Chapter 3 Security Solutions for Intelligent and Complex Systems

Stuart Armstrong *Future of Humanity Institute, UK*

Roman V. Yampolskiy JB Speed School of Engineering, USA

ABSTRACT

Superintelligent systems are likely to present serious safety issues, since such entities would have great power to control the future according to their possibly misaligned goals or motivation systems. Oracle AIs (OAI) are confined AIs that can only answer questions and do not act in the world, represent one particular solution to this problem. However even Oracles are not particularly safe: humans are still vulnerable to traps, social engineering, or simply becoming dependent on the OAI. But OAIs are still strictly safer than general AIs, and there are many extra layers of precautions we can add on top of these. This paper begins with the definition of the OAI Confinement Problem. After analysis of existing solutions and their shortcomings, a protocol is proposed aimed at making a more secure confinement environment which might delay negative effects from a potentially unfriendly superintelligence while allowing for future research and development of superintelligent systems.

DOI: 10.4018/978-1-5225-0741-3.ch003

Copyright ©2017, IGI Global. Copying or distributing in print or electronic forms without written permission of IGI Global is prohibited.

INTRODUCTION

With the likely development of superintelligent programs in the near future, many scientists have raised the issue of safety as it relates to such technology (Bostrom, 2006; Chalmers, 2010; Hall, 2000; Hibbard, 2005; Yampolskiy, 2011a, 2011b; Yampolskiy & Fox, 2012a, 2012b; Yudkowsky, 2008). A common theme in Artificial Intelligence (AI¹) safety research is the possibility of keeping a superintelligent agent in a sealed hardware so as to prevent it from doing any harm to humankind. Such ideas originate with scientific visionaries such as Eric Drexler who has suggested confining transhuman machines so that their outputs could be studied and used safely (Drexler, 1986). Similarly, in 2010 David Chalmers proposed the idea of a "leakproof" singularity (Chalmers, 2010). He suggested that for safety reasons, AI systems first be restricted to simulated virtual worlds until their behavioral tendencies could be fully understood under the controlled conditions.

This chapter is based on combined and extended information from three previously published papers: (Armstrong, 2011; Armstrong, Sandberg, & Bostrom, 2012; Yampolskiy, 2012a)*. We evaluate feasibility of previously presented proposals and suggest a protocol aimed at enhancing safety and security of such methodologies. While it is unlikely, that long-term and secure confinement of AI is possible, we are hopeful that the proposed protocol will give researchers a little more time to find a permanent and satisfactory solution for addressing existential risks associated with appearance of superintelligent machines.

In this chapter we will review specific proposals aimed at creating restricted environments for safely interacting with artificial minds. The key question is: are there strategies that reduce the potential existential risk from a superintelligent AI so much that while implementing it as a free AI would be impermissible a confined implementation would be permissible? The chapter will start by laying out the general design assumptions for the confined AI and formalizing the notion of confinement. Then it will touch upon some of the risks and dangers deriving from the humans running and interaction with the confined AI. The final section looks at some of the other problematic issues concerning the confined AI, such as its ability to simulate human beings within it and its status as a moral agent itself.

Motivation for AI Confinement

There are many motivations to pursue the goal of developing AI. While some motivations are non-instrumental, such as scientific and philosophical curiosity about the nature of thinking or a desire for creating non-human beings, a strong set of motivations is the instrumental utility of AI. Such machines would benefit their owners by being able to do tasks that currently require human intelligence, and possibly tasks 50 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage: <u>www.igi-</u> <u>global.com/chapter/security-solutions-for-intelligent-and-</u> <u>complex-systems/164692</u>

Related Content

Generic Application Security in Current and Future Networks

Silke Holtmannsand Pekka Laitinen (2008). *Handbook of Research on Wireless Security (pp. 379-394).* www.irma-international.org/chapter/generic-application-security-current-future/22059

Impact of Employer Branding on Job Engagement and Organizational Commitment in Indian IT Sector

Geeta Rana, Ravindra Sharma, S.P Singhand Vipul Jain (2019). *International Journal of Risk and Contingency Management (pp. 1-17).* www.irma-international.org/article/impact-of-employer-branding-on-job-engagement-andorganizational-commitment-in-indian-it-sector/228997

An Integrated Machine Learning Framework for Fraud Detection: A Comparative and Comprehensive Approach

Karim Ouazzane, Thekla Polykarpou, Yogesh Pateland Jun Li (2022). International Journal of Information Security and Privacy (pp. 1-17).

www.irma-international.org/article/an-integrated-machine-learning-framework-for-frauddetection/300314

Data Reduction

Yu Wang (2009). Statistical Techniques for Network Security: Modern Statistically-Based Intrusion Detection and Protection (pp. 172-219). www.irma-international.org/chapter/data-reduction/29698

Data Breach on Consumer Behavior

Ruimin Chen (2021). Research Anthology on Privatizing and Securing Data (pp. 1861-1879).

www.irma-international.org/chapter/data-breach-on-consumer-behavior/280260