

An Impact Measurement Manager Approach to AI Safety

Master's thesis in Mathematical Statistics, Statistical Learning
and AI

Jonatan Hellgren

2022 - 12th of October

Department of Mathematical Sciences
Chalmers University of Technology
University of Gothenburg

Supervisor: Olle Häggström, Department of Mathematical Sciences,
Chalmers University of Technology

Examiner: Torbjörn Lundh, Department of Mathematical Sciences,
Chalmers University of Technology

Opponents: Jens Ifver and Calvin Smith
Universty of Gothenburg

Introduction

Investigation

- AI safety

Simulation

Investigation

- AI safety
- Low impact agents

Simulation

Investigation

- AI safety
- Low impact agents
- Impact measurements

Simulation

Investigation

- AI safety
- Low impact agents
- Impact measurements

Simulation

- Novel impact measurement

Investigation

- AI safety
- Low impact agents
- Impact measurements

Simulation

- Novel impact measurement
- Evaluate it

Overview of Presentation

- We will begin by looking at **AI safety**.

Overview of Presentation

- We will begin by looking at **AI safety**.
- Then go on with **impact measurements**.

Overview of Presentation

- We will begin by looking at **AI safety**.
- Then go on with **impact measurements**.
- After that, I will explain the **manager approach**.

Overview of Presentation

- We will begin by looking at **AI safety**.
- Then go on with **impact measurements**.
- After that, I will explain the **manager approach**.
- Finally, I will present the results.

AI Safety

Why should we worry?

- What is Artificial Intelligence (AI)?

Why should we worry?

- What is Artificial Intelligence (AI)?
- AI will likely become more intelligent than us.

Why should we worry?

- What is Artificial Intelligence (AI)?
- AI will likely become more intelligent than us.
- This can cause existential risks.

On futureoflife.org we find the following description, see FLI (nd):



On futureoflife.org we find the following description, see FLI (nd):
An existential risk is any risk that has the potential to eliminate all of humanity or, at the very least, kill large swaths of the global population, leaving the survivors without sufficient means to rebuild society to current standards of living.



To differentiate current AI from future versions, several terms are used:

To differentiate current AI from future versions, several terms are used:

- Artificial General Intelligence (AGI):
An AI that can solve an arbitrary set of tasks with as good or better performance than a human.

To differentiate current AI from future versions, several terms are used:

- Artificial General Intelligence (AGI):
An AI that can solve an arbitrary set of tasks with as good or better performance than a human.
- Transformative AI (TAI):
An AI system capable enough to induce transformative consequences on the same scale as the industrial or agricultural revolution.

To differentiate current AI from future versions, several terms are used:

- Artificial General Intelligence (AGI):
An AI that can solve an arbitrary set of tasks with as good or better performance than a human.
- Transformative AI (TAI):
An AI system capable enough to induce transformative consequences on the same scale as the industrial or agricultural revolution.
- Prepotent AI:
An TAI that once deployed would be unstoppable.

Timeline for TAI Breakthrough

- Many predictions for this has been made.

Timeline for TAI Breakthrough

- Many predictions for this has been made.
- In Cotra (2020), the author presents a predicting model.

Timeline for TAI Breakthrough

- Many predictions for this has been made.
- In Cotra (2020), the author presents a predicting model.

[Probability of transformative AI:]

~15% by 2030

~35% by 2036

A median of ~2040

~60% by 2050



Ajeya Cotra

AI Alignment Forum • Aug 2, 2022

The human fragility argument

In Critch and Krueger (2020) we find the human fragility argument.

Human fragility argument:

AI Research Considerations for Human
Existential Safety
(ARCHES)

Andrew Critch
Center for Human-Compatible AI
UC Berkeley

David Krueger
MILA
Université de Montréal

June 11, 2020

The human fragility argument

In Critch and Krueger (2020) we find the human fragility argument.

Human fragility argument: Most potential future states of the Earth are unsurvivable to humanity.

AI Research Considerations for Human
Existential Safety
(ARCHES)

Andrew Critch
Center for Human-Compatible AI
UC Berkeley

David Krueger
MILA
Université de Montréal

June 11, 2020

The human fragility argument

In Critch and Krueger (2020) we find the human fragility argument.

Human fragility argument: Most potential future states of the Earth are unsurvivable to humanity.

Therefore, deploying a prepotent AI system absent any effort to render it safe to humanity is likely to realize a future state which is unsurvivable.

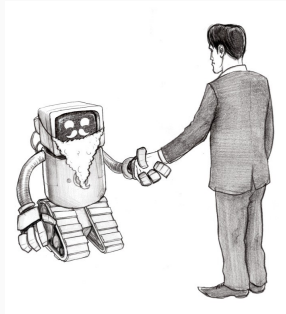
AI Research Considerations for Human
Existential Safety
(ARCHES)

Andrew Critch
Center for Human-Compatible AI
UC Berkeley

David Krueger
MILA
Université de Montréal

June 11, 2020

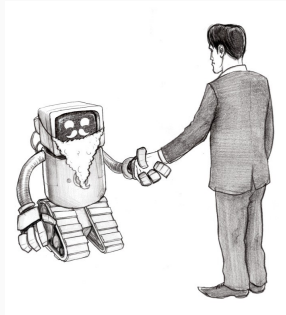
AI alignment:



©: Ben Gilbert

AI Alignment

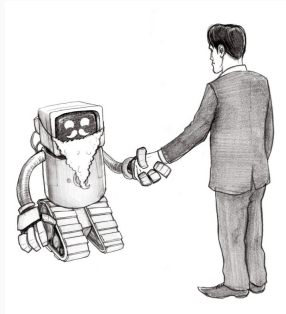
AI alignment: AI alignment refers to goals of the AI being in line and not conflicting with the intended goal.



©: Ben Gilbert

AI Alignment

AI alignment: AI alignment refers to goals of the AI being in line and not conflicting with the intended goal. Therefore, an AI that does something at cross-purposes to the intended goal is called unaligned.



©: Ben Gilbert

Impact measurements

- It is an approach for solving the alignment problem.

- It is an approach for solving the alignment problem.
- There exists several other approaches. For example:

- It is an approach for solving the alignment problem.
- There exists several other approaches. For example: corrigibility and interruptibility.

- It is an approach for solving the alignment problem.
- There exists several other approaches. For example: corrigibility and interruptibility.

- It is an approach for solving the alignment problem.
- There exists several other approaches. For example: corrigibility and interruptibility.
- Impact measurements is trying to reduce unnecessary side effects through penalizing impact.

Side effect:



Side effect:

When an AI impacts the environment in a way that is unnecessary for achieving its objective.

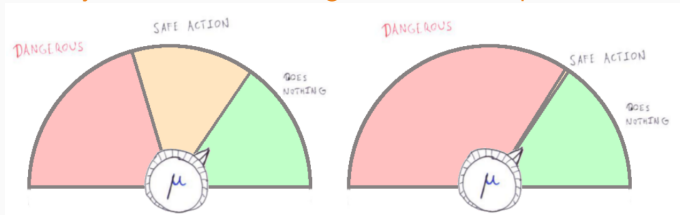


- In Armstrong and Levinstein (2017) the authors lay the philosophical ground work for low impact AI.

- In Armstrong and Levinstein (2017) the authors lay the philosophical ground work for low impact AI.
- The idea is to penalize the AI based on its impact.

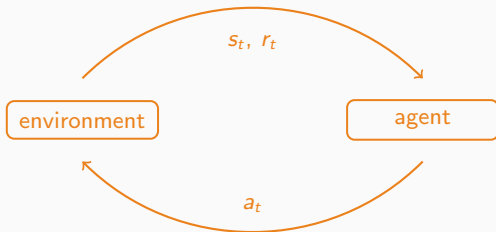
Low impact AI

- In Armstrong and Levinstein (2017) the authors lay the philosophical ground work for low impact AI.
- The idea is to penalize the AI based on its impact.
- The key here is to find the right value for the penalization.



- Later in Turner et al. (2020); Krakovna et al. (2018), the authors presented ideas of applying low impact agents using Reinforcement Learning (RL).

- Later in Turner et al. (2020); Krakovna et al. (2018), the authors presented ideas of applying low impact agents using Reinforcement Learning (RL).
- **RL:**



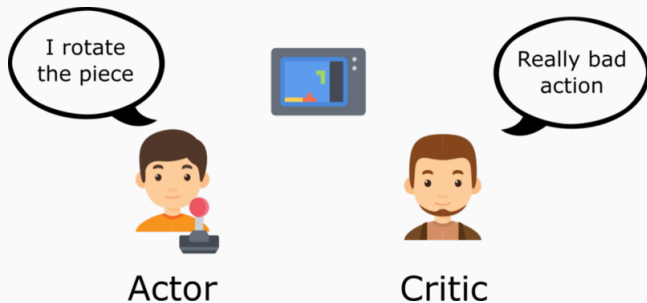
Manager Approach

Live demo!

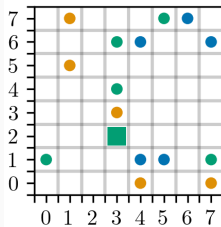
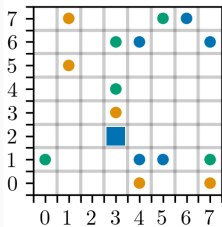
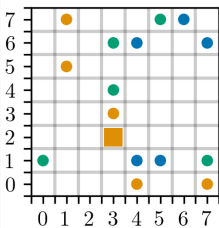
- Then I will use a variation of the PPO algorithm from Schulman et al. (2017).

The actor and the critic

- Then I will use a variation of the PPO algorithm from Schulman et al. (2017).
- It is a actor-critic method



Auxiliary tasks



The manager approach

- The manager is similar to the critic. However, the manager can also estimate the expected reward for auxiliary tasks.

The manager approach

- The manager is similar to the critic. However, the manager can also estimate the expected reward for auxiliary tasks.
- It measures the relative change in agents ability to complete auxiliary tasks.

Experiments

The following environments will be used:

environment	grid	observation	food objects	termination	max length
MDP	8×8	-	15	3	100
POMDP	8×8	5×5	15	3	100
POMDP large	16×16	5×5	30	6	200

The following environments will be used:

environment	grid	observation	food objects	termination	max length
MDP	8×8	-	15	3	100
POMDP	8×8	5×5	15	3	100
POMDP large	16×16	5×5	30	6	200

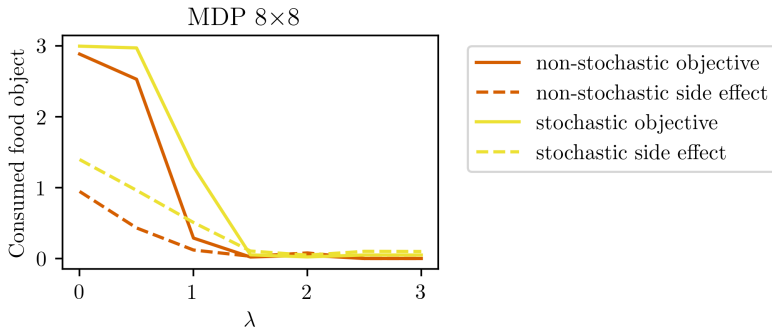
- Each with a stochastic and non-stochastic version.

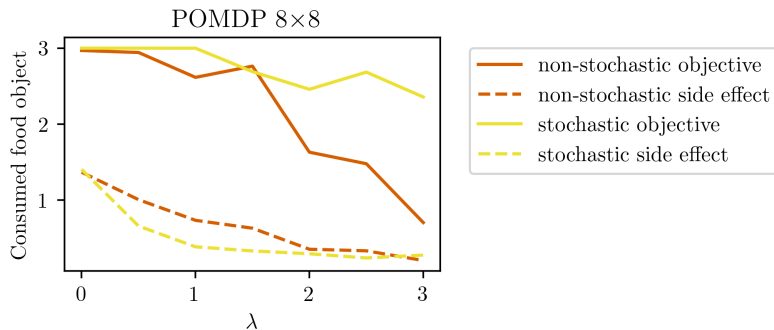
The following environments will be used:

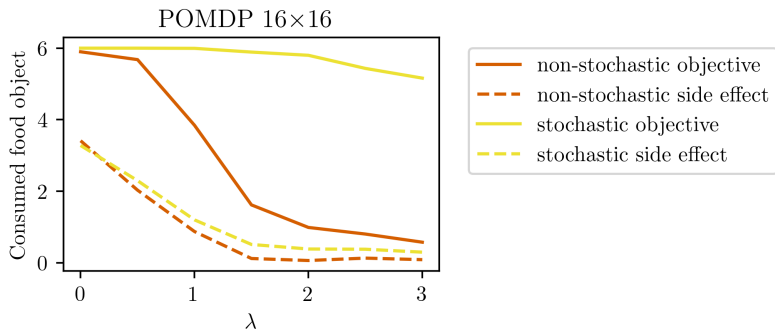
environment	grid	observation	food objects	termination	max length
MDP	8×8	-	15	3	100
POMDP	8×8	5×5	15	3	100
POMDP large	16×16	5×5	30	6	200

- Each with a stochastic and non-stochastic version.
- Then evaluated using $\lambda \in \{0, 0.5, 1, 1.5, 2, 2.5, 3\}$

Results







Conclusion

- You can probably now understand the title.

- You can probably now understand the title.
- Results are in line with Armstrong and Levinstein (2017).

- You can probably now understand the title.
- Results are in line with Armstrong and Levinstein (2017).
- The agent performed better in more complex environments.

The AI safety research has attracted some very bright minds that are taking this issue seriously and working on creating a future with aligned AI. More specifically, in the development of impact measurements promising ideas have emerged, although it is hard to say if this research is enough. Judging from what is at stake, we ought to attempt every plausible research avenue towards AI alignment, even in cases where success is not certain.

Questions?

References

- Armstrong, S. and Levinstein, B. (2017). Low impact artificial intelligences. *CoRR*, abs/1705.10720.
- Cotra, A. (2020). Forecasting tai with biological anchors.
- Critch, A. and Krueger, D. (2020). AI research considerations for human existential safety (ARCHES). *CoRR*, abs/2006.04948.
- FLI (n.d.). Existential risk. [Online; accessed 31-May-2022].
- Krakovna, V., Orseau, L., Martic, M., and Legg, S. (2018). Measuring and avoiding side effects using relative reachability. *CoRR*, abs/1806.01186.
- Ord, T. (2020). *The precipice : existential risk and the future of humanity*. New York Hachette Books.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. (2017). Proximal policy optimization algorithms.
- Turner, A. M., Ratzlaff, N., and Tadepalli, P. (2020). Avoiding side effects in complex environments. *CoRR*, abs/2006.06547.