<u>Obelisk</u> is a research lab focused on <u>brainlike AGI safety</u> and the development of <u>brainlike AGI</u> more generally. They think that there is a good chance that the first AGIs will be <u>neuromorphic</u>. Aligning brainlike AI comes with a different set of opportunities and challenges than other development paths, and Obelisk is studying <u>computational neuroscience</u> and human goal formation to prepare and explore the options.

One threat model they think about is <u>misaligned model-based RL agents</u>, and two possible paths to alignment they're looking at are <u>"Controlled AGI"</u> and <u>"Social-instinct AGI"</u>2.

The main crux for their agenda is whether transformative AI will be brainlike, and whether it's possible to sufficiently align brainlike AI. They are relatively agnostic on the exact failure modes caused by unaligned brainlike AGI.

Interested in helping out?

Making progress on Obelisk's agenda requires advanced knowledge of neuroscience and related disciplines. A strong grounding in machine learning is an asset here, but being a top-level ML engineer or ML research scientist is not required. They hire both within and outside academia.

Related

- What approaches are Al alignment organizations working on?
- B How would we align an AGI whose learning algorithms / cognition look like human ...

¹ Automatically assessing the Al's thoughts and enforcing conservatism in value extrapolation

² Building an Al which deeply feels something like an idealized form of kinship and love for humanity, by understanding the computational structure of these drives in humans.