

Towards Moral Autonomous Systems

Vicky Charisi¹, Louise Dennis², Michael Fisher², Robert Lieck³,
Andreas Matthias⁴, Marija Slavkovik⁵, Janina Loh (Sombetzki)⁶,
Alan F. T. Winfield⁷, and Roman Yampolskiy⁸

¹University of Twente

²University of Liverpool

³Universität Stuttgart

⁴Lingnan University

⁵marija.slavkovik@uib.no, University of Bergen

⁶University of Vienna

⁷University of the West of England, Bristol

⁸University of Louisville

November 1, 2017

Abstract

Both the ethics of autonomous systems and the problems of their technical implementation have by now been studied in some detail. Less attention has been given to the areas in which these two separate concerns meet. This paper, written by both philosophers and engineers of autonomous systems, addresses a number of issues in machine ethics that are located at precisely the intersection between ethics and engineering. We first discuss the main challenges which, in our view, machine ethics poses to moral philosophy. We then consider different approaches towards the conceptual design of autonomous systems and their implications on the ethics implementation in such systems. Then we examine problematic areas regarding the specification and verification of ethical behavior in autonomous systems, particularly with a view towards the requirements of future legislation. We discuss transparency and accountability issues that will be crucial for any future wide deployment of autonomous systems in society. Finally we consider the, often overlooked, possibility of intentional misuse of AI systems and the possible dangers arising out of deliberately unethical design, implementation, and use of autonomous robots.

Keywords: Robot ethics, Machine ethics, Artificial morality, Autonomous systems, Verification, Transparency, Unethical AI

This article is the result of a series of discussions taken within the scope of the Dagstuhl Seminar 6222 Fisher *et al.* [2016]. Dennis, Fisher and Winfield wish to thank EPSRC for their support, particularly via the “Verifiable Autonomy” research project (EP/L024845 and EP/L024861) Corresponding Author: Marija Slavkovic, University of Bergen, P.O.Box 7802, 5020 Bergen, Norway, marija.slavkovic@uib.no

1 Introduction

The so called “trolley problem” is a thought experiment introduced in Foot [1967], whose ethical conundrum continues to fascinate today Wallach and Allen [2008]. Roughly, it can be summarized as follows: there is a runaway trolley on a railroad track and ahead on the track there are five people tied up, unable to escape being killed by the trolley. The track splits in two by a lever you control. The lever can divert the trolley on to a second track where there is one person tied up and unable to move. Is it more ethical to divert the train, or let it run its course?

The emergence of driver-less cars in regular traffic has brought the trolley problem to public attention. Articles such as “Should Your Car Kill You To Save Others?”¹ are flooding popular science media. It is easy, given the same problem context of traffic, to get sidetracked into confusing solving the trolley problem with controlling the impact driver-less cars will have on traffic and our society as a whole. This, of course, is not the case. Enabling machines to exhibit ethical behavior is a very complex and very real time-sensitive issue. The driver-less cars are only the forefront of a whole generation of intelligent systems that can operate autonomously and will operate as part of our society. Which ethical theory to employ for deciding whose death to avoid in a difficult traffic situation is a difficult problem. This, however is not necessarily the most important problem we would need to solve. The goal of this position paper is to discuss what does *engineering machine ethics* entail. Once we as a society have discerned what is the right thing for an artificial autonomous system to do, how can we make sure the machine does it?

The problem of identifying, discerning, and recommending concepts of right and wrong is the domain of moral philosophy. Moral philosophy, together with the law, act as a system of recommendations regarding which

¹<http://www.popularmechanics.com/cars/a21492/the-self-driving-dilemma/>

possible actions are to be considered right or wrong (ignoring, for the moment, ethics systems that do not specifically address the morality of individual *actions*, e.g. character-based ethics, and which are less useful for the problem at hand).

Driver-less vehicles are the most visible of a whole range of technologies. This range includes also the assisted living technologies, as well as the various embedded decision aid software and solutions. We have been using intelligent systems with varying degree of autonomy for the past fifty years: industrial robots, intelligent programming of household appliances, automated trains, etc. However, what all of these machines have in common is that they either operate in a segregated space, a so called *working envelope*, or they have no capabilities to cause damage to their environment. Driver-less vehicles are obviously going to “break” both these restrictions. How can we build intelligent autonomous systems that uphold the ethical values of the society in which they are embedded? This is the main concern of *machine ethics*, a new interdisciplinary area of research within Artificial Intelligence (AI) Moor [2006]; Allen *et al.* [2006]; Wallach *et al.* [2008]; Wallach and Allen [2008]; Anderson *et al.* [2016].

It is clear that choosing the best moral theory to implement in a particular intelligent autonomous system is not a simple question. It is primarily a question for moral philosophy and opens new challenges in this field. In Section 2 we give a brief overview of these challenges.

For a long time, people and societies have been the only intelligent decision-makers. Moral philosophy has been developed with the often implicit assumption that the moral agent is a human. It is not at all clear to what extent existing moral theories extend to non-human decision-makers. Even if it is shown to be easy to replace a human agent with an artificial agent in a moral theory, and some societal decision is made concerning which ethical behavior in machines is desirable or sufficient, we are still faced with a set of problems regarding the *implementation* of ethical reasoning. These are the problems we analyze here: implementation, verification, trust, confidence and transparency and the prevention of intentionally unethical systems.

Moral theories are inherently ambiguous in recommendations of moral behavior, thus reflecting the context dependency of what constitutes a moral choice. We already know that artificial systems, when compared to people, are not good at handling ambiguity. Enabling machines to deal with context-ambiguity in decision-making is a core Artificial Intelligence problem Russell and Norvig [2015]. In Section 3 we give an overview of the most intuitive approaches to implementing ethical behavior in autonomous systems and

discuss the advantages and shortcomings of these approaches.

Human societies have a multitude of means for ensuring its members behave within the socially accepted boundaries of morality. We can say that a person behaves ethically because they have an individual and personal motivation to do so, without going into how this motivation is formed. The motivation for an artificial agent to behave ethically originates not personally from the agent, but from other actors. These actors can broadly be organized into three groups: the designers of the artificial agent, its users, and the various societal regulators whose job it is to make sure that order in society is maintained. It is all of these actors that need to be reassured that their own particular motivations for the AI system to behave ethically are met. Hence, a big concern when engineering machine ethics is not only that ethical behavior is accomplished, but also that the ethical behavior can be verified. This is the issue we discuss in Section 4.

In Section 5 we focus on issues of transparency and accountability for machine ethics implementations. Since there are several actors outside of the ethical agent who are supplying the motivation for ethical behavior, the implementation of this behavior must be transparent to those actors to the extent and in a manner sufficient for their needs. Transparency is a key element in enabling society to have the right amount of trust and confidence in the operations of an AI system.

Lastly in Section 6 we discuss the possible dangers for society that lie in engineering machine ethics. Like all technology, AI systems can also be abused to further criminal activities. AI systems can be deliberately built to behave unethically and illegally, but they also can be induced, deliberately or by accident, into exhibiting socially undesirable behavior.

The main contribution of this position paper is an integral overview of the immediate challenges and open questions faced when pursuing the problem of engineering of machine ethics. This paper is a result of a week long discussion among experts of different fields within the scope of the Dagstuhl Seminar 16222², and incorporates various ideas that arose as a result of discussions among interdisciplinary experts. Position papers that focus on machine ethics as a whole have been produced and they offer interesting insights in the problem as a whole, see for example Moor [2006]; Anderson and Leigh Anderson [2007]; Bonnefon *et al.* [2016]; Anderson *et al.* [2016], but to the best of our knowledge, this is the only document devoted specifically to the problem of engineering machine ethics.

²The report on this seminar is available Fisher *et al.* [2016].

2 Challenges for Machine Moral Philosophy

Normative ethics is the branch of moral philosophy concerned with developing methods for deciding how one ought to act. The purpose of devising a moral theory, within normative ethics, is to develop a consistent and coherent system that can be followed to unambiguously identify an action as good or bad. Numerous specific theories have been developed, the most notable of which are perhaps utilitarianism Harsanyi [1977], Kantianism Ellington [1993], and Ross’s ethical theory Ross [1930].

All of the moral theories so far developed in philosophy have been built around several underlying assumptions. First and foremost is the assumption that the reasoning and decision making when following the moral theory is done by an agent that is a human. The immediate specific question we can pose to moral philosophy is: given a specific normative moral theory, how is the theory affected if an artificial agent replaces a human agent in it?

A more subtle assumption in normative ethics is the assumption that the agent has *de facto* the ability of being a moral agent. This means that the theory is developed for a being that is capable of understanding concepts of right and wrong. The debate of whether an artificially created entity can be a moral agent is still far from settled Etzioni and Etzioni [2017]. Wallach and Allen [Wallach and Allen, 2008, Chapter 2] hint at the idea that how *ethically sensitive* an artificial system can be depends on how able of autonomous action it is. At this point in the development of the field of machine ethics it is fair to summarise that the capacity for moral agency of an artificial system is more of a sliding scale rather than a Boolean value. There are at least two questions to moral philosophy we can pose. The first is: how to define a scale of moral agency to describe the moral abilities of current and future artificial agents. The second question again is about going back to specific normative moral theories and checking how the theory is affected by replacing a full moral agent with an agent that is “lower” on the newly devised moral agency scale. Perhaps the concept of moral agency is altogether inadequate when discussing autonomous intelligent systems and a new concept needs to be devised.

There is a further weaker assumption in normative ethics, which is the assumption that the agents of the theory are able to accept and act upon *considered judgements* Elgin [1996]. Considered judgments are “common sense agreements” about what is good or bad in particular instances or cases. For example, the idea that murder is bad except in exceptional circumstances, is a considered judgment. This, let us call it *considered judgment ability assumption*, is reflected in the numerous moral dilemmas encountered in

the normative ethics literature, the trolley problem being the most popular example.

A dilemma in normative ethics, understood in a very broad sense, is a problem of choosing between two options each of which violating one or more considered judgment. An artificial agent does not possess a common sense understanding unless one is programmed into him. The engineering of common sense, or rather background knowledge, has been a notoriously elusive problem in artificial intelligence. The dilemmas we encounter in machine ethics reflect this. Consider for example the so called Cake or Death problem introduced in Armstrong [2015] which describes a situation in which an agent is unsure if killing people is ethical or baking them a nice cake is ethical.

Solving Cake or Death is not an ethical problem in moral philosophy, it is trivial. However, the question of which are the essential considered judgments that necessarily have to be implemented in machines is not a trivial problem. This would depend on the nature and abilities of a specific class of artificial agents. In contrast to focussing on a general theory of discerning right from wrong, there is a need for normative ethics to identify and develop a minimal such theory.

Building moral machines by implementing human morality is a natural approach. After all, it is human society that those machine are entering and it is human sensibilities and values that they have to uphold. An alternative, or perhaps parallel approach would be to build a normative ethics theory exclusively for artificial agent. An example of such a theory comes to us from science fiction - the three laws of robotics of Prof. Isaac Asimov [1950]. The shortcomings of Asimov's laws of robotics have been extensively argued by the author himself, but they have also been given a serious philosophical consideration Leigh Anderson [2008] and attempts have been made for their implementation Dennis *et al.* [2015]; Vanderelst and Winfield [2017]; Caycedo Alvarez *et al.* [2017].

The development of machine moral theory is an interesting open area for study in normative ethics. The question that has to be addressed first perhaps is what constitutes a viable, desirable, or good theory for artificial agents? We are perhaps primed by the cultural influence of Asimov's laws of robotics to ask: what is it that a robot should never do? But thinking in absolutes is not likely to be viable for machine moral theories any more than it is viable for human normative ethics. Regardless of what comes out of normative ethics in the future, any moral theory developed for artificial agents must be developed to the point of being *implementable*. Prescriptions of good behaviour suffice for people, for machines we need algorithms.

An algorithm necessarily includes a specification of all possible scenarios and context in which an ethical decision can be made. Therefore it is necessary that machine normative ethicists, computer scientists and engineers collaborate closely towards developing machine moral philosophy, with this collaboration perhaps being one of the challenges as well.

3 Different Approaches, their Advantages and Challenges

An intelligent system is one that is capable of communicating with, and reasoning about, its environment and other systems. An autonomous system is one that is capable of, to a certain extent, unsupervised operation and decision-making. Wallach, Allen, and Smit Wallach *et al.* [2008]; Wallach and Allen [2008] argue that very intuitively, ethical behavior in machines can be accomplished in at least two different ways. The first approach is to identify a set of ethical rules, perhaps by choosing a normative ethic theory, around which a decision-making algorithm can be implemented. The second approach is to have a machine evolve or “learn” to discern right from wrong without it having be explicitly guided by any one ethic theory. They refer to these two approaches as the *top-down* and *bottom-up* approach respectively. A hybrid approach in their sense is one in which an agent starts with a set of rules or values and modifies them into a system for discerning right from wrong.

Artificial Intelligence (AI) has grown into a large field that incorporates many approaches, which can be, very tentatively, classified into *soft computing approaches*, which include statistical methods, machine learning and probabilistic reasoning, and *traditional symbolic AI methods*, which includes logic-based reasoning Russell and Norvig [2015]. The question of how to implement machine ethics in an intelligent autonomous system necessarily hinges on the AI methods that system uses. Different AI approaches would be subject to different machine ethics implementations and we need to consider their malleability to machine ethic approaches, as well as their risks and advantages in this respect.

We here roughly classify the current and future machine ethics implementations based on the main AI approach used into *soft machine ethics* and *symbolic machine ethics* mirroring the two largest traditional branches of AI methodology. We discuss both of these approaches and their advantages and challenges. We should note at this point that a hybrid approach here would be one that combines symbolic methods, such as for example rule

based reasoning, with soft methods such as for example supervised learning.

The bottom-up approach of Wallach *et al.* [2008]; Wallach and Allen [2008] naturally lends itself to be approached by using soft computing AI methods, whereas their top-down approach is perhaps best “served” by symbolic AI methods.

3.1 Soft and Symbolic AI Methods for machine ethics

Within engineering, a top-down approach towards solving a task consists in breaking down the task iteratively into smaller sub-tasks until one obtains tasks that can be directly implemented. Problems best solvable by a top-down approach are ones in which the problem, and its context, are fully understood and can be formally specified. This is normally the case when the problem occurs in a controlled environment. These problem properties are also ones required for a successful solution by implementation using symbolic AI methods such as rule based reasoning.

There are numerous ways in which a symbolic AI approach can be taken to develop ethical behaviour in a system. The most frequent in the literature is to constrain the choices of the system using rules derived from an ethical theory. This is the approach taken in Arkin *et al.* [2012], for developing the concept of *ethical governor*, and also in Dennis *et al.* [2016b] where the ethical theory used is a version of Ross’s ethical theory Ross [1930]. The Dennis *et al.* [2016b] work considers a hybrid autonomous system three-layer architecture: a continuous system controlled by a rational software agent, which makes discrete decisions, via a continuous control layer that allows for a dynamic environment interaction and feedback. The rational software agent is provided with an ethical policy, a total order over abstract ethical principles such as “do no harm”, “do not damage property” etc. The agent relies on external entities to identify if, and which, of her possible actions impinges on some of these abstract ethical principles. Having her actions annotated, the agent chooses between possible actions by selecting the one that is minimally unethical with respect to the given ethical policy. In contrast, Bendel [2016] proposes a method for building *annotated decision trees* for making simple ethical choices.

A bottom-up approach to problem solving in engineering starts with describing instances of desired solutions by using adequate parameters and the proceeding to build up a procedure for identifying solutions based on these parameters. Machine learning methods in AI take a bottom-up approach to problem solving. There are several examples of using machine learning to implement machine ethics, such as for example Anderson and

Anderson [2014] and Abel *et al.* [2016]. In Anderson and Anderson [2014] inductive logic programming is used over a corpus of particular cases of ethical dilemma to discover ethical preference principles. Each case relates two actions, one more ethical than the other. The preference between the actions depends on ethically relevant features that actions involve such as harm, benefit, respect for autonomy, etc. Each feature is represented as an integer that specifies the degree of its presence (positive value) or absence (negative value) in a given action. The system is able to extract an ethical rule from the cases it is presented with and thus to a certain extent is able to learn to discern right from wrong. In Abel *et al.* [2016] reinforcement learning is used to learn what the most moral of two actions is, by rewarding the “correct” decisions an agent makes and “punishing” the bad “wrong” ones.

Whether a soft or symbolic AI method is used depends on the nature of the problem that needs to be solved. The two families of methods tackle problem from different sides and are not mutually exclusive. Each of the methods comes with its own advantages and challenges with respect to building ethical behaviour in an intelligent autonomous system.

3.2 Advantages and challenges of using Symbolic AI Methods

Symbolic AI methods are best suited for ethical reasoning in limited domains when the context of the decision-making problems can be predicted. To use symbolic based reasoning, an ethical theory needs to be chosen before constructing and deploying the system and this theory does not change throughout the system’s usage. This allows for a thorough and well-informed process of decision-making and the verification of the system prior its practical application. Different ethical principles and theories may explicitly be encoded into the system giving clear options to decide upon. Any parameters that are left open for definition by the customer or to be learned from interaction with the environment have a clear function and it is possible to verify that they do not violate more general ethical principles.

Symbolic AI methods also come with their set of challenges and limitations. General ethical guidelines are typically formulated on a very abstract level. Much philosophical discourse on ethics is concerned with problems occurring when applying such general guidelines to concrete situations. It is thus not clear if, and how, general ethical guidelines can be leveraged to solve concrete problems of decision making in practice. For a real-world system the connection to the non-discrete sensory-motor level must be made. There are many ways to transform continuous sensor values into discrete

propositions that can be used in reasoning. General guidelines or even single terms and concepts are only (if at all) implementable in a reduced way, i.e. restricted to one preferably clear interpretation. Due to their context sensitive definition it is not possible to consider every possible interpretation of an abstract guideline or term in implementing them in an artificial system Matthias [2011].

Furthermore, symbolic AI approaches risk conflicts between the implemented ethical theories and principles. If only one theory is implemented, e.g. Kant's Categorical Imperative Ellington [1993] or Isaac Asimov's first Law of Robotics³, then this theory would determine the specific maxims that are to be defined situationally by the artificial system. Winfield Winfield *et al.* [2014] describes experimental trials of a minimally ethical robot which implements Asimov's three laws of robotics. The chosen theory must be such as to allow implementable rules to be derived from it. Such a monistic approach assumes that there exist no moral dilemmas, i.e. that the implemented theory is able to give a conflict-free rule to make a decision in every context (Winfield *et al.* [2014] experimentally shows how a single ethical rule performs when faced with a balanced ethical dilemma).

Deciding on a specific set of ethical principles involves settling long-standing philosophical disputes in an ad-hoc way. It is possible that governmental bodies might take the lead in outlining high-level ethical guidance to designers. For example, Germany's ministry of transport recently announced the intention to set out a basic ethical policy to be followed by car designers stipulating that property damage takes always precedence over personal injury, that there must be no classification of people, for example, on size, age and the like, and that – ultimately – it is the manufacturer who is liable⁴.

Symbolic AI approaches require that any kind of “common sense” is explicitly coded using the formal language of the method. This severely impacts the scalability of the machine ethics solution. As it is well understood in the AI sub-discipline of knowledge representation and reasoning, the more expressive the formal language for encoding knowledge is, the more computationally expensive reasoning becomes. It is not a problem of having or symbolically encoding a large amount of information, but a problem of computing logical entailment or consistency, which are in the core of deep reasoning and are known to be of non-deterministic computational com-

³The laws can be found quoted in Wikipedia, at http://en.wikipedia.org/wiki/Three_Laws_of_Robotics

⁴<http://www.wiwo.de/politik/europa/selbstfahrende-autos-dobrindt-gruendet-ethikkommission-fuer-automatisiertes-fahren/14513384.html>

plexity. Therefore, it is unsurprising that the most recent major AI breakthroughs have been achieved using statistical processing of information and shallow reasoning. That being said, some symbolic-methods, such as model checking, are scalable within reasonable parameters, and have been vastly deployed in the information processing industry.

3.3 Advantages and challenges of using Soft AI Methods

Soft AI methods are best applicable when we know the kind of data the AI system receives from interacting with the environment and, while the overall objective might not be well known or specified, we still have an idea how to process these data in a useful manner. For instance, the vision pipeline of a household robot is designed to extract obstacles (walls, tables, etc.), objects of interest (books on a shelf etc.), and its own position from the sensory data because this information is useful for a wide range of tasks it will be required to perform. Most real-world AI systems will be partly designed using a soft AI method at least on the lower sensory-motor level. Because soft AI methods are based on synthesis of actions and choices, with respect to the task of building an ethical AI the major question here is: how do components designed in a bottom-up fashion affect the overall ethical properties of the system?

Soft AI methods do not require predetermines moral principles , ethical theories or sets of rules, but instead formulate basal parameters and intend to implement competences whereby an artificial system acts autonomously. This can be done, for example, via trial and error or other modes of learning such as imitation, induction and deduction, exploration, learning through reward, association and conditioning Cangelosi and Schlesinger [2014]. Soft AI methods can be separated into models of evolution Froese and Di Paolo [2010] and models of human socialization Fong *et al.* [2003]; Breazeal and Scassellati [2002]. The former simulate evolutionary moral learning, by assessing slightly different programs in an artificial system to evaluate an ethical case. Those programs that can solve the ethical task sufficiently go through to a “next round” where they are (re)combined to solve further ethical tasks. Evolutionary approaches can be used in earlier stages of moral development before considering models of human socialization.

Models of human socialization consider the role of empathy and emotion for moral learning. They assume that a robot learns morality via empathy Slote [2007]. What is controversial in the philosophical discourse is that there exist two types of empathy Stüeber [2006]: perceptual empathy, when an emotion triggers an equivalent or congruent reaction in the ob-

server Misselhorn [2009], and imaginative empathy that requires a change in perspective in the form of empathising with the other, putting oneself in the observed other's shoes. Perceptual empathy is explicable with the help of specific *theories of mind* or neuronal resonance and mirror neurons and has been implemented in a rudimentary fashion in artificial systems Balconi and Bortolotti [2012]; Rizzolatti and Fabbri-Destro [2008]; Mataric [2000]. Ekman Ekman [1992] implements perceptual empathy in the form of a basal affect program as an autonomous reaction scheme as a route to the implementation of morality in robots. Young children and chimpanzees are equipped with this fundamental form of empathy which forms the basis for pre-social behavior Warneken and Tomasello [2009]; Hoffman [2001]. Imaginative empathy is much more complex and develops on the basis of perceptual empathy only. It is exhibited only in human socialisation, not in non-human primates. This form of empathy is cognitively more ambitious and is involved in more complex moral reasoning and acting Gallagher [2012]. We are not aware of any attempt to implement imaginative empathy in artificial systems.

Since, by means of a soft AI solution the AI system becomes a moral agent (if only in the narrowest sense of the word) one might ask whether it is likely to be more adaptable to making ethical choices in situations that are not pre-determined (which is a strong limitation to using the symbolic AI methods). Since the system learns its own ethical rules, it circumvents, to an extent (one could argue), the need to choose one particular ethical theory to implement. But this seems at least questionable. Every self-learning system must still be configured to pay attention to particular *features* of the data set, and to ignore others. Looking at the *consequences* of an action, instead of the agent's *motivation* (for example) is such a choice of features that essentially determines the choice of moral theory. It seems difficult to judge at this point whether we can hope to create ethical-theory-agnostic AI systems, since every choice of relevant data features is already, to some extent, a choice of moral theory.

A major challenge with using soft AI methods is that it is hard to certify whether the system fulfils any requirements one might want to impose. Indeed this is a challenge for all machine learning systems. A machine learning solution virtually behaves as a black box - the approach solves a problem successfully most of the time, but it is unclear whether a solution can be expected for sure, or why a particular solution was learned or developed for a particular problem. Soft AI methods, and machine learning in particular, have had a dramatic success recently, with machine learning methods being used in a variety of problems and contexts. This success has prompted

for calls to ensure that some level of explainability for the choices of the system is required, which in turn have given rise for the Explainable AI (XAI) DARPA programme⁵. In Anderson and Leigh Anderson [2015] we find one of the earliest specific implementations of XAI in machine ethics. Their system extracts a moral rule from a collection of cases and is able to explain why a particular decision is identified as more ethical referring to the learning data.

Nonetheless, the black box nature of soft AI methods is likely to mean that these solutions are unsuitable for implementation in critical systems. This fundamental problem occurs irrespective of whether the ethical system itself or only low-level sub-systems are built using a soft AI solution.

3.4 Modular and Hybrid approaches

Both the soft and symbolic AI methods come with advantages and challenges, but they also can complement each-other. A system is an entity comprised of several entities, thus in principle an AI system can be built using components that exploit both solution approaches. We are unaware of any implemented hybrid ethical reasoning system⁶, but we can very briefly discuss some recommendations for how such a system can be created.

One approach would be to separate decision-making by the ethical principles it involves. For example, decisions involving the possibility of human death should be made using a pre-programmed ethical policy, while decisions involving violation of autonomy can be based on ethical preferences learned through interaction with the system’s owner. Another approach would be to separate decision-making in different contexts, with soft AI methods being allowed as the default ethical decision-making method, while symbolic AI approaches being required to be implemented for certain specific pre-determined contexts. Alternatively a system can be designed so it first learns to recognize the ethical implications of its actions and then those implications can be used to follow an implemented ethical theory when choosing some specific course of action.

Implementing ethical reasoning within a system is not sufficient, we must execute such implementation in a way that allows for verification of the quality of the resulting ethical behavior. The designers and manufacturers

⁵<https://www.darpa.mil/program/explainable-artificial-intelligence>

⁶ Although the title of van Rysewyk and Pontier [2015] claims a hybrid implemented system upon closer inspection it is not clear in which sense is the solution not a pure bottom-up approach in the sense of Wallach and Allen [2008], while the authors themselves do not offer an analysis of this type in their paper.

of AI systems necessarily have to offer reasons for their users to trust the artificial ethical system, and they also need to foresee possible malfunctions and provide means to deal with them.

4 Specification and Verification of Ethical Behavior

Within our society, entities that are in a position to do us harm, be it a complex machine production tool, the surgeon operating on our unconscious body, the other drivers on the highway, or a chainsaw, are subject to licensing and certification. Certification informs consumers and experts of the properties of a product, a system, or a person in a position of responsibility. Knowing that a standard has been met allows individuals to have confidence in using machinery and to trust the decisions and actions of professionals. Tools and systems are certified to operate within designated parameters, while under (well defined) proper care. Certification confirms that the manufacturer has taken all steps necessary to avoid or minimize foreseeable risks that arise in relation to the usage of the tool. Certification for persons in position of responsibility is more complex because it involves a (possibly continuous) examination to demonstrate that the certified person has the understanding and skills necessary to perform his/her duties. Typically, this involves regulations prescribing *expected* behavior — often, humans must pass an examination concerning these regulations. Once we move to an autonomous system, with no human directly in control, what are our means to ensure that a systems actually matches the relevant criteria?

In order to be confident in a system’s behavior we need to *specify* what we can expect the system to do in a particular circumstance, *verify* that the system does actually achieve this, and *validate* that our requirements are actually what the end-users want. There exist a vast range of different techniques, for example developed over many years within the field of *Software Engineering* Sommerville [2001]. These techniques range from the *formal*, such as proof, through *structured*, such as testing, to *informal*, such as user validation. All these approaches can, in principle, be applied across the range of autonomous systems, including robotics Fisher *et al.* [2013].

4.1 Who is the confirmation of ethical behavior for?

What constitutes an appropriate specification and verification methodology for ethical behavior depends on who is to use the results. In the case

of intelligent autonomous systems at least three interested parties can be discerned: the designers including developers and engineers working on developing and maintaining the systems, the end-users, owners or customers, and lastly various government and trade regulatory bodies and insurance agents. Although these three categories are the evident interested parties, this issue of interest discernment is an open problem in its own right, and as some preliminary investigations show⁷ finer discernment may be required.

Although those actually constructing the AI system may have an intimate knowledge of its internal workings, it is still important that developers and engineers not only have confidence in their prototypes but have techniques for highlighting where issues still remain. The technology itself should not be a black box, but should be open to maintenance and analysis, and must be flexible enough to be improved dynamically.

For end-users, customers and owners, the primary concern is that the AI system they interact with is safe and behaves ethically with respect to the ethical norms they themselves follow, as long as these are within the scope of what is considered ethical and legal within their society. *Trust* is a key issue and, in order to have trust extended to AI systems, the user needs to be informed of its range of capabilities. The future of AI systems and their proper integration within our society is subject, paradoxically, to undue levels of both optimism and pessimism in terms of the extent to which people can trust such systems. Close attention must be paid to nurturing the appropriate level of trust.

AI systems are an exciting technological development that have long been anticipated as part of the future in various works of fiction and there is the temptation to play-up their apparent capabilities, particularly by early marketing when the producers are still seeking financiers for their products. This could lead to the customers placing an unwarranted level of trust in some technology, even when adequate disclaimers and use guidelines are outlined by the manufacturer, which can in turn lead to disastrous consequences⁸. Such misplacement of trust is dangerous for users in the present, and may cause society to over-react in order to limit integration of technologies which given proper time to adequately develop would have been advantageous to the same society.

The appearance of trustworthiness is similarly an issue when people interact with an AI system. For example, a robot might *appear* “experienced,”

⁷<http://robohub.org/should-a-carebot-bring-an-alcoholic-a-drink-poll-says-it-depends-on-who-owns-the-robot/>

⁸<https://www.theguardian.com/technology/2016/jun/30/tesla-autopilot-death-self-driving-car-elon-musk>

“benevolent,” or “sympathetic”. Such appearances are of particular concern for AI systems that are integrated in assisted living technologies. Concerns have been raised with respect to the impact assisted living technologies can have on the elderly Sharkey and Sharkey [2012]. Similarly, it has been shown that children who interact with robots derive expectations of them and ascribe abilities to them. We need to develop an understanding of the potential long-term effects of robots on child development Matthias [2015].

Trust should play a considerable role in choosing an ethical theory to implement in AI systems. The ethical theory that is easiest to implement may not necessarily be the one that is most trusted by society. This was demonstrated in the case of utilitarianism and driver-less cars Bonnefon *et al.* [2016].

It is important to note that *trust* Arkin *et al.* [2012] is not equal to ethics. Trust is a social construct intimately concerned with how each individual views the behavior of a robot or system. There may well be some varieties of *objective* trustworthiness, but there will remain many varieties of *subjective* trustworthiness. Many items affect users’ level of trust Salem *et al.* [2015], for example, the relationship between *trust* and *harm*. If you could show that robot causes no harm, would you trust it more?

Those who must regulate AI systems and their integration within society also need *confidence* in the system. In addition, the Insurance industry needs to be clear where responsibility Sombetzki [2015] lies and so where *liability* lies. The concept of liability is likely to be complex and may be split over several actors, such as the manufacturers/designers, the operators and the environment. *Regulation* is crucial and first steps have been taken to go beyond safety and reliability regulations Bryson and Winfield [2017]; International Organization for Standardization (ISO) [2014] into considering the ethical aspects that should be taken into account British Standards Institution (BSI) [2016].

Any system operating in the real world would eventually find itself in a situation in which it will malfunction and AI systems are no exception. The question is thus how certain can one be in the verified ethical behavior of an AI system and what measures can be taken to mitigate the consequences of, and learn from, a system’s potential failure. This is the issue of having confidence in the system.

In terms of safety standards, the “gold standard” is currently that of aircraft autopilot software where safety is measured as the number of accidents per miles flown. We might think that, for AI systems, at least as much confidence is needed. But, is the standard too high or even achievable? There are of course, noticeable differences between aircraft and other autonomous

systems. While the operational environment of an aircraft is very controlled and limited, as the aircraft must adhere to a strictly defined flying corridor, the severity of accidents is very high. *E.g.*, any malfunction in the air is certainly fatal for all of the aircraft passengers measuring in the hundreds, whereas a miscalculation on the road does not need to be fatal since cars carry fewer passengers than airplanes. It is possible that the hours of operation per accident alone is not always the best measure to assess safety of an autonomous system, but that the severity of the damage caused and the number of individuals involved in an accident should also be taken into account Patchett *et al.* [2015]; Kelly and McDermid [2001]; Denney and Pai [2014]; Webster *et al.* [2014].

4.2 What do we want the system to do?

A key problem is specifying *what* our expectations of an AI system are. Although this is beginning to be codified where safety is considered, for example through *robot safety* standards⁹, it is less clear where the ethical/moral requirements should come from and in what form should they be represented? The BS8611 standard British Standards Institution (BSI) [2016], for example, does not prescribe *what* the ethical requirements should be, but maps out the issues over which ethical decisions should be considered.

An obvious route for ethical and legal requirements is through regulatory or standards bodies. These entities have the ability to set overall standards, potentially with the help of domain experts. In addition, designers may well have built in specific ethical codes that go beyond (though do not contradict) those prescribed by regulations. Finally, the user herself may wish to input her ethical preferences, ensuring that the AI acts in a way that is personally acceptable. Since there are multiple actors that need to define and refine the ethical requirements of the system, each with varying levels of technical expertise, the issue arises of how the ethical requirements are represented for the machine and for concerned actors. No one clear methodology emerges. One possibility is to have them represented in the form of a set of legal or formal rules, as argued in Saptawijaya and Moniz Pereira [2016]. Another possibility is to use a set of example scenarios developed to test specific ethical choices, as in Anderson and Anderson [2014]. A third, but by no means final, possibility is as a statistical envelope around a large (and possibly random) set of test cases, against which the AI system must be exhaustively assessed.

⁹See International Organization for Standardization (ISO) [2016] for a range of robotic safety standards.

4.3 How do we show that the AI systems meets the expectations?

There is a well-established body of work tackling the *Verification and Validation* (V&V) of systems, both hardware-centred and software-rich. The aim of *Verification* is to ensure that a system meets its requirements; *Formal Verification* takes this further, not only having precise formal requirements, but carrying out a comprehensive mathematical analysis of the system to ‘prove’ whether it corresponds to these formal requirements. There are many varieties of formal verification, the most popular being *model checking* Clarke *et al.* [1999]; Armstrong *et al.* [2012], whereby formal requirements are checked (usually automatically) against *all* possible executions of the system. Verification, via model checking, is widely used especially for the analysis of the *safety* and *reliability* of robotic systems both in terms of physical navigation Mitsch *et al.* [2013] and in terms of internal decision-making Dennis *et al.* [2016a]. What is being verified is that the behaviour of a particular system conforms to defined expectations. In terms of ethical/moral verification, it seems clear that if an AI system acts by following mathematically specified rules, we can potentially formally verify its high-level behavior. Only recently, however, has the use of formal verification for *ethical* or *moral* issues begun to be addressed Dennis *et al.* [2016b, 2015].

A practical alternative to fully formal verification is to use sophisticated *coverage-driven analysis* methods, appealing to Monte-Carlo techniques and dynamic test refinement in order to systematically “cover” a wide range of practical situations. Especially where real-world interactions and devices are involved, testing is likely to be crucial. Indeed, testing for safety and reliability of robotic systems is well-established Mossige *et al.* [2015]. Such model-based testing is a well-developed technology but, as we move to more complex (ethical) issues sophisticated extensions may well be required. Though such approaches are typically used *before* deployment, related techniques provide a basis for run-time verification and compliance testing Rosu and Havelund [2005]. Testing is not as exhaustive as formal proof, but can cover many more scenarios.

Validation is the process of confirming that the final system has the intended behavior once it is active in its target environment, and is often concerned with satisfying external stakeholders. For example, does our system match ethical standards or legal rules set by regulators? Does our system perform acceptably from a customer point of view, and how well do users feel that it works Lehmann *et al.* [2013]? There are many approaches to carrying out validation, typically involving the assessment of accuracy,

repeatability, trust, usability, resilience, etc. All must be extended to cope with ethical and moral concerns.

It is clear that the strength and breadth of V&V research should allow us to extend and develop this towards ethical and moral concerns. However, a number of issues remain, as follows.

- If the core software is not purely rule-based, for example involving some sub-symbolic learning procedures, then we will need a symbolic representation of the learned content if we are to carry out formal verification of the above form. One of the limitations of both formal verification and testing is likely to be in verifying learning procedures, especially where new ethical principles and preferences of behavior are learned.
- Fully formal verification is likely to be unrealistic for complete, complex systems both because of non-symbolic components (as mentioned above) and because of practical complexity limits.

However, we can formally verify *parts* of the system under particular circumstances. There are things that can be proved about core parts of *the system* and about the system's *outputs*. Consequently, formal verification techniques can provide *some* evidence. In assessing how much confidence we need in the V&V of AI system ethics, it may be possible to leave the burden of this decision to the regulator, manufacturer or end-user as appropriate. So long as a clear indication of the extent of the V&V of a system exists a user or other interested may take the decision about the risk involved in using the system. Note that we can potentially separate regulation from verification and so allow a variety of different V&V techniques to be applied.

Lastly, we would like to include here the existing efforts of validating a system that uses soft AI methods, which is the discussion of *Ethical Turing Tests*. Ethical Turing Tests were introduced in Allen *et al.* [2000]. In Anderson and Anderson [2014] this idea is further fleshed out and implemented. Under an ethical Turing test, both the AI system and an ethicist resolve the same dilemmas. The system passes the test if its choices are sufficiently similar to the ones of the ethicist. Whether a variant of a Turing test is a sufficient indicator of a certain type of “human-like” behaviour from a machine is a topic that has been argued as long as the artificial intelligence field exists. All the issues that have been raised, and exhaustively discussed in artificial intelligence, against the Turing test original can be argued to hold for an Ethical Turing test.

5 Transparency and accountability

The opacity and transparency of deep neural network algorithms has become a major research subject area in recent years. Without knowing how the algorithm functions concerns about inherent biases in data and algorithm itself create difficulties in either discerning forensically or through explainability how decisions were made. Furthermore, there are questions as to whether the actual logic used by learning systems can be explained to people, or whether any explanation would be created after the fact.

The problem with interpretability/explainability is that in some cases it may be impossible to provide a complete and accurate explanation for how a black box system has arrived at its decision. As complexity of systems increases, it is not unusual for the algorithm to extract a million-dimensional feature vector and assign unique weights to each feature. Any human-readable explanation for the decision will include some top N most important features and completely ignore the rest. A human-comprehensible explanation cant be too long or too complex. A good metaphor for this is how we explain things to children if they are not old enough to fully appreciate nuances of the problem. Where do kids come from? You buy them at the store! Any human-friendly explanation from a sufficiently complex system has to be a partially inaccurate simplification or a complete lie.

The choice of the relevant criteria for an AI system to be deemed ethical will eventually need to be taken by society as a whole. Therefore *transparency* is of utmost importance and thus ensuring transparency is a major challenge. To this end it is necessary to identify *what* has to be transparent to *whom*, and *how* this can be realized.

Transparency is a key requirement for ethical machines. Important attributes flow from transparency including *trust*, because it is hard to trust a machine unless you have some understanding of what it is doing and why, and *accountability*, because without transparency it becomes very difficult to understand who is responsible when a machine does not behave as we expect it to¹⁰. An ethical machine will need to be transparent to different stakeholders in different ways – each suited to that particular stakeholder. In this section we consider the transparency needs of a range of stakeholders before considering aspects of transparency common to all. This section outlines how and why transparency is important to four different groups of stakeholders: users, regulators (including accident investigators), ethicists/lawyers and society at large. Each group has different transparency

¹⁰Although it is important to note that transparency is not the same as accountability.

needs, some of which will have to be met by allowing an AI system’s ethics, and ethical logic, to be human readable, or through public engagement. Other needs will require new human-robot interfaces.

Some literature exists on the topic of transparency in AI and autonomous systems. Owotoki and Mayer-Lindenberg Owotoki and F. [2007] proposes a theoretical framework for providing transparency in computational intelligence (CI) in order to expose the underlying reasoning process of an agent embodying CI models. In a recent book Taylor and Kelsey Taylor and Kelsey [2016] make the case for the importance of transparency in AI systems to an open society. For autonomous robots Wortham *et al.* [2016] describes early results showing that building transparency into robot action-selection can help users build a more accurate understanding of the robot. There is also no doubt that transparency is high on policy agenda: the 2016 UK Parliamentary Select Committee on Science and Technology’s final report on Robotics and AI expresses concerns over both decision making transparency and accountability and liability¹¹. Indeed the EU’s new General Data Protection Regulation, due to take effect as law in 2018, creates a “right to explanation” such that a user will be able to ask for an explanation of an algorithmic decision that was made about them Goodman and Flaxman [2016].

5.1 Transparency to the user

Although the critical importance of the human-machine interface is well understood, what is not yet clear is the extent to which an ethical machine’s ethics should be transparent to its user. It would seem to be unwise to rely on a user to discover a machine’s ethics by trial and error, but at the same time a machine that requires its user to undergo a laborious process of familiarisation may well be unworkable.

For care robots for instance it may be appropriate for the user to configure the “ethics” settings (perhaps expressing the user’s preference for more or less privacy) or, at the very least, allowing the user to choose between a small number of “preset” ethics options. There is of course always some danger that many users will rely on the default setting. What is clear is that how these options are presented to the user is very important Matthias [2015]; they should for instance help and guide the user in thinking about their ‘value hierarchy’. The robot might for instance explain to the user

¹¹<http://www.publications.parliament.uk/pa/cm201617/cmselect/cmsctech/145/145.pdf>

“what would happen” in different situations and hence guide their preferences Theodorou *et al.* [2016].

For other robot types, driverless cars for instance, the ethics settings may be fixed (perhaps by law) and therefore not user configurable. However, the need for the user to understand how the car would behave in certain situations remains critically important – especially if the car’s design (or the law) requires her to act as a safety driver and assume manual control when the autopilot cannot cope. Even for fully autonomous cars in which the user is only ever a passenger, the person with legal responsibility for the car should be aware of the car’s ethics settings. For fully autonomous cars there should still be some user interface so that the passenger can discover, or perhaps ask for help, if the vehicle become unexpectedly immobile or starts behaving erratically.

5.2 Transparency to regulatory bodies

It is clear that the ethics of ethical robots needs to be transparent to those responsible for (i) certifying the safety of ethical machines, and (ii) accident investigators. Both regulators and accident investigators will be working within a governance framework which includes standards and protocols. The role of the protocols is to set out how robots are certified against those standards, and – following an accident – how the accident is investigated.

Regulators will need the ability to determine that a machine’s ethics comply with the appropriate standards¹², and making such a determination will require those ethics to be coded and embedded into the robot in a readable way. We might imagine something like a standard Ethics Markup Language (EML – perhaps based on XML) which codes the ethics. The EML script would be embedded in the robot in a way that is accessible to the regulator, noting that the script will need to be secured to prevent attack from hackers.

Accident Investigators. When serious accidents happen, as they inevitably will (see Section 4.1), they will need to be investigated. To allow such investigation, data must to be recorded, suggesting the need for a robot equivalent of the flight data recorder. Therefore an **Ethical Black Box** (EBB) is proposed Winfield and Jirotko [2017]– a device that records all relevant data including, crucially, internal state data on the robot’s ethical governor. Although the data stored by the EBB would be vital for investigating all aspects of an accident, including causes unrelated to the robot’s

¹²Standards Which do not yet exist.

ethics, here we are interested in accidents which might have been caused by a fault or deficiency in the robot’s ethics programming. By recording the sequence of internal states of the ethical reasoning in the moments before the accident, the EBB would allow an investigator to discover exactly why the robot made an incorrect decision. Information that would be important both in determining accountability, and to make recommendations for upgrading the robot’s ethics and prevent the same accident happening again.

Specifying the EBB is beyond the scope of this paper, this work has been carried further in Winfield and Jirotko [2017]. It is however clear that research is needed to determine what data the EBB must record, the frequency and time window of that data, and how the privacy of that data is maintained. One thing we can be sure of however is the need for an industry standard EBB (as in the aviation industry). Different EBBs and EBB standards will of course be needed for different applications, but for driverless cars for instance, a single standard EBB should be mandated. Such an EBB would itself require an industry standard, and protocols for certification, fitting and maintenance of EBBs.

5.3 Transparency to ethicists / lawyers

A third group of stakeholders includes lawyers, who might be required to advocate for AI systems’ owners, or on behalf of anyone who makes a claim against an AI system’s owner, or ethicists who might, for instance, be required to act as expert witnesses, in a court of law. If we consider an accident in which a robot’s ethics are implicated (see Section 5.2 above), it is clear that both lawyers and ethicists will need to understand (i) a robot’s ethics, (ii) the process the robot uses to make an ethical decision (in other words how its ethical reasoning works), and (iii) the data captured by the ethical black box. Providing this kind of transparency to lawyers and ethicists will not only be necessary, but is also likely to be challenging, as robot manufacturers and designers may regard such details, especially (ii), as proprietary IP.

Another category of expert stakeholder includes psychologists, who might be required to either evaluate robots for their potential to cause psychological harm to the user, or as expert witnesses, in providing an investigation with an expert evaluation of the psychological harm caused to robot user(s) in a particular case.

5.4 Transparency to the whole of society

AI systems – and especially ethical AI systems – are a disruptive technology, with potentially significant societal and economic impact, thus an easily overlooked but important stakeholder is society as a whole. We only need to consider driverless cars and trucks to appreciate the level of potential disruption, to jobs and transport policy for instance, as already reflected in the level of public and press interest in this technology.

It is therefore very important that the ethics of ethical AI systems should be transparent to society at large, for two reasons. First, because citizens should be able to make informed judgments about the kinds of AI system they wish to have in their lives, and even more importantly those they do not want in their lives, so that they can lobby their elected representatives and ensure that government policy properly reflects those views. And second, if society is to have confidence in the ethics of a class of ethical AI system (driverless cars, for example) then it should accept a degree of collective responsibility for those ethics.

5.5 Technical means to bring about transparency

It is clear that the different stakeholders outlined above have very different transparency needs. Some of those needs are met through making the ethical rules and logic readable (for instance for regulators, ethicists or lawyers), but for others transparency can only be met through technical means. Here we briefly outline several approaches to meeting those needs.

- Assisted living AI systems would benefit from a “Why did you do that?” button which, when pressed, causes the robot to explain – perhaps using speech synthesized text – why it carried out the previous action. We could call the system behind this an “explanation module”. For an AI system with a fixed set of responses the explanation module should be relatively easy to implement, but for an AI system which learns its ethics such an implementation could be challenging; in either case the explanation module and its user interface would need very careful design in order to meet the needs of a non-technical user.
- An ethical AI system which makes use of simulation based internal models as part of some ethical governor (for example Winfield *et al.* [2014]) might allow us to go further than the “Why did you do that?” button, by making the robot’s internal simulation accessible to the user. This would enable the user to ask the robot “What would you do?” in a given situation. Clearly such a facility would need a much

more sophisticated user interface than a button press, but through visualisation tools we can imagine the user watching the robot’s internal simulation running through various scenarios on a connected laptop or tablet device. Note that a similar visualisation interface would be of great value to accident investigators (Section 5.2), and expert witnesses or lawyers (section 5.3) to *play back* a robot’s internal simulation in the moments leading up to an accident, and what the alternatives open to the robot at the time might have been.

- The technical requirements for an ethical back box (EBB) were already outlined in Section 5.2 above.

6 Dangerous and Deliberately Unethical AI

Finally, it is important to be aware of the ways people may abuse or manipulate AI systems. As with all technology AI systems can also be deliberately abused for malice or to further one’s illegal goals Yampolskiy [2016]. While our primary concern is to contribute towards designing AI systems that behave ethically within a human society Sotola and Yampolskiy [2015]; Yampolskiy [2015a] and promote human and animal welfare, some concern also needs to be raised about how that AI system can protect itself against abuse Yampolskiy and Spellchecker [2016]. By abuse we, of course, do not mean mistreating the AI system in the sense in which a person or an animal can be mistreated, but taking advantage of the capabilities and opportunities offered by the AI system to commit criminal acts.

The abuse of an AI system can be achieved by hacking an existing system or by deliberately creating an unethical AI system Pistono and Yampolskiy [2016]; Vanderelst and Winfield [2016]. Hacking itself can be accomplished in several ways. The code of the AI system might be directly hacked. But a system can also be manipulated by interaction and such manipulation does not necessarily require technical knowledge. This is illustrated by the short-lived Tay experiment. Tay was an artificial intelligence chatter-bot released by Microsoft Corporation on March 23, 2016 and taken offline 16 hours after launch¹³. Tay was programmed to learn from conversation, however it took the netizens a very short time to “train” it into making morally questionable statements.

Manipulation by interaction can be accomplished both deliberately and by accident. A learning based system can be led Yampolskiy [2014] into eliciting bad conclusions through crafted case descriptions, etc. By this

¹³<http://www.bbc.com/news/technology-35890188>

means one can slowly train systems away from moral behavior. As an example of accidental manipulation consider the example of children learning that driverless cars slow down in their presence, they might choose to make a game out of it. Children playing with car's reactions might annoy passengers by causing delay; and might ultimately lead to the disabling of safeguards.

Purposeful creation of Malevolent AI can be attempted by a number of diverse agents with varying degrees of competence and success. Each such agent would bring its own goals/resources into the equation, but what is important to understand here is just how prevalent such attempts will be and how numerous such agents can be. For example, we should be concerned about: the *military* developing cyber-weapons and robot soldiers to achieve dominance; *governments* attempting to use AI to establish hegemony, control people, or take down other governments; *corporations* trying to achieve monopoly, destroying the competition through illegal means; *villains* trying to take over the world and using AI as a dominance tool; *black hats* attempting to steal information, resources or destroy cyber infrastructure targets; *doomsday cults* attempting to bring the end of the world by any means; the *depressed* looking to commit suicide by AI; *psychopaths* trying to add their name to history books in any way possible; *criminals* attempting to develop proxy systems to avoid risk and responsibility; AI *risk deniers* attempting to demonstrate that AI is not a risk factor and so ignoring caution; and even AI *safety researchers*, if unethical, attempting to justify funding and secure jobs by purposefully developing problematic AI.

The ethical and unethical behaviors of an AI system are not necessarily symmetrical. Existing systems define only a small part of the problem space Yampolskiy [2015b]. Apart from ethical and unethical behavior, an AI system can also exhibit a behavior that has neither been programmed nor predicted as a particular combination of otherwise ethical rules and choices.

Lastly we must mention the potential for “cultural imperialism” when designing the ethical behavior of an AI system. With globalisation, a product's production and consumers are diverse. What constitutes ethical behavior in one region may even be considered unethical in another. All the involved actors, the designers, users and society, both on the supplier and on the demand end of the AI system need to be aware of the reality that the supplier society ethics influences the ethical behavior of the AI system, which in turn influences the ethics of the society in which the AI system operates.

7 Summary

Moral philosophy has a very rich history of studying how to discern right from wrong in a systematic, consistent and coherent way. Today we have a real need for a *functional* system of ethical reasoning as AI systems that function as part of our society are ready to be deployed. Building an AI system that behaves ethically is a multifaceted challenge. The questions of which ethical theory should be used to govern the AI system's behavior has received most of the attention. Here, we focus on the problem that comes next, after what is right or wrong for a machine to do is decided – how to implement the ethical behavior.

The problem of engineering ethical behavior is made complex because of the prime motivators for such behavior. For humans, the motivation for behaving ethically is primarily internal. Without falling into difficult philosophical arguments on the existence of free will, we accept that people are capable of behaving ethically because they choose to do so, although, of course, they too can be motivated towards ethical behaviour by incentives, punishments and assignment of liability. For AI systems, however, the motivation towards ethical behaviour is exclusively external because it can always be traced back to their design and it cannot be reinforced in the same way as it can be done with people. This motivation furthermore comes from several stakeholders. We cannot claim that the full list of these stakeholders can even be known before the AI systems are fully deployed, but we can discern between the three most evident groups of stakeholders: the designers, the users and the various regulatory organs of society. Each of these stakeholders needs to play their role in deciding what is the best ethical behavior for a given AI system, but they also need to be convinced in an adequate way that the implemented behavior actually yields the desired results. A moral AI system needs to be adequately transparent and accountable to each group of stakeholders.

Unlike people, who more or less share the same “hardware” and reasoning capabilities, machines and AI systems can be built using many different approaches. The implementation of ethical reasoning will depend not only on what the stakeholders need and desire, but also on what is possible given the chosen problem-solving implementation. We discussed the two basic implementation approaches reflecting two large families of AI methods: the soft AI and the symbolic AI families, and identify the challenges and advantages of each.

An AI system capable of ethical behavior is necessarily a complex system. With complex systems two things are evident: that they will malfunction

and that they can be used to attain criminal goals. We discuss methods of verifying that an AI system behaves as designed within specified parameters, but we also discuss how the engineering of the ethical behavior impacts available options once a system malfunctions. Lastly we discuss in broad strokes what the stakeholders need to be aware of in terms of abuse of an AI system with ethical behavior capabilities, both when that abuse is intentional and accidental.

References

- D. Abel, J. MacGlashan, and M. L. Littman. Reinforcement Learning As a Framework for Ethical Decision Making. In *Workshops at the Thirtieth AAAI Conference on Artificial Intelligence*, 2016.
- C. Allen, G. Varner, and J. Zinser. Prolegomena to any future artificial moral agent. *Journal of Experimental and Theoretical Artificial Intelligence*, 12:251–261, 2000.
- C. Allen, W. Wallach, and I. Smit. Why machine ethics? *IEEE Intelligent Systems*, 21(4):12–17, 2006.
- M. Anderson and S. Leigh Anderson. Geneth: A general ethical dilemma analyzer. In *Proceedings of the Twenty-Eighth AAAI Conference on Artificial Intelligence, July 27 -31, 2014, Québec City, Québec, Canada.*, pages 253–261, 2014.
- M. Anderson and S. Leigh Anderson. Machine ethics: Creating an ethical intelligent agent. *AI Magazine*, 28(4):15–26, 2007.
- M. Anderson and S. Leigh Anderson. Toward ensuring ethical behavior from autonomous systems: a case-supported principle-based paradigm. *Industrial Robot*, 42(4):324–331, 2015.
- M. Anderson, S. Leigh Anderson, and Vincent Berenz. Ensuring ethical behavior from autonomous systems. In *Artificial Intelligence Applied to Assistive Technologies and Smart Environments, Papers from the 2016 AAAI Workshop, Phoenix, Arizona, USA, February 12, 2016*, 2016.
- R.C. Arkin, P. Ulam, and A. R. Wagner. Moral Decision Making in Autonomous Systems: Enforcement, Moral Emotions, Dignity, Trust, and Deception. *Proceedings of the IEEE*, 100(3):571–589, 2012.

- P.J. Armstrong, M. Goldsmith, G. Lowe, J. Ouaknine, H. Palikareva, A. W. Roscoe, and J. Worrell. Recent Developments in FDR. In P. Madhusudan and S. A. Seshia, editors, *Computer Aided Verification: 24th International Conference, CAV 2012, Berkeley, CA, USA, July 7-13, 2012 Proceedings*, volume 7358 of *LNCS*, pages 699–704. Springer Berlin Heidelberg, 2012.
- S. Armstrong. Motivated value selection for artificial agents. In *Artificial Intelligence and Ethics, Papers from the 2015 AAAI Workshop, Austin, Texas, USA, January 25, 2015.*, 2015. <http://aaai.org/ocs/index.php/WS/AAAIW15/paper/view/10183>.
- I. Asimov. *I, Robot*. Gnome Press, 1950.
- M. Balconi and A. Bortolotti. Detection of the facial expression of emotion and self-report measures in empathic situations are influenced by sensorimotor circuit inhibition by low-frequency rtms. *Brain Stimulation*, 5(3):330 – 336, 2012.
- O. Bendel. Annotated decision trees for simple moral machines. 2016.
- J.F. Bonnefon, A. Shariff, and I. Rahwan. The social dilemma of autonomous vehicles. *Science*, 352(6293):1573–1576, 2016.
- C. Breazeal and B. Scassellati. Robots that imitate humans. *Trends in Cognitive Sciences*, 6(11):481–487, 2002.
- British Standards Institution (BSI). BS 8611 Robots and Robotic Devices — Guide to the ethical design and application. <http://www.bsigroup.com>, 2016.
- J. Bryson and A.F.T. Winfield. Standardizing ethical design for artificial intelligence and autonomous systems. *IEEE Computer*, 50(5):116–119, 2017.
- A. Cangelosi and M. Schlesinger. *Developmental Robotics: From Babies to Robots*. The MIT Press, 2014.
- M. Caycedo Alvarez, Ø. S. Berge, A. S. Berget, E. S. Bjørknes, D.V.K. Johnsen, F. O. Madsen, and M. Slavkovik. Implementing Asimov’s first law of robotics. In *30th Norsk Informatikkonferanse, NIK 2017, Westerdals Oslo ACT, November, 27-29, 2017*, 2017. forthcoming.
- E. M. Clarke, O. Grumberg, and D. Peled. *Model Checking*. MIT Press, 1999.

- E. Denney and J.P. Pai. Automating the assembly of aviation safety cases. *IEEE Trans. Reliability*, 63(4):830–849, 2014.
- L. A. Dennis, M. Fisher, and A. F. T. Winfield. Towards Verifiably Ethical Robot Behaviour. In *Proc. AAAI Workshop on AI and Ethics*, 2015. <http://aaai.org/ocs/index.php/WS/AAAIW15/paper/view/10119>.
- L. A. Dennis, M. Fisher, N. K. Lincoln, A. Lisitsa, and S. M. Veres. Practical Verification of Decision-Making in Agent-Based Autonomous Systems. *Automated Software Engineering*, 23(3):305–359, 2016.
- L. A. Dennis, M. Fisher, M. Slavkovik, and M. P. Webster. Formal Verification of Ethical Choices in Autonomous Systems. *Robotics and Autonomous Systems*, 77:1–14, 2016.
- P. Ekman. Are there basic emotions? *Psychological Review*, 99(3):550–553, 1992.
- C. Elgin. *Considered Judgment*. Princeton: New Jersey: Princeton University Press, 1996.
- J. W. Ellington. *Translation of: Grounding for the Metaphysics of Morals: with On a Supposed Right to Lie because of Philanthropic Concerns by Kant, I. [1785]*. Hackett Publishing Company, 1993.
- A. Etzioni and O. Etzioni. Incorporating ethics into artificial intelligence. *The Journal of Ethics*, pages 1–16, 2017.
- M. Fisher, L. A. Dennis, and M. Webster. Verifying Autonomous Systems. *ACM Communications*, 56(9):84–93, 2013.
- M. Fisher, C. List, M. Slavkovik, and A. F. T. Winfield. Engineering moral agents - from human morality to artificial morality (dagstuhl seminar 16222). *Dagstuhl Reports*, 6(5):114–137, 2016.
- T. Fong, I. Nourbakhsh, and K. Dautenhahn. A survey of socially interactive robots. *Robotics and Autonomous Systems*, 42(3–4):143 – 166, 2003. Socially Interactive Robots.
- P. Foot. The problem of abortion and the doctrine of double effect. *Oxford Review*, 5:5–15, 1967.
- T. Froese and E. A. Di Paolo. Modelling social interaction as perceptual crossing: An investigation into the dynamics of the interaction process. *Connection Science*, 22(1):43–68, March 2010.

- S. Gallagher. Empathy, simulation, and narrative. *Science in Context*, 25(3):355–381, 009 2012.
- B. Goodman and S. Flaxman. Eu regulations on algorithmic decision-making and a “right to explanation”. *arXiv preprint arXiv:1606.08813*, 2016.
- J.C. Harsanyi. Rule utilitarianism and decision theory. *Erkenntnis (1975-)*, 11(1):25–53, 1977.
- M. Hoffman. *Empathy and Moral Development: Implications for Caring and Justice*. Cambridge University Press, 2001.
- International Organization for Standardization (ISO). ISO 13482: Robots and robotic devices — Safety requirements for Personal Care Robots. <http://www.iso.org>, 2014.
- International Organization for Standardization (ISO). TC299 — Robotics, 2016.
- T.P. Kelly and J.A. McDermid. A systematic approach to safety case maintenance. *Rel. Eng. & Sys. Safety*, 71(3):271–284, 2001.
- H. Lehmann, D. S. Syrdal, K. Dautenhahn, G.J. Gelderblom, S. Bedaf, and F. Amirabdollahian. What Should a Robot do for you? Evaluating the Needs of the Elderly in the UK. In *Proc. 6th Int. Conf. on Advances in Computer-Human Interactions*, pages 83–88, 2013.
- S. Leigh Anderson. Asimov’s ‘three laws of robotics’ and machine metaethics. *AI & Society*, 22(4):477–493, 2008.
- M. J. Mataric. Getting humanoids to move and imitate. *IEEE Intelligent Systems*, 15(4):18–24, July 2000.
- A. Matthias. Algorithmic moral control of war robots: Philosophical questions. *Law, Innovation and Technology*, 3(2):279–301, 2011.
- A. Matthias. Robot lies in health care: when is deception morally permissible? *Kennedy Institute of Ethics Journal*, 25(2):169–162, 2015.
- C. Misselhorn. Empathy with inanimate objects and the uncanny valley. *Minds and Machines*, 19(3):345, 2009.
- S. Mitsch, K. Ghorbal, and A. Platzer. On Provably Safe Obstacle Avoidance for Autonomous Robotic Ground Vehicles. In *Robotics: Science and Systems IX*, 2013.

- J. H. Moor. The nature, importance, and difficulty of machine ethics. *IEEE Intelligent Systems*, 21(4):18–21, July 2006.
- M. Mossige, A. Gotlieb, and H. Meling. Testing Robot Controllers using Constraint Programming and Continuous Integration. *Information & Software Technology*, 57:169–185, 2015.
- P. Owotoki and Mayer-Lindenberg F. Transparency of computational intelligence models. In *Research and Development in Intelligent Systems XXIII, The 26th SGAI International Conference on Innovative Techniques and Applications of Artificial Intelligence, Proceedings*, pages 387–392. Springer, 2007.
- C. Patchett, M. Jump, and M. Fisher. Safety and Certification of Unmanned Air Systems. In *Engineering and Technology Reference*. 2015.
- F. Pistono and R. V Yampolskiy. Unethical research: How to create a malevolent artificial intelligence. In *25th International Joint Conference on Artificial Intelligence (IJCAI-16). Ethics for Artificial Intelligence Workshop (AI-Ethics-2016)*, 2016.
- G. Rizzolatti and M. Fabbri-Destro. The mirror system and its role in social cognition. *Current Opinion in Neurobiology*, 18(2):179 – 184, 2008. Cognitive neuroscience.
- W.D. Ross. *The Right and the Good*. Oxford University Press, 1930.
- G. Rosu and K. Havelund. Rewriting-Based Techniques for Runtime Verification. *Automated Software Engineering*, 12(2):151–197, 2005.
- S.J. Russell and P. Norvig. *Artificial Intelligence: A Modern Approach*. Pearson Education, 3 edition, 2015.
- M. Salem, G. Lakatos, F. Amirabdollahian, and K. Dautenhahn. Would You Trust a (Faulty) Robot?: Effects of Error, Task Type and Personality on Human-Robot Cooperation and Trust. In *Proc. 10th Annual ACM/IEEE International Conference on Human-Robot Interaction (HRI)*, pages 141–148. ACM, 2015.
- A. Saptawijaya and L. Moniz Pereira. Logic programming for modeling morality. *Logic Journal of the IGPL*, 24(4):510–525, 2016.
- A. Sharkey and N. Sharkey. Granny and the robots: ethical issues in robot care for the elderly. *Ethics and Information Technology*, 14(1):27–40, 2012.

- M. A. Slote. *The Ethics of Care and Empathy*. Routledge, 2007.
- J. Sombetzki. Responsibility in Crisis — From the Traditional Concept of Responsibility to Systems Responsibility, 2015.
- I. Sommerville. *Software Engineering*. Pearson Studium, 2001.
- K. Sotala and R. V Yampolskiy. Responses to catastrophic agi risk: A survey. *Physica Scripta*, 90(1), 2015.
- K. Stüeber. *Rediscovering Empathy: Agency, Folk Psychology, and the Human Sciences*. MIT Press, 2006.
- R. Taylor and T. Kelsey. *Transparency and the open society: Practical lessons for effective policy*. Policy Press, Bristol UK, 2016.
- A. Theodorou, R. Wortham, and J.J. Bryson. Why is my robot behaving like that? designing transparency for real time inspection of autonomous robots. In *AISB Workshop on Principles of Robotics, April 2016, Sheffield UK, Proceedings*, 2016.
- S. P. van Rysewyk and M. Pontier. *A Hybrid Bottom-Up and Top-Down Approach to Machine Medical Ethics: Theory and Data*, pages 93–110. Springer International Publishing, Cham, 2015.
- D. Vanderelst and A.F. Winfield. The dark side of ethical robots. *arXiv preprint arXiv:1606.02583*, 2016.
- D. Vanderelst and A. Winfield. An architecture for ethical robots inspired by the simulation theory of cognition. *Cognitive Systems Research*, 2017.
- W. Wallach and C. Allen. *Moral Machines: Teaching Robots Right from Wrong*. Oxford University Press, 2008.
- W. Wallach, C.n Allen, and I. Smit. Machine morality: Bottom-up and top-down approaches for modelling human moral faculties. *AI Soc.*, 22(4):565–582, March 2008.
- F. Warneken and M. Tomasello. Varieties of altruism in children and chimpanzees. *Trends in Cognitive Sciences*, 13(9):397 – 402, 2009.
- M. Webster, N. Cameron, M. Fisher, and M. Jump. Generating Certification Evidence for Autonomous Unmanned Aircraft Using Model Checking and Simulation. *Journal of Aerospace Information Systems*, 11(5):258–279, 2014.

- A.F.T. Winfield and M. Jirotko. The case for an ethical black box. In *Towards Autonomous Robotic Systems - 18th Annual Conference, TAROS 2017, Guildford, UK, July 19-21, 2017, Proceedings*, pages 262–273, 2017.
- A. F. T. Winfield, C. Blum, and W. Liu. *Towards an Ethical Robot: Internal Models, Consequences and Ethical Action Selection*, pages 85–96. Springer International Publishing, 2014.
- R.H. Wortham, A. Theodorou, and J.J. Bryson. What does the robot think? transparency as a fundamental design requirement for intelligent systems. In *IJCAI-2016 Ethics for Artificial Intelligence Workshop, July 2016, New York USA, Proceedings*, 2016.
- R. V Yampolskiy and MS Spellchecker. Artificial intelligence safety and cybersecurity: a timeline of ai failures. *arXiv preprint arXiv:1610.07997*, 2016.
- R.V. Yampolskiy. Utility function security in artificially intelligent agents. *Journal of Experimental & Theoretical Artificial Intelligence*, 26(3):373–389, 2014.
- R. V. Yampolskiy. Artificial superintelligence: A futuristic approach, 2015.
- R. V Yampolskiy. The space of possible mind designs. In *Artificial General Intelligence: 8th International Conference, AGI 2015, AGI 2015, Berlin, Germany, July 22-25, 2015, Proceedings*, volume 9205, page 218. Springer, 2015.
- R. V Yampolskiy. Taxonomy of pathways to dangerous artificial intelligence. In *Workshops at the Thirtieth AAAI Conference on Artificial Intelligence*, 2016.