

MIRI Briefing on Extinction Risk from AI

September 2024

I. The development of artificial superintelligence poses an imminent risk of human extinction.

"Artificial superintelligence" (ASI) refers to AI that can substantially surpass humans in all strategically relevant activities (economic, scientific, military, etc.).

The timeline to ASI is uncertain, but probably not long. On the present trajectory, MIRI would be uncomfortable ruling out the possibility that ASI is developed in the next year or two, and we'd be surprised if it was still several decades away.¹

AI labs are aggressively scaling up systems they don't understand. The Deep Learning techniques behind the rapid AI progress of the last few years create massive neural networks automatically. The resulting models are vast human-unreadable tangles of machine-written operations more "grown" than designed.² Labs basically discovered a "cheat": Engineers can't tell you why a modern AI makes a given choice, but have nevertheless released increasingly capable systems year after year.³

Sufficiently intelligent AIs will likely develop persistent goals of their own. Humans have wants, and make long-term plans, for reasons that we expect to also apply to sufficiently smart mechanically-grown AIs. (The computer science of this prediction does not fit into a paragraph; inquire further if interested.) We are only barely starting to see this phenomenon in today's AIs, which require a long training process to hammer them into the apparently-obedient form the public is allowed to see.⁴

We expect the the ASI's goals to be hollow and lifeless in the end. Imbuing a superhumanly intelligent AI with a deep, persistent care for worthwhile objectives is much more difficult than training it to answer the right way on an ethics test.⁵ Having spent two decades on the serious version of this problem, our informed view is that the field is nowhere near a solution.⁶

ASI that doesn't value us will end us. Unless it has worthwhile goals,⁷ ASI will put our planet to uses incompatible with our continued survival, just as we fail to concern ourselves with the crabgrass growing on the site of a planned parking lot.⁸ No malice, resentment, or misunderstanding is needed to precipitate our extinction.⁹

II. Human survival likely depends on delaying the creation of ASI as soon as we can for as long as necessary.

A “wait and see” approach to ASI is probably not survivable. A superintelligent adversary will not reveal its full capabilities and telegraph its intentions.¹⁰ It will not offer a fair fight. It will make itself indispensable¹¹ or undetectable until it can strike decisively and/or seize an unassailable strategic position.¹²

MIRI doesn’t see any viable quick fixes or workarounds to misaligned ASI. OpenAI admits that today’s most important methods of steering AI won’t scale to the superhuman regime.¹³ Attempts to restrain¹⁴ or deceive¹⁵ a superior intelligence are prone to fail for reasons both foreseeable and unforeseeable.¹⁶ Our own theoretical work suggests that plans to align ASI using unaligned AIs are similarly unsound.¹⁷ We also don’t think a well-funded crash program to solve alignment would be able to correctly identify solutions that won’t kill us.¹⁸ Our current view is that a safe way forward will likely require ASI to be delayed for a long time.¹⁹

Delaying ASI requires an effective worldwide ban on its development, and tight control over the factors of its production. This is a large ask,²⁰ but domestic oversight, mirrored by a few close allies, will not suffice. This is not a case where we just need the “right” people to build it before the “wrong” people do, as ASI is not a national weapon; it is a global suicide bomb.²¹ If anyone builds it, everyone dies.

To preserve the option of shutting down ASI development if or when the will is found, MIRI advocates promptly building the off-switch.²² The “off-switch” refers to the systems and infrastructure needed for the eventual enactment of a ban.²³ It starts with identifying the relevant parties, tracking the relevant hardware, and requiring that advanced AI work take place within a limited number of monitored and secured locations. It extends to building out the protocols, plans, and chain of command to be followed in the event of a shutdown decision. As the off-switch could also provide resilience to more limited AI mishaps, we hope it will find broader near-term support than a full ban.²⁴

An off-switch can only prevent our extinction from ASI if it has sufficient reach and is actually used to shut down development in time.²⁵ If humanity is to survive this dangerous period, it will have to stop treating AI as a domain for international rivalry and demonstrate a collective resolve equal to the threat.

Endnotes

