

# Inner Optimization

Evan Hubinger, Chris van Merwijk, Vladimir Mikulik, and Joar Skalse

*with special thanks to Paul Christiano, Scott Garrabrant,*

*and the MIRI agent foundations team for originating many of the ideas discussed in this paper.*

(Dated: February 16, 2019)

ABSTRACT.

## I. INTRODUCTION

BODY. [1]

---

- [1] Eliezer Yudkowsky, “Optimization daemons,” .
- [2] Riceissa, “Optimization daemons,” (2018).
- [3] K. E. Drexler, “Reframing superintelligence: Comprehensive ai services as general intelligence,” Technical Report #2019-1, Future of Humanity Institute, University of Oxford (2019).
- [4] Daniel Filan, “Bottle caps aren’t optimisers,” (2018).
- [5] Jiawei Su, Danilo Vasconcellos Vargas, and Kouichi Sakurai, “One pixel attack for fooling deep neural networks,” CoRR (2017).
- [6] Jan Leike, David Krueger, Tom Everitt, Miljan Martic, Vishal Maini, and Shane Legg, “Scalable agent alignment via reward modeling: a research direction,” CoRR (2018).
- [7] Paul Christiano, “What does the universal prior actually look like?” (2016).
- [8] Alex Graves, Greg Wayne, and Ivo Danihelka, “Neural turing machines,” CoRR (2014).
- [9] Paul Christiano, “Open question: are minimal circuits daemon-free?” (2018).
- [10] Borja Ibarz, Jan Leike, Tobias Pohlen, Geoffrey Irving, Shane Legg, and Dario Amodei, “Reward learning from human preferences and demonstrations in atari,” CoRR (2018).
- [11] Kareem Amin and Satinder Singh, “Towards resolving unidentifiability in inverse reinforcement learning,” CoRR (2016).
- [12] Nick Bostrom, *Superintelligence: Paths, Dangers, Strategies* (Oxford University Press, 2014).
- [13] David Manheim and Scott Garrabrant, “Categorizing variants of goodhart’s law,” CoRR (2018).
- [14] Paul Christiano, “Worst-case guarantees,” (2019).
- [15] Jane X Wang, Zeb Kurth-Nelson, Dhruva Tirumala, Hubert Soyer, Joel Z Leibo, Remi Munos, Charles Blundell, Dhharshan Kumaran, and Matt Botvinick, “Learning to reinforcement learn,” CoRR (2016).
- [16] Yan Duan, John Schulman, Xi Chen, Peter L. Bartlett, Ilya Sutskever, and Pieter Abbeel, “RL<sup>2</sup>: Fast reinforcement learning via slow reinforcement learning,” CoRR (2016).

- Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané, “Concrete problems in ai safety,” CoRR (2016).
- Stuart Armstrong and Sören Mindermann, “Occam’s razor is insufficient to infer the preferences of irrational agents,” CoRR (2017).
- Xiaowei Huang, Marta Kwiatkowska, Sen Wang, and Min Wu, “Safety verification of deep neural networks,” CoRR (2016).
- Guy Katz, Clark Barrett, David Dill, Kyle Julian, and Mykel Kochenderfer, “Reluplex: An efficient smt solver for verifying deep neural networks,” CoRR (2017).
- Kexin Pei, Yinzhi Cao, Junfeng Yang, and Suman Jana, “Towards practical verification of machine learning: The case of computer vision systems,” CoRR (2017).
- Nate Soares, Benja Fallenstein, Eliezer Yudkowsky, and Stuart Armstrong, “Corrigibility,” CoRR (2015).