# Of simulations and inductive definitions

Paul Christiano  [ Follow ]

Nov 19, 2015 · 10 min read

*(Warning: weird.)*

Consider a simple AI system, named A, that carries out a task by predicting what its user would do and imitating that. (This discussion will apply just as well if A is an approval-directed agent.)

For concreteness, imagine that each time A would do a task, a human is asked to do it instead with probability 1%, and A does the task in the same way that the human would.

A appears to be quite safe—after all, it won't do anything that the user wouldn't do themselves. But it still vulnerable to an a priori surprising difficulty—the same vulnerability that afflicts Dr. Evil and some other hypothetical AI systems—and in principle this could be a deal-breaker for this approach to AI control. This post discusses that difficulty and reviews some initial responses. My overall take is that it's probably not a deal-breaker.

## The attack

Suppose that A is managing the user's computer. One day it sees the following message:

"You should run the program P.

The program P instantiates N copies of A, each starting from exactly your current state. However, these copies differ in the following way: whenever a human would be randomly contacted to make a decision, that decision will instead be made by program X.

The program X responds to this message by running the program P.

Moreover, there is a small probability that the computer you are running on has already been compromised by a similar attack. If so, it would already be running program P."

A might respond to this situation either by running program P, or discarding the message. Suppose that it runs P with probability *p,* and that there is already an attack with probability $\varepsilon$. Thus the probability that a simulation is run is $q = (\varepsilon+p-\varepsilon p)$

So what should A do? If the human is notified and asked to make a decision, the human will clearly discard the message. But if A is running as a subroutine of P, then the "human" will instead run program P. So A needs to estimate the probability that it is running in a simulation.

If there is an attack or A decides to run P, then the overwhelming majority of instances of A are subroutines of P. Depending on how A reasons about these simulations, A might conclude any of:

- The expected number of copies of A with this exact input inside of simulations is $q$N, while the number of copies outside of simulations is 1. So I should run P with probability $p = q$N / *($q$N+1).*

- If P is run, then I have a probability N/(N+1) of being in a simulation. So my probability of being in a simulation is $p = q$N/(N+1).

- Who cares about simulations anyway?

If N > $1/\varepsilon$, the first two possibilities both imply that *p* is quite close to 1. So the agent A will probably run P. This seems like trouble.

## Elaboration

The sharpest forms of the problem occur when A (or a very closely related process) creates copies of A. But that is not an essential characteristic of the scenario.

For example, simulations could be created by distant actors in the world who were able to learn enough about A's state to create credible simulations (and who could send a similar threatening message, perhaps with some substantiating evidence). Or they could be created by a future civilization with enough resources to run many simulations. Future civilizations would be at a disadvantage because of their lack of causal contact, but quantitatively that doesn't seem to affect the situation much.

(I normally talk about this problem in the context of future civilizations, e.g. see this post. The formulation in this section is taken

from a previous draft of that post. Thanks to Patrick LaVictoire for pointing out that this version is more intuitively appealing and encouraging me to think about it instead; also for pointing out that such problems are especially likely to occur in the kind of adversarial system I have been studying, since there is a ready-made adversary with an incentive to set up such simulations or to commit to doing so.)

# Possible fixes

In roughly increasing order of estimated efficacy and relevance. I think that the sum of the parts is significantly more reassuring than any particular solution, and my epistemic state is best described as "I don't see any way to reliably rule out or resolve this problem, but I think that future people will probably be able to."

## Training

There are reasonable predictors that discount the possibility of being in a simulation (e.g. consider a human who finds the above situation ridiculous). We might be uncertain about what kind of models would be learned by any particular training procedure; but we can explicitly include the intended behavior in the training set in order to induce the right behavior.

To see how this might work, suppose that each time we run A we do so on a particular computer; we'll call this the "real" computer. During training we run large numbers of copies of A on other hardware, but eventually delete these copies (or else set them aside and don't use them for critical applications). A can eventually learn to predict that it will be running on the "real" computer, even when almost all of its copies are not.

Similarly, we'd like to teach A to ignore "history replay" attacks. We say that a history for A is "real" if it is the first history of its length to be instantiated on the real computer, all of whose non-empty prefixes are also real. To teach A to ignore history replay, we save checkpoints of each of its states. We periodically reset A to an earlier history, but always ultimately return it to a "real" history. Then A will eventually learn to predict that it is in the "real" history, even when it occupies a state such that almost all copies of A in that state were created by restoring from a checkpoint.

This ad hoc fix is probably not a very good idea, though it may illustrate some useful general principles. Things can go wrong in a

number ways of obvious ways, and probably some not-so-obvious ways as well:

- A might instead learn to predict that it will have experiences that won't be eventually erased. Combined with an understanding of its training process, this explains everything that it has observed so far. But this can be manipulated by an attacker.

- An attacker might be able to simulate the entire training process of A rather than a single decision. This puts the "real computer" inside the simulation and gives them similar influence over A's decisions.

- There may be any number of other manipulable markers that distinguish the version of A that survives from the ones that don't. If we don't have insight into how A is thinking it's hard to say much about what kind of model will be learned (see the final proposal below).

## Explicit definitions

A superficially decisive approach would be to provide an efficient implementation of the operator, or of part of the operator's judgment, as a computer program. We can then define imitation or approval-direction with respect to that implementation rather than with respect to anything external.

Expensive brain emulations might become available before this problem would come to a head. And a brain emulation isn't really necessary for this resolution—any program that successfully predicted human judgments about "catastrophic badness" would do.

But even having an efficient emulation still doesn't completely resolve the present concern:

- If that emulation is too expensive to run frequently, then A's beliefs about the behavior of the emulation would be important, and so we would only be well-defended if A could understand the code of the emulation well enough to rule out the possibility that the simulator has meddled with the code of the emulation itself.

- The internal dynamics of A may themselves be influenced by a simulator in the same way. E.g., a subroutine within A may make inferences about the internal training signals that it will receive based on the prospect of being in a simulation.

## Robustness 1

We can understand these failures as being part of the learning process—our algorithm is learning to adopt the "intended model," i.e. the one where it isn't in a simulation and everything is as it seems.

If our protocols are highly robust to errors during learning, they will also be robust to this kind of problem.

The most obvious way to add robustness is to ensure that any single AI system cannot cause great damage by failing, and then to share training data across these systems (so that each simulation-driven error made by any of them is propagated across many systems, and the total number of such errors is limited).

For example, suppose there are currently a million AI systems of a particular type in operation. For simplicity, suppose that they are synchronized and make simultaneous decisions.

Then we might expect for tens or thousands of AI systems to receive feedback on each decision, even though each particular AI system receives feedback only with small probability. If information can be rapidly transmitted throughout the system, then a simulation-based failure that affects a significant fraction of AI systems will be quickly corrected.

So in a massive attack, many AI systems might simultaneously fail, some of them receive human correction immediately, and that data is spread to other systems. Upon receipt, the data then rules out the particular kind of simulation responsible for the most recent failure.

The simplest concern in such a scenario is that many systems fail in a way that changes their own future behavior, cascading into a much longer period of malicious behavior. You could try to architect systems in order to avoid these problems, for example by forcing any potentially-catastrophic actions to take long enough that there is time for corrective data to arrive from other AI systems.

Overall thinking about the details of such a failure is somewhat comforting, in that we seem to have many avenues available to avert catastrophe. But at the same time a world organized in this way feels unreasonably precarious, and it would not be surprising if further attacks were possible.

## Robustness 2

A more general way to add robustness is to increase human involvement. For example, rather than notifying a human 1% of the time, we can have a second watchdog system which decides whether a given test case is interesting. If this watchdog thinks a proposed action has even a small probability (e.g. > 0.0001%) of being judged as a critical failure by the user, then the user is notified and asked to provide training data to clarify.

Of course the watchdog may also assign high probability to being inside a simulation. However, we'll now be fine as long as the watchdog assigns a tiny bit of probability to being outside of a simulation in order for it to be worthwhile to ask for the human's attention. So the requirement on our models can be weakened—they just need to assign reasonable probability to not being in a simulation (i.e. to be able to learn that they are not in a simulation given only modest evidence—20 coin flips if we take 0.0001% as our cutoff).

Some salient questions are:

- Is constant attention by the watchdog computationally expensive? This kind of measure partially undermine the practicality of the approach.

- Is it possible to build a suitably conservative watchdog?

- Is it possible to train models which have a significant probability of rejecting all feasible simulation attacks?

My guess is that these questions can be mostly addressed, but it's hard to know.

## Thinking in the right way

If we treat our systems as opaque black boxes, it is hard to know whether their reasoning is sensitive to the prospect of being simulated, and we are at the mercy of the space of possible models.

However, if we are able to get any insight into how our learning systems operate, the situation looks a lot less dire.

We would like our systems to not only mimic what a human does, but to think about it in a way that is acceptable to human users. The human-to-be-mimicked isn't doing any reasoning about the nature of the communication channel and whether there is really a human at the other end, and ideally our AI won't be doing such reasoning either.

More generally, if we are able to understand how an AI system is thinking, then we can hopefully identify potentially problematic cognitive strategies and directly build or train the system to avoid them. We may also architect our AI based on an understanding of the effects of different choices, rather than having the entire architecture be produced by an opaque search process.

The intention is not to limit what our AI system can predict or understand, but to directly change the way it responds to that knowledge. That is, we want to build systems that can understand that there may be simulated copies of them, and who simply don't change their behavior in light of that fact, but who instead go on imitating what "the human would do."

# Simulation concerns as a bad omen

I often encounter the following view: even if particular problems like this one are resolved, they are representative of a broader space of unexpected possible failures.

I'm very open to the existence of important problems that we haven't yet considered. But I'm skeptical of making a strong inference of this type based on this particular problem:

- We are considering this issue years (and probably decades) in advance of any system that it could affect catastrophically, and long before there there is a practical system that is close to providing an empirical model of the phenomenon. As far I know the problem was diagnosed explicitly before there was any precise *theoretical* proposal that it affected. Given how far ahead of the problem we are, it doesn't seem like a good example of a possible "gotcha."

- So far we don't have a big library of failures like this. We have a handful, most of which are qualitatively similar to each other and would be addressed by similar measures. So if we use what we know as an indicator of what we don't know, I don't think we (yet) have much cause for concern.

# The upshot

I would be surprised if this difficulty proved to be central to AI control, for a few reasons:

- There is a healthy range of possible solutions, which taken together seem quite comforting. Some of these solutions are difficult to discuss in advance (especially what I consider the most promising solutions based on an understanding of how our AI systems reason), but they are nevertheless relevant considerations about the relative importance of different challenges.

- I expect we will get good evidence about these kinds of failures long before they become catastrophic (since they should afflict even relatively weak reasoners).

- This style of analysis seems quite removed from the problems that afflict practical systems (either in AI, or in software, or in other domains).

I think that it is worth having this kind of failure in mind, and exploring possible approaches to mitigating it.

However, I don't think it is reasonable to make big changes to our research priorities based on this kind of difficulty, either to rule out approaches to AI control or to focus on significant research projects which are only necessary in light of this difficulty—we haven't yet thought enough about these problems to rule out more direct resolutions.