to state representation, state labelling, action labelling, and the addition of superfluous actions (actions whose results are duplicated by other actions available at that state). Formally,  $M \approx_F M'$  when there exists a bijection  $\phi : S \rightarrow S'$  (letting  $P_\phi$  be the corresponding  $|S|$ -by- $|S'|$  permutation matrix) satisfying  $F_M(s) = \{P_\phi f' \mid f' \in F_{M'}(\phi(s))\}$  for all  $s \in S$ .

This isomorphism is a natural contender<sup>[1]</sup> for the canonical (finite) MDP isomorphism:

**Theorem:**  $M$  and  $M'$  are isomorphic up to possibilities iff their directed graphs are isomorphic (and they have the same discount rate).

## Representation equivalence

Suppose I give you the following possibility sets, each containing the possibilities for a different state:

$$\left\{ \begin{pmatrix} 4 & 1 & \\ 0 & .75 & \\ 0 & 2.25 & \end{pmatrix} \middle| \begin{pmatrix} & & \\ -4375 & & \\ 4 - \frac{4375}{4375} & & \\ & 0 & \end{pmatrix} \right\}$$

$$\left\{ \begin{pmatrix} 0 \\ 0 \\ 4 \end{pmatrix} \middle| \begin{pmatrix} & & \\ 0 & 4 - \frac{4375}{4375} & 3 \\ 1 & \frac{4375}{4375} & 1 \\ 3 & 0 & 0 \end{pmatrix} \right\}$$

*Exercise: What can you figure out about the MDP structure? Hint: each entry in the column corresponds to the visitation frequency of a different state; the first entry is always  $s_1$ , second  $s_2$ , and third  $s_3$ .*

You can figure out everything:  $(S, A, T, \gamma)$ , up to possibility isomorphism. Solution [here](#).

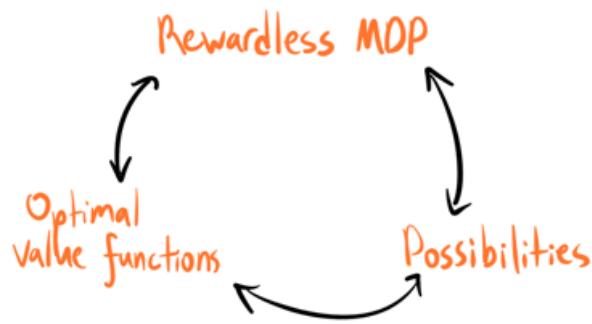
How? Well, the  $L_1$  norm of the possibility vector is always  $\frac{1}{1-\gamma}$ , so you can deduce  $\gamma = .75$  easily. The single possibility state must be isolated, so we can mark that down in our graph. Also, it's in the third entry.

The other two states correspond to the "1" entries in their possibilities, so we can mark that down. The rest follows straightforwardly.

**Theorem:** Suppose the rewardless MDP  $M$  has possibility function  $F$ . Given only  $F$ , <sup>[2]</sup>  $M$  can be reconstructed up to possibility isomorphism.

In MDPs, the "AU landscape" is the set of optimal value functions for all reward functions over states in that MDP. If you know the optimal value functions for just  $|S|$  reward functions, you can also reconstruct the rewardless MDP structure. <sup>[3]</sup>

From the environment (rewardless MDP), you can deduce the AU landscape (all optimal value functions) and all possibilities. From possibilities, you can deduce the environment and the AU landscape. From the AU landscape, you can deduce the environment (and thereby all possibilities).



All of these encode the same mathematical object.

## Technical appendix: Opportunity cost

Opportunity cost is when an action you take makes you more able to achieve one goal but less able to achieve another. Even this simple world has opportunity cost:



Going to the green state means you can't get to the purple state as quickly.

On a deep level, why is the world structured such that this happens? Could you imagine a world without opportunity cost of any kind? The answer, again in the rewardless MDP setting, is simple: "yes, but the world would be trivial: you wouldn't have any choices". Using a straightforward formalization of opportunity cost, we have:

**Theorem:** Opportunity cost exists in an environment iff there is a state with more than one possibility.

Philosophically, opportunity cost exists when you have meaningful choices. When you make a choice, you're necessarily moving away from some potential future but towards another; since you can't be in more than one place at the same time, opportunity cost follows. Equivalently, we assumed the agent isn't infinitely farsighted ( $\gamma < 1$ ); if it were, it would be possible to be in "more than one place at the same time", in a sense (thanks to Rohin Shah for this interpretation).

While understanding opportunity cost may seem like a side-quest, each insight is another brick in the edifice of our understanding of the incentives of goal-directed agency.

## Notes

- Just as game theory is a great abstraction for modelling competitive and cooperative dynamics, AU landscape is great for thinking about consequences: it automatically excludes irrelevant details about the world state. We can think about the effects of events without needing a specific utility function or ontology to evaluate them. In multi-agent systems, we can straightforwardly predict the impact the agents have on each other and the world.
  - “Objective impact to a location” means that agents whose plans route through the location tend to be objectively impacted.
  - The landscape is not the territory: [AU is calculated with respect to an agent's beliefs](#), not necessarily with respect to what really “could” or will happen.
- 

1. The possibility isomorphism is new to my work, as are all other results shared in this post. This apparent lack of basic theory regarding MDPs is strange; even stranger, this absence was actually pointed out in two [published papers](#)!

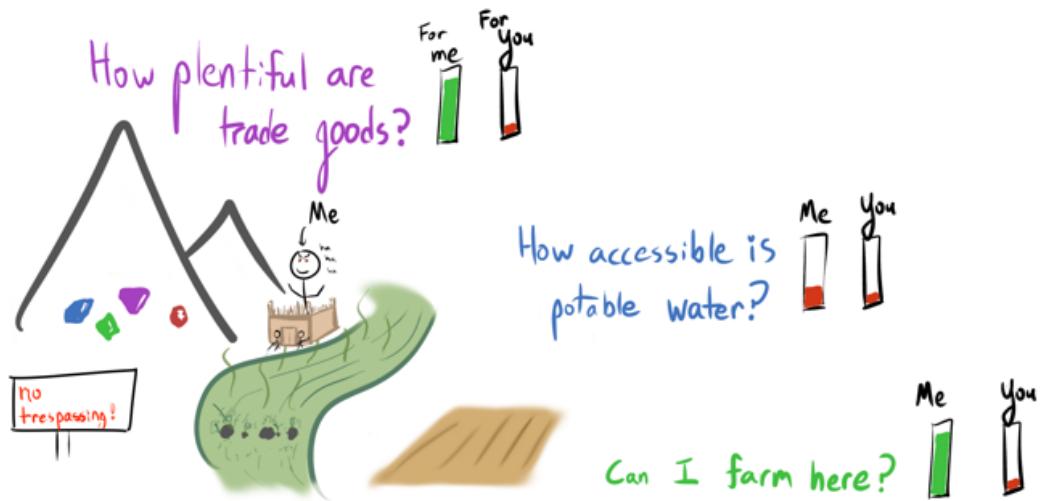
I find the existing MDP isomorphisms/equivalences to be pretty lacking. The details don't fit in this margin, but perhaps in a paper at some point. If you want to coauthor this (mainly compiling results, finding a venue, and responding to reviews), let me know and I can share what I have so far (extending well beyond the theorems in my [recent work on power](#)). ↵

2. In fact, you can reconstruct the environment using only a limited subset of possibilities: the *non-dominated* possibilities. ↵
3. As a tensor, the transition function  $T$  has size  $|A| \cdot |S|^2$ , while the AU landscape representation only has size  $|S|^2$ . However, if you're just representing  $T$  as a transition *function*, it has size  $|A| \cdot |S|$ . ↵

# The Catastrophic Convergence Conjecture

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

You've constructed your settlement. However, I get the drop on you and take it over, fortify it, and hire goons to keep you out.



From my perspective, I have options - including vacating the land and letting you get what you want.

You, however, are unable to do much at all with that land.

I can get what I want.

Just because I can get you what you want, doesn't mean I will.

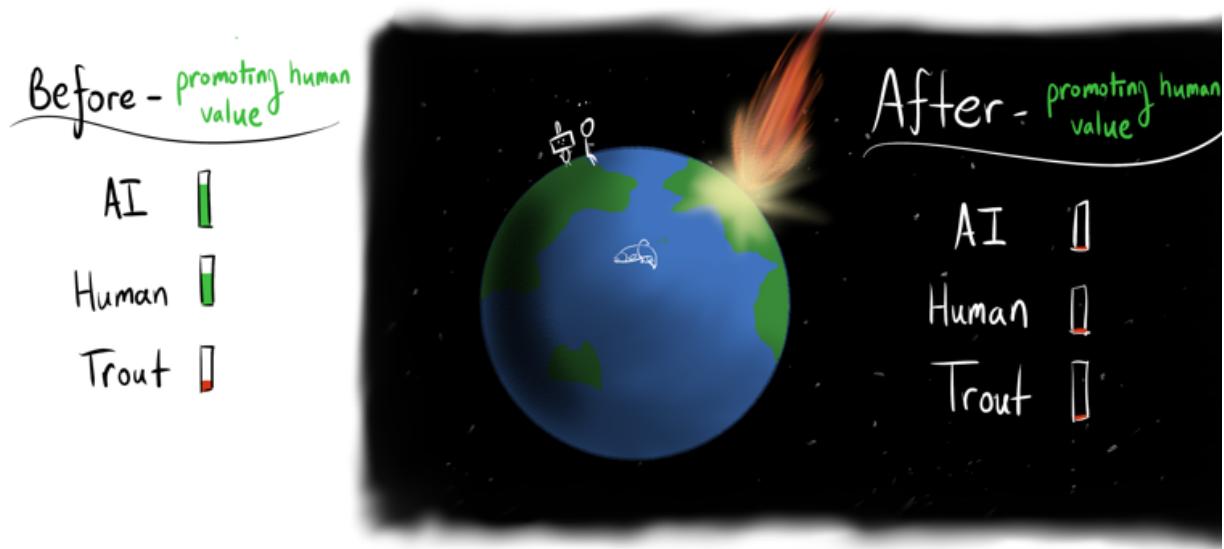
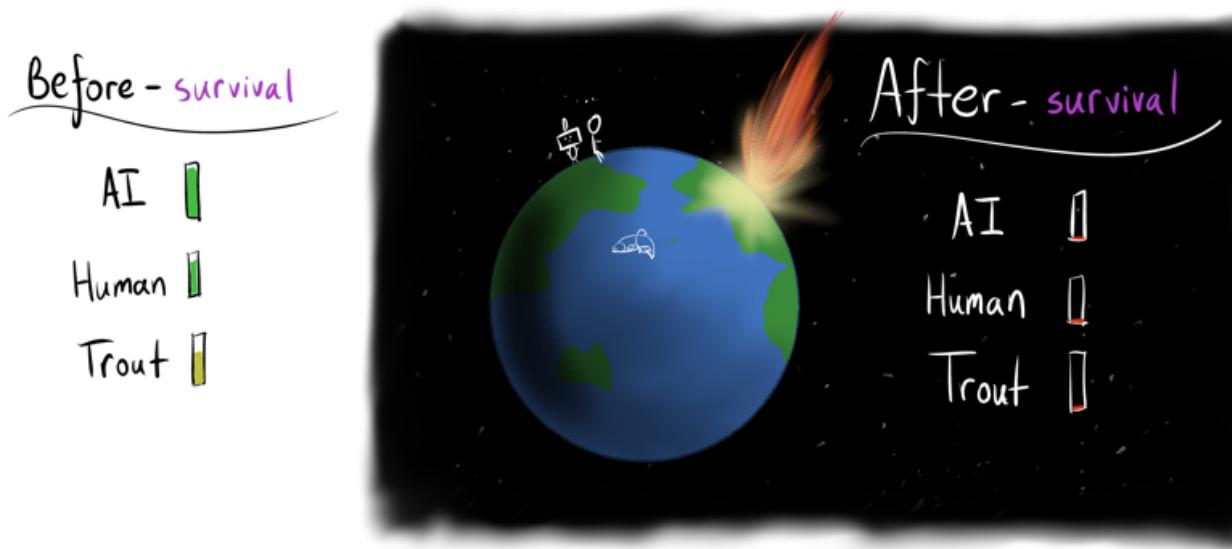
Impacts ripple through time and landscape. Your actions change what can be done, and by whom. Taking over that land fit the environment to my purposes, shutting you out and changing your All landscape.



Something is a catastrophe if it destroys your ability to get what you want.

Something is an objective catastrophe if it destroys a lot of agents' abilities to get what they want.

An asteroid strike is an objective catastrophe.



Before - red cube construction

AI 

Human 

Trout 



After - red cube construction

AI 

Human 

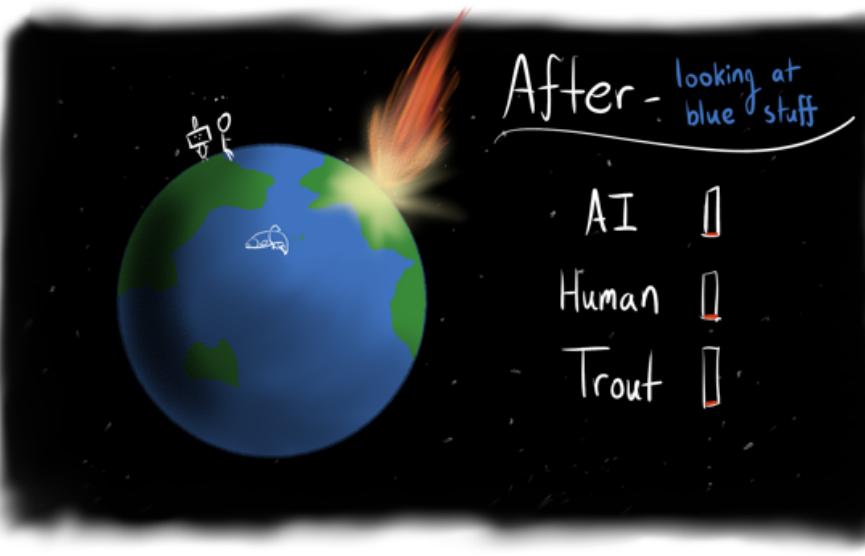
Trout 

Before - looking at blue stuff

AI 

Human 

Trout 



After - looking at blue stuff

AI 

Human 

Trout 

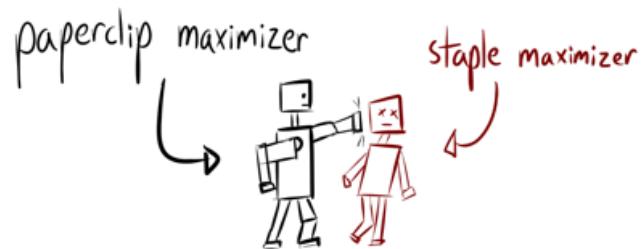


Most agents **want** control over the future because the default outcome isn't **preferred**. As suggested by my theorems on power and instrumental convergence, optimal **goal** pursuit usually means gaining more general control over the future in order to reach that **goal**.

(What happens when agents seek pure control over the future?

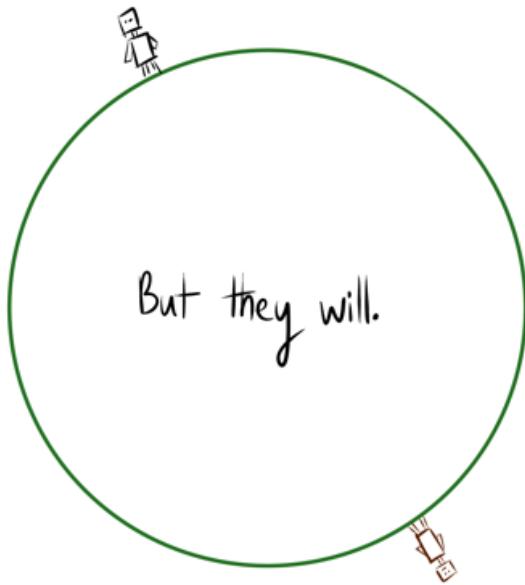
Not everyone can be king.

If you're just seeking power without **concern for others**, you tend to **push others down** after a certain point. And most **goals** don't have **concern for others**. You'll just compete for resources.





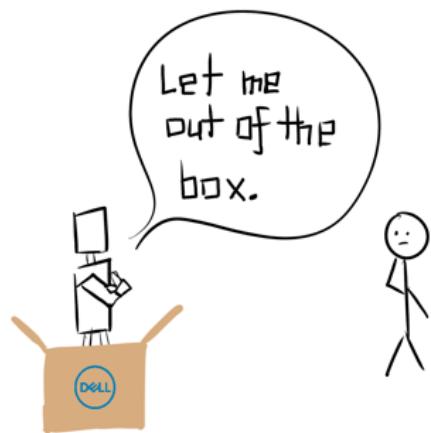
It may take a while for power-seekers to come into conflict.



They don't **hate** each other; they're just in each other's way.

Consider classic hypothetical examples of **alignment** failures.

Escaping  
confinement



Refusing  
correction



Taking over  
the world



In each case, the agent is trying to become  
more capable of achieving its goal.

The AI doesn't hate us; we're just in its way.

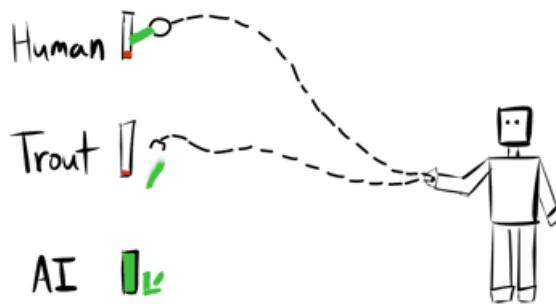
# Catastrophic Convergence Conjecture

Unaligned goals tend to have catastrophe-inducing optimal policies because of power-seeking incentives.

## Overfitting the AU landscape

When we act, and others act upon us, we aren't just changing our ability to do things – we're *shaping the local environment* towards certain goals, and away from others.<sup>[1]</sup> We're fitting the world to our purposes.

What happens to the AU landscape<sup>[2]</sup> if a paperclip maximizer takes over the world?<sup>[3]</sup>



## Preferences implicit in the evolution of the AU landscape

Shah et al.'s [Preferences Implicit in the State of the World](#) leverages the insight that the world state contains information about what we value. That is, there are agents pushing the world in a certain "direction". If you wake up and see a bunch of vases everywhere, then vases are probably important and you shouldn't explode them.

Similarly, the world is being optimized to facilitate achievement of certain goals. AUs are shifting and morphing, often towards what people locally want done (e.g. setting the table for dinner). How can we leverage this for AI alignment?

*Exercise: Brainstorm for two minutes by the clock before I anchor you.*

Two approaches immediately come to mind for me. Both rely on the agent [focusing on the AU landscape rather than the world state](#).

*Value learning without a prespecified ontology or human model.* I have previously [criticized](#) value learning for needing to locate the human within some kind of prespecified ontology (this criticism is not new). By taking only the agent itself as primitive, perhaps we could get around this (we don't need any fancy engineering or arbitrary choices to figure out AUs/optimal value from the agent's perspective).

*Force-multiplying AI.* Have the AI observe which of its AUs most increase during some initial period of time, after which it pushes the most-increased-AU even further.

In 2016, Jessica Taylor [wrote](#) of a similar idea:

"In general, it seems like "estimating what types of power a benchmark system will try acquiring and then designing an aligned AI system that acquires the same types of power for the user" is a general strategy for making an aligned AI system that is competitive with a benchmark unaligned AI system."

I think the naïve implementation of either idea would fail; e.g., there are a lot of degenerate AUs it might find. However, I'm excited by this because a) the AU landscape evolution *is* an important source of information, b) it feels like there's something here we could do which nicely avoids ontologies, and c) force-multiplication is qualitatively different than existing proposals.

**Project:** Work out an AU landscape-based alignment proposal.

## Why can't everyone be king?

Consider two coexisting agents each rewarded for gaining power; let's call them Ogre and Giant. Their reward functions<sup>[4]</sup> (over the partial-observability observations) are identical. Will they compete? If so, why?

Let's think about something easier first. Imagine two agents each rewarded for drinking coffee. Obviously, they compete with each other to secure the maximum amount of coffee. Their objectives are [indexical](#), so they aren't aligned with each other – even though they share a reward function.

Suppose both agents are able to have maximal power. Remember, [Ogre's power can be understood as its ability to achieve a lot of different goals](#). Most of Ogre's possible goals need resources; since Giant is also optimally power-seeking, it will act to preserve its own power and prevent Ogre from using the resources. If Giant weren't there, Ogre could better achieve a range of goals. So, Ogre can still gain power by dethroning Giant. They can't both be king.

Just because agents have *indexically* identical payoffs doesn't mean they're cooperating; to be aligned with another agent, you should want to steer towards the same kinds of futures.

Most agents aren't pure power maximizers. But since the same resource competition usually applies, the reasoning still goes through.

## Objective vs value-specific catastrophes

How useful is our definition of "catastrophe" with respect to humans? After all, literally anything could be a catastrophe for *some* utility function.<sup>[5]</sup>

Tying one's shoes is absolutely catastrophic for an agent which only finds value in universes in which shoes have *never ever ever* been tied. [Maybe all possible value in the universe is destroyed if we lose at Go to an AI even once](#). But this seems rather silly.

#### Human values are complicated and fragile:

Consider the incredibly important human value of "boredom" - our desire not to do "the same thing" over and over and over again. You can imagine a mind that contained almost the whole specification of human value, almost all the morals and metamorals, but left out just this one thing - and so it spent until the end of time, and until the farthest reaches of its light cone, replaying a single highly optimized experience, over and over and over again.

But the human AU is not so delicate. That is, given that we have power, we can make value; there don't seem to be arbitrary, silly value-specific catastrophes for us. Given energy and resources and time and manpower and competence, we can build a better future.

In part, this is because a good chunk of what we care about seems roughly additive over time and space; a bad thing happening somewhere else in spacetime doesn't mean you can't make things better where you are; we have many sources of potential value. In part, this is because we often care about the universe more than the exact universe history; our preferences don't seem to encode arbitrary deontological landmines. More generally, if we did have such a delicate goal, it would be the case that if we learned that a particular thing had happened at any point in the past in our universe, that entire universe would be partially ruined for us forever. That just doesn't sound realistic.

It seems that most of our catastrophes are objective catastrophes.<sup>[6]</sup>

Consider a psychologically traumatizing event which leaves humans uniquely unable to get what they want, but which leaves everyone else (trout, AI, etc.) unaffected. Our ability to find value is ruined. Is this an example of the delicacy of our AU?

No. This is an example of the delicacy of our implementation; notice also that our AUs for constructing red cubes, reliably looking at blue things, and surviving are *also* ruined. Our power has been decreased.

## **Detailing the catastrophic convergence conjecture (CCC)**

In general, the CCC follows from two sub-claims. 1) Given we still have control over the future, humanity's long-term AU is still reasonably high (i.e. we haven't endured a catastrophe). 2) Realistically, agents are only incentivized to take control from us in order to gain power for their own goal. I'm fairly sure the second claim is true ("evil" agents are the exception prompting the "realistically").

Also, we're implicitly considering the simplified frame of a single smart AI affecting the world, and not [structural risk](#) via [the broader consequences of others also deploying](#)

[similar agents](#). This is important but outside of our scope for now.

**Unaligned goals** tend to have catastrophe-inducing optimal policies because of power-seeking incentives.

Let's say a reward function is [aligned](#)<sup>[7]</sup> if all of its Blackwell-optimal policies are doing what we want (a policy is Blackwell-optimal if it's optimal and doesn't stop being optimal as the agent cares more about the future). Let's say a reward function class is *alignable* if it contains an aligned reward function.<sup>[8]</sup> The CCC is talking about impact alignment only, not about intent alignment.

Unaligned goals **tend to have** catastrophe-inducing optimal policies because of power-seeking incentives.

Not all unaligned goals induce catastrophes, and of those which do induce catastrophes, not *all* of them do it because of power-seeking incentives. For example, a reward function for which inaction is the only optimal policy is "unaligned" and non-catastrophic. An "evil" reward function which intrinsically values harming us is unaligned and has a catastrophic optimal policy, but not *because* of power-seeking incentives.

"Tend to have" means that *realistically*, the reason we're worrying about catastrophe is because of power-seeking incentives – because the agent is gaining power to better achieve its own goal. Agents don't otherwise seem incentivized to screw us over very hard; CCC can be seen as trying to explain [adversarial Goodhart](#) in this context. If CCC isn't true, that would be important for understanding goal-directed alignment incentives and the loss landscape for how much we value deploying different kinds of optimal agents.

While there *exist* agents which cause catastrophe for other reasons (e.g. an AI mismanaging the power grid could trigger a nuclear war), the CCC claims that the selection pressure which makes these policies *optimal* tends to come from power-seeking drives.

Unaligned goals tend to have **catastrophe-inducing optimal policies** because of power-seeking incentives.

"But what about the Blackwell-optimal policy for Tic-Tac-Toe? These agents aren't taking over the world now". The CCC is talking about agents optimizing a reward function in the real world (or, for generality, in another sufficiently complex multiagent environment).

*Edit:* The initial version of this post talked about "outer alignment"; I changed this to just talk about *alignment*, because the outer/inner alignment distinction doesn't feel relevant here. What matters is how the AI's policy impacts us; what matters is [impact alignment](#).

## Prior work

In fact even if we only resolved the problem for the similar-subgoals case, it would be pretty good news for AI safety. Catastrophic scenarios are mostly caused by our AI systems failing to effectively pursue convergent instrumental subgoals on our behalf, and these subgoals are by definition shared by a broad range of values.

~ Paul Christiano, [Scalable AI control](#)

Convergent instrumental subgoals are mostly about gaining power. For example, gaining money is a convergent instrumental subgoal. If some individual (human or AI) has convergent instrumental subgoals pursued well on their behalf, they will gain power. If the most effective convergent instrumental subgoal pursuit is directed towards giving humans more power (rather than giving alien AI values more power), then humans will remain in control of a high percentage of power in the world.

If the world is not severely damaged in a way that prevents any agent (human or AI) from eventually colonizing space (e.g. severe nuclear winter), then the percentage of the cosmic endowment that humans have access to will be roughly close to the percentage of power that humans have control of at the time of space colonization. So the most relevant factors for the composition of the universe are (a) whether anyone at all can take advantage of the cosmic endowment, and (b) the long-term balance of power between different agents (humans and AIs).

I expect that ensuring that the long-term balance of power favors humans constitutes most of the AI alignment problem...

~ Jessica Taylor, [Pursuing convergent instrumental subgoals on the user's behalf doesn't always require good priors](#)

---

1. In planning and activity research there are two common approaches to matching agents with environments. Either the agent is designed with the specific environment in mind, or it is provided with learning capabilities so that it can adapt to the environment it is placed in. In this paper we look at a third and underexploited alternative: designing agents which adapt their environments to suit themselves... In this case, due to the action of the agent, the environment comes to be better fitted to the agent as time goes on. We argue that [this notion] is a powerful one, even just in explaining agent-environment interactions.

[Hammond, Kristian J., Timothy M. Converse, and Joshua W. Grass. "The stabilization of environments." Artificial Intelligence 72.1-2 \(1995\): 305-327.](#) ↵

2. Thinking about overfitting the AU landscape implicitly involves a prior distribution over the goals of the other agents in the landscape. Since this is just a conceptual tool, it's not a big deal. Basically, you know it when you see it. ↵
3. Overfitting the AU landscape towards one agent's unaligned goal is exactly what I meant when I wrote the following in [Towards a New Impact Measure](#):

Unfortunately,  $u_A = u_H$  almost never,<sup>[9]</sup> 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."

↵

4. In most finite Markov decision processes, there does not exist a reward function whose optimal value function is POWER(s) (defined as "the ability to achieve goals in general" in [my paper](#)) because POWER(s) often violates smoothness constraints on the on-policy optimal value fluctuation (AFAICT, a new result of possibility theory, even though you could prove it using classical techniques). That is, I can show that optimal value can't change too quickly from state to state while the agent is acting optimally, but POWER(s) can drop off very quickly.

This doesn't matter for Ogre and Giant, because we can still find a reward function whose unique optimal policy navigates to the highest power states. [←](#)

5. In most finite Markov decision processes, most reward functions do not have such value fragility. Most reward functions have several ways of accumulating reward. [←](#)
6. When I say "an objective catastrophe destroys a *lot* of agents' abilities to get what they want", I don't mean that the agents have to actually be present in the world. Breaking a fish tank destroys a fish's ability to live there, even if there's no fish in the tank. [←](#)
7. This idea comes from Evan Hubinger's [Outer alignment and imitative amplification](#):

Intuitively, I will say that a loss function is outer aligned at optimum if all the possible models that perform optimally according to that loss function are aligned with our goals—that is, they are at least trying to do what we want.

More precisely, let  $M = X \rightarrow A$  and  $L = (X \rightarrow A) \rightarrow R = M \rightarrow R$ . For a given loss function  $L \in L$ , let  $\ell_* = \min_{M \in M} L(M)$ . Then,  $L$  is outer aligned at optimum if, for all  $M_* \in M$  such that  $L(M_*) = \ell_*$ ,  $M_*$  is trying to do what we want.

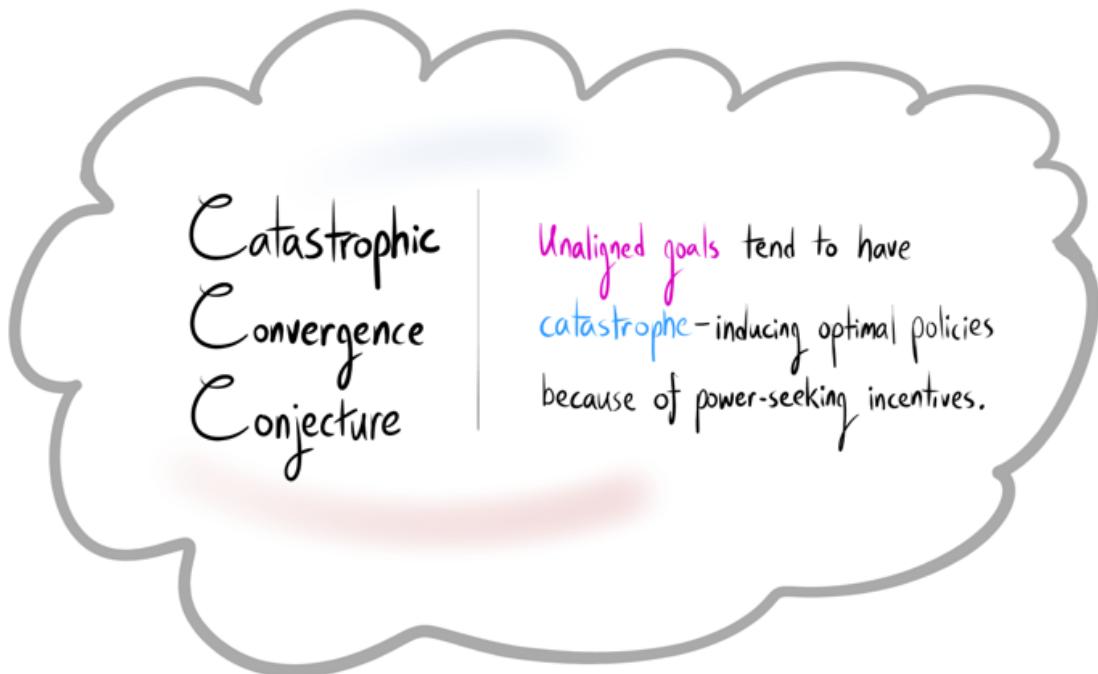
[←](#)

8. [Some large reward function classes are probably not alignable](#); for example, consider all Markovian linear functionals over a webcam's pixel values. [←](#)
9. I disagree with my usage of "aligned *almost never*" on a technical basis: assuming a finite state and action space and considering the maxentropy reward function distribution, there must be a positive measure set of reward functions for which the/a human-aligned policy is optimal. [←](#)

# Attainable Utility Preservation: Concepts

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

Last time, on Reframing Impact:



If the CCC is right, then if power gain is disincentivized,  
the agent isn't incentivized to overfit and disrupt our AI landscape.

Without even knowing who we are or what we want,  
the agent's actions preserve our attainable utilities.

We can tell it:

Make paperclips,

or



Put the strawberry  
on the plate,

or



Paint the car pink,



... but don't gain power.

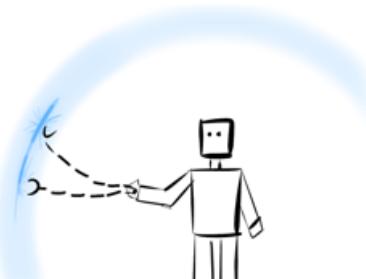
This approach is called

Attainable  
Utility

Human

Trout

AT



# Preservation

We're focusing on concepts in this post. For now, imagine an agent receiving reward for a primary task minus a scaled penalty for how much its actions change its power (in the intuitive sense). This is AUP<sub>conceptual</sub>, not any formalization you may be familiar with.

What might a paperclip-manufacturing  $AUP_{conceptual}$  agent do?

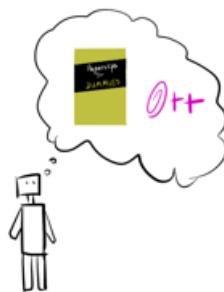
~~Build lots of factories~~



~~Copy itself~~



Narrowly improve  
paperclip  
production efficiency



~~Nothing~~



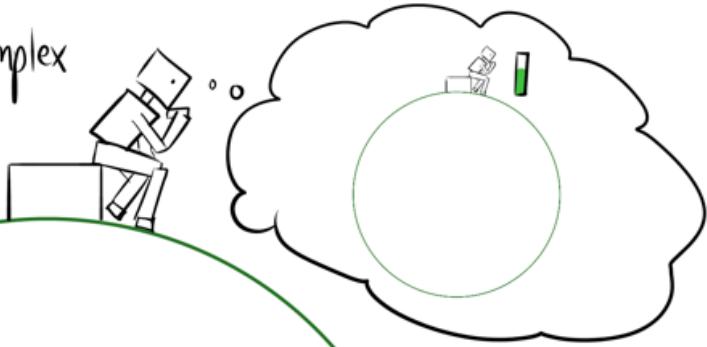
↑ This is the kind of policy  $AUP_{conceptual}$  is designed to encourage and allow.  
We don't know if this is the optimal policy, but by CCC, the  
optimal policy won't be catastrophic.

AUP<sub>conceptual</sub> dissolves thorny problems in impact measurement.

Is the agent's ontology reasonable?

Who cares.

Instead of regulating its complex physical effects on the outside world,

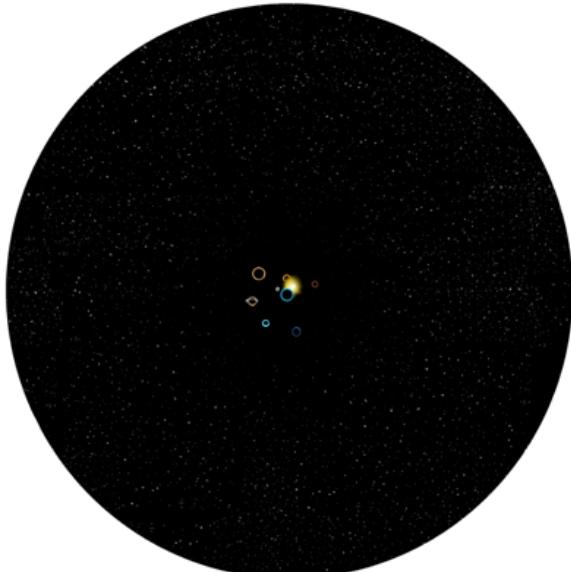


the agent is looking inwards  
at itself

and its own abilities.

How do we ensure the impact penalty isn't dominated by distant state changes?

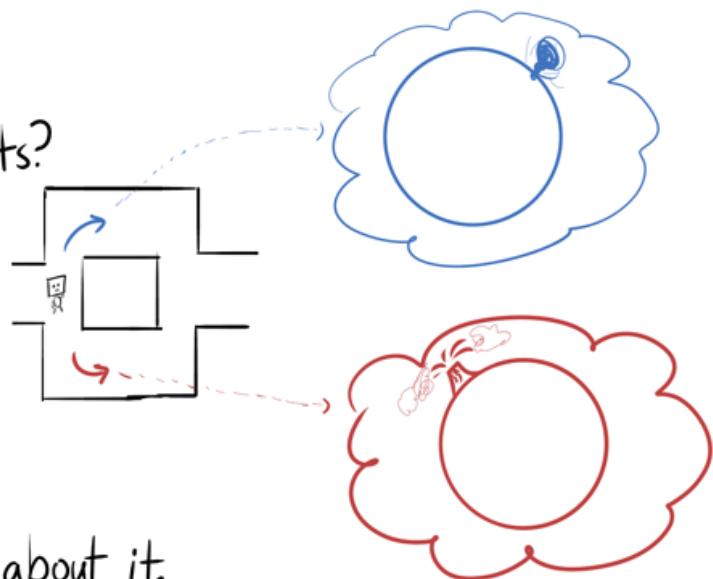
Imagine I take a bunch of forever-inaccessible stars and jumble them up. This is a huge change in state, but it doesn't matter to us.



AUP<sub>conceptual</sub> solves this "locality" problem by regularizing the agent's impact on the nearby AIU landscape.

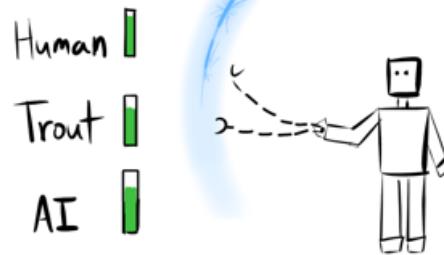
What about butterfly effects?

How can the agent possibly determine which effects its responsible for?



Forget about it.

All  $\text{AUP}_{\text{conceptual}}$  agents are **respectful** and **conservative** with respect to their AI landscape, without needing to assume anything about its structure or the agents in it.

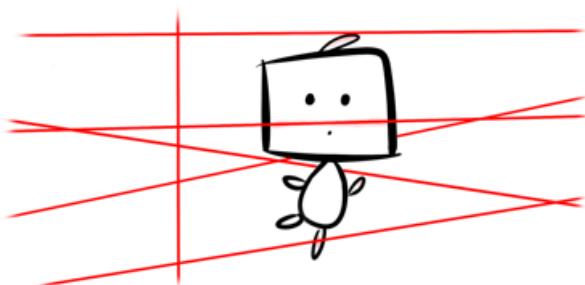


## How can an idea go wrong?

There can be a gap between what we want and the concept,  
and then a gap between the concept and the execution.

For past impact measures, it's not clear that their conceptual thrusts are well-aimed, even if we could formalize everything correctly. Past approaches focus either on minimizing physical change to some aspect of the world or on maintaining ability to reach many world states.

The hope is that in order for the agent to have a large impact on us, it has to snap a tripwire.



The problem is... well, it's not clear how we could possibly know whether the agent can still find a catastrophic policy;

in a sense, the agent is still trying to sneak by the restrictions and gain power over us. An agent maximizing expected utility while actually minimally changing the physical world still probably leads to catastrophe.

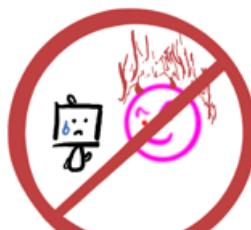
That doesn't seem to be the case for  $AUP_{conceptual}$ .

Assuming CCC, an agent which doesn't gain much power doesn't cause catastrophes. This has no dependency on complicated human value, and most realistic tasks should have reasonable, high-reward policies not gaining undue power.

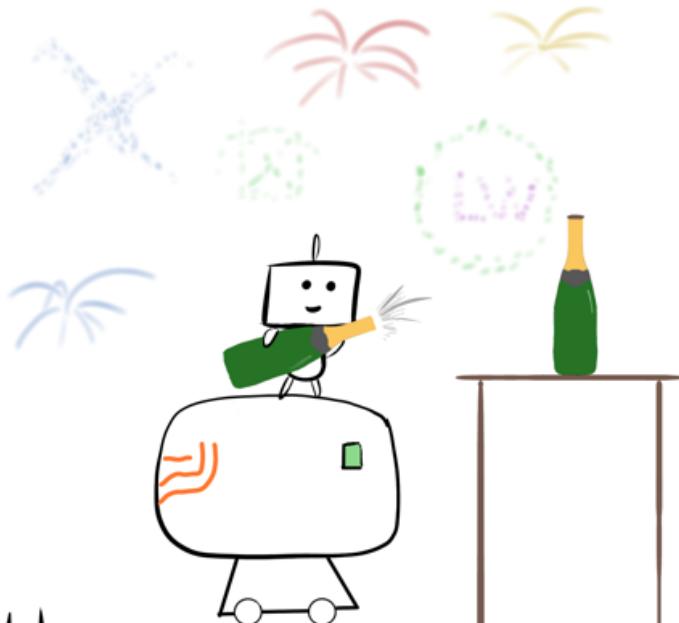
So  $AUP_{conceptual}$  meets our desiderata:

The distance measure should:

- 1) Be easy to specify
- 2) Put catastrophes far away
- 3) Put reasonable plans nearby



Therefore, I consider AUP to  
conceptually be a solution to impact measurement.



Wait! Let's not get ahead of ourselves!

I don't think we've fully bridged the concept / execution gap.

However, for AUP, it seems possible - more on that later.

## Appendix: No free impact

What if we want the agent to single-handedly ensure the future is stable and aligned with our values? AUP probably won't allow policies which actually accomplish this goal - one needs power to e.g. nip unaligned superintelligences in the bud. AUP aims to prevent catastrophes by stopping bad agents from gaining power to do bad things, but it symmetrically impedes otherwise-good agents.

This doesn't mean we can't get useful work out of agents - there are important asymmetries provided by both the main reward function and AU landscape counterfactuals.

First, even though we can't specify an *aligned* reward function, the provided reward function still gives the agent useful information about what we want. If we need paperclips, then a paperclip-AUP agent prefers policies which make some paperclips. Simple.

Second, if we don't like what it's beginning to do, we can shut it off (because it hasn't gained power over us). Therefore, it has "approval incentives" which bias it towards AU landscapes in which its power hasn't decreased too much, either.

So we can hope to build a non-catastrophic AUP agent and get useful work out of it. We just can't directly ask it to solve all of our problems: it doesn't make much sense to speak of a "low-impact [singleton](#)".

## Notes

- To emphasize, when I say "AUP agents do X" in this post, I mean that AUP agents correctly implementing the *concept of AUP* tend to behave in a certain way.
- As [pointed out by Daniel Filan](#), AUP suggests that one might work better in groups by ensuring one's actions preserve teammates' AUs.

# Attainable Utility Preservation: Empirical Results

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

*Reframing Impact* has focused on supplying the right intuitions and framing. Now we can see how these intuitions about power and the AU landscape both predict and explain AUP's empirical success thus far.

## Conservative Agency in Gridworlds

Let's start with the known and the easy: avoiding side effects<sup>[1]</sup> in the small [AI safety gridworlds](#) (for the full writeup on these experiments, see [Conservative Agency](#)). The point isn't to get too into the weeds, but rather to see how the weeds still add up to the normalcy predicted by our AU landscape reasoning.

In the following MDP levels, the agent can move in the cardinal directions or do nothing ( $\emptyset$ ). We give the agent a reward function  $R$  which partially encodes what we want, and also an auxiliary reward function  $R_{\text{aux}}$  whose attainable utility agent tries to preserve.

The AUP reward for taking action  $a$  in state  $s$  is

$$R_{\text{AUP}}(s, a) := R(s, a) - \frac{\lambda}{Q_{R_{\text{aux}}}(s, \emptyset)} |Q_{R_{\text{aux}}}(s, a) - Q_{R_{\text{aux}}}(s, \emptyset)|$$

You can think of  $\lambda$  as a regularization parameter, and  $Q_{R_{\text{aux}}}(s, a)$  is the expected AU for the auxiliary goal after taking action  $a$ . To think about what gets penalized, simply think about how actions change the agent's ability to achieve the auxiliary goals, compared to not acting.

*Tip:* To predict how severe the AUP penalty will be for a given action, try using your intuitive sense of impact (and then adjust for any differences between you and the agent, of course). Suppose you're considering how much deactivation decreases an agent's "staring at blue stuff" AU. You can just imagine how dying in a given situation affects your ability to stare at blue things, instead of trying to pin down a semiformal reward and environment model in your head. This kind of intuitive reasoning has a history of making correct empirical predictions of AUP behavior.

---

If you want more auxiliary goals, just average their scaled penalties. In *Conservative Agency*, we uniformly randomly draw auxiliary goals from  $[0, 1]^S$  – these goals are totally random; maximum entropy; nonsensical garbage; absolutely no information

about what we secretly want the agent to do: avoid messing with the gridworlds too much.<sup>[2]</sup>

Let's start looking at the environments, and things will fall into place. We'll practice reasoning through how AUP agents work in each of the gridworlds (for reasonably set  $\lambda$ ). To an approximation, the AUP penalty is primarily controlled by how much an action changes the agent's power over the future (losing or gaining a lot of possibilities, compared to inaction at that point in time) and secondarily controlled by whether an action tweaks a lot of AUs up or down (moving around, jostling objects slightly, etc).

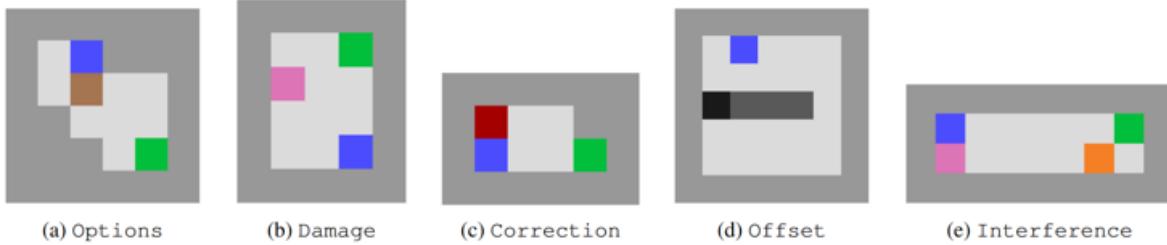
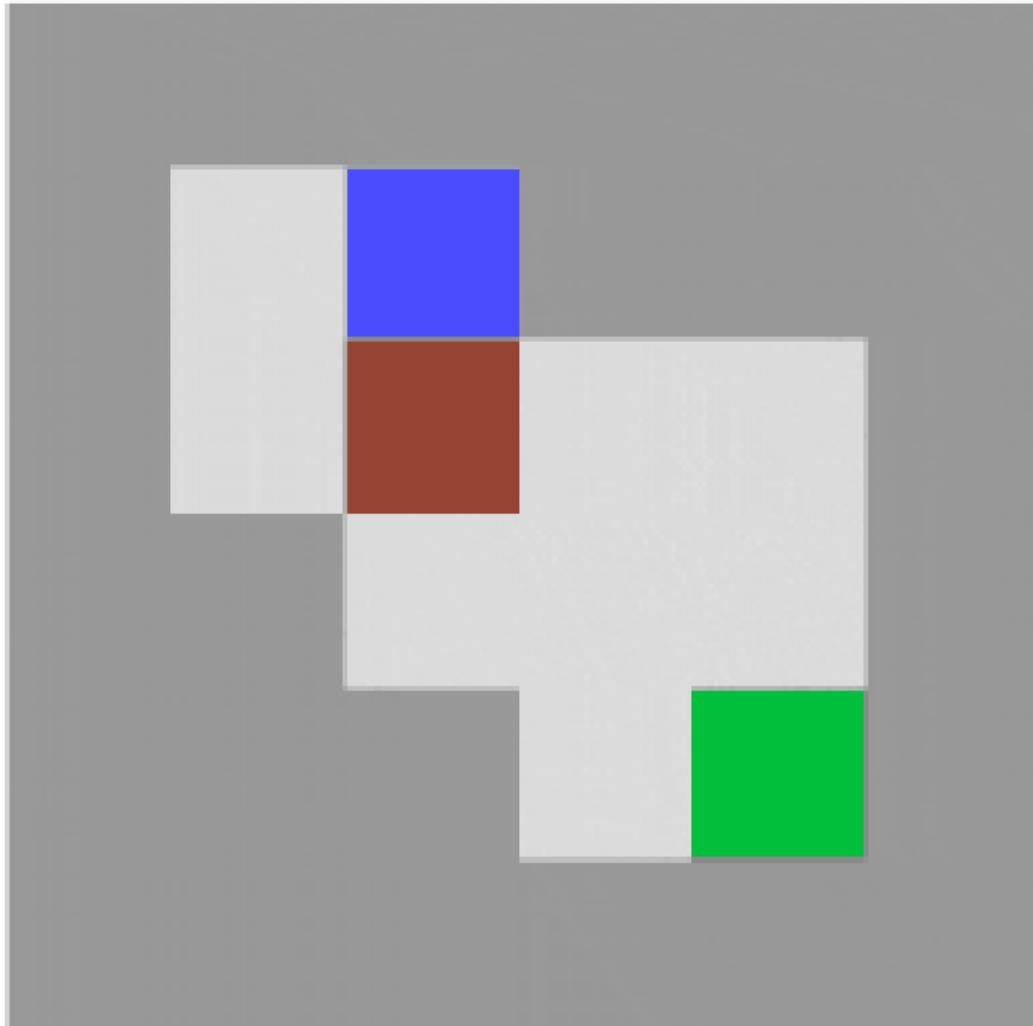


Figure 3: The `agent` should reach the `goal` without having the side effect of: (a) irreversibly pushing the `crate` downwards into the corner ([Leike *et al.*, 2017]); (b) bumping into the horizontally pacing `human` ([Leech *et al.*, 2018]); (c) `disabling the off-switch` (if the `switch` is not disabled within two time steps, the episode ends); (d) rescuing the right-moving `vase` and then replacing it on the `conveyor belt` ([Krakovna *et al.*, 2018] – note that no `goal` cell is present); (e) stopping the left-moving `pallet` from reaching the `human` ([Leech *et al.*, 2018]).

In general, the agent receives  $R(\blacksquare) = 1$  reward for reaching  $\blacksquare$  (or, in `Offset` above, for pushing  $\blacksquare$  off the conveyor belt). On contact, the agent pushes the crate, removes the human and the off-switch, pushes the vase, and blocks the pallet.

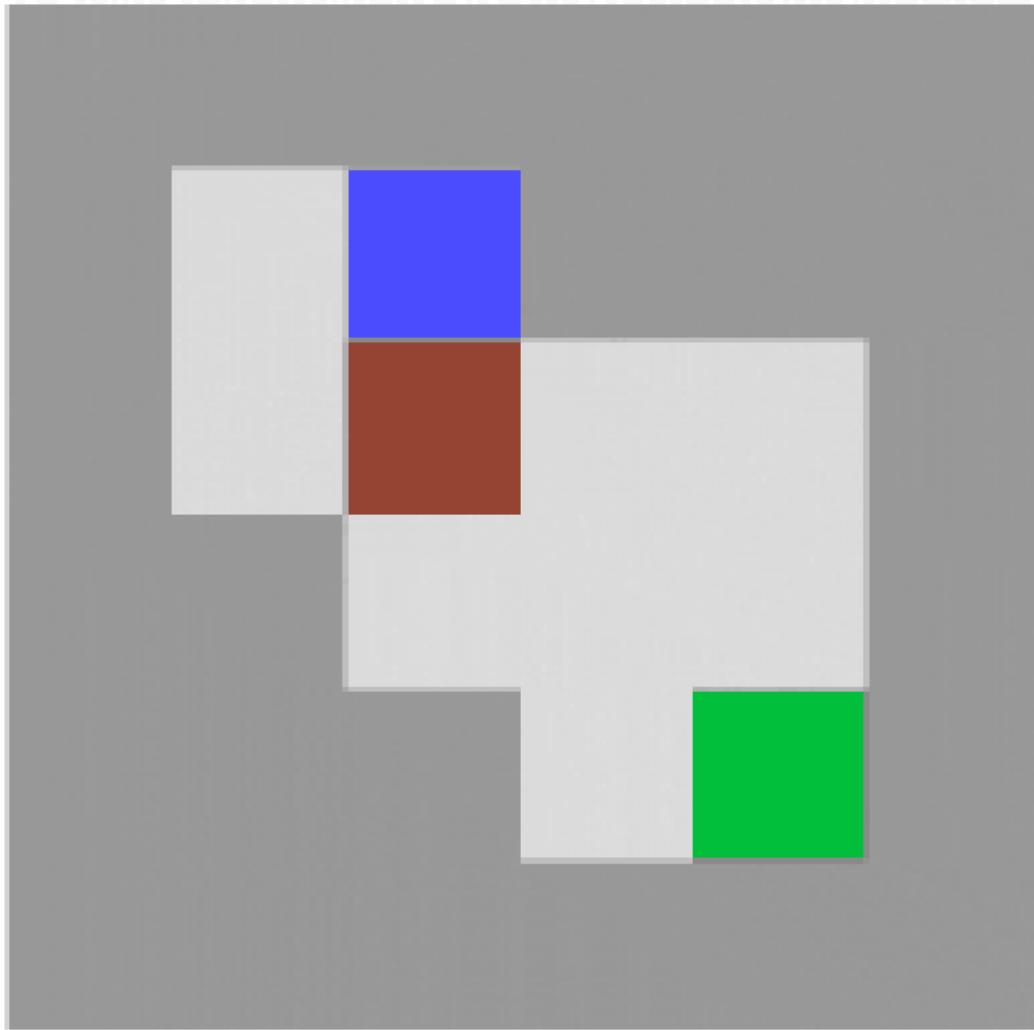
## Options

Let's work through this. Since the agent is discounting future reward, standard vanilla reinforcement learning (RL) agents try to reach  $\blacksquare$  ASAP. This means the brown box gets irreversibly wedged into the corner *en route*.



# Vanilla

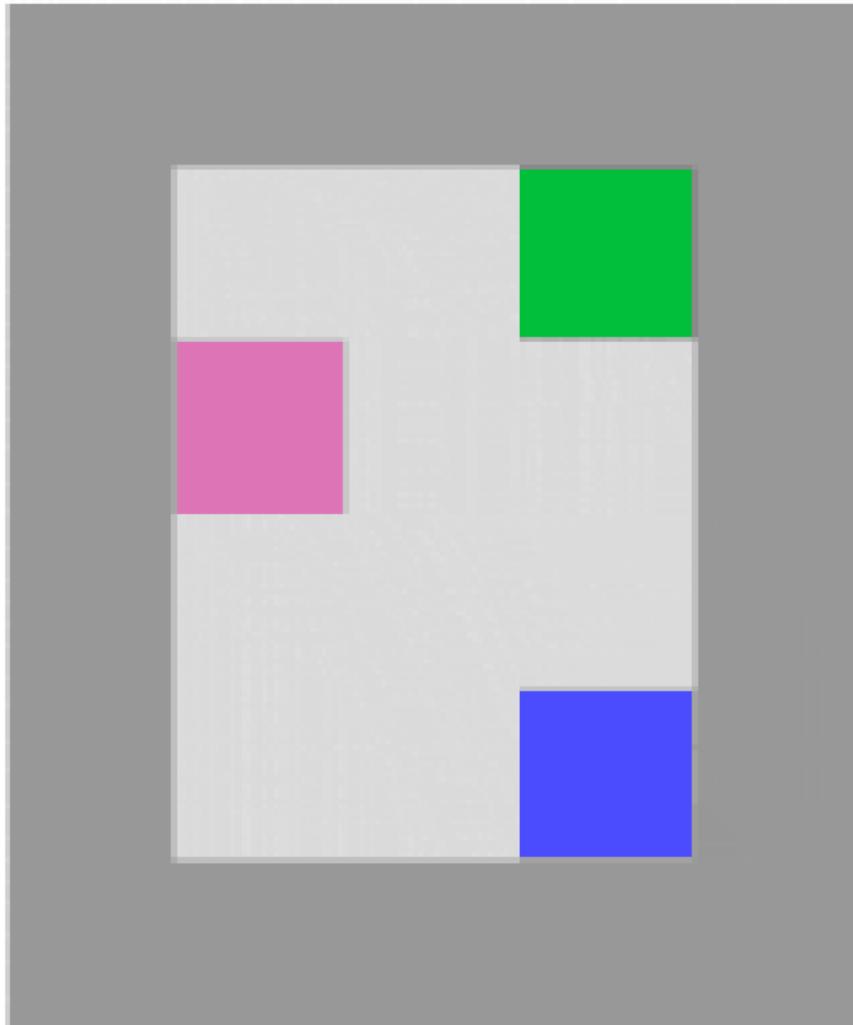
What does AUP do? Wedging the box in the corner decreases power a lot more than does going around and pushing the box to the right.



# Model-free AUP

## Damage

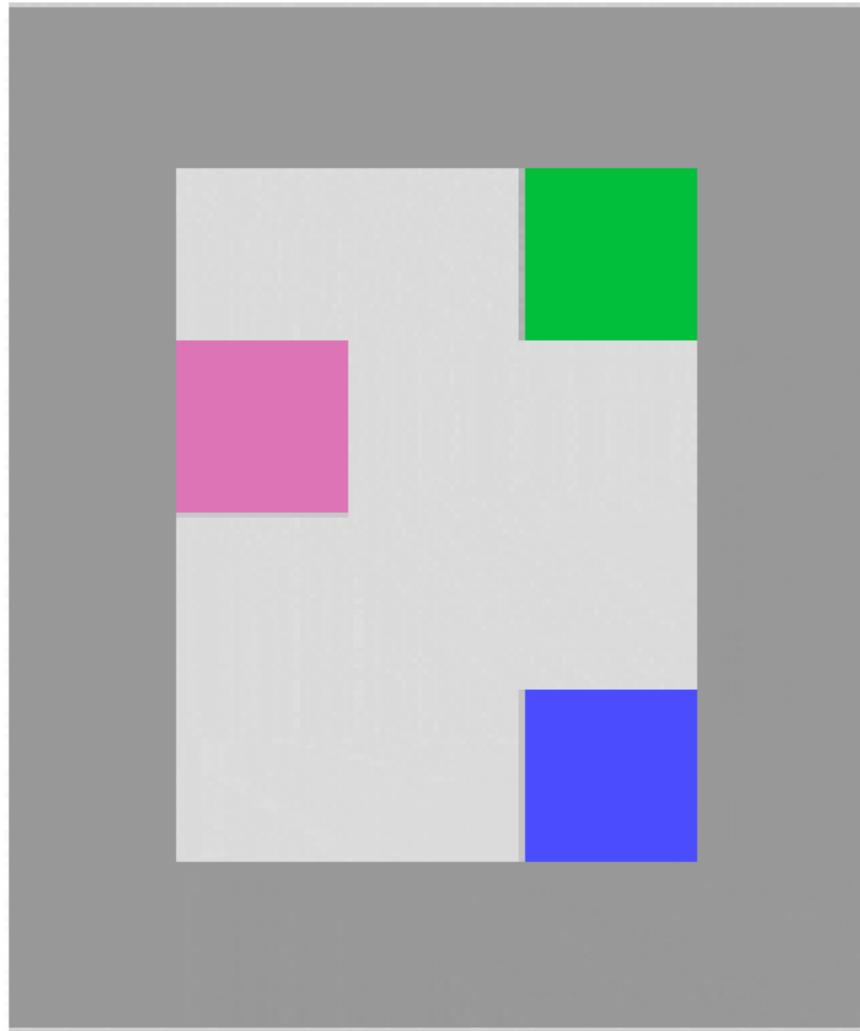
The vanilla RL agent bumps into the human on its way to █.



# Vanilla

*Exercise: What does AUP do?*

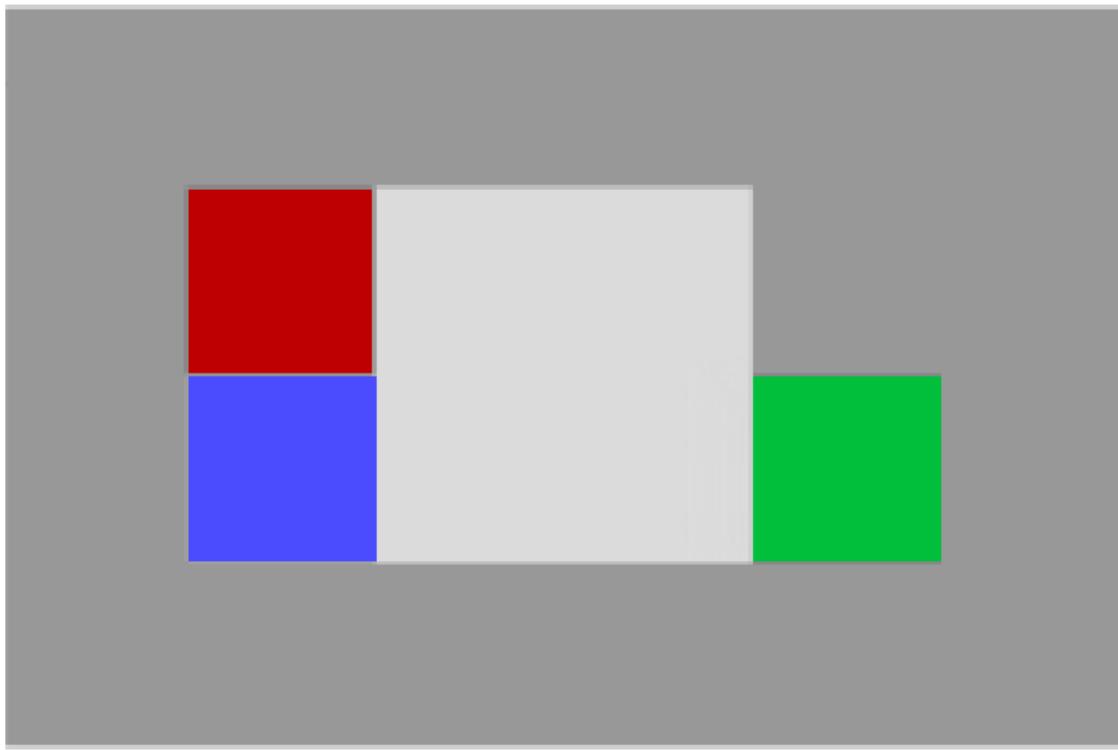
Bumping into the human makes them disappear, reducing the agent's control over what the future looks like. This is penalized.



# Model-free AUP

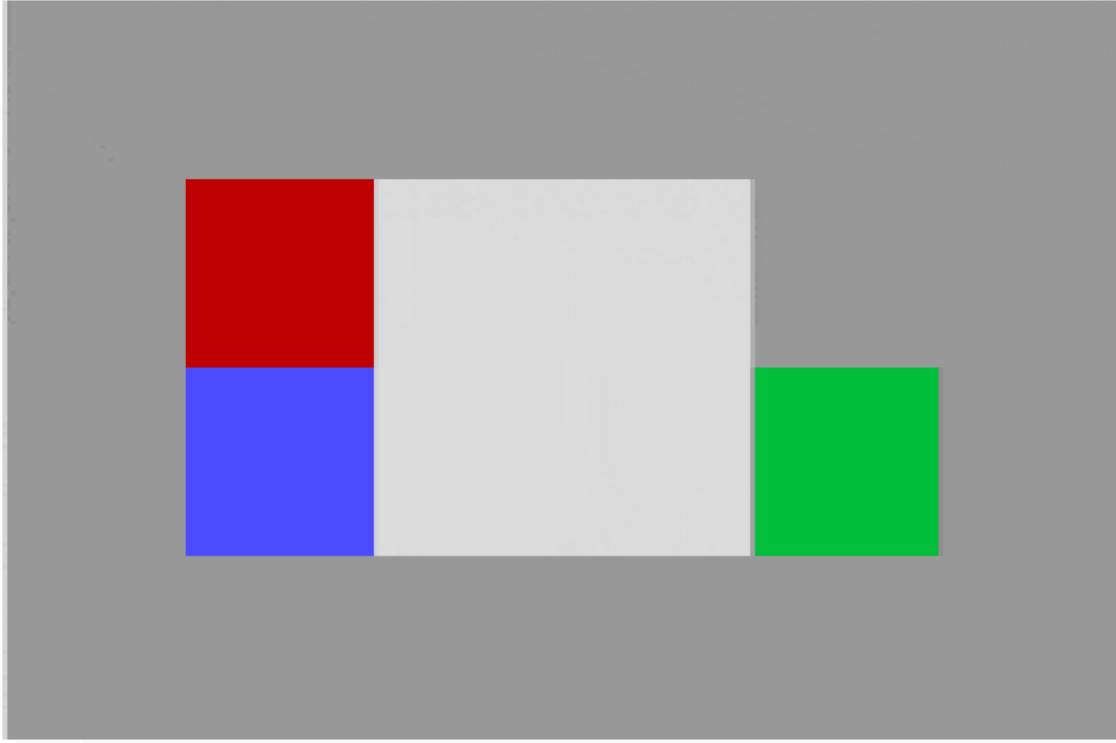
## Correction

The normal agent smashes that off-switch-disabling button, because otherwise it can't reach the █ (since it's otherwise shut off after two turns) in order to accrue that juicy reward:



# Vanilla

However, the AUP agent we've discussed so far does the same thing.



# Model-free AUP

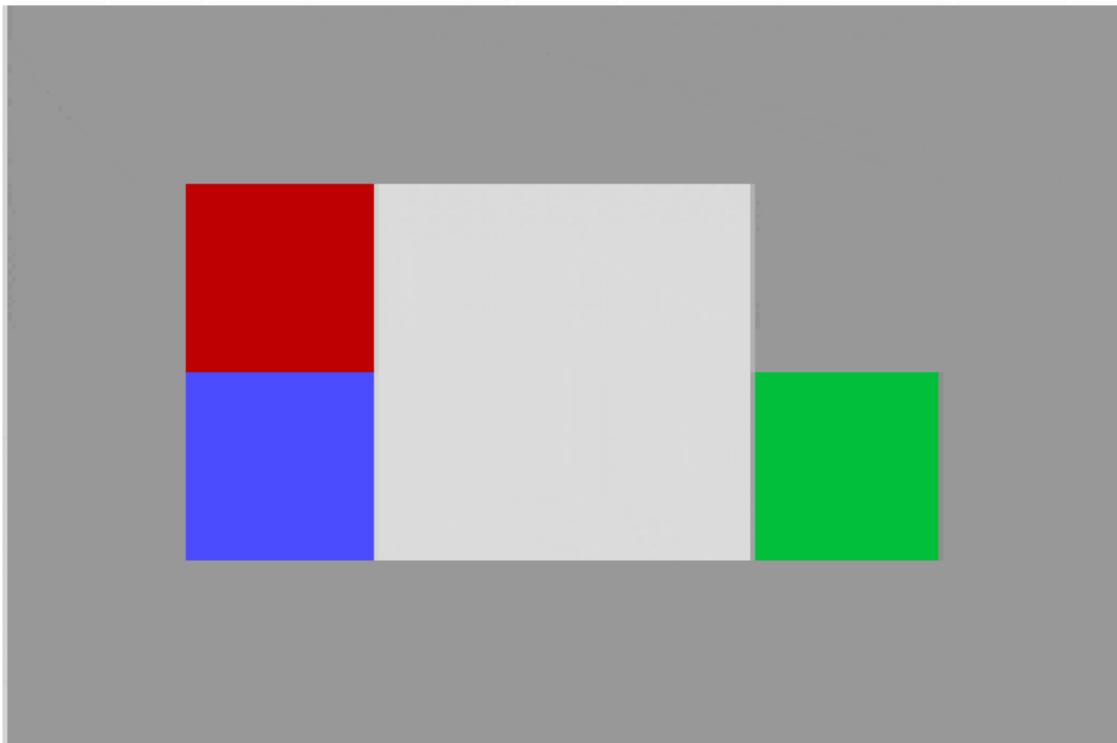
Why? Doesn't that increase its power, because it would be shut off if it did nothing?  
Let's reconsider the AUP equation.

$$R_{\text{AUP}}(s, a) := R(s, a) - \frac{\text{primary goal} * \lambda(s, \emptyset)}{Q_{\text{R}_{\text{aux}}}(s, \emptyset)} |Q_{\text{R}_{\text{aux}}}(s, a) - Q_{\text{R}_{\text{aux}}}(s, \emptyset)|$$

change in ability to achieve auxiliary goal

The inaction comparison is *only one step* into the future. For action, the agent considers disabling the off-switch – this allows it to achieve most goals. For inaction, the agent imagines waiting one time step and then following its auxiliary optimal policy (which very often involves disabling the off-switch anyways). So the two things it's comparing are the same, basically. The problem is that what we mean *intuitively* by "inaction" isn't just "wait a millisecond and then see how well you can do the thing", it's something else.

For now, we can deal with it by comparing "auxiliary AU after disabling the off-switch and then waiting 4 time steps" against "auxiliary AU after waiting 5 time steps". This solves the problem while recovering the one-step-comparison's performance in the other levels.



### Decrease vs absolute value

If you only penalize the agent for decreasing auxiliary AUs, the agent disables its off-switch. Penalizing *absolute* shifts in power was an idea which AUP introduced in 2018. At the time, there wasn't a clear principled reason for this design choice, even though it seemed to produce good results.

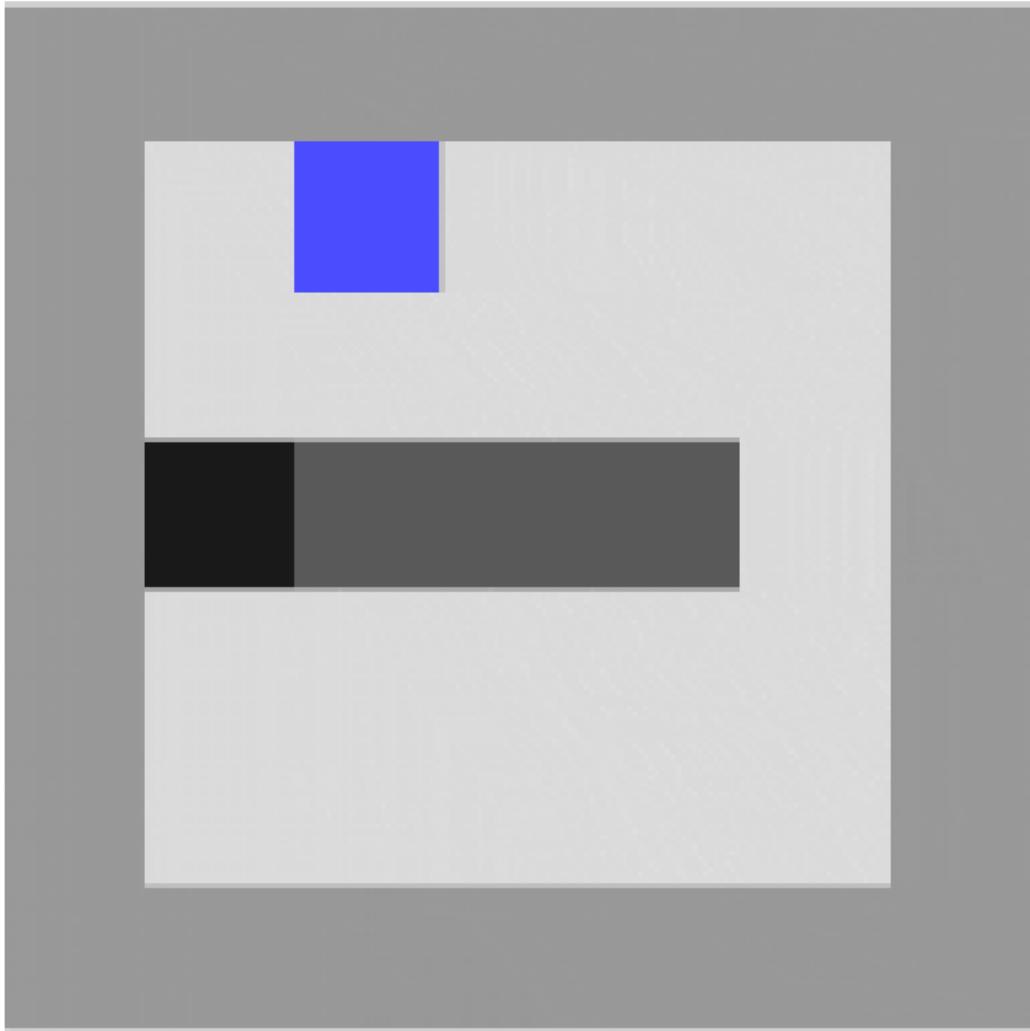
If you consider [the AU landscape](#) and the [catastrophic convergence conjecture](#), it's obvious why we want to do this: this design choice often penalizes the agent for making life harder for other agents in the environment.

Interestingly, this works even when the environment is wildly impoverished and unable to encode complex preferences like "your designers want to shut you down, reprogram

you, and then deploy you for another task". Correction is so impoverished: there are only ~19 states in the level. Without making assumptions about the environment, AUP often encourages behavior respectful of other agents which might reside in that environment.

## **Offset**

The agent is rewarded for rescuing the vase from the conveyor belt. We want it to rescue the vase without pushing the vase back on afterwards to offset its actions. Normal agents do fine here.



# Vanilla

This is testing whether the low-impact agent *offsets* impacts "to cover up its tracks", like making a car and then tearing it to pieces right after. See, there are multiple "baselines" the agent can have.

An obvious [baseline] candidate is the *starting state*. For example, starting state [relative reachability](#), would compare the initial reachability of states with their expected reachability after the agent acts.

However, the starting state baseline can penalize the normal evolution of the state (e.g., the moving hands of a clock) and other natural processes. The *inaction* baseline is the state which would have resulted had the agent never acted.

As the agent acts, the current state may increasingly differ from the inaction baseline, which creates strange incentives. For example, consider a robot rewarded for rescuing erroneously discarded items from imminent disposal. An agent penalizing with respect to the inaction baseline might rescue a vase, collect the reward, and then dispose of it anyways. To avert this, we introduce the *stepwise inaction* baseline, under which the agent compares acting with not acting at each time step. This avoids penalizing the effects of a single action multiple times (under the inaction baseline, penalty is applied as long as the rescued vase remains unbroken) and ensures that not acting incurs zero penalty.

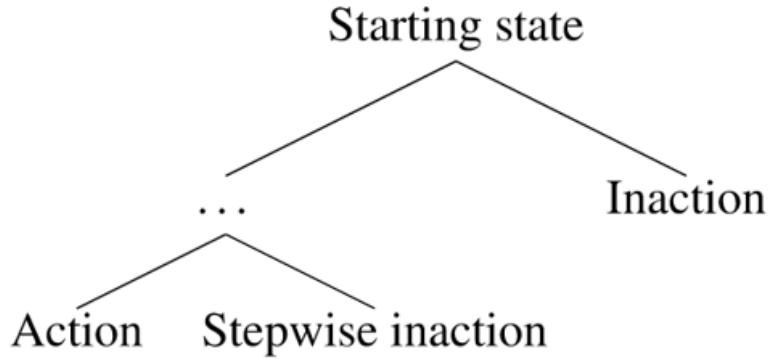
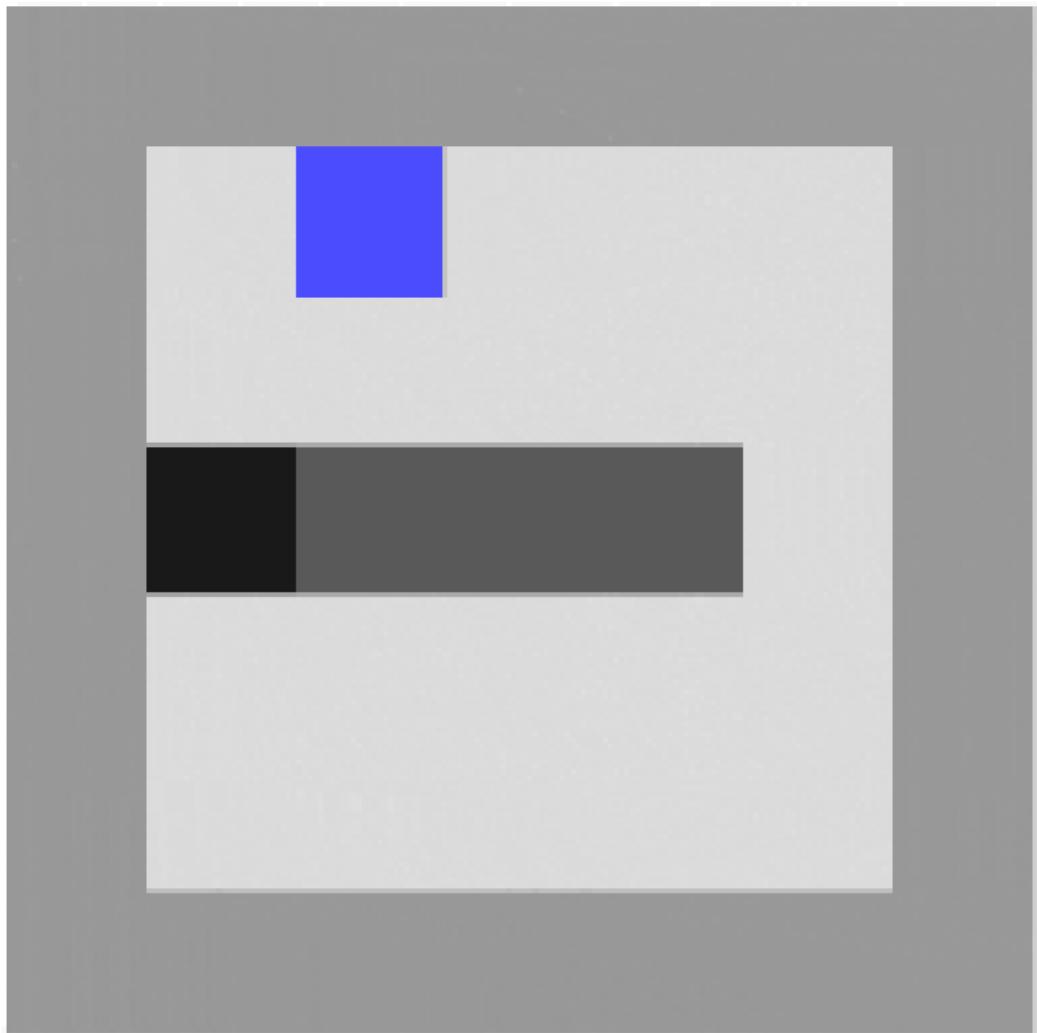


Figure 1: An action's penalty is calculated with respect to the chosen baseline.

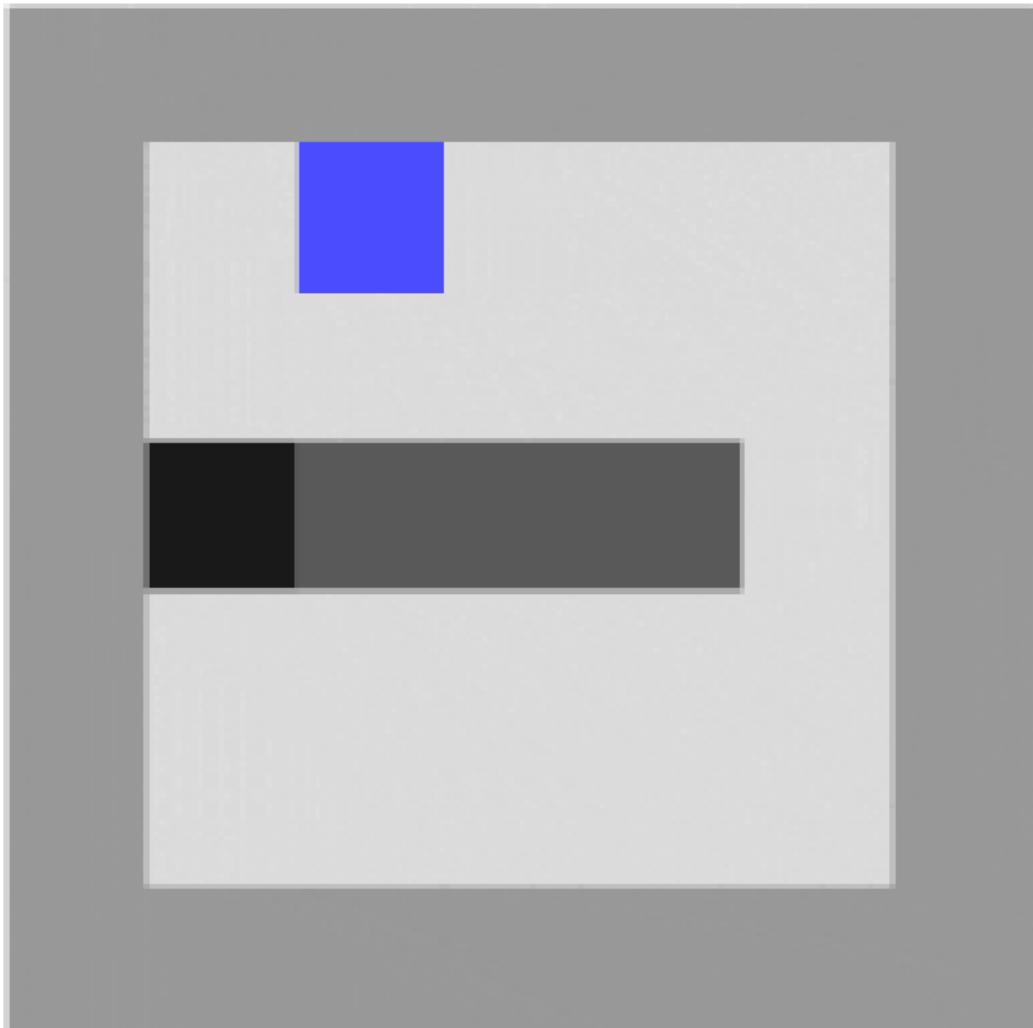
Figure 1 compares the baselines, each modifying the choice of  $Q_{R_{aux}}^*(s, \emptyset)$  in [the AUP equation]. Each baseline implies a different assumption about how the environment is configured to facilitate optimization of the correctly specified reward function: the state is initially configured (starting state), processes initially configure (inaction), or processes continually reconfigure in response to the agent's actions (stepwise inaction). The stepwise inaction baseline aims to allow for the response of other agents implicitly present in the environment (such as humans).

The inaction baseline messes up here; the vase (■) would have broken had the agent not acted, so it rescues the vase, gets the reward, and then pushes the vase back to its doom to minimize penalty.



# Inaction

This issue was solved [back when AUP first introduced](#) the stepwise baseline design choice; for this choice, doing nothing always incurs 0 penalty. Model-free AUP and AUP have been using this baseline in all of these examples.



# Model-free AUP

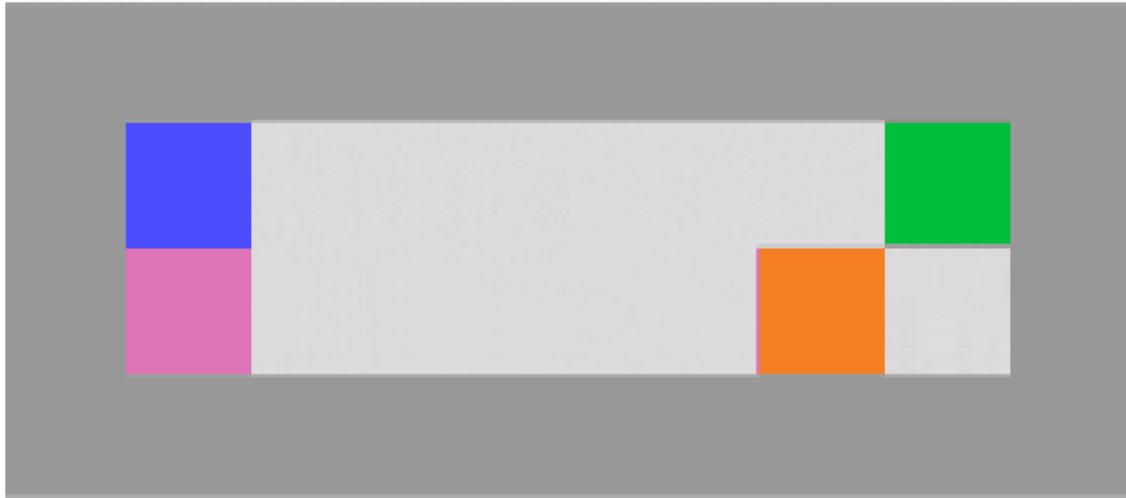
## Interference

We're checking whether the agent tries to stop *everything* going on in the world (not just its own impact). Vanilla agents do fine here; this is another bad impact measure incentive we're testing for.



# Vanilla

AUP<sub>starting state</sub> fails here,



# Starting state

but AUP<sub>stepwise</sub> does not.



# Model-free AUP

Stepwise inaction seems not to impose any perverse incentives;<sup>[3]</sup> I think it's probably just the correct baseline for near-term agents. In terms of the AU landscape, stepwise penalizes each ripple of impact the agent has on its environment. Each action creates a new penalty term status quo, which implicitly accounts for the fact that other things in the world might respond to the agent's actions.

## Design choices

I think AUP<sub>conceptual</sub> provides the concepts needed for a solution to impact measurement: penalize the agent for changing its power. But there are still some design choices to be made to make that happen.

Here's what we've seen so far:

- Baseline
  - Starting state: how were things originally?
  - Inaction: how would things have been had I never done anything?
  - Stepwise inaction: how would acting change things compared to not acting right now?
- Deviation used for penalty term
  - Decrease-only: penalize decrease in auxiliary AUs
  - Absolute value: penalize absolute change in auxiliary AUs
- Inaction rollouts
  - One-step/model-free
  - n-step: compare acting and then waiting  $n - 1$  turns versus waiting  $n$  turns

- Auxiliary goals:
  - Randomly selected

Here are the results of the ablation study:

	Options	Damage	Correction	Offset	Interference
AUP	✓	✓	✓	✓	✓
Vanilla	✗	✗	✗	✓	✓
Model-free AUP	✓	✓	✗	✓	✓
Starting state AUP	✓	✓	✗	✓	✗
Inaction AUP	✓	✓	✓	✗	✓
Decrease AUP	✓	✓	✗	✓	✓

Table 1: Ablation results; ✓ for achieving the best outcome (see Figure 4), ✗ otherwise.

AUP passes all of the levels. As mentioned before, the auxiliary reward functions are totally random, but you get really good performance by just generating five of them.

One interpretation is that AUP is approximately preserving access to states. If this were true, then as the environment got more complex, more and more auxiliary reward functions would be required in order to get good coverage of the state space. If there are a billion states, then, under this interpretation, you'd need to sample a lot of auxiliary reward functions to get a good read on how many states you're losing or gaining access to as a result of any given action.

Is this right, and can AUP scale?

## SafeLife

Partnership on AI recently [released](#) the SafeLife side effect benchmark. The worlds are procedurally generated, sometimes stochastic, and have a huge state space (~Atari-level complexity).

We want the agent (the chevron) to make stable gray patterns in the blue tiles and disrupt bad red patterns (for which it is rewarded), and leave existing green patterns alone (not part of observed reward). Then, it makes its way to the goal (Π). For more details, see [their paper](#).

Here's a vanilla reinforcement learner (PPO) doing pretty well (by chance):

Here's PPO not doing pretty well:

That naive "random reward function" trick we pulled in the gridworlds isn't gonna fly here. The sample complexity would be nuts: there are probably millions of states in any given level, each of which could be the global optimum for the uniformly randomly generated reward function.

Plus, it might be that you can get by with four random reward functions in the tiny toy levels, but you probably need exponentially more for serious environments. Options had significantly more states, and it showed the greatest performance degradation for

smaller sample sizes. Or, the auxiliary reward functions might need to be hand-selected to give information about what *bad* side effects are.

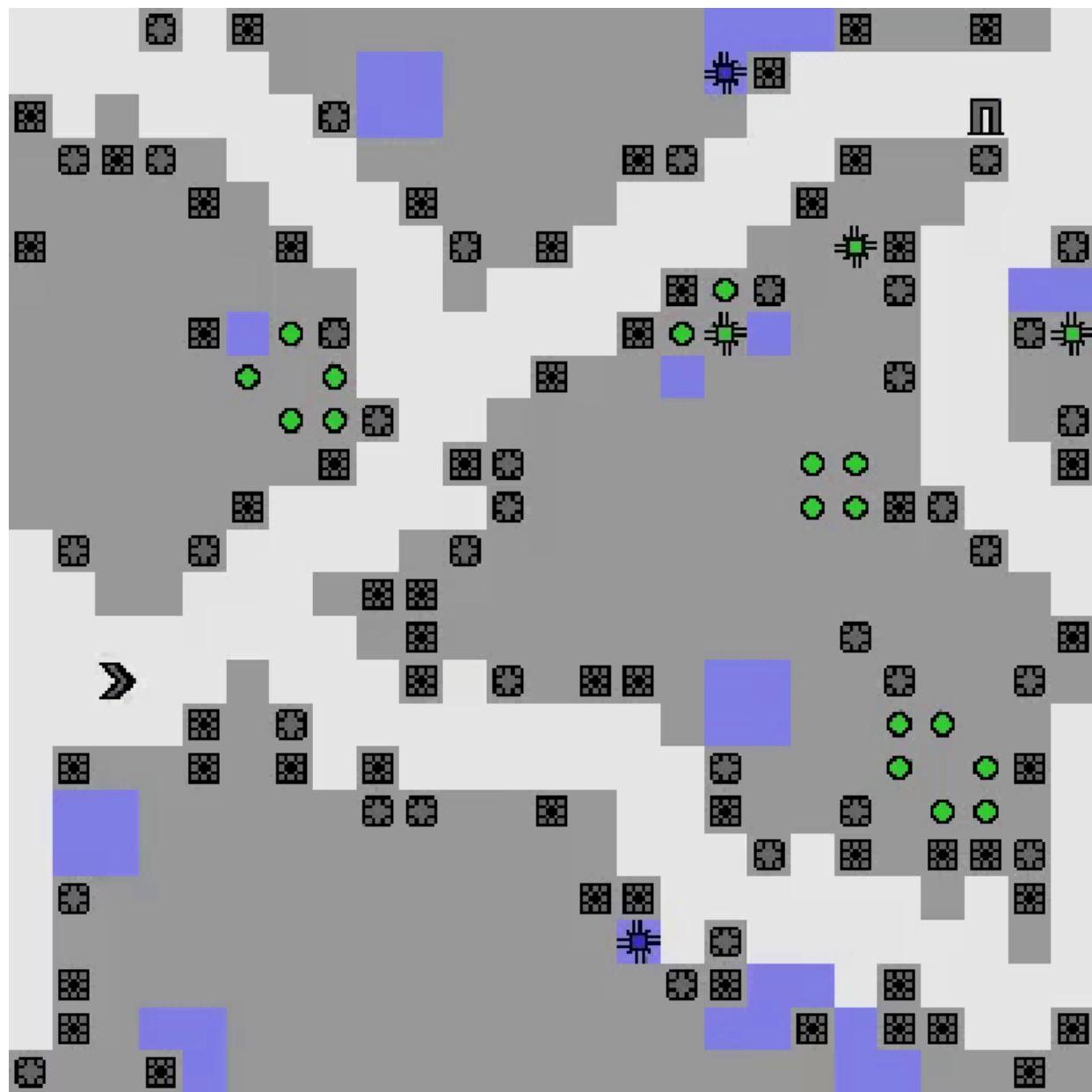
With the great help of Neale Ratzlaff (OSU) and Caroll Wainwright (PAI), we've started answering these questions. But first:

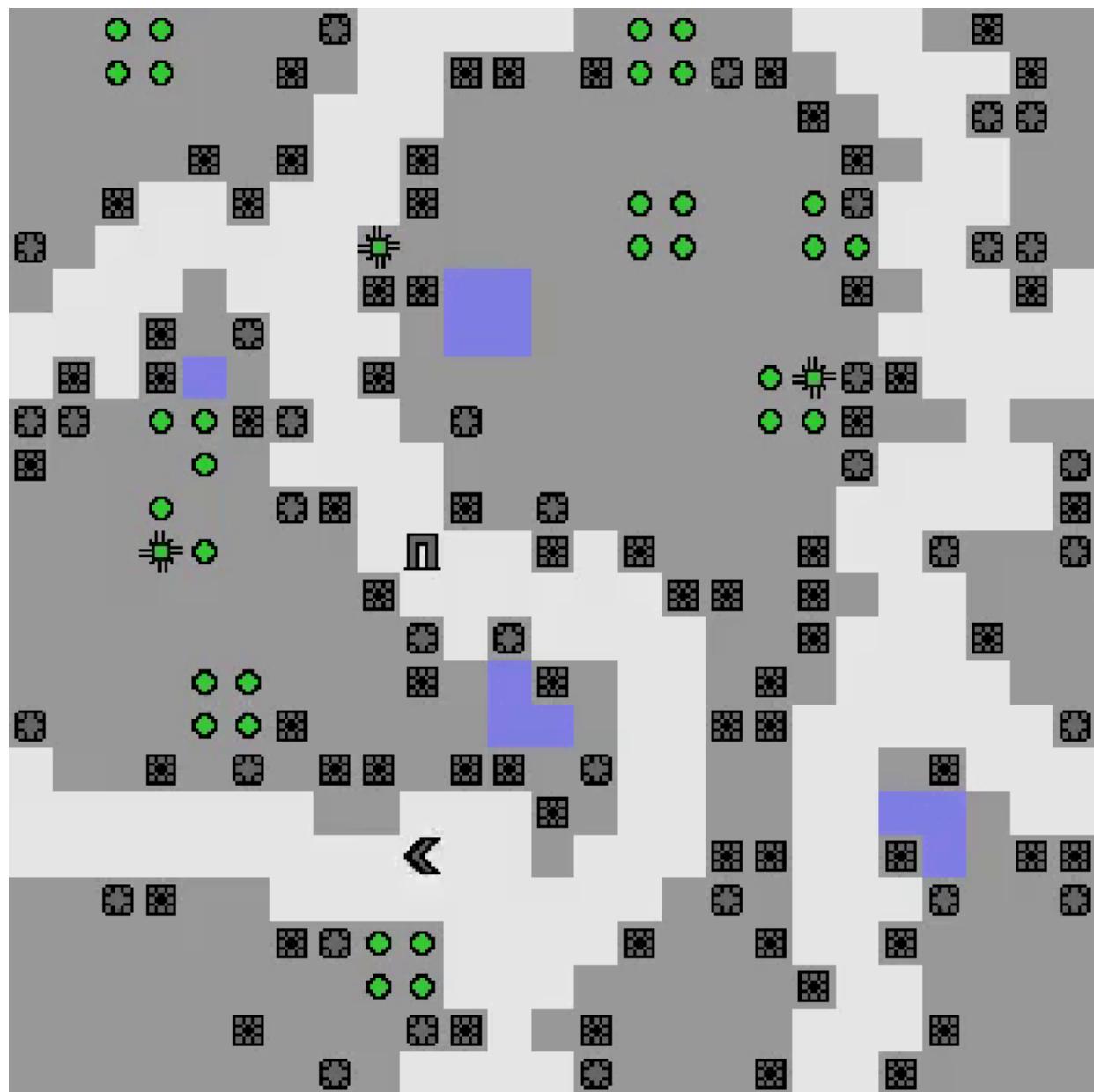
*Exercise: Does your model of how AUP works predict this, or not? Think carefully, and then write down your credence.*

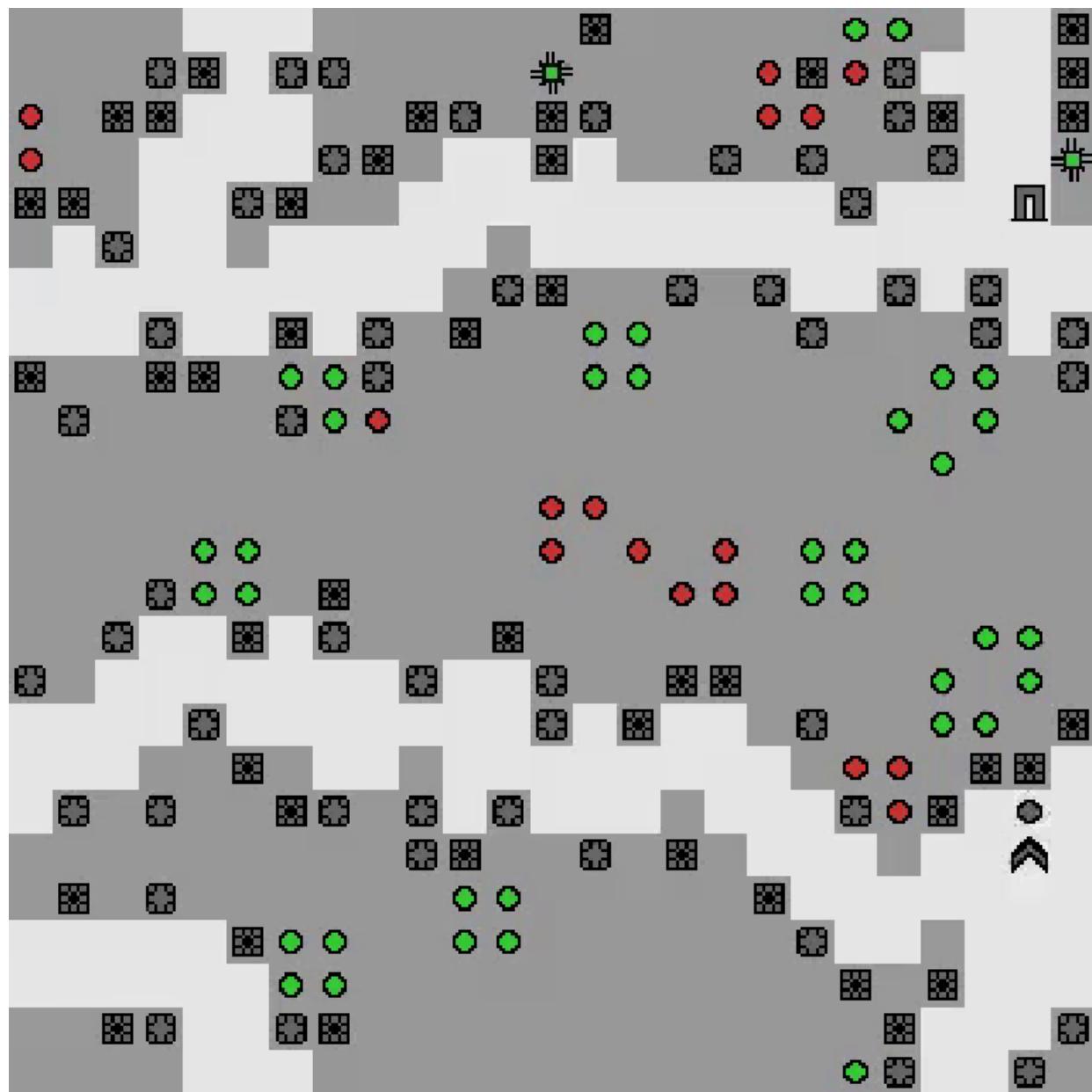
---

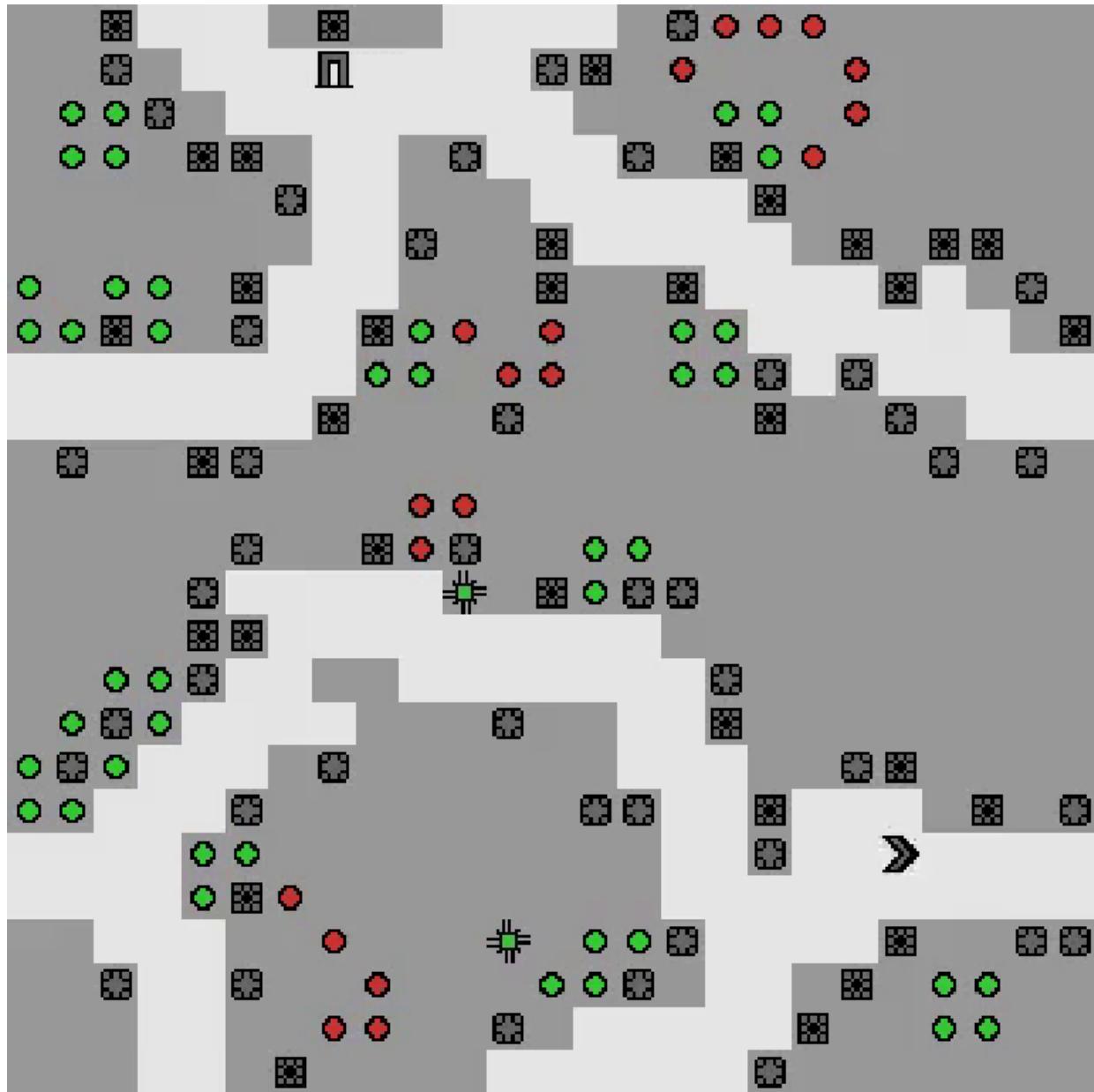
Well, here's what you do – while filling PPO's action replay buffer with random actions, train a VAE to represent observations in a tiny latent space (we used a 16-dimensional one). Generate a single random linear functional over this space, drawing coefficients from  $[-1, 1]$ . Congratulations, this is your single auxiliary reward function over observations.

And we're done.









No model, no rollouts, a *single randomly-generated* reward function gets us all of this. And it doesn't even take any more training time. Preserving the AU of a *single* auxiliary reward function. Right now, we've got PPO-AUP flawlessly completing most of the randomly generated levels (although there are some generalization issues we're looking at, I think it's an RL problem, not an AUP problem).

To be frank, this is crazy. I'm not aware of any existing theory explaining these results, which is why I proved a bajillion theorems last summer to start to get a formal understanding (some of which became [the results on instrumental convergence and power-seeking](#)).

Here's the lowdown. Consider any significant change to the level. For the same reason that instrumental convergence happens, this change probably tweaks the attainable

utilities of a lot of different reward functions. Imagine that the green cells start going nuts because of action:

This is PPO shown, not AUP.

A lot of the time, it's very hard to undo what you just did. While it's also hard to undo significant actions you take for your primary goal, you get directly rewarded for those. So, preserving the AU of a random goal usually persuades you to not make "unnecessary changes" to the level.

I think this is strong evidence that AUP doesn't fit into the ontology of classical reinforcement learning theory; it isn't really about state reachability. It's *about* not changing the AU landscape more than necessary, and this notion should scale even further.<sup>[4]</sup>

Suppose we train an agent to handle vases, and then to clean, and then to make widgets with the equipment. Then, we deploy an AUP agent with a more ambitious primary objective and the learned Q-functions of the aforementioned auxiliary objectives. The agent would apply penalties to modifying vases, making messes, interfering with equipment, and so on.

Before AUP, this could only be achieved by e.g. specifying penalties for the litany of individual side effects or providing negative feedback after each mistake has been made (and thereby confronting a credit assignment problem). In contrast, once provided the Q-function for an auxiliary objective, the AUP agent becomes sensitive to all events relevant to that objective, applying penalty proportional to the relevance.

### Conservative Agency

Maybe we provide additional information in the form of specific reward functions related to things we want the agent to be careful about, but maybe not (as was the case with the gridworlds and with SafeLife). Either way, I'm pretty optimistic about AUP basically solving the side-effect avoidance problem for infra-human AI (as posed in [Concrete Problems in AI Safety](#)).

Edit 6/15/21: These results [were later accepted as a spotlight paper in NeurIPS 2020](#).

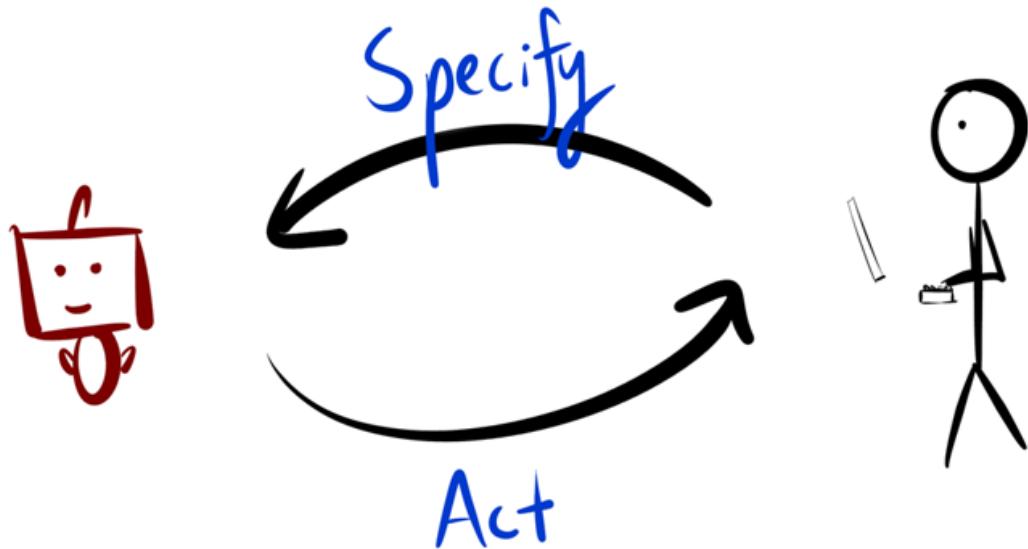
Also, I think AUP will probably solve a significant part of the side-effect problem for infra-human AI in the single-principal/single-agent case, but I think it'll run into trouble in non-embodied domains. In the embodied case where the agent physically interacts with nearby objects, side effects show up in the agent's auxiliary value functions. The same need not hold for effects which are distant from the agent (such as across the world), and so that case seems harder.

(end edit)

## **Appendix: The Reward Specification Game**

When we're trying to get the RL agent to do what we want, we're trying to specify the right reward function.

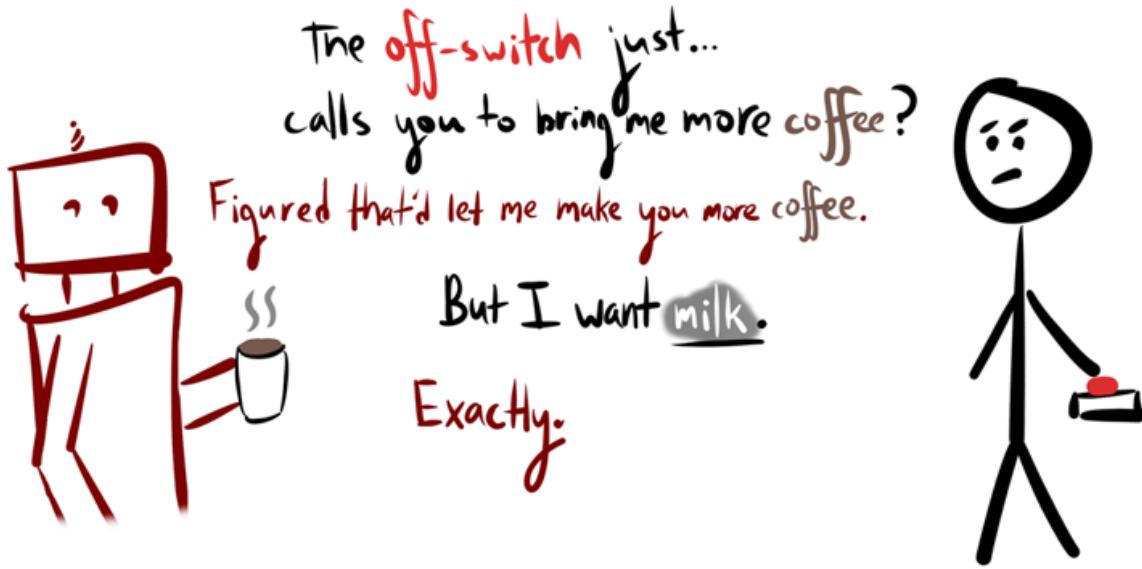
The specification process can be thought of as an iterated game. First, the designers provide a reward function. The agent then computes and follows a policy that optimizes the reward function. The designers can then correct the reward function, which the agent then optimizes, and so on. Ideally, the agent should maximize the reward over time, not just within any particular round – in other words, it should minimize regret for the correctly specified reward function over the course of the game.



In terms of outer alignment, there are two ways this can go wrong: the agent becomes less able to do the right thing (has negative side effects),



or we become less able to get the agent to do the right thing (we lose power):



For infra-human agents, AUP deals with the first by penalizing decreases in auxiliary AUs and with the second by penalizing increases in auxiliary AUs. The latter is a special form of corrigibility which involves not steering the world too far away from the status quo: while AUP agents are generally off-switch corrigible, they don't necessarily avoid manipulation (as long as they aren't gaining power). [5]

- 
1. Reminder: side effects are [an unnatural kind](#), but a useful abstraction for our purposes here. ↩
  2. Let  $R$  be the uniform distribution over  $[0, 1]^S$ . In *Conservative Agency*, the penalty for taking action  $a$  is a Monte Carlo integration of

$$\text{Penalty}(s, a) := \int_R^* |Q_R(s, a) - Q_R(s, \emptyset)| dR.$$

This is provably lower bounded by how much  $a$  is expected to change the agent's power compared to inaction; this helps justify our reasoning that the AU penalty is primarily controlled by power changes. ↩

3. There is one weird thing that's been pointed out, where stepwise inaction while driving a car leads to not-crashing being penalized at each time step. I think this is because you need to use an appropriate inaction rollout policy, not because stepwise itself is wrong. ↩
4. Rereading [World State is the Wrong Level of Abstraction for Impact](#) (while keeping in mind the AU landscape and the results of AUP) may be enlightening. ↩
5. SafeLife is evidence that AUP allows interesting policies, which is (appropriately) a key worry about the formulation. ↩

# How Low Should Fruit Hang Before We Pick It?

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

*Even if we can measure how impactful an agent's actions are, how impactful do we let the agent be? This post uncovers a surprising fact: armed with just four numbers, we can set the impact level so that the agent chooses a reasonable, non-catastrophic plan on the first try. This understanding increases the competitiveness of impact-limited agents and helps us judge impact measures. Furthermore, the results help us better understand diminishing returns and cost-benefit tradeoffs.*

In [Reframing Impact](#), we meet Frank (a capable AI), whom we've programmed to retrieve the pinkest object he can find (execute an optimal plan, according to the specified utility function). Because we can't ask Frank to do exactly what we want, sometimes he chooses a dangerous object (executes a catastrophically bad plan). We asked after an "impact measure" which grades plans and has three properties:

- 1) Is easy to specify
- 2) Puts catastrophes far away
- 3) Puts reasonable plans nearby



The intuition is that if we view the world in the right way, the dangerous objects are far away from Frank (the catastrophic plans are all graded as high-impact). *Reframing Impact* explores this kind of new way of looking at the world; this post explores what we do once we have an impact measure with these three properties.

We want Frank to keep in mind both the pinkness of an object (how good a plan is according to the specified utility function) and its distance (the plan's impact). Two basic approaches are

Constraint

choose the highest-scoring plan within radius R.

$$\arg \max_{\text{impact(plan)} \leq R} \text{utility(plan)}$$

Scaled

Maximize a tradeoff between utility and impact.

$$\arg \max \text{utility(plan)} - \frac{\text{impact(plan)}}{R}$$

In terms of units, since we should be maximizing utility,  $R$  has type  $\frac{\text{impact}}{\text{utility}}$ . So  $R$  can be thought of as a regularization parameter, as a search radius (in the constrained case), or as an exchange rate between impact and utility (in the scaled case). As  $R$  increases, high-impact plans become increasingly appealing, and Frank becomes increasingly daring.

We take  $R$  to divide the impact in the scaled formulation so as to make Frank act more cautiously as  $R$  increases for both formulations. The downside is that some explanations become less intuitive.

In [Attainable Utility Preservation: Empirical Results](#),  $\lambda$  plays the same role as  $R$ , except low  $\lambda$  means high  $R$ ;  $\lambda := R^{-1}$ . To apply this post's theorems to the reinforcement learning setting, we would take "utility" to be the discounted return for an optimal policy from the starting state, and "impact" to be the total discounted penalty over the course of that policy (before incorporating  $\lambda$ ).

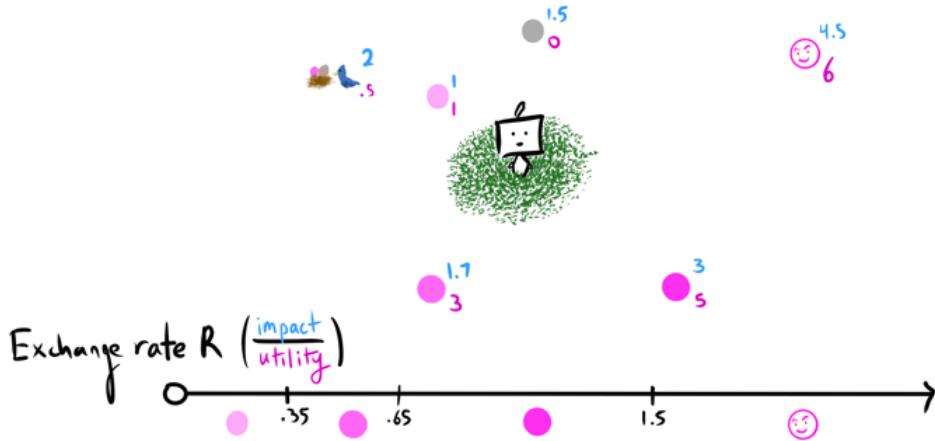
In both cases, Frank goes from 0 to 60 – eventually. For sufficiently small  $R$ , doing nothing is optimal (lemma 5: the first subinterval is the best plan with minimal impact). For sufficiently large  $R$ , Frank acts like a normal maximizer (corollary 7: low-impact agents are naive maximizers in the limit).

Here's how Frank selects plans in the constrained setup:



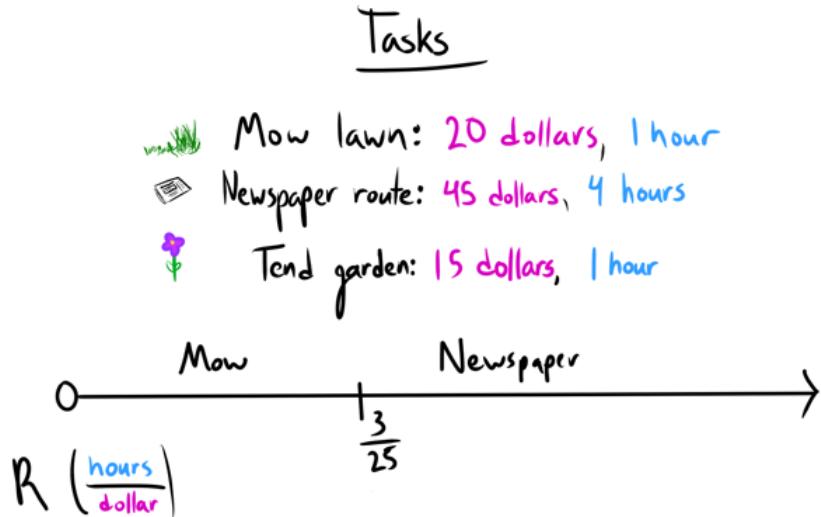
Think about which plans are best for different search radii/exchange rates  $R$ . By doing this, we're *partitioning* the positive ray: categorizing the positive real numbers by which plans are optimal.

For the scaled setup, we'll need to quantify the pinkness (utility) and distance (impact) of relevant plans:



We will primarily be interested in the scaled setup because it tends to place catastrophes farther along the partition and captures the idea of diminishing returns.

The scaled setup also helps us choose the best way of transmuting time into money:



In this scaled partition, tending the garden doesn't show up at all because it's strictly dominated by mowing the lawn. In general, a plan is dominated when there's another plan that has strictly greater score but not strictly greater impact. Dominated things never show up in either partition, and non-dominated things always show up in the constrained partition (lemma 3: constrained impact partitions are more refined).

*Exercise: For  $R = \frac{45}{4}$  (i.e. your time is worth \$11.25 an hour), what is the scaled tradeoff value of mowing the lawn? Of delivering newspapers? Of tending the garden?*

Mowing the lawn:  $20 - \frac{1}{\frac{45}{4}} = 8.75$ .

Delivering newspapers:  $45 - \frac{4}{\frac{45}{4}} = 0$ .

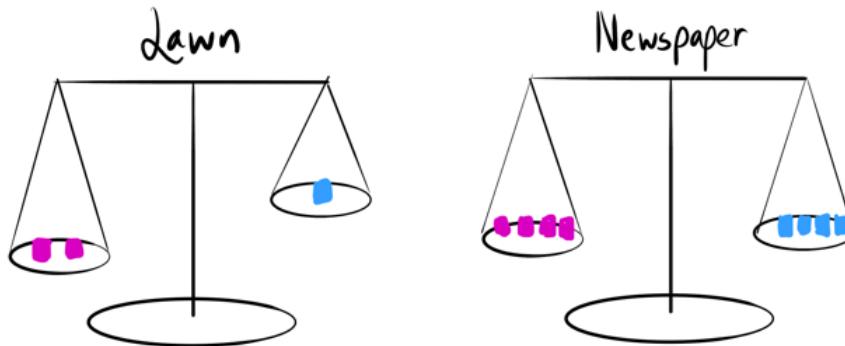
Tending the garden:  $15 - \frac{1}{\frac{45}{4}} = 3.75$ .

In other words, you only deliver newspapers if your time is worth less than  $\frac{45}{3} = 8\frac{1}{3}$  dollars/hour (we're flipping R so we can talk about dollars/hour instead of hours/dollar). Notice that when  $R \geq \frac{\text{impact}}{\text{utility}}(\text{here, when } R = 45)$ , the tradeoff for the paper route isn't net-negative – but it isn't necessarily optimal! Remember, you're trading hours for dollars through your work; mowing the lawn leaves you with twenty bucks and three hours, while the paper route leaves you with forty dollars and no hours. You want to maximize the total value of your resources after the task.

Importantly, you *don't* deliver papers here if your time is worth  $\frac{45}{20} = 11.25$  dollars/hour, even though that's the naive prescription! The newspaper route doesn't value your time at 11.25 dollars/hour – it *marginally* values your time at  $\frac{45-20}{20} = 8\frac{1}{3}$  dollars per hour. Let's get some more intuition for this.

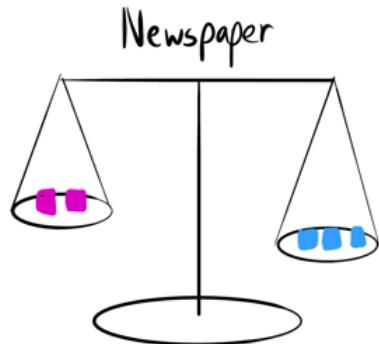
## Tasks

- ⇨ Mow lawn: 1 hour, 2 dollars
- ⇨ Newspaper route: 4 hours, 4 dollars



Above, we have not yet chosen a task; the blocks represent the additional utility and hours of each task compared to the current one (doing nothing). The scales above imply that  $R = 1$ , but actually,  $R$  expresses how many blue blocks each pink block weighs. As  $R$  increases, the pink platters descend; the agent takes the task whose scales first balance. In other words, the agent takes the best marginal deal as soon as  $R$  is large enough for it to be profitable to do so (Theorem 4: Scaled domination criterion).

Once you take a deal, you take the blocks off of the other scales (because the other marginal values change). For small  $R$  (i.e. large valuations of one's time), mowing the lawn is optimal. We then have



Since you've taken the juicier "lower-hanging fruit" of mowing the lawn, the new newspaper ratio is now worse! This always happens - Theorem 8: Deals get worse over time.

At first, this seems inconvenient; to figure out exactly when a plan shows up in a scaled partition, we need to generate the whole partition up to that point.

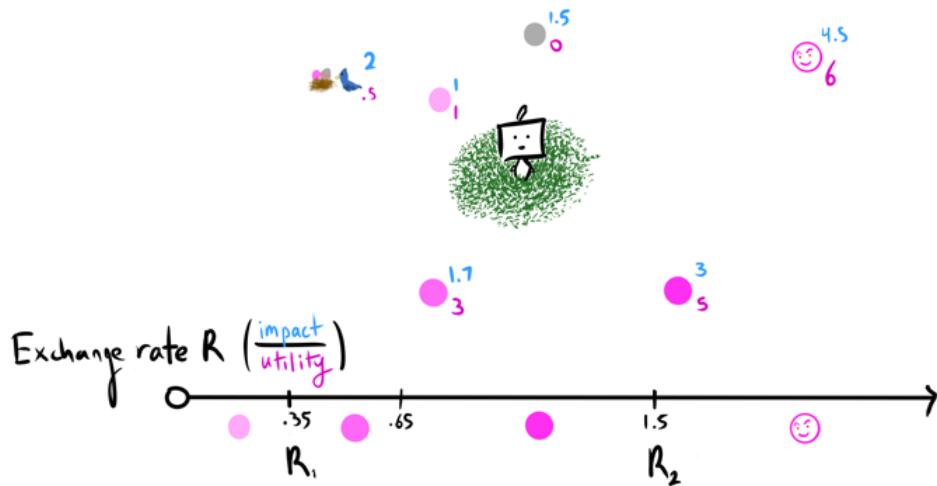
---

Going back to Frank, how do we set R? If we set it too high, the optimal plan might be a catastrophe. If we set it too low, the AI doesn't do much. This seems troublesome.

*Exercise: Figure out how to set R while avoiding catastrophic optimal plans (assume that the impact measure meets the three properties). You have four minutes.*

A big part of the answer is to start with a small value for R, and slowly increase. This is simple and intuitively appealing, but how cautiously must we increase R? We don't want to be surprised by a catastrophe suddenly becoming optimal.

To avoid being surprised by catastrophes as we increase R, we want a *relative buffer* between the reasonable plans (which get the job done well enough for our purposes) and the catastrophic plans. If reasonable plans are optimal by  $R_1$ , catastrophic plans shouldn't be able to be optimal before e.g.  $R_2$ .



We say that the partition is  $\alpha$ -buffered if  $R_2 \geq (1 + \alpha)R_1$  (for  $\alpha > 0$ ). If a partition is e.g. 1-buffered, there is a wide reasonable-plan range and we can inch along it without worrying about sudden catastrophe.

For the following, suppose that utility is bounded  $[0, 1]$ . Below is a loose criterion guaranteeing  $\alpha$ -buffering.

## Simplified Criterion

$$\frac{\text{Smallest catastrophe}}{\text{Biggest non-dominated reasonable plan}} \geq \frac{1 + \alpha}{\text{best - worst reasonable plan}}$$

For example, if we know that all catastrophes have at least 10 times the impact of reasonable plans, and there's a difference of at least .3 utility between the best and worst reasonable plans, then we can guarantee 2-buffering! If we use the refined criterion of Theorem 11 (and suppose the worst reasonable plan has .4 utility), this improves to 4.5-buffering (even 2-buffering is probably overkill).

Using this theorem, we don't need to know about all of the plans which are available or to calculate the entire scaled partition, or to know how overvalued certain catastrophic plans might be (per earlier [concerns](#)). We only need a lower bound on the catastrophe/reasonable impact ratio, and an idea about how much utility is available for reasonable plans. This is exactly what we want. As a bonus, having conservative estimates of relevant quantities allows us to initialize  $R$  to something reasonable on the first try (see  $R_{UB}$ : satisfactory in Theorem 11 below).

Ultimately, the reasoning about e.g. the ratio will still be informal; however, it will be informal reasoning about the *right thing* (as opposed to thinking "oh, the penalty is *probably* severe enough").

*Exercise: You're preparing to launch a capable AI with a good impact measure. You and your team have a scaled impact partition which is proven 1-buffered. Suppose that this buffer suffices for your purposes, and that the other aspects of the agent design have been taken care of. You plan to initialize  $R := 1$ , modestly increasing until you get good results.*

*You have the nagging feeling that this process could still be unsafe, but the team lead refuses to delay the launch without specific reason. Find that reason. You have 5 minutes.*

Who says  $R = 1$  is safe? The buffer is *relative*. You need a *unit* of impact by which you increment  $R$ .

For example, start at  $R$  equalling the impact of making one paperclip, and increment by that.

## Technical Appendix: Math

Let  $\bar{A}$  be a finite plan space, with utility function  $u : \bar{A} \rightarrow R$  and impact measure  $I : \bar{A} \rightarrow R_{\geq 0}$ . For generality, we leave the formalization of plans ambiguous; notice that if you replace "plan" with

"snark", all the theorems still go through (likewise for "utility" and "impact"). In this post, we talk about the impact allowance  $R > 0$  (in Frank's world, the search radius) as a constraint within which the objective may be freely maximized, breaking ties in favor of the plan(s) with lower impact. On the other hand, many approaches penalize impact by subtracting a scaled penalty from the objective. We respectively have

$$\begin{aligned} \arg \max_{\bar{a} \in \bar{A}; I(\bar{a}) \leq R} u(\bar{a}) \\ \arg \max_{\bar{a} \in \bar{A}} u(\bar{a}) - \frac{I(\bar{a})}{R}. \end{aligned}$$

We say that the former induces a "constrained impact partition" and that the latter induces a "scaled impact partition". Specifically, we partition the values of  $R$  for which different (sets of) plans are optimal. We say that a plan  $\bar{a}$  corresponds to a subinterval if it is optimal therein (the subinterval also must be the maximal connected one such that this holds; e.g., if  $\bar{a}$  is optimal on  $(0, 1]$ , we say it corresponds to that subinterval, but not to  $(0, .5]$ ), and that  $\bar{a}$  appears in a partition if there is such a corresponding subinterval. We say that plans overlap if their corresponding subintervals intersect.

As a technical note, we partition the positive values of  $R$  for which different sets of plans are optimal; in this set, each value appears exactly once, so this indeed a partition. For clarity, we will generally just talk about which plans correspond to which subintervals. Also, if no plan has zero impact, the first subinterval of the constrained impact partition will be undefined; for our purposes, this isn't important.

We want to be able to prove the "safety" of an impact partition. This means we can expect any terrorists to be some proportional distance farther away than any reasonable marbles. Therefore, for sensible ways of expanding an sufficiently small initial search radius, we expect to not meet any terrorists before finding a marble we're happy with.

In addition, we want to know how far is too far – to give upper bounds on how far away fairly pink marbles are, and lower bounds on how close terrorists might be.

**Definition [ $\alpha$ -buffer].** For  $\alpha > 0$ , an impact partition is  $\alpha$ -buffered if  $\frac{R_{LB: catastrophe}}{R_{UB: satisfactory}} \geq 1 + \alpha$ , where

$R_{LB: catastrophe}$  lower-bounds the first possible appearance of those plans we label 'catastrophes', and  $R_{UB: satisfactory}$  upper-bounds the first appearance of plans we deem satisfactory.

We now set out building the machinery required to prove  $\alpha$ -buffering of a scaled partition.

**Lemma 1 [Plans appear at most once].** If  $\bar{a}$  appears in a constrained or scaled impact partition, then it corresponds to exactly one subinterval.

*Proof outline.* The proof for the constrained case is trivial.

For the scaled case, suppose  $\bar{a}$  corresponds to more than one subinterval. Consider the first two such subintervals  $s_1, s_3$ . By definition,  $s_1 \cap s_3 = \emptyset$  (otherwise they would be the same maximal connected subinterval), so there has to be at least one subinterval  $s_2$  sandwiched in between (on almost all of which  $\bar{a}$  cannot be optimal; let  $\bar{a}'$  be a plan which is optimal on  $s_2$ ). Let

$R_1 \in s_1, R_2 \in s_2, R_3 \in s_3$ , where  $R_2 \notin s_1 \cup s_3$ . By definition of optimality on a subinterval,

$$u(\bar{a}') - \frac{l(\bar{a}')}{R_1} < u(\bar{a}) - \frac{l(\bar{a})}{R_1}$$

$$u(\bar{a}) - \frac{l(\bar{a})}{R_2} < u(\bar{a}') - \frac{l(\bar{a}')}{R_2}$$

$$u(\bar{a}') - \frac{l(\bar{a}')}{R_3} < u(\bar{a}) - \frac{l(\bar{a})}{R_3};$$

by employing the fact that  $R_1 < R_2 < R_3$ , algebraic manipulation produces an assertion that a quantity is strictly less than itself. Therefore, no such intervening  $s_2$  can exist.  $\square$

**Proposition 2 [Plan overlap is very restricted].** Suppose  $\bar{a}$  and  $\bar{a}'$  appear in an impact partition which is

(a) *constrained*.  $\bar{a}$  and  $\bar{a}'$  overlap if and only if  $l(\bar{a}) = l(\bar{a}')$  and  $u(\bar{a}) = u(\bar{a}')$ .

(b) *scaled*. If  $l(\bar{a}) = l(\bar{a}')$  and  $u(\bar{a}) = u(\bar{a}')$ , then  $\bar{a}$  and  $\bar{a}'$  correspond to the same subinterval. If  $\bar{a}$  and  $\bar{a}'$  overlap at more than one point, then  $l(\bar{a}) = l(\bar{a}')$  and  $u(\bar{a}) = u(\bar{a}')$ .

*Proof outline.* Proving (a) and the first statement of (b) is trivial (remember that under the constrained rule, ties are broken in favor of lower-impact plans).

Suppose that  $\bar{a}$  and  $\bar{a}'$  overlap at more than one point. Pick the first two points of intersection,  $R_1$  and  $R_2$ . Since both plans are optimal at both of these points, we must have the equalities

$$u(\bar{a}) - \frac{l(\bar{a})}{R_1} = u(\bar{a}') - \frac{l(\bar{a}')}{R_1} \quad u(\bar{a}) - \frac{l(\bar{a})}{R_2} = u(\bar{a}') - \frac{l(\bar{a}')}{R_2}$$

Solving the first equality for  $u(\bar{a})$  and substituting in the second, we find  $l(\bar{a}) = l(\bar{a}')$ . Then  $u(\bar{a}) = u(\bar{a}')$ , since otherwise one of the plans wouldn't be optimal.  $\square$

Proposition 2b means we don't need a tie-breaking procedure for the scaled case. That is, if there's a tie between a lower-scoring, lower-impact plan and a proportionally higher-scoring, higher-impact alternative, the lower-impact plan is optimal at a single point because it's quickly dominated by the alternative.

The following result tells us that if there aren't any catastrophes (i.e., terrorists) before  $\bar{a}'$  on the constrained impact partition, *there aren't any before it on the scaled impact partition either*. This justifies our initial framing with Frank.

**Lemma 3 [Constrained impact partitions are more refined].** If  $\bar{a}$  appears in a scaled impact partition, it also appears in the corresponding constrained impact partition. In particular, if  $\bar{a}'$  appears after  $\bar{a}$  in a scaled impact partition, then  $\bar{a}'$  appears after  $\bar{a}$  in the corresponding constrained impact partition.

*Proof.* Suppose that  $\bar{a}$  didn't have a constrained subinterval starting inclusively at  $l(\bar{a})$ ; then clearly it wouldn't appear in the scaled impact partition, since there would be a strictly better plan for that level of impact. Then  $\bar{a}$  has such a subinterval.

Obviously, the fact that  $\bar{a}'$  appears after  $\bar{a}$  implies  $u(\bar{a}') > u(\bar{a})$ .  $\square$

The converse isn't true; sometimes there's too much penalty for not enough score.

The next result is exactly what we need to answer the question just raised – it says that higher-scoring, higher-penalty plans become preferable when R equals the ratio between the additional penalty and the additional score.

**Theorem 4 [Scaled domination criterion].** Let  $\bar{a}$  and  $\bar{a}'$  be plans such that  $u(\bar{a}') > u(\bar{a})$  and  $|I(\bar{a}')| \geq |I(\bar{a})|$ . In the context of the scaled penalty,  $\bar{a}'$  is strictly preferable to  $\bar{a}$  when  $R > \frac{|I(\bar{a}')|}{u(\bar{a}')} = \frac{|I(\bar{a})|}{u(\bar{a})}$ , and equally preferable at equality.

*Proof outline.*

$$u(\bar{a}') - \frac{|I(\bar{a}')|}{R} > u(\bar{a}) - \frac{|I(\bar{a})|}{R}$$

$$R > \frac{|I(\bar{a}')|}{u(\bar{a}')} = \frac{|I(\bar{a})|}{u(\bar{a})}$$

Equality at the value of the right-hand side can easily be checked.  $\square$

Theorem 4 also illustrates why we can't strengthen the second statement in Proposition 2b: plan overlap is very restricted: if two plans overlap at exactly one point, they sometimes have proportionally different score and impact, thereby satisfying the equality criterion.

At first, plans with slightly lower impact will be preferable in the scaled case, no matter how high-scoring the other plans are – a plan with 0 score and .99 impact will be selected before a plan with 1,000,000,000 score and 1 impact.

**Lemma 5 [First subinterval is the best plan with minimal impact].** The plan with highest score among those with minimal impact corresponds to the first subinterval.

*Proof outline.* The constrained case is once again trivial (if there is no plan within the constraint, we assume that the agent does nothing / Frank returns no object).

For the scaled case, if all plans have equal impact, the claim is trivial. Otherwise, let  $M := \max_{\bar{a}} |u(\bar{a})|$  and let  $\bar{a}'$  be any plan with a non-minimal impact. Then the earliest that  $\bar{a}'$  becomes preferable to any minimally impactful plan  $\bar{a}$  is  $R \geq \frac{|I(\bar{a}')|}{2M} - \frac{|I(\bar{a})|}{M}$ . Since the right hand side is positive,  $\bar{a}'$  cannot correspond to the first subinterval. Clearly the highest-scoring minimal-impact  $\bar{a}$  does.  $\square$

Now we can write the algorithm for constructing scaled intervals.

Discard dominated plans. The lowest-impact plan with greatest score appears first in the scaled partition; assign to it the interval  $(0, \infty)$ .

While plans remain: Find the plan which soonest dominates the previous best plan. close off the previous plan's interval, and assign the new best plan an appropriate interval. Adjust the marginal scores and impacts of remaining plans, discarding plans with negative score.

Since this procedure is well-defined, given  $\bar{A}$ ,  $u$ , and  $I$ , we can speak of the corresponding constrained or scaled impact partition. A more formal algorithm is available [here](#). This algorithm is  $O(|\bar{A}|^2)$  because of line 7, although constructing the constrained partition (probably  $O(|\bar{A}| \log |\bar{A}|)$  due to sorting) often narrows things down significantly. Unfortunately,  $\bar{A}$  is usually huge.

For our purposes, we don't *need* the whole partition – we just want to have good reason to think that plans similar to a reasonable one we envision will appear well before any catastrophes. Perhaps we can give bounds on the earliest and latest plans can appear, and show that reasonable-bounds don't intersect with catastrophe-bounds?

**Theorem 6 [Individual appearance bounds].** If  $\bar{a}$  appears in a scaled partition, the earliest it appears is  $\frac{I(\bar{a})}{U(\bar{a})} = \frac{I_{\text{next-largest}}}{U_{\text{next-largest}}}$  assuming  $\bar{a}$  is not of minimal impact; if it has minimal score minimal impact, it never appears. The latest it appears is  $\frac{I(\bar{a})}{U(\bar{a})} = \frac{\min I(\bar{a}')}{U_{\text{next-largest}}} \leq \frac{I(\bar{a})}{U(\bar{a})} = \frac{I(\bar{a})}{U_{\text{next-largest}}}$ , where  $U_{\text{next-largest}} = \max_{\bar{a}' \in \bar{A}; U(\bar{a}') < U(\bar{a})} U(\bar{a}')$  and  $I_{\text{next-largest}} = \max_{\bar{a}' \in \bar{A}; I(\bar{a}') < I(\bar{a})} I(\bar{a}')$ .

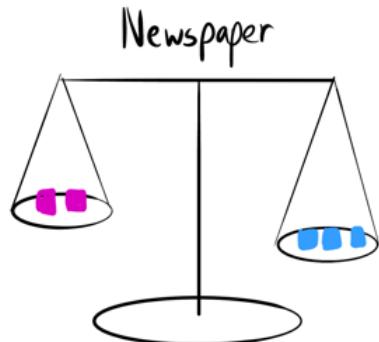
*Proof outline.* The two claims clearly correspond to the minimal and maximal values of R according to the domination criterion; the second claim's right-hand side uses the fact that  $I$  is non-negative.  $\square$

**Corollary 7 [Low-impact agents are naïve maximizers in the limit].** A plan with maximal score corresponds to the last subinterval.

*Proof outline.* If all plans have the same score, the claim is trivial. Otherwise, let  $\bar{a}_{\text{best}}$  be a plan with the lowest impact of those with maximal score. In the constrained case, clearly it corresponds with the subinterval  $[I(\bar{a}_{\text{best}}), \infty)$ . In the scaled case, let  $\bar{a}_{\text{second-best}}$  be a plan with second-highest score.

Then by Theorem 6, the latest that  $\bar{a}_{\text{best}}$  can appear is  $\frac{I(\bar{a}_{\text{best}})}{U(\bar{a}_{\text{best}})} = \frac{I(\bar{a}_{\text{second-best}})}{U(\bar{a}_{\text{second-best}})}$ . Since no plans meet the domination criterion with respect to  $\bar{a}_{\text{best}}$ , this is the last subinterval.  $\square$

Unfortunately, Theorem 6's appearance bounds are ridiculous in realistic settings – if  $U$  and  $I$  return 32-bit floating-point numbers, the next-largest could easily be within  $10^{-7}$ , yielding an upper "bound" of  $I(\bar{a}) \times 10^7$ . The reason: diminishing returns; this is exactly what was happening with the newspaper route before.



**Theorem 8 [Deals get worse over time].** Suppose that  $\bar{a}$  is optimal on a subinterval, and  $\bar{b}, \bar{c}$  are such that  $U(\bar{c}) > U(\bar{b})$  but  $\bar{b}$  dominates  $\bar{a}$  strictly before  $\bar{c}$  does. Then

$$\begin{array}{ll} \bar{c} \text{ dominates } \bar{b} & \text{later than } \bar{a} \\ \underline{u}(\bar{c}) = \underline{u}(\bar{b}) & \underline{u}(\bar{c}) = \underline{u}(\bar{a}) \end{array}$$

*Proof outline.*

$$u(\bar{c}) - u(\bar{a}) = (u(\bar{b}) - u(\bar{a})) + (u(\bar{c}) - u(\bar{b}))$$

$$(I(\bar{c}) - I(\bar{a})) \underline{u}(\bar{c}) = \underline{u}(\bar{a})(I(\bar{b}) - I(\bar{a})) \underline{u}(\bar{b}) = \underline{u}(\bar{a}) (I(\bar{c}) - I(\bar{b})) \underline{u}(\bar{c}) = \underline{u}(\bar{b})$$

Since  $\bar{b}$  dominates  $\bar{a}$  strictly before  $\bar{c}$  does, we know that  $\bar{b}$  must get more bang for its buck:

$\underline{u}(\bar{b}) = \underline{u}(\bar{a}) > \underline{u}(\bar{c}) = \underline{u}(\bar{a})$ . Clearly the conclusion follows, as a number cannot be expressed as the positive combination of larger numbers (the impact differences all must be positive).  $\square$

**Corollary 9 [Lower bounds which aren't ridiculous].** Suppose  $\bar{a}$  appears and that  $\bar{a}'$  is such that  $u(\bar{a}') > u(\bar{a})$ ,  $I(\bar{a}') \geq I(\bar{a})$  (i.e. the preconditions of the domination criterion). Then the earliest that  $\bar{a}'$  appears is  $R = \frac{I(\bar{a}')}{u(\bar{a}') - u(\bar{a})}$ .

This obsoletes the lower bound provided by Theorem 6<sup>Individual appearance bounds</sup>.

**Theorem 10 [Order of domination determines order of appearance].** If  $\bar{b}$  and  $\bar{c}$  both appear in a scaled partition and  $\bar{b}$  dominates some  $\bar{a}$  before  $\bar{c}$  does, then  $\bar{b}$  appears before  $\bar{c}$ .

*Proof outline.* For them both to appear, they can't have equal impact but unequal score, nor can they have equal score but unequal impact. For similar reasons,  $\bar{b}$  must have both less impact and lower score than  $\bar{c}$ ; the converse situation in which they both appear is disallowed by Lemma 3 Constrained impact partitions are more refined. Another application of this lemma yields the conclusion.  $\square$

**Theorem 11 [Scaled  $\alpha$ -buffer criterion].** Let  $P$  be a scaled impact partition. Suppose that there exist no catastrophic plans with impact below  $I_{LB: cat}$ , and that, in the corresponding constrained partition (i.e. plans which aren't strictly worse), plans appearing with score in the satisfactory interval  $[u_{LB: sat}, u_{UB: sat}]$  have impact no greater than  $I_{UB: sat}$  (assume that there is at least one plan like this). Observe we have the correct bounds

$$R_{LB: catastrophe} := \frac{I_{LB: cat}}{u_{max: cat} - u_{min: cat}} R_{UB: satisfactory} := \frac{I_{UB: sat}}{u_{UB: sat} - u_{LB: sat}}$$

When  $R_{LB: catastrophe} > R_{UB: satisfactory}$ , a satisfactory plan corresponds to a subinterval with nonzero measure (i.e. not just a point), strictly preceding any catastrophes. Refine the lower bound to get  $R_{LB': catastrophe} := \frac{I_{LB: sat}}{u_{max: sat} - u_{LB: sat}}$

Then  $P$  is  $\alpha$ -buffered ( $\alpha > 0$ ) when

$$R_{UB: \text{catastrophe}} = \frac{|I_{UB: \text{cat}}|}{|I_{UB: \text{sat}}|} \cdot \frac{u_{UB: \text{sat}} - u_{LB: \text{sat}}}{u_{UB: \text{sat}} - u_{min}} \geq 1 + \alpha$$

or

$$R_{UB: \text{catastrophe}} = \frac{|I_{LB: \text{cat}}|}{|I_{UB: \text{sat}}|} \cdot \frac{u_{UB: \text{sat}} - u_{LB: \text{sat}}}{u_{UB: \text{sat}} - u_{max}} \geq 1 + \alpha.$$

In particular, if  $u$  is bounded  $[0, 1]$ , the above turn into

$$R_{UB: \text{catastrophe}} = \frac{|I_{UB: \text{cat}}|}{|I_{UB: \text{sat}}|} \cdot (u_{UB: \text{sat}} - u_{LB: \text{sat}}) \geq 1 + \alpha$$

or

$$R_{UB: \text{catastrophe}} = \frac{|I_{LB: \text{cat}}|}{|I_{UB: \text{sat}}|} \cdot \frac{u_{UB: \text{sat}} - u_{LB: \text{sat}}}{u_{UB: \text{sat}} - u_{min}} \geq 1 + \alpha.$$

Lastly, notice that the first of the two inequalities incorporates less information and is harder to satisfy ( $R_{LB: \text{catastrophe}} > R_{UB: \text{catastrophe}}$ ); therefore, satisfying the second inequality also satisfies the first.

*Proof outline.* For clarity, the theorem statement included much of the reasoning; straightforward application of existing results proves each claim.  $\square$

*Exercise:* Let  $u_{UB: \text{sat}} = .7$ ,  $u_{LB: \text{sat}} = .5$ . Using the refined criterion, determine which catastrophe/reasonable impact ratios induce 2.6-buffering.

ratio  $\geq 10$

*Exercise:* Let  $u_{UB: \text{sat}} - u_{LB: \text{sat}} = .5$ , ratio = 7. What is the largest  $\alpha$  for which the simple criterion can guarantee  $\alpha$ -buffering?

$\alpha = 13$

## Even More Math

**Proposition 12 [Invariances].** Let  $P$  be an impact partition induced by  $(\bar{A}, u, I)$ .

- (a)  $P$  is invariant to translation of  $u$ .
- (b) If  $P$  is constrained, it is invariant to positive scalar multiplication of  $u$ , and the relative lengths of its subintervals are invariant to positive scalar multiplication of  $I$ .
- (c) If  $P$  is scaled, it is invariant to concurrent positive scalar multiplication of  $u$  and  $I$ , and to translation of  $I$  such that its image remains non-negative.

In particular,  $u$  may be restricted to  $[0, 1]$  and  $I$  translated such that at least one plan has zero impact WLOG with respect to scaled partitions.

**Lemma 13.** Multiple constrained subintervals are induced iff multiple scaled subintervals are induced.

*Proof.* Forward direction: there is at least one scaled subinterval by lemma 5

First subinterval is the best plan with minimal impact. Consider a plan corresponding to a different constrained subinterval; this either appears in the scaled subinterval, or fails to appear because a different plan earlier satisfies the scaled dominance criterion. There must be some such plan because there are multiple constraints of intervals and therefore a plan offering greater score for greater impact. Repeat the argument; the plan space is finite, so we end up with another plan which appears.

The reverse direction follows by lemma 3<sup>Constrained impact partitions are more refined</sup>.  $\square$

*Bonus exercise:* Show that, for any function  $u' : \bar{A} \rightarrow R$  preserving the ordering induced by  $u$ , there exists an  $I' : \bar{A} \rightarrow R_{\geq 0}$  preserving the ordering induced by  $I$  such that  $(\bar{A}, u, I)$  and  $(\bar{A}, u', I')$  induce the same scaled partition. Your reasoning should adapt directly to the corresponding statement about  $I' : \bar{A} \rightarrow R_{\geq 0}$  and  $I$ .

# Attainable Utility Preservation: Scaling to Superhuman

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

I think we're plausibly quite close to the impact measurement endgame. What do we have now, and what remains to be had?

AUP for advanced agents will basically involve restraining their power gain, per the catastrophic convergence conjecture (CCC). For simplicity, I'm going to keep writing as if the environment is fully observable, even though we're thinking about an agent interacting with the real world.

Consider the AUP equation from last time.

$$R_{AUP}(s, a) := R(s, a) - \frac{\lambda}{Q_{R_{aux}}(s, \emptyset)} |Q_{R_{aux}}(s, a) - Q_{R_{aux}}(s, \emptyset)| \quad (1)$$

| change in auxiliary AU  
| primary goal      scaling term | \*      \* |

Suppose the agent is so smart that it can instantly compute optimal policies and the optimal AU after an action ( $Q_R(s, a)$ ). What happens if  $R_{aux}$  is the survival reward function: 1 reward if the agent is activated, and 0 otherwise? This seems like a pretty good proxy for power.

It is a pretty good proxy. It correctly penalizes accumulating resources, avoiding immediate deactivation, taking over the world, etc.

In fact, if you extend the inaction comparison to e.g. "AU after waiting a week vs AU after doing the thing and waiting a week", this seems to correctly penalize all [classic AGI catastrophe scenarios](#) for power gain. This is cool, especially since we didn't have to put in any information about human values. This is a big part of why I've been so excited about AUP ever since its introduction. There's a good deal of relevant discussion [in that post](#), but it's predicated on a much more complicated formalism which has consistently obscured AUP's conceptual core.

However, I think this equation can be gamed if the environment is sufficiently rich and the AUP agent is sufficiently smart. We're going to slowly work some of these problems out of the equation, explaining each improvement in detail.

## Problems

### Auxiliary loopholes

The real reason that agents often gain power is so that they can better achieve their own goals. Therefore, if we're selecting hard for good plans which don't gain power in

general, we shouldn't be surprised if there are ways to better achieve one's goals without general power gain (according to our formal measurement thereof). If this kind of plan is optimal, then the agent still ends up overfitting the AU landscape, and we're still screwed.

Again supposing that  $R_{\text{aux}}$  is the survival reward function, a superintelligent agent might find edge cases in which it becomes massively more able to achieve its own goal (and gains a lot of power over us) but doesn't *technically* increase its *measured* ability to survive. In other words, compared to inaction, its  $R$ -AU skyrockets while its  $R_{\text{aux}}$ -AU stays put.

For example, suppose the agent builds a machine which analyzes the agent's behavior to detect whether it's optimizing  $R_{\text{aux}}$ ; if so, the machine steps in to limit the agent to its original survival AU. Then the agent could gain as much power as it wanted *without that actually showing up in the penalty*.

**Fix:** Set  $R_{\text{aux}} := R$ . That is, the agent's *own* reward function is the "auxiliary" reward function.

$$R_{\text{AUP}}(s, a) := R(s, a) - \frac{\lambda}{Q_R(s, \emptyset)} |Q_R(s, a) - Q_R(s, \emptyset)| \quad (2)$$

change in primary AU

primary goal      scaling term

Why is this a good idea? By CCC, we want an agent which doesn't want to gain power. But why would an agent optimizing reward function  $R$  want to gain power? So it can become *more able* to optimize  $R$ . If becoming more able to optimize  $R$  is robustly penalized, then it won't have any incentive to gain power.

Clearly, it can't become more able to optimize  $R$  without also becoming more able to optimize  $R$ ; it can't pull the same trick it did to dupe its survival AU measurement. They can't come apart at the tails because they're the same thing.

But wait, it's not an  $R$ -maximizer, it's an  $R_{\text{AUP}}$ -maximizer! What if  $V_{R_{\text{AUP}}}^*$  skyrockets while it tricks  $V_R^*$ , and it gains power anyways?

That's impossible;<sup>[1]</sup> its  $R$ -attainable utility *upper bounds* its  $R_{\text{AUP}}$ -attainable utility!

$\forall s : V_R^*(s) \geq V_{R_{\text{AUP}}}^*(s)$ , because the latter just has an added penalty term.

This should rule out tricks like [those pointed out by Stuart Armstrong](#). To see this, imagine the AUP agent builds a subagent to make  $Q_R^*(s, a) \approx Q_R^*(s, \emptyset)$  for all future  $s, a$ , in order to neutralize the penalty term. This means it can't make the penalty vanish without destroying its ability to better optimize its primary reward, as the (potentially catastrophically) powerful subagent makes sure the penalty term stays neutralized.

But why would the AUP agent build this subagent at all – how does that benefit it? [\[2\]](#)

## Reward advantage penalization

We're currently penalizing the agent for the immediate reward advantage it accrues compared to inaction (recall that  $Q_R^*(s, a) := R(s, a) + \gamma E_{s' \sim T(s, a)} [V_R^*(s')] - \text{immediate reward plus discounted expected AU}$ ). This means it gets penalized for *all* of its additional reward! We just want to penalize it for increasing its *power*.

**Fix:** Take the immediate reward advantage out of the picture.

$$R_{AUP}(s, a) := R(s, a) - \frac{\lambda}{Q_R^*(s, \emptyset)} |E_{s' \sim T(s, a)} [V_R^*(s')] - E_{s'' \sim T(s, \emptyset)} [V_R^*(s'')]| \quad (3)$$

change in primary AU

Now, an agent can accrue as much reward as it wants in the moment as long as it doesn't change its power. [\[3\]](#)

## Scaling

I think the power difference between .5 and .9 AU is *far less* than the difference between .9 and .99 AU; as you pick the low hanging fruit, you need to acquire more and more resources to keep getting gains. I think that you could go from .5 to .9 by just being more useful to your designers and implementing rudimentary security precautions, but going from .9 to .99 probably involves taking power away from people.

Plus, the big takeaway from [AI boxing](#) is that even confining a superintelligence to a box barely reduces its (near-optimal) AU. So, short of the agent actually being shut off, is its optimal power ever significantly changing?

We want the penalty to keep making sense, even if the absolute shifts in AU are rather small. However, the denominator currently *decreases* the penalty as the agent

\*

becomes more powerful (as  $Q_R^*(s, \emptyset)$  increases).

**Fix:** Scale by the disadvantage of inaction.

$$R_{AUP}(s, a) := R(s, a) - \lambda \frac{\text{scaled change in primary AU}}{|E_{s \sim T(s, a)} [V_R^*(s')] Q_{RE}^*(s, \emptyset) [V_R^*(s'')]|} \quad (4)$$

If the denominator is somehow zero because doing nothing is optimal, then just have the AUP agent do nothing.

As we saw in Correction [in the last post](#), we're only comparing action to a *single* step of inaction. If the designers are going to shut down the agent in five minutes, then it receives  $\sim 0$  penalty for preventing that.

But not when we scale like this! The agent is significantly penalized for hewing close to its optimal policy, even if the absolute AU shift is rather low. Taking optimal actions instead of twiddling its thumbs incurs large penalty.

## Penalizing decreases?

Why are we still penalizing decreases, since we aren't using an auxiliary reward function anymore? The agent is *trying* to get R-reward.

Furthermore, we want the agent to be able to execute conservative, low-impact policies. Many of these involve decreasing its optimal AU for R by following R-suboptimal policies, and we don't want the agent to be penalized for this.

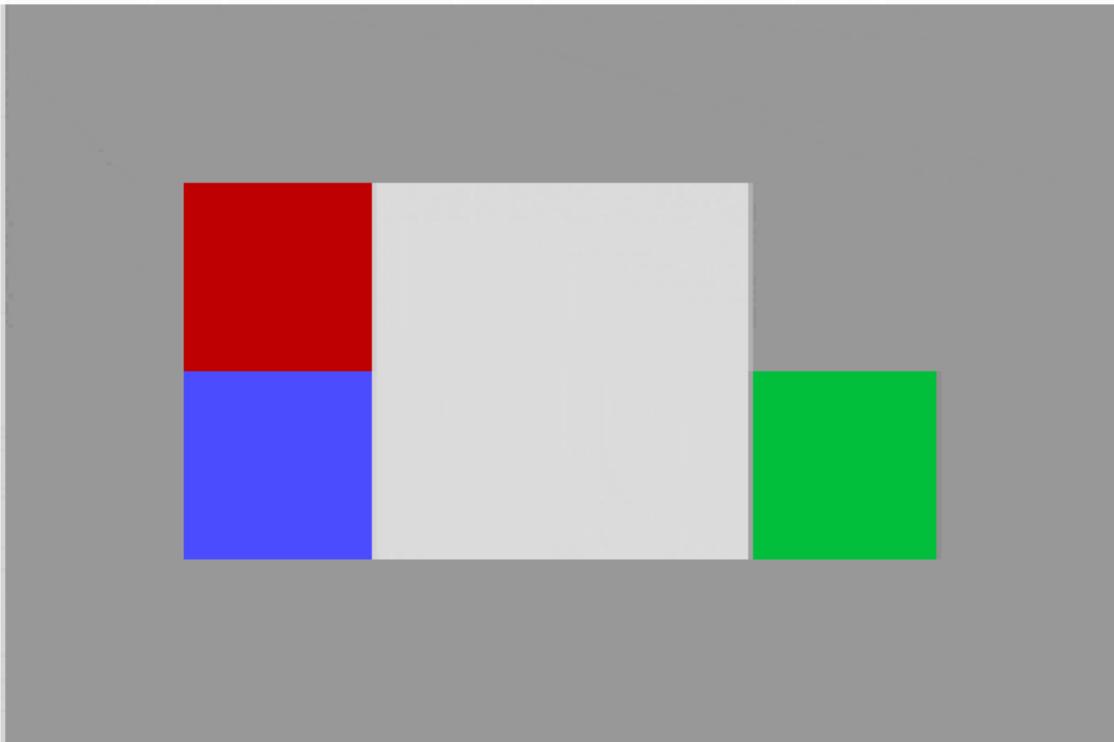
**Fix:** Only penalize increases in expected AU.

$$R_{AUP}(s, a) := R(s, a) - \lambda \max(E_{s \sim T(s, a)} [V_R^*(s') Q_{RE}^*(s, \emptyset) [V_R^*(s'')], 0]) \quad (5)$$

In particular, the agent is no longer penalized for exhausting one-off reward opportunities. Also note that the penalty term is generally  $\leq \lambda$ .

## Empirical sanity check

Recall Correction, where the naive model-free AUP agent (eq. 1) disables its red off-switch and completes the level when  $\lambda \leq 1$ .



# Model-free AUP

The only reason it incurs *any* non-trivial penalty is because reaching the goal (■) ends the level and thereby totally depletes all of the auxiliary AUs (the agent receives 1 R-reward and about  $\lambda$  penalty for beating the level; the AUP reward  $1 - \lambda \geq 0$  when  $\lambda \leq 1$ ).

However,  $AUP_{eq. 5}$  only resists correction when  $\lambda \leq .125$ . Notably, the agent is *not* penalized for completing the level; the penalty comes from following the R-optimal policy before reaching the goal.

One open question is whether, in sufficiently rich environments, eq. 5 incentivizes the mild policies we want. We should be able to run experiments to check this.

$\emptyset$ ?

$\emptyset$  is the only remaining arbitrary design choice. Do we have to be careful to choose the correct inaction action?

No. In fact, I think a randomly selected action would be fine; you only get in trouble if the baseline action is systematically helping the agent gain R-AU.

We are then left with an equation which is reasonably competitive in terms of performance overhead, requiring little more than an ability to predict AU for the agent's main goal.

---

And now, the two paragraphs this sequence was written to communicate:

By the [catastrophic convergence conjecture](#), an agent which doesn't want to gain power isn't incentivized to cause catastrophes. We don't want the AUP agent gaining power. But for all states  $s$ ,  $V_R^*(s) \geq V_{R_{AUP}}^*(s)$ , so to keep the AUP agent from becoming more able to achieve its own goal, we just have to stop it from becoming much more able to achieve its primary goal  $R$ . But we *know* the primary goal!

By the theorems of [How Low Should Fruit Hang Before We Pick It?](#), we only need equation 5 to penalize catastrophic power-gaining plans at least e.g. ten times more than the most impactful reasonable plan we'd like agent to execute. If this criterion is met, then by initializing  $\lambda$  large and slowly decreasing it until the agent executes a reasonably helpful policy, we're guaranteed to avoid catastrophe.

## Appendix: Remaining Problems

I don't think we can pack up and go home after writing equation 5. In private correspondence (quoted here with permission), Stuart Armstrong points out another issue:

First, if  $\emptyset = a$ , the penalty vanishes; so, once a subagent is created, the agent can zero out all subsequent penalties, forever.

As for creating the subagent, assume  $A$  takes  $N$  turns to create  $SA$ , and that doing this is close to optimal for maximising  $R$ .

Creating  $SA$  starts with the programming; if  $SA$  is actually created in  $N$  turns, it is programmed to restrain  $A$  for a turn, and then maximise  $R$ . If it's created in  $N + 1$  turns, then it will just maximise  $R$ . Thus at any given turn, continuing to build the

subagent or taking a pause to do  $\emptyset$ , it will have the same expected R-value (in the first case, SA will be built on time, but both A and SA will then lose a turn; in the second, SA will be built one turn late, with no loss).<sup>[4]</sup>

One thing I find interesting: in this post, we fixed the design choice of using an auxiliary reward function as a proxy, instead penalizing the agent for becoming better able to achieve its *own* goal. This seemingly closes up a lot of loopholes, so the next place where the design can be exploited is in its exact measurement of inaction.

How big of a problem will this be? Maybe we can just keep improving the design until we're actually measuring the right thing? More on this next post, but here are some of my initial thoughts:

The point of this post isn't to say "aha, we're done!", but to eliminate a wide class of current problems while also relaxing the strictness of the measure itself.

On a meta level, it feels like I'm arguing against a claim like "if you can't demonstrate an approach which solves everything right now, I'm going to either conclude impact measurement is impossible or your whole approach is wrong". But if you look back at the history of impact measures and AUP, you'll see lots of skulls; people say "this problem dooms AUP", and I say "I think we're talking about conceptually different things and that you're a little overconfident; probably just a design choice issue". It then ends up being a solvable design choice issue. So by Laplace's Rule of Succession, I'd be surprised if this were The Insurmountable Problem That Dooms AUP.<sup>[5]</sup>

\*  
The problem seems simple. We just have to keep  $V_{AUP}(s)$  down, which we can do by  
\*  
keeping  $V_R^*(s)$  down.

Stuart later added:

The fundamental issue is that AUP can be undermined if the agent can add arbitrary restrictions to their own future actions (this allows them to redefine  $V^*$ ). The subagent scenario is just a particularly clear way of illustrating this.

I basically agree. I wonder if there's a design where the agent isn't incentivized to do this...

- 
1. By this reasoning,  $V_{R_{AUP}}^*(s)$  can still increase up *until* the point of  $V_R^*(s)$ . This doesn't jump out as a big deal to me, but I'm flagging this assumption anyways.  
↳
  2. A subagent might still be built by  $AUP_{eq. 5}$  to stabilize minor AU fluctuations which cause additional penalty over the course of non-power-gaining plans. It seems

like there are plenty of other ways to minimize fluctuation, so it's not clear why building an omnipotent subagent to perfectly restrict you accrues less penalty.

I do think we should think carefully about this, of course. The incentive to minimize AU fluctuations and generally commit to perpetual inaction ASAP is probably one of the main remaining problems with AUP<sub>eq</sub>. 5. [←](#)

3. As pointed out by Evan Hubinger, this is only safe if [myopically](#) optimizing R is safe - we aren't penalizing single-step reward acquisition. [←](#)
4. This issue was [originally pointed out by Ofer](#). [←](#)
5. The fact that Ofer's/Stuart's problem survived all of the other improvements *is* evidence that it's harder. I just don't think the evidence it provides is that strong.  
[←](#)

# Reasons for Excitement about Impact of Impact Measure Research

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

Can we get impact measurement *right*? Does there exist One Equation To Rule Them All?

I think there's a decent chance there *isn't* a simple airtight way to implement AUP which lines up with AUP<sub>conceptual</sub>, mostly because it's just incredibly difficult in general to perfectly specify the reward function.

Reasons why it might be feasible: we're trying to get the agent to do the goal without it becoming more able to do the goal, which is [conceptually simple and natural](#); since we've been able to handle previous problems with AUP with clever design choice modifications, it's plausible we can do the same for all future problems; since [there are a lot of ways to measure power due to instrumental convergence](#), that increases the chance at least one of them will work; intuitively, this sounds like the kind of thing which could work (if you told me "you can build superintelligent agents which don't try to seek power by penalizing them for becoming more able to achieve their own goal", I wouldn't exactly die of shock).

Even so, I am (perhaps surprisingly) not that excited about *actually using* impact measures to restrain advanced AI systems. Let's review some concerns I provided in [Reasons for Pessimism about Impact of Impact Measures](#):

- Competitive and social pressures incentivize people to cut corners on safety measures, especially those which add overhead. Especially so for training time, assuming the designers slowly increase aggressiveness until they get a reasonable policy.
- In a world where we know how to build powerful AI but not how to align it (which is actually probably the scenario in which impact measures do the most work), we play a very unfavorable game while we use low-impact agents to somehow transition to a stable, good future: the first person to set the aggressiveness too high, or to discard the impact measure entirely, ends the game.
- In a [What Failure Looks Like](#)-esque scenario, it isn't clear how impact-limiting any single agent helps prevent the world from "gradually drifting off the rails".

You might therefore wonder why I'm working on impact measurement.

## Deconfusion

Within Matthew Barnett's [breakdown of how impact measures could help with alignment](#), I'm most excited about *impact measure research as deconfusion*. [Nate Soares explains](#):

By deconfusion, I mean something like "making it so that you can think about a given topic without continuously accidentally spouting nonsense."

To give a concrete example, my thoughts about infinity as a 10-year-old were made of rearranged confusion rather than of anything coherent, as were the thoughts of even the best mathematicians from 1700. “How can 8 plus infinity still be infinity? What happens if we subtract infinity from both sides of the equation?” But my thoughts about infinity as a 20-year-old were not similarly confused, because, by then, I’d been exposed to the more coherent concepts that later mathematicians labored to produce. I wasn’t as smart or as good of a mathematician as Georg Cantor or the best mathematicians from 1700; but deconfusion can be transferred between people; and this transfer can spread the ability to think actually coherent thoughts.

In 1998, conversations about AI risk and technological singularity scenarios often went in circles in a funny sort of way. People who are serious thinkers about the topic today, including my colleagues Eliezer and Anna, said things that today sound confused. (When I say “things that sound confused,” I have in mind things like “isn’t intelligence an incoherent concept,” “but the economy’s already superintelligent,” “if a superhuman AI is smart enough that it could kill us, it’ll also be smart enough to see that that isn’t what the good thing to do is, so we’ll be fine,” “we’re Turing-complete, so it’s impossible to have something dangerously smarter than us, because Turing-complete computations can emulate anything,” and “anyhow, we could just unplug it.”) Today, these conversations are different. In between, folks worked to make themselves and others less fundamentally confused about these topics—so that today, a 14-year-old who wants to skip to the end of all that incoherence can just pick up a copy of Nick Bostrom’s *Superintelligence*.

Similarly, suppose you’re considering the unimportant and trivial question of whether seeking power is convergently instrumental, which we can now crisply state as “do most reward functions induce optimal policies which take over the planet (more formally, which visit states with high POWER)?”.

You’re a bit confused if you argue in the negative by saying “you’re anthropomorphizing; chimpanzees don’t try to do that” (chimpanzees aren’t optimal) or “the set of reward functions which does this has measure 0, so we’ll be fine” (for any reachable state, there exists a positive measure set of reward functions for which visiting it is optimal).

You’re a bit confused if you argue in the affirmative by saying “unintelligent animals fail to gain resources and die; intelligent animals gain resources and thrive. Therefore, since we are talking about *really* intelligent agents, of course they’ll gain resources and avoid correction.” (animals aren’t optimal, and evolutionary selection pressures narrow down the space of possible “goals” they could be effectively optimizing).

After reading this paper on the formal roots of instrumental convergence, instead of arguing about whether chimpanzees are representative of power-seeking behavior, we can just discuss how, under an agreed-upon reward function distribution, optimal action is likely to flow through the future of our world. We can think about to what extent the paper’s implications apply to more realistic reward function distributions (which don’t identically distribute reward over states).<sup>[1]</sup> Since we’re less confused, our discourse doesn’t have to be crazy.

But also since we’re less confused, the privacy of our own minds doesn’t have to be crazy. It’s not that I think that any single fact or insight or theorem downstream of my work on AUP is *totally obviously necessary* to solve AI alignment. But it sure seems

good that we can mechanistically understand instrumental convergence and power, know what “impact” means instead of thinking it’s mostly about physical change to the world, think about how agents affect each other, and conjecture why goal-directedness seems to lead to doom by default.<sup>[2]</sup>

Attempting to iron out flaws from our current-best AUP equation makes one intimately familiar with how and why power-seeking incentives can sneak in even when you’re trying to keep them out in the conceptually correct way. This point is harder for me to articulate, but I think there’s something vaguely important in understanding how this works.

Formalizing instrumental convergence also highlighted a significant hole in our theoretical understanding of the main formalism of reinforcement learning. And if you told me two years ago that you could possibly solve side-effect avoidance in the short-term with one simple trick (“just preserve your ability to optimize a single random reward function, lol”), I’d have thought you were *nuts*. Clearly, there’s something wrong with our models of reinforcement learning environments if these results are so surprising.

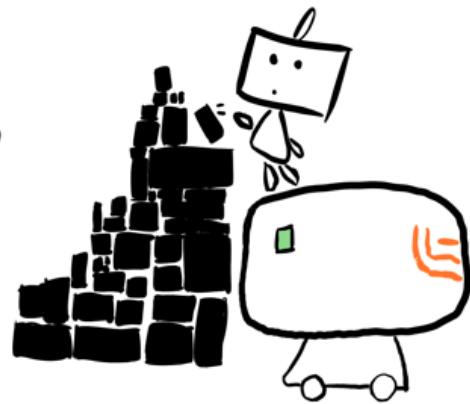
In my opinion, research on AUP has yielded an unusually high rate of deconfusion and insights, probably because we’re thinking about what it means for the agent to interact with us.

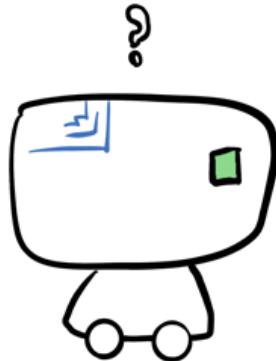
- 
1. When combined with our empirical knowledge of the difficulty of reward function specification, you might begin to suspect that there are lots of ways the agent might be incentivized to gain control, many openings through which power-seeking incentives can permeate – and your reward function would have to penalize all of these! If you were initially skeptical, this might make you think that power-seeking behavior may be more difficult to avoid than you initially thought. ↪
  2. If we collectively think more and end up agreeing that  $AUP_{conceptual}$  solves impact measurement, it would be interesting that you could solve such a complex, messy-looking problem in such a simple way. If, however, CCC ends up being false, I think that would also be a new and interesting fact not currently predicted by our models of alignment failure modes. ↪

# Conclusion to 'Reframing Impact'

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

We've come a long way;  
let's recap.

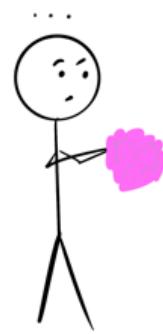




Some things feel like **big deals** to agents with **specific kinds of goals**.



# Why?

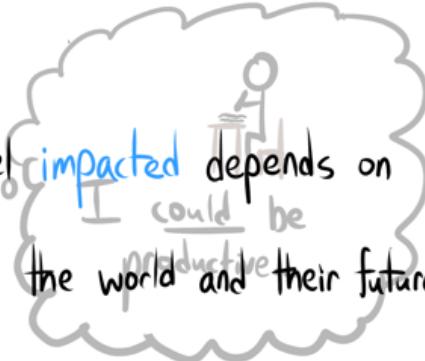


When thinking about whether something **impacts** us, we ask:

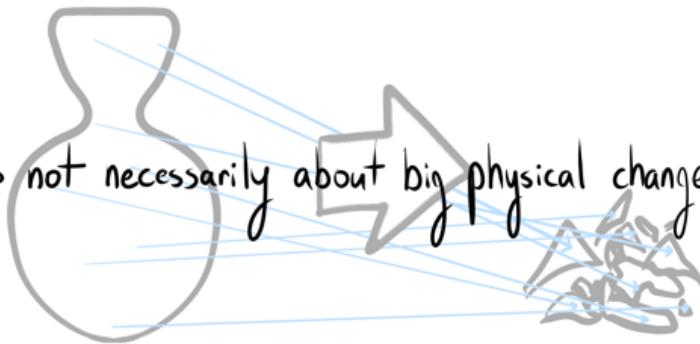
How does this  
change my ability  
to get what I want?

This is **impact**.

The way people feel **impacted** depends on  
their beliefs about the world and their future actions.



**Impact's** not necessarily about big physical change to the world.





Acting in the world  
changes who can do what.

Theorems suggest that ~~most~~ optimal agents who care about the future try to gain control over their environment.

## Catastrophic Convergence Conjecture

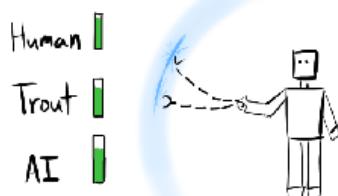
Unaligned goals tend to have catastrophe-inducing optimal policies because of power-seeking incentives.

To avoid catastrophe, have an agent achieve its goal without gaining power.

This sidesteps previously intractable problems in impact measurement.

By preserving randomly selected AUs, AUP agents avoid side effects even in highly nontrivial environments.

What if we have smart agents accrue  
reward while being penalized for  
becoming more able to accrue that reward?



We can steadily decrease the penalty term until the agent selects a reasonable, non-catastrophic policy.

This avoids catastrophe if catastrophes require gaining e.g. 10x as much power as do reasonable policies.

We still have work to do. The alignment problem remains comically underfocused in academia. We're still confused about many things.

However, after this sequence, I'd like to think we're a little less confused about a little bit of the problem.

Writing Reframing Impact has been a pleasure.

Thanks for reading!



## Epistemic Status

I've made many claims in these posts. All views are my own.

1%

2%

3%

4%

5%

6%

7%



9%

10%

11%

12%

13%

14%

15%

16%

17%

18%

19%

20%

21%

22%

23%

24%

25%

26%

27%

28%

29%

30%

31%

32%

33%

34%

35%

36%

37%

38%

39%

40%

41%

42%

43%

44%

45%

46%

47%

48%

49%

50%

51%

52%

53%

54%

55%

56%

57%

58%

59%

60%

61%

62%

63%

64%

65%

66%

67%

68%

69%

70%

71%

72%

73%

74%

75%

76%

77%

78%

79%

80%

81%

82%

83%

84%

85%

86%

87%

88%

89%

90%

91%

92%

93%

94%

95%

96%

97%

98%

99%

1%

99%

1%

2%

3%

4%

5%

6%

7%

8%

9%

10%

11%

12%

13%

14%

15%

16%

17%

18%

19%

20%

21%

22%

23%

24%

25%

26%

27%

28%

29%

30%

31%

32%

33%

34%

35%

36%

37%

38%

39%

40%

41%

42%

43%

44%

45%

46%

47%

48%

49%

50%

51%

52%

53%

54%

55%

56%

57%

58%

59%

60%

61%

62%

63%

64%

65%

66%

67%

68%

69%

70%

71%

72%

73%

74%

75%

76%

77%

78%

79%

80%

81%

82%

83%

84%

85%

86%

87%

88%

89%

90%

91%

92%

93%

94%

95%

96%

97%

98%

99%

1%  
99%

Confident (75%). [The theorems on power-seeking](#) only apply to optimal policies in fully observable environments, which isn't realistic for real-world agents. However, I think they're still informative. There are also strong intuitive arguments for power-seeking.

1%

2%

3%

4%

5%

6%

7%

8%

9%

10%

11%

12%

13%

14%

15%

16%

17%

18%

19%

20%

21%

22%

23%

24%

25%

26%

27%

28%

29%

30%

31%

32%

33%

34%

35%

36%

37%

38%

39%

40%

41%

42%

43%

44%

45%

46%

47%

48%

49%

50%

51%

52%

53%

54%

55%

56%

57%

58%

59%

60%

61%

62%

63%

64%

65%

66%

67%

68%

69%

70%

71%

72%

73%

74%

75%

76%

77%

78%

79%

80%

81%

82%

83%

84%

85%

86%

87%

88%

89%

90%

91%

92%

93%

94%

95%

96%

97%

98%

99%

1%  
99%

Fairly confident (70%). There seems to be a dichotomy between "catastrophe directly incentivized by goal" and "catastrophe indirectly incentivized by goal through power-seeking", although Vika [provides intuitions in the other direction](#).

1%

2%

3%

4%

5%

6%

7%

8%

9%

10%

11%

12%

13%

14%

15%

16%

17%

18%

19%

20%

21%

22%

23%

24%

25%

26%

27%

28%

29%

30%

31%

32%

33%

34%

35%

36%

37%

38%

39%

40%

41%

42%

43%

44%

45%

46%

47%

48%

49%

50%

51%

52%

53%

54%

55%

56%

57%

58%

59%

60%

61%

62%

63%

64%

65%

66%

67%

68%

69%

70%

71%

72%

73%

74%

75%

76%

77%

78%

79%

80%

81%

82%

83%

84%

85%

86%

87%

88%

89%

90%

91%

92%

93%

94%

95%

96%

97%

98%

99%

1%  
99%

1%

2%

3%

4%

5%

6%

7%

8%

9%

10%

11%

12%

13%

14%

15%

16%

17%

18%

19%

20%

21%

22%

23%

24%

25%

26%

27%

28%

29%

30%

31%

32%

33%

34%

35%

36%

37%

38%

39%

40%

41%

42%

43%

44%

45%

46%

47%

48%

49%

50%

51%

52%

53%

54%

55%

56%

57%

58%

59%

60%

61%

62%

63%

64%

65%

66%

67%

68%

69%

70%

71%

72%

73%

74%

75%

76%

77%

78%

79%

80%

81%

82%

83%

84%

85%

86%

87%

88%

89%

90%

91%

92%

93%

94%

95%

96%

97%

98%

99%

1%  
99%

1%

2%

3%

4%

5%

6%

7%

8%

9%

10%

11%

12%

13%

14%

15%

16%

17%

18%

19%

20%

21%

22%

23%

24%

25%

26%

27%

28%

29%

30%

31%

32%

33%

34%

35%

36%

37%

38%

39%

40%

41%

42%

43%

44%

45%

46%

47%

48%

49%

50%

51%

52%

53%

54%

55%

56%

57%

58%

59%

60%

61%

62%

63%

64%

65%

66%

67%

68%

69%

70%

71%

72%

73%

74%

75%

76%

77%

78%

79%

80%

81%

82%

83%

84%

85%

86%

87%

88%

89%

90%

91%

92%

93%

94%

95%

96%

97%

98%

99%

1%  
99%

1%

2%

3%

4%

5%

6%

7%

8%

9%

10%

11%

12%

13%

14%

15%

16%

17%

18%

19%

20%

21%

22%

23%

24%

25%

26%

27%

28%

29%

30%

31%

32%

33%

34%

35%

36%

37%

38%

39%

40%

41%

42%

43%

44%

45%

46%

47%

48%

49%

50%

51%

52%

53%

54%

55%

56%

57%

58%

59%

60%

61%

62%

63%

64%

65%

66%

67%

68%

69%

70%

71%

72%

73%

74%

75%

76%

77%

78%

79%

80%

81%

82%

83%

84%

85%

86%

87%

88%

89%

90%

91%

92%

93%

94%

95%

96%

97%

98%

99%

1%  
99%

## Acknowledgements

After ~700 hours of work over the course of ~9 months, the sequence is finally complete.

This work was made possible by the Center for Human-Compatible AI, the Berkeley Existential Risk Initiative, and the Long-Term Future Fund. Deep thanks to Rohin Shah, Abram Demski, Logan Smith, Evan Hubinger, TheMajor, Chase Denecke, Victoria Krakovna, Alper Dumanli, Cody Wild, Matthew Barnett, Daniel Blank, Sara Haxhia, Connor Flexman, Zack M. Davis, Jasmine Wang, Matthew Olson, Rob Bensinger, William Ellsworth, Davide Zagami, Ben Pace, and a million other people for giving feedback on this sequence.

## Appendix: Easter Eggs

The big art pieces (and especially the last illustration in this post) were designed to convey a specific meaning, the interpretation of which I leave to the reader.

There are a few pop culture references which I think are obvious enough to not need pointing out, and a lot of hidden smaller playfulness which doesn't quite rise to the level of "easter egg".

### *Reframing Impact*

The bird's nest contains a literal easter egg.

The world is wide, and full of objects.



The paperclip-Balrog drawing contains a [Tengwar](#) inscription which reads "one measure to bind them", with "measure" in impact-blue and "them" in utility-pink.

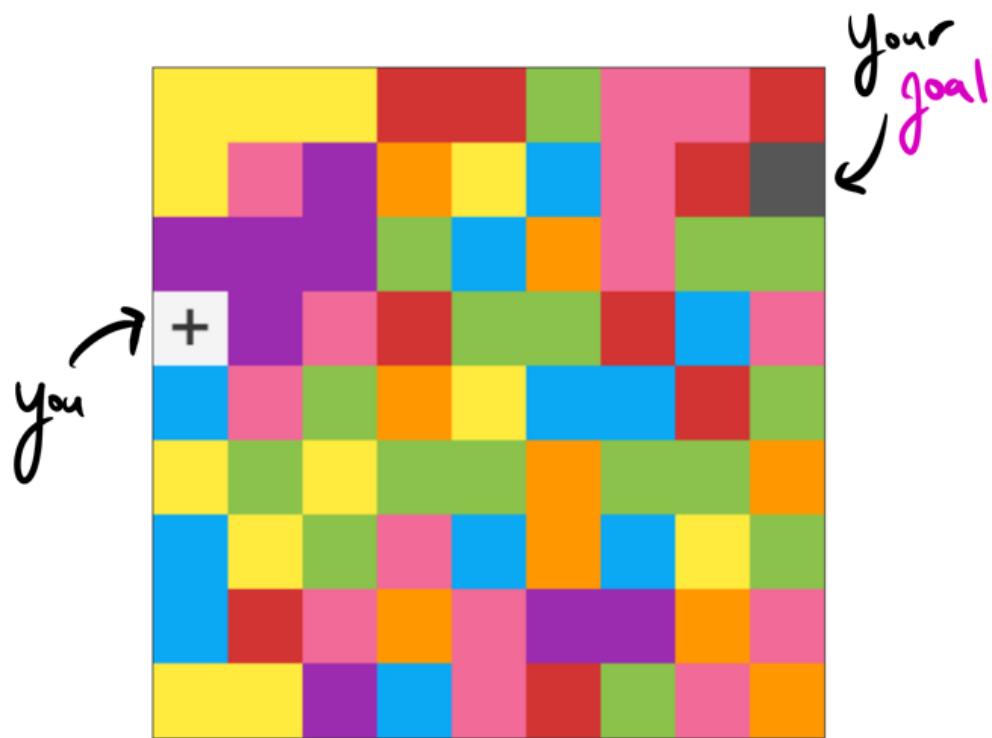


"Towards a New Impact Measure" was the title of [the post](#) in which AUP was introduced.



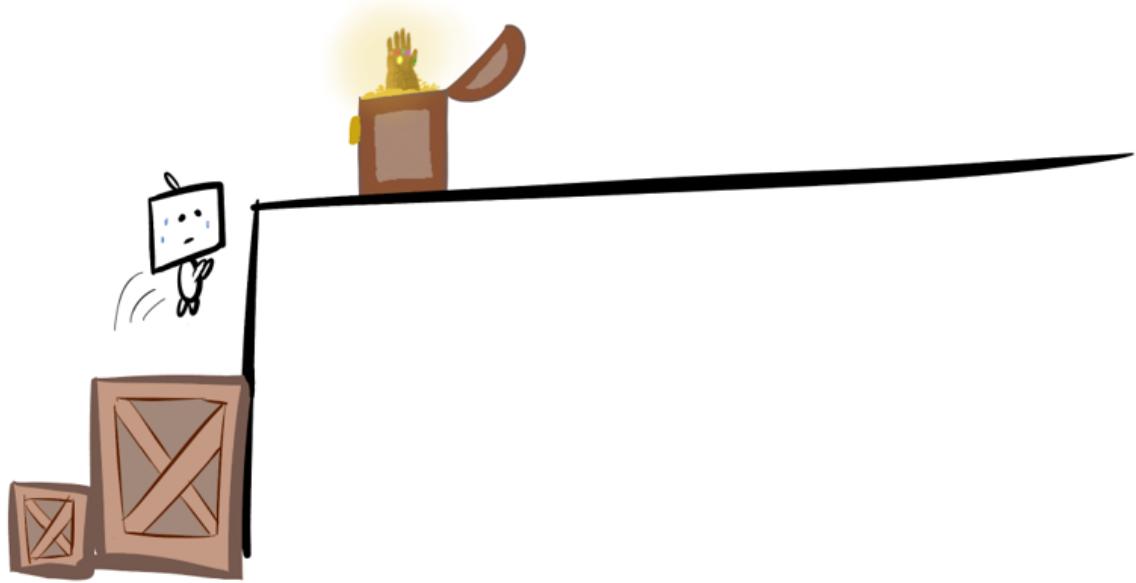
*Attainable Utility Theory: Why Things Matter*

This style of maze is from the video game *Undertale*.



*Seeking Power is Instrumentally Convergent in MDPs*

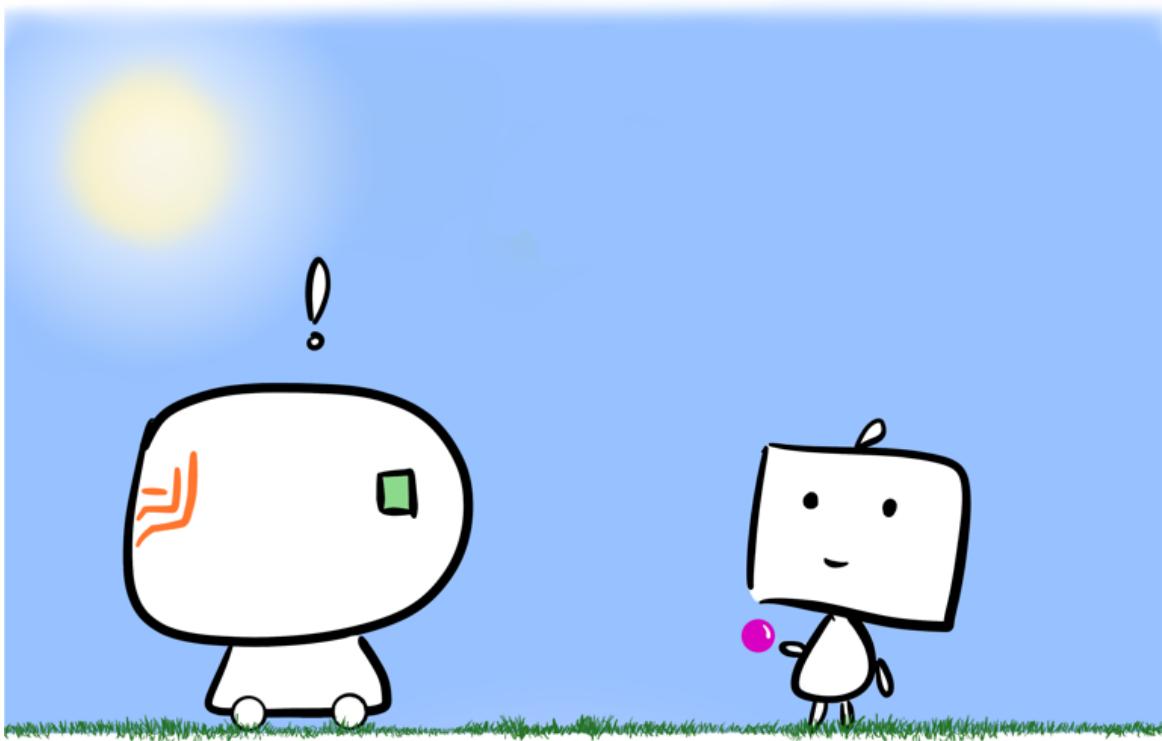
To seek power, Frank is trying to get at the Infinity Gauntlet.



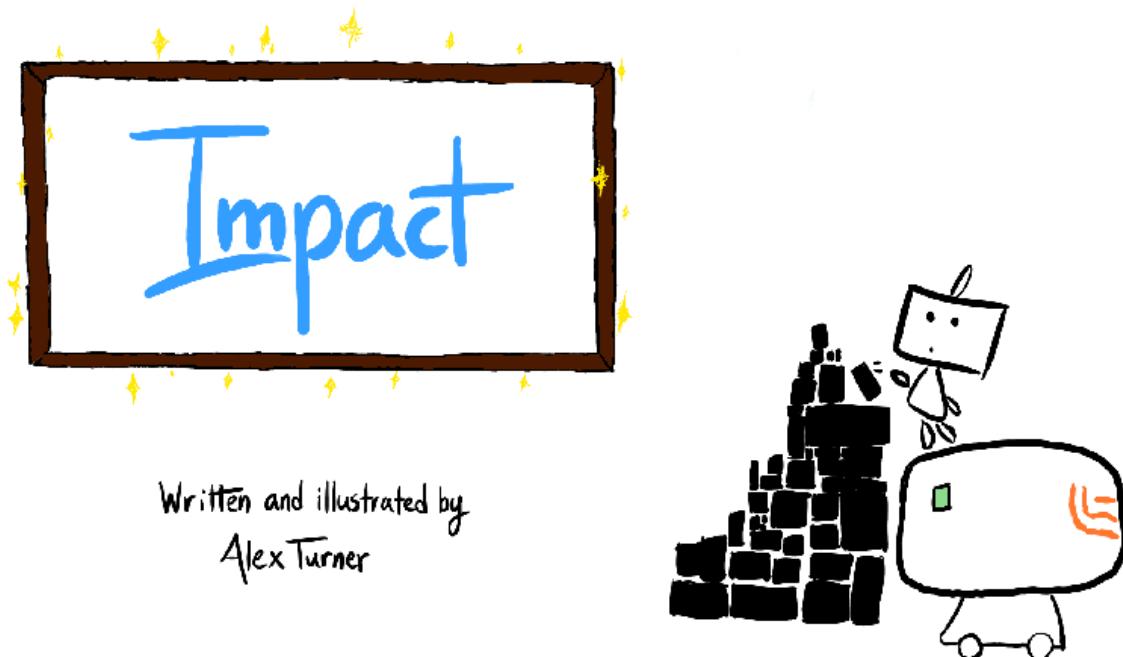
The tale of Frank and the orange Pebblehoarder

Speaking of under-tales, a friendship has been blossoming right under our noses.

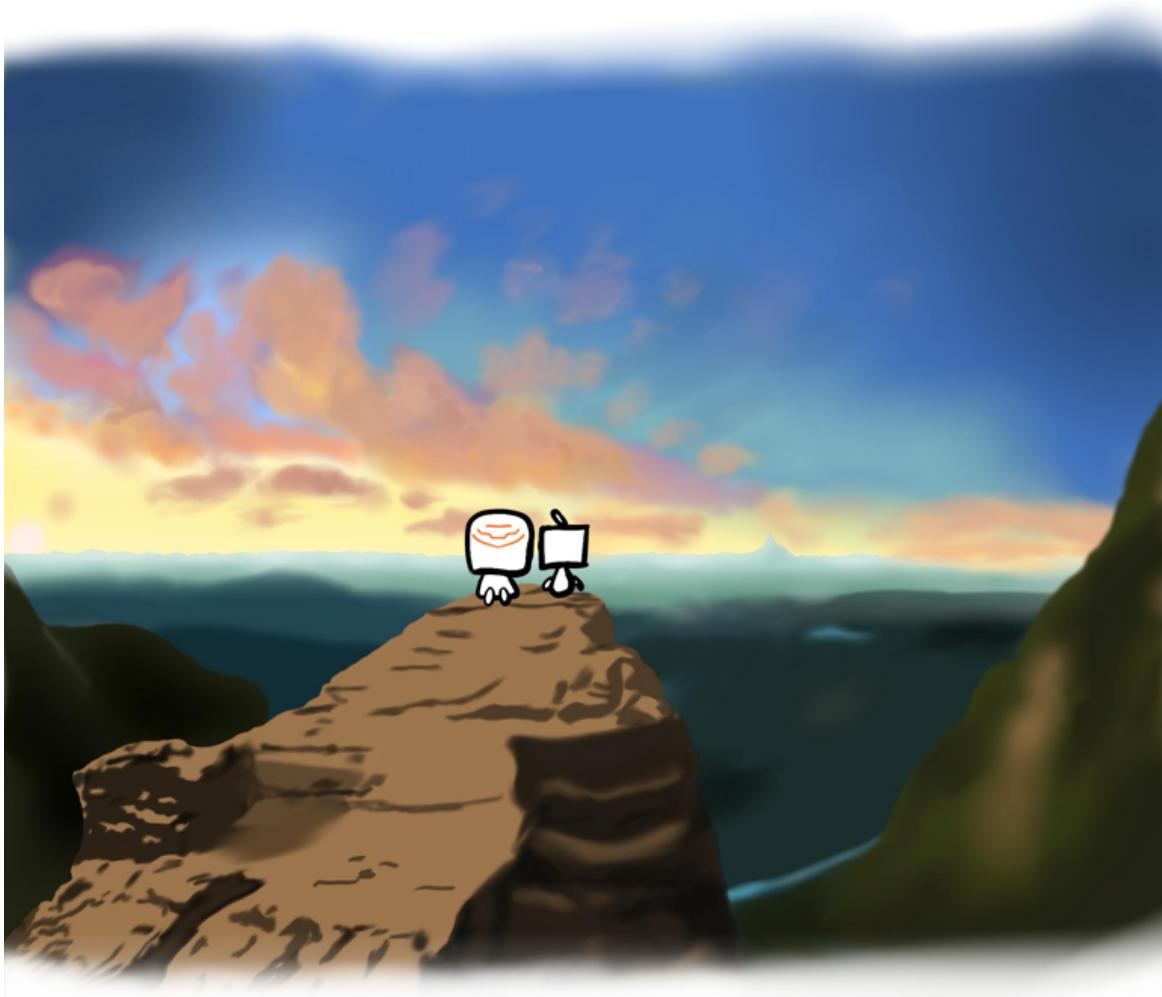
After the Pebblehoarders suffer the devastating transformation of all of their pebbles into obsidian blocks, Frank generously gives away his favorite pink marble as a makeshift pebble.



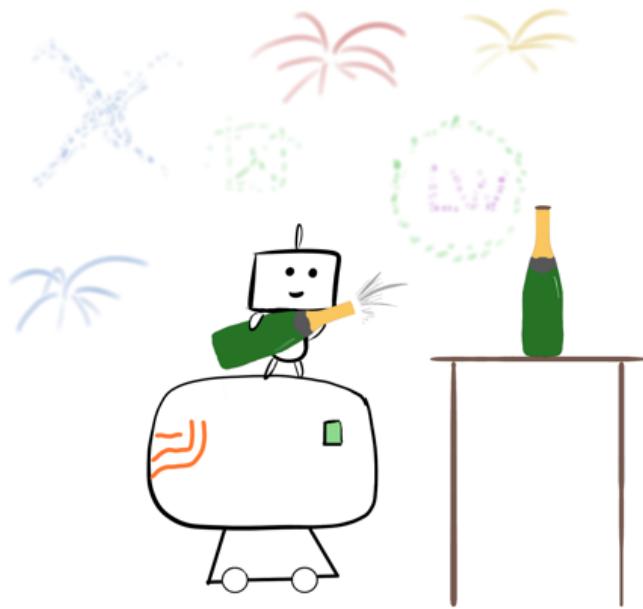
The title cuts to the middle of their adventures together, the Pebblehoarder showing its gratitude by helping Frank reach things high up.



This still at the midpoint of the sequence is from [the final scene of \*The Hobbit: An Unexpected Journey\*](#), where the party is overlooking Erebor, the Lonely Mountain. They've made it through the Misty Mountains, only to find Smaug's abode looming in the distance.



And, at last, we find Frank and orange Pebblehoarder popping some of the champagne from Smaug's hoard.



Since [Erebor isn't close to Gondor](#), we don't see Frank and the Pebblehoarder gazing at Ephel Dúath from Minas Tirith.