

# *Anthropomorphic Algorithms for AI Safety*

Tom Wallis

## *Anthropomorphic algorithms*

FOR A LONG TIME, sociotechnical systems analysis within computing science has been developing formalisms of human-like traits, such as trust and comfort. These allow an intelligent agent to interact with other agents in its environment in measured, cautious ways; they might be used, for example, to decide whether it should accept information from another agent if its behaviour is becoming erratic (or to discard previous data which is no longer “trustworthy”).

These anthropomorphic algorithms are undergoing continual improvement<sup>1,2</sup>, but two problems remain untouched:

1. Various different anthropomorphic algorithms have been developed, but none have been combined into a system with several traits.

For example, an algorithm might be designed where an agent’s ratings of trust and comfort in a given scenario influence each other — not unlike a human’s lack of trust in an agent making it less comfortable with certain situations.

2. Lots of philosophical questions arise when developing anthropomorphic agents. There are questions in roboethics, such as the morality of creating an intelligent agent which might exhibit racial bias after its training. There are also questions in machine ethics, involving ethical decisions that an agent might make, and how they can be affected by trust and comfort.

HOWEVER, one exciting unexplored problem is that of AI safety: could anthropomorphic algorithms give humanity an edge in developing friendly AI? If we can limit its space of mind<sup>3</sup>, perhaps we can limit potential damage from an artificial agent. Perhaps developing anthropomorphic agents allows us to reason better about how an artificial agent can be unfriendly, allowing us to better predict catastrophes related to the agent turning malicious. Perhaps such a method is provably inadequate for solving the problem of AI safety. In any case, further research is required to determine the method’s efficacy.

<sup>1</sup> Seifeddine Kramdi. A modal approach to model computational trust. *PhD Thesis*, 2015. URL <https://tel.archives-ouvertes.fr/tel-01328169>

<sup>2</sup> Joana Urbano, Ana Paula Rocha, and Eugénio Oliveira. An approach to computational social trust. 27:113–131, 2014. DOI: 10.3233/AIC-130587

<sup>3</sup> Murray Shanahan. *Embodiment and the Inner Life: Cognition and Consciousness in the Space of Possible Minds*. Oxford University Press, 2010; and Aaron Sloman. The Structure of the Space of Possible Minds. pages 35–42, 1984

## *Responsibility and its implications*

AS A MASTERS STUDENT at Glasgow University, my current research is in developing a computational formalism of responsibility. This formalism would be the first algorithmic definition of an agent's responsibility, and fits the above problems perfectly. Responsibility has a wealth of philosophical literature already at its disposal, making it an ideal platform for beginning to address problems of machine ethics and roboethics. The formalism is similar to current trust and comfort models, making it easy to integrate with existing frameworks into an agent with several anthropomorphic traits. Most interestingly, an intelligent agent with a concept of responsibility is useful to analyse from the perspective of AI safety in a way that trust and comfort models are less suitable.

As a result of the wealth of literature on responsibility for human agents, much work can be done to teach artificial agents to act in responsible ways; either in developing machine learning algorithms which tune the parameters of an agent's feeling of responsibility, or in imposing a strict sense of responsibility on that agent. That computational responsibility might improve agent corrigibility<sup>4</sup> is an exciting prospect. Computational responsibility also provides an analogous framework to existing human responsibility to begin reasoning about how intelligent agents might relate to traditional moral responsibility. Plenty of work relating to existing philosophical literature becomes applicable to intelligent agents on the advent of computational responsibility.

I PROPOSE that the breadth and practicality of this work represents a substantial addition to the current literature on intelligent agents, and that the introduction of computational responsibility to the growing arsenal of anthropomorphic algorithms represents a turning point in the relevance of anthropomorphic algorithms to philosophical literature. Moreover, it offers a rich and exciting opportunity for collaboration between Computing Science and Philosophy. At a stretch, anthropomorphic algorithms may even represent a new area in the study of artificial intelligence safety, which promises to advance literature for both computing science and philosophy.

One interesting question which springs to mind is that of space of possible minds. Earlier it was alluded to that anthropomorphic algorithms might limit the space of possible minds that an intelligent agent can occupy; computational responsibility, on the other hand, might broaden the space of possible responsibilities that an intelligent agent might inhabit. These agents are not bound by human instincts like self preservation and can be expected to "value" different things, making a theory of moral philosophy for a broader range of minds than biological ones possible and practical.

<sup>4</sup> Stuart Armstrong, Nate Soares, Benja Fallenstein, and Eliezer Yudkowsky. Corrigibility. *AAAI Workshops: Workshops at the Twenty-Ninth AAAI Conference on Artificial Intelligence*, 2015

## *My suitability*

Given my experience developing this computational responsibility formalism, I am uniquely equipped to begin the research to be done. The formalism I have designed has been purposefully created with a philosophical foundation in mind, drawing from work by P.F. Strawson<sup>5</sup> and Ben Colburn. In addition, the breadth of the work to be done makes it ideal for pursuit as a PhD project.

MY OTHER EXPERIENCE in research falls into two fields.

THE FIRST, sociotechnical systems, involves analysis of complex systems of people and technology they interact with; my work was to create a modelling system which allowed simulation of sociotechnical workflows without barriers to entry like understanding of a domain-specific modelling language. These models were then dynamically altered during runtime by a code fuzzer I designed and implemented, which injected human-like variance into the model, allowing for accurate simulations of human behaviour using simple modelling techniques. This research won an award for “Best Software Product” for my year.

OTHER RESEARCH I pursue lies in the field of experimental storytelling. Project Albert explores the possibility that improvisation of children’s bedtime stories might be made easier and more accessible by use of design patterns, which generalise complex ideas by turning them into interrelated rules which are defined semantically. A full suite of design patterns have been developed, with example stories which follow the patterns. The intention is to provide a framework for parents to share stories with their children which have some sentimental value as a result of the improvisation and random aspect, as well as to provide a way for parents to connect through stories when they cannot necessarily afford books to read from.

<sup>5</sup> P.F. Strawson. Freedom and resentment. *Proceedings of the British Academy*, Vol. 48, 1960

A paper on this work is currently being developed, which can be found at <http://bit.ly/2gi4GDo>. My original dissertation can be found at <http://bit.ly/2fYbcgy>.

Information on Project Albert can be found at <http://projectalbert.net/>.

An essay on the work so far and its efficacy is currently under development, which can be found at <http://bit.ly/2fZvggr>.