

INVITED COMMENT • **OPEN ACCESS**

# Responses to catastrophic AGI risk: a survey

To cite this article: Kaj Sotala and Roman V Yampolskiy 2015 *Phys. Scr.* **90** 018001

View the [article online](#) for updates and enhancements.

## Related content

- [How feasible is the rapid development of artificial superintelligence?](#)  
Kaj Sotala
- [‘The concept of information in physics’: an interdisciplinary topical lecture](#)  
T Dittrich
- [The great downside dilemma for risky emerging technologies](#)  
Seth D Baum

## Recent citations

- [James Babcock \*et al\*](#)
- [Serap Uur and Gulsun Kurubacak](#)
- [Why We Do Not Evolve Software? Analysis of Evolutionary Algorithms](#)  
Roman V Yampolskiy



# Corrigendum: Responses to catastrophic AGI risk: a survey (2015 *Phys. Scr.* 90 018001)

Kaj Sotala<sup>1</sup> and Roman V Yampolskiy<sup>2</sup>

<sup>1</sup> Machine Intelligence Research Institute, Berkeley, CA, USA

<sup>2</sup> University of Louisville, KY, USA

Parts of the reference list are corrected to the following:

- [105] Goertzel B 2002 Thoughts on AI morality *Dynamical Psychology* ([www.goertzel.org/dynapsyc/2002/AIMorality.htm](http://www.goertzel.org/dynapsyc/2002/AIMorality.htm))
- [106] Goertzel B 2004 Encouraging a positive transcension *Dynamical Psychology* ([www.goertzel.org/dynapsyc/2004/PositiveTranscension.htm](http://www.goertzel.org/dynapsyc/2004/PositiveTranscension.htm))
- [107] Goertzel B 2004 Growth, choice and joy *Dynamical Psychology* ([www.goertzel.org/dynapsyc/2004/GrowthChoiceJoy.htm](http://www.goertzel.org/dynapsyc/2004/GrowthChoiceJoy.htm))
- [108] Goertzel B 2006 Apparent limitations on the 'AI friendliness' and related concepts imposed by the complexity of the world ([www.goertzel.org/papers/LimitationsOnFriendliness.pdf](http://www.goertzel.org/papers/LimitationsOnFriendliness.pdf))
- [109] Goertzel B 2010 Coherent aggregated volition *The Multiverse According to Ben* (<http://multiverseaccordingtoben.blogspot.ca/2010/03/coherent-aggregated-volition-toward.html>)
- [110] Goertzel B 2010 GOLEM (<http://goertzel.org/GOLEM.pdf>)
- [111] Goertzel B 2012 Should humanity build a global AI nanny to delay the singularity until it's better understood? *J. Consciousness Stud.* **19** 96–111
- [112] Goertzel B 2012 When should two minds be considered versions of one another? *Int. J. Mach. Consciousness* **4** 177–85
- [113] Goertzel B 2012 CogPrime ([http://wiki.opencog.org/w/CogPrime\\_Overview](http://wiki.opencog.org/w/CogPrime_Overview))
- [114] Goertzel B and Bugaj S V 2008 Stages of ethical development in artificial general intelligence systems *Artificial General Intelligence (Frontiers in Artificial Intelligence and Applications no. 171)* (Amsterdam: IOS) pp 448–59
- [115] Goertzel B and Pitt J 2012 Nine ways to bias open-source AGI toward friendliness *J. Evol. Technol.* **22** 116–31
- [133] Hanson R 1994 If uploads come first *Extropy* **6** 10–15
- [134] Hanson R 1998 Economic growth given machine intelligence (<http://hanson.gmu.edu/aigrow.pdf>)
- [135] Hanson R 2007 Shall we vote on values, but bet on beliefs? (<http://hanson.gmu.edu/futarchy.pdf>)
- [136] Hanson R 2008 Economics of the singularity *IEEE Spectr.* **45** 45–50
- [137] Hanson R 2009 Prefer law to values *Overcoming Bias* ([www.overcomingbias.com/2009/10/prefer-law-to-values.html](http://www.overcomingbias.com/2009/10/prefer-law-to-values.html))
- [138] Hanson R 2012 Meet the new conflict, same as the old conflict *J. Consciousness Stud.* **19** 119–25
- [146] Hibbard B 2001 Super-intelligent machines ACM SIGGRAPH *Comput. Graph.* **35** 13–5
- [147] Hibbard B 2005 The ethics and politics of super-intelligent machines ([https://sites.google.com/site/whibbard/g/SI\\_ethics\\_politics.doc](https://sites.google.com/site/whibbard/g/SI_ethics_politics.doc))
- [148] Hibbard B 2005 Critique of the SIAI collective volition theory ([www.ssec.wisc.edu/~billh/g/SIAI\\_CV\\_critique.html](http://www.ssec.wisc.edu/~billh/g/SIAI_CV_critique.html))
- [149] Hibbard B 2008 Open source AI *Artificial General Intelligence Frontiers (Artificial Intelligence and Applications no. 171)* ed P Wang, B Goertzel and S Franklin (Amsterdam: IOS) pp 473–7
- [150] Hibbard B 2012 Model-based utility functions *J. Artificial Gen. Intell.* **3** 1–24
- [151] Hibbard B 2012 Decision support for safe AI design *Artificial General Intelligence (Lecture Notes in Artificial Intelligence no. 7716)* ed J Bach, B Goertzel and M Ikl (New York: Springer) pp 117–25
- [152] Hibbard B 2012 The error in my 2001 VisFiles column ([www.ssec.wisc.edu/~billh/g/visfiles\\_error.html](http://www.ssec.wisc.edu/~billh/g/visfiles_error.html))
- [153] Hibbard B 2012 Avoiding unintended AI behaviors *Artificial General Intelligence (Lecture Notes in Artificial Intelligence no. 7716)* ed J Bach, B Goertzel and M Ikl (New York: Springer) pp 107–16
- [306] Yudkowsky E 1996 Staring into the singularity (<http://yudkowsky.net/obsolete/singularity.html>)
- [307] Yudkowsky E 2001 Creating friendly AI 1.0 (<http://intelligence.org/files/CFAI.pdf>)
- [308] Yudkowsky E 2004 Coherent extrapolated volition (<http://intelligence.org/files/CEV.pdf>)
- [309] Yudkowsky E 2011 Artificial intelligence as a positive and negative factor in global risk *Global Catastrophic Risks* ed N Bostrom and M M Cirkovic (Oxford: Oxford University Press) pp 308–45
- [310] Yudkowsky E 2008 Hard takeoff *Less Wrong* ([http://lesswrong.com/lw/wf/hard\\_takeoff/](http://lesswrong.com/lw/wf/hard_takeoff/))
- [311] Yudkowsky E 2009 Value is fragile *Less Wrong* ([http://lesswrong.com/lw/y3/value\\_is\\_fragile/](http://lesswrong.com/lw/y3/value_is_fragile/))
- [312] Yudkowsky E 2011 Complex value systems are required to realize valuable futures (<http://intelligence.org/files/ComplexValues.pdf>)
- [313] Yudkowsky E 2012 Reply to Holden on tool AI *Less Wrong* ([http://lesswrong.com/lw/cze/reply\\_to\\_holden\\_on\\_tool\\_ai/](http://lesswrong.com/lw/cze/reply_to_holden_on_tool_ai/))

## Invited Comment

# Responses to catastrophic AGI risk: a survey

Kaj Sotala<sup>1</sup> and Roman V Yampolskiy<sup>2</sup>

<sup>1</sup> Machine Intelligence Research Institute

<sup>2</sup> University of Louisville

Received 15 April 2014

Accepted for publication 13 November 2014

Published 19 December 2014



## Abstract

Many researchers have argued that humanity will create artificial general intelligence (AGI) within the next twenty to one hundred years. It has been suggested that AGI may inflict serious damage to human well-being on a global scale ('catastrophic risk'). After summarizing the arguments for why AGI may pose such a risk, we review the field's proposed responses to AGI risk. We consider societal proposals, proposals for external constraints on AGI behaviors and proposals for creating AGIs that are safe due to their internal design.

Keywords: artificial general intelligence, existential risk, catastrophic risk, AI risk, artificial intelligence, friendly AI, machine ethics

## 1. Introduction<sup>3</sup>

Many have argued that in the next twenty to one hundred years we will create artificial general intelligences (AGIs) [39, 46, 170, 193, 196, 200, 235, 288]<sup>4</sup>. Unlike current 'narrow' AI systems, AGIs would perform at or above the human level not merely in particular domains (e.g., chess or arithmetic), but in a wide variety of domains, including novel ones<sup>5</sup>. They would have a robust understanding of natural language and be capable of general problem solving.

<sup>3</sup> This paper elaborates and expands upon some of the discussion originally found in Yampolskiy [304], as well as reviewing a large amount of additional material.

<sup>4</sup> For a preliminary 'AGI roadmap', see Adams *et al* [10]. For a variety of views, see Eden [5]. The term 'AGI' was introduced by Gubrud [124]. For overviews of AGI approaches, see Wang *et al* [288] and Adams *et al* [10]. Some closely related terms are 'strong AI' (e.g., [170]) and 'human-level AI' (e.g., [67, 190]). Unlike the term 'human-level AI', the term 'artificial general intelligence' does not necessarily presume that the intelligence will be human-like.

<sup>5</sup> For this paper, we use a binary distinction between narrow AI and AGI. This is merely for the sake of simplicity: we do not assume the actual difference between the two categories to necessarily be so clean-cut.

The creation of AGI could pose challenges and risks of varied severity for society, such as the possibility of AGIs out competing humans in the job market [60, 189]. This article, however, focuses on the suggestion that AGIs may come to act in ways not intended by their creators, and in this way pose a *catastrophic* [52] or even an *existential* [47] risk to humanity<sup>6</sup>. We will organize and summarize the proposals that have been made so far for responding to catastrophic AGI risk, so as to provide a map of the field to newcomers and veterans alike<sup>7</sup>.

Section 2 explains why AGI may pose a catastrophic risk. Sections 3–5 review three categories of proposals for dealing with AGI risk: societal proposals, proposals for external constraints on AGI behaviors, and proposals for creating AGIs that are safe due to their internal design. Although the main purpose of this paper is to provide a summary of existing work, we briefly provide commentary on

<sup>6</sup> A catastrophic risk is something that might inflict serious damage to human well-being on a global scale and cause ten million or more fatalities [52]. An existential risk is one that threatens human extinction [47]. Many writers argue that AGI might be a risk of such magnitude [3, 44, 47, 64, 65, 68, 72, 82, 117, 126, 129, 160, 189, 193, 278, 289, 300, 309].

<sup>7</sup> One important work in this field is Bostrom [3], published after this paper was originally written. It introduces some additional proposals as well as discussing many of the ones reviewed in this work, but time constraints involved in the publication of this paper did not allow us to update this work to take it properly into account.



Content from this work may be used under the terms of the Creative Commons Attribution 3.0 licence. Any further distribution of this work must maintain attribution to the author(s) and the title of the work, journal citation and DOI.

the proposals in each major subsection of sections 3–5 and highlight some of the proposals we consider the most promising in section 6.

Attempting to influence the course of developing technologies involves a great deal of uncertainty concerning the nature of those technologies and the likely long-term impacts of societal decisions [50]. While we aim to provide a preliminary analysis of many proposals, we do not have final answers. The purpose of this article is to act as a review and to highlight some considerations that will be useful starting points for further research.

In the hopes of fostering further debate, we also highlight some of the proposals that we consider the most promising. In the medium term, these are regulation (section 3.3), merging with machines (section 3.4), AGI confinement (section 4.1), Oracle AI (section 5.1) and motivational weaknesses (section 5.6). In the long term, the most promising approaches seem to be value learning (section 5.2.5) and human-like architectures (section 5.3.4). Section 6 provides an extended discussion of the various merits and problems of these proposals.

## 2. Catastrophic AGI risk

We begin with a brief sketch of the argument that AGI poses a catastrophic risk to humanity. At least two separate lines of argument seem to support this conclusion.

First, AI has already made it possible to automate many jobs [60] and AGIs, when they are created, should be capable of performing *most* jobs better than humans [130, 136, 189, 289]. As humanity grows increasingly reliant on AGIs, these AGIs will begin to wield more and more influence and power. Even if AGIs initially function as subservient tools, an increasing number of decisions will be made by autonomous AGIs rather than by humans. Over time it would become ever more difficult to replace the AGIs, even if they no longer remained subservient.

Second, there may be a sudden discontinuity in which AGIs rapidly become far more numerous or intelligent [72, 117, 244, 251, 306]. This could happen due to (1) a conceptual breakthrough which makes it easier to run AGIs using far less hardware, (2) AGIs using fast computing hardware to develop ever-faster hardware, or (3) AGIs crossing a threshold in intelligence that allows them to carry out increasingly fast software self-improvement. Even if the AGIs were expensive to develop at first, they could be cheaply copied and could thus spread quickly once created.

Once they become powerful enough, AGIs might be a threat to humanity even if they are not actively malevolent or hostile. Mere indifference to human values—including human survival—could be sufficient for AGIs to pose an existential threat [211, 212, 309, 312].

We will now lay out the above reasoning in more detail.

### 2.1. Most tasks will be automated

Ever since the Industrial Revolution, society has become increasingly automated. Brynjolfsson [60] argue that the current high unemployment rate in the United States is partially due to rapid advances in information technology, which has made it possible to replace human workers with computers faster than human workers can be trained in jobs that computers cannot yet perform. Vending machines are replacing shop attendants, automated discovery programs which locate relevant legal documents are replacing lawyers and legal aides, and automated virtual assistants are replacing customer service representatives.

Labor is becoming automated for reasons of cost, efficiency and quality. Once a machine becomes capable of performing a task as well as (or almost as well as) a human, the cost of purchasing and maintaining it may be less than the cost of having a salaried human perform the same task. In many cases, machines are also capable of doing the same job faster, for longer periods and with fewer errors. In addition to replacing workers entirely, machines may also take over aspects of jobs that were once the sole domain of highly trained professionals, making the job easier to perform by less-skilled employees [298].

If workers can be affordably replaced by developing more sophisticated AI, there is a strong economic incentive to do so. This is already happening with narrow AI, which often requires major modifications or even a complete redesign in order to be adapted for new tasks. ‘A roadmap for US robotics’ [154] calls for major investments into automation, citing the potential for considerable improvements in the fields of manufacturing, logistics, health care and services. Similarly, the US Air Force Chief Scientist’s [78] ‘Technology horizons’ report mentions ‘increased use of autonomy and autonomous systems’ as a key area of research to focus on in the next decade, and also notes that reducing the need for manpower provides the greatest potential for cutting costs. In 2000, the US Congress instructed the armed forces to have one third of their deep strike force aircraft be unmanned by 2010, and one third of their ground combat vehicles be unmanned by 2015 [4].

To the extent that an AGI could learn to do many kinds of tasks—or even *any* kind of task—without needing an extensive re-engineering effort, the AGI could make the replacement of humans by machines much cheaper and more profitable. As more tasks become automated, the bottlenecks for further automation will require adaptability and flexibility that narrow-AI systems are incapable of. These will then make up an increasing portion of the economy, further strengthening the incentive to develop AGI.

Increasingly sophisticated AI may eventually lead to AGI, possibly within the next several decades [39, 200]. Eventually it will make economic sense to automate all or nearly all jobs [130, 136, 289]. As AGIs will possess many advantages over humans [200, 253], a greater and greater proportion of the workforce will consist of intelligent machines.

## 2.2. AGIs might harm humans

AGIs might bestow overwhelming military, economic, or political power on the groups that control them [47]. For example, automation could lead to an ever-increasing transfer of wealth and power to the owners of the AGIs [54, 60]. AGIs could also be used to develop advanced weapons and plans for military operations or political takeovers [47, 124, 163]. Some of these scenarios could lead to catastrophic risks, depending on the capabilities of the AGIs and other factors.

Our focus is on the risk from the possibility that AGIs could behave in unexpected and harmful ways, even if the intentions of their owners were benign. Even modern-day narrow-AI systems are becoming autonomous and powerful enough that they sometimes take unanticipated and harmful actions before a human supervisor has a chance to react. To take one example, rapid automated trading was found to have contributed to the 2010 stock market ‘flash crash’ [70]<sup>8</sup>. Autonomous systems may also cause people difficulties in more mundane situations, such as when a credit card is automatically flagged as possibly stolen due to an unusual usage pattern [16], or when automatic defense systems malfunction and cause deaths [240].

As machines become more autonomous, humans will have fewer opportunities to intervene in time and will be forced to rely on machines making good choices. This has prompted the creation of the field of ‘machine ethics’ [1, 16, 282], concerned with creating AI systems designed to make appropriate moral choices. Compared to narrow-AI systems, AGIs will be even more autonomous and capable, and will thus require even more robust solutions for governing their behavior<sup>9</sup>.

If some AGIs were both powerful and indifferent to human values, the consequences could be disastrous. At one extreme, powerful AGIs indifferent to human survival could bring about human extinction. As Yudkowsky [309] writes, ‘The AI does not hate you, nor does it love you, but you are made out of atoms which it can use for something else’.

Omohundro [211, 212] and Bostrom [51] argue that standard microeconomic theory prescribes particular instrumental behaviors which are useful for the achievement of almost any set of goals. Furthermore, any agents which do not follow certain axioms of rational behavior will possess vulnerabilities which some other agent may exploit to their own benefit. Thus AGIs which understand these principles and wish to act efficiently will modify themselves so that their

behavior more closely resembles rational economic behavior [213]. Extra resources are useful in the pursuit of nearly any set of goals and self-preservation behaviors will increase the probability that the agent can continue to further its goals. AGI systems which follow rational economic theory will then exhibit tendencies toward behaviors such as self-replicating, breaking into other machines and acquiring resources without regard for anyone else’s safety. They will also attempt to improve themselves in order to more effectively achieve these and other goals, which could lead to rapid improvement even if the designers did not intend the agent to self-improve.

Even AGIs that were explicitly designed to behave ethically might end up acting at cross-purposes to humanity, because it is difficult to precisely capture the complexity of human values in machine goal systems [30, 199, 309, 312].

Muehlhauser [199] caution that moral philosophy has found no satisfactory formalization of human values. All moral theories proposed so far would lead to undesirable consequences if implemented by superintelligent machines. For example, a machine programmed to maximize the satisfaction of human (or sentient) preferences might simply modify people’s brains to give them desires that are maximally easy to satisfy.

Intuitively, one might say that current moral theories are all *too simple*—even if they seem correct at first glance, they do not actually take into account all the things that we value and this leads to a catastrophic outcome. This could be referred to as the *complexity of value thesis*. Recent psychological and neuroscientific experiments confirm that human values are highly complex [199], that the pursuit of pleasure is not the only human value [7] and that humans are often unaware of their own values [96, 197, 302].

Still, perhaps powerful AGIs would have desirable consequences so long as they were programmed to respect *most* human values. If so, then our inability to perfectly specify human values in AGI designs need not pose a catastrophic risk. Different cultures and generations have historically had very different values from each other and it seems likely that over time our values would become considerably different from current-day ones. It could be enough to maintain some small set of core values, though what exactly would constitute a core value is unclear<sup>10</sup>.

Yudkowsky [312] argues that, due to the fragility of value, the basic problem remains. He argues that, even if an AGI implemented *most* human values, the outcome might still be unacceptable. For example, an AGI which failed to incorporate the value of novelty could create a solar system filled with countless minds experiencing one highly optimal and satisfying experience over and over again, never doing or feeling anything else [311]<sup>11</sup>.

<sup>8</sup> On the less serious front, see Eisen [94] for an amusing example of automated trading going awry.

<sup>9</sup> In practice, there have been two separate communities doing research on automated moral decision-making [13, 199, 246]. The ‘AGI ethics’ community has concentrated specifically on advanced AGIs (e.g., [110, 309]), while the ‘machine ethics’ community typically has concentrated on more immediate applications for current-day AI (e.g., [1, 284]). In this paper, we have cited the machine ethics literature only where it seemed relevant, leaving out papers that seemed to be too focused on narrow-AI systems for our purposes. In particular, we have left out most discussions of military machine ethics [26], which focus primarily on the constrained special case of creating systems that are safe for battlefield usage. Note that while the term ‘machine ethics’ is relatively established, ‘AGI ethics’ is not. One proposed alternative name for the ‘AGI ethics’ discipline is ‘AI safety engineering’ [304, 305].

<sup>10</sup> For example, different people may disagree over whether freedom or well-being is a more important value.

<sup>11</sup> Miller [189] similarly notes that, despite a common belief to the contrary, it is impossible to write laws in a manner that would match our stated moral principles without a judge needing to use a large amount of implicit common-sense knowledge to correctly interpret them: ‘Laws shouldn’t always be interpreted literally because legislators can’t anticipate all possible contingencies. Also, humans’ intuitive feel for what constitutes murder goes beyond anything we can commit to paper. The same applies to friendliness’ [189].

In this paper, we will frequently refer to the problem of ‘AGI safety’ or ‘safe AGI’, by which we mean the problem of ensuring that AGIs respect human values, or perhaps some extrapolation or idealization of human values<sup>12</sup>. We do not seek to imply that current human values would be the best possible ones, that AGIs could not help us in developing our values further, or that the values of other sentient beings would be irrelevant. Rather, by ‘human values’ we refer to the kinds of basic values that nearly all humans would agree upon, such as that AGIs forcibly reprogramming people’s brains, or destroying humanity, would be a bad outcome. In cases where proposals related to AGI risk might change human values in some major but not as obviously catastrophic way, we will mention the possibility of these changes but remain agnostic on whether they are desirable or undesirable.

We conclude this section with one frequently forgotten point: in order to avoid catastrophic risks or worse, it is not enough to ensure that only some AGIs are safe. Proposals which seek to solve the issue of catastrophic AGI risk need to also provide some mechanism for ensuring that *most* (or perhaps even ‘nearly all’) AGIs are either created safe or prevented from doing considerable harm.

### 2.3. AGIs may become powerful quickly

There are several reasons why AGIs may quickly come to wield unprecedented power in society. ‘Wielding power’ may mean having direct decision-making power, or it may mean carrying out human decisions in a way that makes the decision maker reliant on the AGI. For example, in a corporate context an AGI could be acting as the executive of the company, or it could be carrying out countless low-level tasks which the corporation needs to perform as part of its daily operations.

Bugaj [63] consider three kinds of AGI scenarios: capped intelligence, soft takeoff and hard takeoff. In a *capped intelligence* scenario, all AGIs are prevented from exceeding a predetermined level of intelligence and remain at a level roughly comparable with humans. In a *soft takeoff* scenario, AGIs become far more powerful than humans, but on a timescale which permits ongoing human interaction during the ascent. Time is not of the essence, and learning proceeds at a relatively human-like pace. In a *hard takeoff* scenario, an AGI will undergo an extraordinarily fast increase in power, taking effective control of the world within a few years or less<sup>13</sup>. In this scenario, there is little time for error correction or a gradual tuning of the AGI’s goals.

<sup>12</sup> Within the AGI ethics literature, safe autonomous AGI is sometimes called ‘friendly AI’ [115, 183, 189, 307, 309, 312]. Yudkowsky [307] defines ‘friendly AI’ as ‘the production of human-benefiting, non-human-harming actions in Artificial Intelligence systems that have advanced to the point of making real-world plans in pursuit of goals’. However, some papers (e.g., Goertzel [105, 108, 294]) use ‘friendly AI’ as a narrower term to refer to safe AGI designs as advocated by Yudkowsky [307, 309]. These designs have the goal of benefiting humans as their overarching value, from which all the other goals and values of the system are derived. In this paper we use the term ‘friendly AI’ to refer to Yudkowsky’s proposal and ‘safe AGI’ as our more general term.

<sup>13</sup> Bugaj and Goertzel defined hard takeoff to refer to a period of months or less. We have chosen a somewhat longer time period, as even a few years might easily turn out to be too little time for society to properly react.

The viability of many proposed approaches depends on the hardness of a takeoff. The more time there is to react and adapt to developing AGIs, the easier it is to control them. A soft takeoff might allow for an approach of incremental machine ethics [223], which would not require us to have a complete philosophical theory of ethics and values, but would rather allow us to solve problems in a gradual manner. A soft takeoff might however present its own problems, such as there being a larger number of AGIs distributed throughout the economy, making it harder to contain an eventual takeoff.

Hard takeoff scenarios can be roughly divided into those involving the quantity of hardware (the *hardware overhang* scenario), the quality of hardware (the *speed explosion* scenario) and the quality of software (the *intelligence explosion* scenario). Although we discuss them separately, it seems plausible that several of them could happen simultaneously and feed into each other.

**2.3.1. Hardware overhang.** Hardware progress may outpace AGI software progress. Contemporary supercomputers already rival or even exceed some estimates of the computational capacity of the human brain, while no software seems to have both the brain’s general learning capacity and its scalability.

Bostrom [46] estimates that the effective computing capacity of the human brain might be somewhere around  $10^{17}$  operations per second (OPS) and Moravec [195] estimates it at  $10^{14}$  OPS. As of November 2012, the fastest supercomputer in the world had achieved a top capacity of  $10^{16}$  floating-point operations per second (FLOPS) and the five-hundredth fastest a top capacity of  $10^{13}$  FLOPS [188]. Note however that OPS and FLOPS are not directly comparable and there is no reliable way of inter-converting the two. Sandberg and Bostrom [234] estimate that OPS and FLOPS grow at a roughly comparable rate.

If such trends continue, then by the time the software for AGI is invented there may be a *computing overhang*—an abundance of cheap hardware available for running thousands or millions of AGIs, possibly with a speed of thought much faster than that of humans [244, 253, 310].

As increasingly sophisticated AGI software becomes available, it would be possible to rapidly copy improvements to millions of servers, each new version being capable of doing more kinds of work or being run with less hardware. Thus, the AGI software could replace an increasingly large fraction of the workforce<sup>14</sup>. The need for AGI systems to be trained for some jobs would slow the rate of adoption, but powerful computers could allow for fast training. If AGIs end

<sup>14</sup> The speed that would allow AGIs to take over most jobs would depend on the cost of the hardware and the granularity of the software upgrades. A series of upgrades over an extended period, each producing a 1% improvement, would lead to a more gradual transition than a single upgrade that brought the software from the capability level of a chimpanzee to a rough human equivalence. Note also that several companies, including Amazon and Google, offer vast amounts of computing power for rent on an hourly basis. An AGI that acquired money and then invested all of it in renting a large amount of computing resources for a brief period could temporarily achieve a much larger boost than its budget would otherwise suggest.

up doing the vast majority of work in society, humans could become dependent on them.

AGIs could also plausibly take control of Internet-connected machines in order to harness their computing power [253]; Internet-connected machines are regularly compromised<sup>15</sup>.

**2.3.2. Speed explosion.** Another possibility is a *speed explosion* [72, 158, 251, 306], in which intelligent machines design increasingly faster machines. A hardware overhang might contribute to a speed explosion, but is not required for it. An AGI running at the pace of a human could develop a second generation of hardware on which it could run at a rate faster than human thought. It would then require a shorter time to develop a third generation of hardware, allowing it to run faster than on the previous generation, and so on. At some point, the process would hit physical limits and stop, but by that time AGIs might come to accomplish most tasks at far faster rates than humans, thereby achieving dominance. (In principle, the same process could also be achieved via improved software.)

The extent to which the AGI needs humans in order to produce better hardware will limit the pace of the speed explosion, so a rapid speed explosion requires the ability to automate a large proportion of the hardware manufacturing process. However, this kind of automation may already be achieved by the time that AGI is developed<sup>16</sup>.

**2.3.3. Intelligence explosion.** Third, there could be an *intelligence explosion*, in which one AGI figures out how to create a qualitatively smarter AGI and that AGI uses its increased intelligence to create still more intelligent AGIs, and so on<sup>17</sup>, such that the intelligence of humankind is quickly left far behind and the machines achieve dominance [72, 117, 177, 200]<sup>18</sup>.

<sup>15</sup> Botnets are networks of computers that have been compromised by outside attackers and are used for illegitimate purposes. Rajab *et al* [226] review several studies which estimate the sizes of the largest botnets as being between a few thousand to 350 000 bots. Modern-day malware could theoretically infect any susceptible Internet-connected machine within tens of seconds of its initial release [257]. The Slammer worm successfully infected more than 90% of vulnerable hosts within ten minutes and had infected at least 75 000 machines by the thirty minute mark [192]. The previous record holder in speed, the Code Red worm, took fourteen hours to infect more than 359 000 machines [191].

<sup>16</sup> Loosemore [177] also suggest that current companies carrying out research and development are more constrained by a lack of capable researchers than by the ability to carry out physical experiments.

<sup>17</sup> Most accounts of this scenario do not give exact definitions for ‘intelligence’ or explain what a ‘superintelligent’ AGI would be like, instead using informal characterizations such as ‘a machine that can surpass the intellectual activities of any man however clever’ [117] or ‘an intellect that is much smarter than the best human brains in practically every field, including scientific creativity, general wisdom and social skills’ [46]. Yudkowsky [309] defines intelligence in relation to ‘optimization power’, the ability to reliably hit small targets in large search spaces, such as by finding the *a priori* exceedingly unlikely organization of atoms which makes up a car. A more mathematical definition of machine intelligence is offered by Legg [173]. Sotala [253] discusses some of the functional routes to actually achieving superintelligence.

<sup>18</sup> One example of an AGI framework designed specifically for repeated self-improvement is offered by Schmidhuber [236].

Yudkowsky [309, 310] argues that an intelligence explosion is likely. So far, natural selection has been improving human intelligence and human intelligence has to some extent been able to improve itself. However, the core process by which natural selection improves humanity has been essentially unchanged and humans have been unable to deeply affect the cognitive algorithms which produce their own intelligence. Yudkowsky suggests that if a mind became capable of directly editing itself, this could spark a rapid increase in intelligence, as the actual process causing increases in intelligence could itself be improved upon. (This requires that there exist powerful improvements which, when implemented, considerably increase the rate at which such minds can improve themselves.)

Hall [130] argues that, based on standard economic considerations, it would not make sense for an AGI to focus its resources on solitary self-improvement. Rather, in order not to be left behind by society at large, it should focus its resources on doing the things that it is good at and trade for the things it is not good at. However, once there exists a community of AGIs that can trade with one another, this community could collectively undergo rapid improvement and leave humans behind.

A number of formal growth models have been developed which are relevant to predicting the speed of a takeoff; an overview of these can be found in Sandberg’s [231] paper. Many of them suggest rapid growth. For instance, Hanson [134] suggests that AGI might lead to the economy doubling in months rather than years. However, Hanson is skeptical about whether this would prove a major risk to humanity and considers it mainly an economic transition similar to the Industrial Revolution.

To some extent, the soft/hard takeoff distinction may be a false dichotomy: a takeoff may be soft for a while, and then become hard. Two of the main factors influencing the speed of a takeoff are the pace at which computing hardware is developed and the ease of modifying minds [253]. This allows for scenarios in which AGI is developed and there seems to be a soft takeoff for, say, the initial ten years, causing a false sense of security until a breakthrough in hardware development causes a hard takeoff.

Another factor that might cause a false sense of security is the possibility that AGIs can be developed by a combination of insights from humans and AGIs themselves. As AGIs become more intelligent and it becomes possible to automate portions of the development effort, those parts accelerate and the parts requiring human effort become bottlenecks. Reducing the amount of human insight required could dramatically accelerate the speed of improvement. Halving the amount of human involvement required might at most double the speed of development, possibly giving an impression of relative safety, but going from 50% human insight required to 1% human insight required could cause the development to become ninety-nine times faster<sup>19</sup>.

<sup>19</sup> The relationship in question is similar to that described by Amdahl’s [17] law.

From a safety viewpoint, the conservative assumption is to presume the worst [307]. Yudkowsky argues that the worst outcome would be a hard takeoff, as it would give us the least time to prepare and correct errors. On the other hand, it can also be argued that a soft takeoff would be just as bad, as it would allow the creation of multiple competing AGIs, allowing the AGIs that were the least burdened with goals such as ‘respect human values’ to prevail. We would ideally like a solution, or a combination of solutions, which would work effectively for both a soft and a hard takeoff.

### 3. Societal proposals

The notion of catastrophic AGI risk is not new and this concern was expressed by early thinkers in the field [64, 117, 269, 300]. Hence, there have also been many proposals concerning what to do about it. The proposals we review are neither exhaustive nor mutually exclusive: the best way of achieving a desirable outcome may involve pursuing several proposals simultaneously.

Proposals can be divided into three general categories: proposals for societal action, design proposals for external constraints on AGI behavior and design recommendations for internal constraints on AGI behavior. In this section we briefly review societal proposals. These include doing nothing, integrating AGIs with society, regulating research, merging with machines and relinquishing research into AGI.

#### 3.1. Do nothing

**3.1.1. AI is too distant to be worth our attention.** One response is that, although AGI is possible in principle, there is no reason to expect it in the near future. Typically, this response arises from the belief that, although there have been great strides in narrow AI, researchers are still very far from understanding how to build AGI. Distinguished computer scientists such as Gordon Bell and Gordon Moore, as well as cognitive scientists such as Douglas Hofstadter and Steven Pinker, have expressed the opinion that the advent of AGI is remote [6]. Davis [80] reviews some of the ways in which computers are still far from human capabilities. Bringsjord [58] even claim that a belief in AGI this century is fideistic, appropriate within the realm of religion but not within science or engineering.

Some writers also actively criticize any discussion of AGI risk in the first place. The philosopher Alfred Nordmann [208, 209] holds the view that ethical concern is a scarce resource, not to be wasted on unlikely future scenarios such as AGI. Likewise, Dennett [86] considers AGI risk an ‘imprudent pastime’ because it distracts our attention from more immediate threats.

Others think that AGI is far off and not yet a major concern, but admit that it might be valuable to give the issue some attention. A presidential panel of the Association for the Advancement of Artificial Intelligence considering the long-term future of AI concluded that there was overall skepticism about AGI risk, but that additional research into the topic and

related subjects would be valuable [156]. Posner [220] writes that dedicated efforts for addressing the problem can wait, but that we should gather more information about the problem in the meanwhile.

The potential negative consequences of AGI are enormous, ranging from economic instability to human extinction. ‘Do nothing’ could be a reasonable course of action if near-term AGI seemed extremely unlikely, if it seemed too early for any proposals to be effective in reducing risk, or if those proposals seemed too expensive to implement.

As a comparison, asteroid impact prevention is generally considered a topic worth studying, even though the probability of a civilization-threatening asteroid impact in the near future is not considered high. Napier [202] discusses several ways of estimating the frequency of such impacts. Many models produce a rate of one civilization-threatening impact per five hundred thousand or more years, though some models suggest that rates of one such impact per hundred thousand years cannot be excluded.

An estimate of one impact per hundred thousand years would suggest less than a 0.1% chance of a civilization-threatening impact within the next hundred years. The probability of AGI being developed within the same period seems considerably higher [200] and there is likewise a reasonable chance of a hard takeoff after it has been developed [309, 310], suggesting that the topic is at the very least worth studying. Even without a hard takeoff society is becoming increasingly automated and even narrow AI is starting to require ethical guidelines [26, 282].

We know neither which fields of science will be needed nor how much progress in them will be necessary for safe AGI. If much progress is needed and we believe effective progress to be possible this early on, it becomes reasonable to start studying the topic even before AGI is near. Muehlhauser [199] suggest that, for one safe AGI approach alone (value learning, discussed further in section 5), efforts by AGI researchers, economists, mathematicians and philosophers may be needed. Safe AI may require the solutions for some of these problems to come well before AGI is developed.

**3.1.2. Little risk, no action needed.** Some authors accept that a form of AGI will probably be developed but do not consider autonomous AGI to be a risk, or consider the possible negative consequences acceptable. Bryson [61] argue that, although AGI will require us to consider ethical and social dangers, the dangers will be no worse than those of other technologies. Whitby [298] writes that there has historically been no consistent trend of the most intelligent people acquiring the most authority and that computers will augment humans rather than replace them. Whitby [299] further argue that AGIs will not have any particular motivation to act against us. Jenkins [159] agrees with these points to the extent of saying that a machine will only act against humans if it is programmed to value itself over humans, although she does find AGI to be a real concern.

Another kind of ‘no action needed’ response argues that AGI development will take a long time [59], implying that

there will be plenty of time to deal with the issue later on. This can also be taken as an argument for later efforts being more effective, as they will be better tuned to AGI as it develops.

Others argue that superintelligence will not be possible at all<sup>20</sup>. McDermott [182] points out that there are no good examples of algorithms which could be improved upon *indefinitely*. Deutsch [87] argues that there will never be superintelligent AGIs, because human minds are already universal reasoners and computers can at best speed up the experimental work that is required for testing and fine-tuning theories. He also suggests that even as the speed of technological development increases, so will our ability to deal with change. Anderson [22] likewise suggests that the inherent unpredictability of the world will place upper limits on an entity's effective intelligence.

Heylighen [145] argues that a single, stand-alone computer is exceedingly unlikely to become superintelligent and that individual intelligences are always outmatched by the distributed intelligence found in social systems of many minds. Superintelligence will be achieved by building systems that integrate and improve the 'global brain', the collective intelligence of everyone on Earth. Heylighen does acknowledge that this kind of a transition will pose its own challenges, but not of the kind usually evoked in discussions of AGI risk.

The idea of AGIs not having a motivation to act against humans is intuitively appealing, but there seem to be strong theoretical arguments against it. As mentioned earlier, Omohundro [211, 212] and Bostrom [51] argue that self-replication and the acquisition of resources are useful in the pursuit of many different kinds of goals and that many types of AI systems will therefore exhibit tendencies toward behaviors such as breaking into other machines, self-replicating and acquiring resources without regard for anyone else's safety. The right design might make it possible to partially work around these behaviors [243, 287], but they still need to be taken into account. Furthermore, we might not foresee all the complex interactions of different AGI mechanisms in the systems that we build and they may end up with very different goals than the ones we intended [88, 229, 309, 312].

Can AGIs become superintelligent? First, we note that AGIs do not necessarily need to be much more intelligent than humans in order to be dangerous. AGIs already enjoy advantages such as the ability to rapidly expand their population by having themselves copied [133, 136, 253], which may confer on them considerable economic and political influence even if they were not superintelligent. A better-than-human ability to coordinate their actions, which AGIs of a similar design could plausibly have [253], might then be enough to tilt the odds in their favor.

Another consideration is that AGIs do not necessarily need to be qualitatively more intelligent than humans in order to outperform humans. An AGI that merely thought twice as fast as any single human could still defeat him at intellectual tasks that had a time constraint, all else equal. Here an 'intellectual' task should be interpreted broadly to refer not only to 'book smarts' but to any task that animals cannot perform due to their mental limitations—including tasks involving social skills [309]. Straightforward improvements in computing power could provide AGIs with a considerable advantage in speed [72, 158, 251, 253, 306], which the AGI could then use to study and accumulate experiences that improved its skills.

As for Heylighen's [145] 'global brain' argument, there does not seem to be a reason to presume that powerful AGIs could not be geographically distributed, or that they could not seize control of much of the Internet. Even if individual minds were not very smart and needed a society to make progress, for minds that are capable of copying themselves and communicating perfectly with each other, individual instances of the mind might be better understood as parts of a whole than as separate individuals [253, 254]. In general, the distinction between an individual and a community might not be meaningful for AGIs [115]. If there were enough AGIs, they might be able to form a community sufficient to take control of the rest of the Earth. Heylighen [144] himself has argued that many of the features of the Internet are virtually identical to the mechanisms used by the human brain. If the AGI is not carefully controlled, it might end up in a position where it made up the majority of the 'global brain' and could undertake actions which the remaining parts of the organism did not agree with.

**3.1.3. Let them kill us.** Dietrich [89] argues that humanity frequently harms other species and that people have also evolved to hurt other people by engaging in behaviors such as child abuse, sexism, rape and racism. Therefore, human extinction would not matter, as long as the machines implemented only the positive aspects of humanity.

De Garis [82] suggests that AGIs destroying humanity might not matter. He writes that on a cosmic scale, with hundreds of billions of stars in our galaxy alone, the survival of the inhabitants of a single planet is irrelevant. As AGIs would be more intelligent than us in every way, it would be better if they replaced humanity.

AGIs being more intelligent and therefore more valuable than humans equates intelligence with value, but Bostrom [49] suggests ways by which a civilization of highly intelligent entities might lack things which we thought to have value. For example, such entities might not be conscious in the first place. Alternatively, there are many things which we consider valuable for their own sake, such as humor, love, game-playing, art, sex, dancing, social conversation, philosophy, literature, scientific discovery, food and drink, friendship, parenting, and sport. We value these due to the fact that we have dispositions and preferences which have been evolutionarily adaptive in the past, but for a future civilization

<sup>20</sup> The opposite argument is that superior intelligence will inevitably lead to more moral behavior. Some of the arguments related to this position are discussed in the context of evolutionary invariants (section 5.3.1), although the authors advocating the use of evolutionary invariants do believe AGI risk to be worth our concern.

few or none of them might be, creating a world with very little of value. Bostrom [51] proposes an orthogonality thesis, by which an artificial intelligence can have any combination of intelligence level and goal, including goals that humans would intuitively deem to be of no value.

**3.1.4. 'Do nothing' proposals: our view.** As discussed above, completely ignoring the possibility of AGI risk at this stage would seem to require a confident belief in at least one of the following propositions.

1. AGI is very remote.
2. There is no major risk from AGI even if it is created.
3. Very little effective work can be done at this stage.
4. AGIs destroying humanity would not matter.

In the beginning of this paper, we mentioned several experts who considered it plausible that AGI might be created in the next twenty to one hundred years; in this section we have covered experts who disagree.

In general, there is a great deal of disagreement among people who have made AGI predictions, and no clear consensus even among experts in the field of artificial intelligence. The lack of expert agreement suggests that expertise in the field does not contribute to an ability to make reliable predictions<sup>21</sup>. If the judgment of experts is not reliable, then, probably, neither is anyone else's. This suggests that it is unjustified to be highly certain of AGI being near, but also of it *not* being near. We thus consider it unreasonable to have a confident belief in the first proposition.

The second proposition also seems questionable: as discussed in sections 2.1, 2.3 and 3.1.2, AGIs seem very likely to obtain great power, possibly very quickly. Furthermore, as discussed in section 2.2, the complexity and fragility of value theses imply that it could be very difficult to create AGIs which would not cause immense amounts of damage if they had enough power.

It also does not seem like it is too early to work on the problem: as we summarize in section 6, there seem to be a number of promising research directions which can already be pursued. We also agree with Yudkowsky [309], who points out that research on the philosophical and technical requirements of safe AGI might show that broad classes of possible AGI architectures are fundamentally unsafe, suggesting that such architectures should be avoided. If this is the case, it seems better to have that knowledge as early as possible, before there has been a great deal of investment into unsafe AGI designs.

In response to the suggestion that humanity being destroyed would not matter, we certainly agree that there is much to be improved in today's humanity, and that our future descendants might have very little resemblance to ourselves.

<sup>21</sup> Armstrong [29] point out that many of the task properties which have been found to be conducive for developing reliable and useful expertise are missing in AGI timeline forecasting. In particular, one of the most important factors is whether experts get rapid (preferably immediate) feedback, while a timeline prediction that is set many decades in the future might have been entirely forgotten by the time that its correctness could be evaluated.

Regardless, we think that much about today's humans is valuable and worth preserving, and that we should be able to preserve it without involving the death of present humans.

### 3.2. Integrate with society

Integration proposals hold that AGI might be created in the next several decades and that there are indeed risks involved. These proposals argue that the best way to deal with the problem is to make sure that our societal structures are equipped to handle AGIs once they are created.

There has been some initial work toward integrating AGIs with existing legal and social frameworks, such as considering questions of their legal position [31, 62, 132, 174, 252, 296, 297] and moral rights [62, 125, 175, 260, 261, 291].

**3.2.1. Legal and economic controls.** Hanson [138] writes that the values of older and younger generations have often been in conflict with each other and he compares this to a conflict between humans and AGIs. He believes that the best way to control AGI risk is to create a legal framework such that it is in the interest of both humans and AGIs to uphold it. Hanson [137] suggests that if the best way for AGIs to get what they want is via mutually agreeable exchanges, then humans would need to care less about what the AGIs wanted. According to him, we should be primarily concerned with ensuring that the AGIs will be law-abiding enough to respect our property rights. Miller [189] summarizes Hanson's argument and the idea that humanity could be content with a small fraction of the world's overall wealth and let the AGIs have the rest. An analogy to this idea is that humans do not kill people who become old enough to no longer contribute to production, even though younger people could in principle join together and take the wealth of the older people. Instead, old people are allowed to keep their wealth even while in retirement. If things went well, AGIs might similarly allow humanity to 'retire' and keep its accumulated wealth, even if humans were no longer otherwise useful for AGIs.

Hall [128] also says that we should ensure that the interactions between ourselves and machines are economic, 'based on universal rules of property and reciprocity'. Moravec [196] likewise writes that governmental controls should be used to ensure that humans benefit from AGIs. Without government intervention, humans would be squeezed out of existence by more efficient robots, but taxation could be used to support human populations for a long time. He also recommends laws which would require any AGIs to incorporate programming that made them safe and subservient to human desires. Sandberg [232] writes that relying only on legal and economic controls would be problematic, but that a strategy which also incorporated them in addition to other approaches would be more robust than a strategy which did not.

However, even if AGIs were integrated with human institutions, it does not guarantee that human values would survive. If humans were reduced to a position of negligible power, AGIs might not have any reason to keep us around.

Economic arguments, such as the principle of comparative advantage, are sometimes invoked to argue that AGI would find it more beneficial to trade with us than to do us harm. However, technological progress can drive the wages of workers below the level needed for survival [60, 74, 100, 189] and there is already a possible threat of technological unemployment [60]. AGIs keeping humans around due to gains from trade implicitly presumes that they would not have the will or the opportunity to simply eliminate humans in order to replace them with a better trading partner and then trade with the new partner instead.

Humans already eliminate species with low economic value in order to make room for more humans, such as when clearing a forest in order to build new homes. Clark uses the example of horses in Britain: their population peaked in 1901, with 3.25 million horses doing work such as plowing fields, hauling wagons and carriages short distances, and carrying armies into battle. The internal combustion engine replaced so many of them that by 1924 there were fewer than two million. Clark writes:

There was always a wage at which all these horses could have remained employed. But that wage was so low that it did not pay for their feed, and it certainly did not pay enough to breed fresh generations of horses to replace them. Horses were thus an early casualty of industrialization [74].

There are also ways to harm humans while still respecting their property rights, such as by manipulating them into making bad decisions, or selling them addictive substances. If AGIs were sufficiently smarter than humans, humans could be tricked into making a series of trades that respected their property rights but left them with negligible assets and caused considerable damage to their well-being.

A related issue is that AGIs might become more capable of changing our values than we are capable of changing AGI values. Mass media already convey values that have a negative impact on human well-being, such as idealization of rare body types, which causes dissatisfaction among people who do not have those kinds of bodies [12, 122]. AGIs with a deep understanding of human psychology could engineer the spread of values which shifted more power to them, regardless of their effect on human well-being.

Yet another problem is ensuring that the AGIs have indeed adopted the right values. Making intelligent beings adopt specific values is a difficult process which often fails. There could be an AGI with the wrong goals that would pretend to behave correctly in society throughout the whole socialization process. AGIs could conceivably preserve and conceal their goals far better than humans could.

Society does not know of any methods which would reliably instill our chosen values in *human* minds, despite a long history of trying to develop them. Our attempts to make AGIs adopt human values would be hampered by our lack of experience and understanding of the AGI's thought processes, with even tried-and-true methods for instilling positive values in humans possibly being ineffective. The limited success that

we do have with humans is often backed up by various incentives as well as threats of punishment, both of which might fail in the case of an AGI developing to become vastly more powerful than us.

Additionally, the values which a being is likely to adopt, or is even capable of adopting, will depend on its mental architecture. We will demonstrate these claims with examples from humans, who are not blank slates on whom arbitrary values can be imposed with the right education [218]. Although the challenge of instilling specific values in humans is very different from the challenge of instilling them in AGIs, our examples are meant to demonstrate the fact that the existing properties of a mind will affect the process of acquiring values. Just as it is difficult to make humans permanently adopt some kinds of values, the kind of mental architecture that an AGI has will affect its inclination to adopt various values.

Psychopathy is a risk factor for violence and psychopathic criminals are much more likely to re-offend than non-psychopaths [139]. Harris [140] argue that therapy for psychopaths is ineffective<sup>22</sup> and may even make them more dangerous, as they use their improved social skills to manipulate others more effectively. Furthermore, 'cult brainwashing' is generally ineffective and most cult members will eventually leave [25] and large-scale social engineering efforts often face widespread resistance, even in dictatorships with few scruples about which methods to use [chapters 6 and 7 [238]]. Thus, while one can try to make humans adopt values, this will only work to the extent that the individuals in question are actually disposed toward adopting them.

**3.2.2. Foster positive values.** Kurzweil [170], considering the possible effects of many future technologies, notes that AGI may be a catastrophic risk. He generally supports regulation and partial relinquishment of dangerous technologies, as well as research into their defensive applications. However, he believes that with AGI this may be insufficient and that, at the present time, it may be infeasible to develop strategies that would guarantee safe AGI. He argues that machine intelligences will be tightly integrated into our society and that, for the time being, the best chance of avoiding AGI risk is to foster positive values in our society. This will increase the likelihood that any AGIs that are created will reflect such positive values.

One possible way of achieving such a goal is moral enhancement [91], the use of technology to instill people with better motives. Persson [215, 216] argue that, as technology improves, we become more capable of damaging humanity, and that we need to carry out moral enhancement in order to lessen our destructive impulses.

**3.2.3. 'Integrate with society' proposals: our view.** Proposals to incorporate AGIs into society suffer from the issue that some AGIs may never adopt benevolent and cooperative values, no matter what the environment. Neither does the intelligence of the AGIs necessarily affect their values [51].

<sup>22</sup> Salekin [230] offers a more optimistic opinion.

Sufficiently intelligent AGIs could certainly come to eventually understand human values, but humans can also come to understand others' values while continuing to disagree with them.

Thus, in order for these kinds of proposals to work, they need to incorporate strong enforcement mechanisms to keep non-safe AGIs in line and to prevent them from acquiring significant power. This requires an ability to create value-conforming AGIs in the first place, to implement the enforcement. Even a soft takeoff would eventually lead to AGIs wielding great power, so the enforcement could not be left to just humans or narrow AIs<sup>23</sup>. In practice, this means that integration proposals must be combined with some proposal for internal constraints which is capable of reliably creating value-conforming AGIs. Integration proposals also require there to be a soft takeoff in order to work, as having a small group of AGIs which rapidly acquired enough power to take control of the world would prevent any gradual integration schemes from working.

Therefore, because any effective integration strategy would require creating safe AGIs, and the right safe AGI design could lead to a positive outcome even if there were a hard takeoff, we believe that it is currently better to focus on proposals which are aimed at furthering the creation of safe AGIs.

### 3.3. Regulate research

Integrating AGIs into society may require explicit regulation. Calls for regulation are often agnostic about long-term outcomes but nonetheless recommend caution as a reasonable approach. For example, Hibbard [147] calls for international regulation to ensure that AGIs will value the long-term well-being of humans, but does not go into much detail. Daley [79] calls for a government panel for AGI issues. Hughes [157] argues that AGI should be regulated using the same mechanisms as previous technologies, creating state agencies responsible for the task and fostering global cooperation in the regulation effort<sup>24</sup>.

Current mainstream academic opinion does not consider AGI a serious threat [156], so AGI regulation seems unlikely in the near future. On the other hand, many AI systems are becoming increasingly autonomous and a number of authors are arguing that even narrow-AI applications should be equipped with an understanding of ethics [282]. Currently there are calls to regulate AI in the form of high-frequency trading (HFT) [250] and AI applications that have a major impact on society might become increasingly regulated. At the same time, legislation has a well-known tendency to lag behind technology and regulating AI applications will probably not translate into regulating basic research into AGI.

<sup>23</sup> For proposals which suggest that humans could use technology to remain competitive with AGIs and thus prevent them from acquiring excessive amounts of power, see section 3.4.

<sup>24</sup> Some of these proposals are written in the context of the USA and refer to the actions that the US government should take, but the general logic behind the proposals is not US-specific.

**3.3.1. Review boards.** Yampolskiy [305] note that university research programs in the social and medical sciences are overseen by institutional review boards. They propose setting up analogous review boards to evaluate potential AGI research. Research that was found to be AGI related would be restricted with measures ranging from supervision and funding limits to partial or complete bans. At the same time, research focusing on safety measures would be encouraged.

Posner [p 221, 220] suggests the enactment of a law which would require scientific research projects in dangerous areas to be reviewed by a federal catastrophic risks assessment board and forbidden if the board found that the project would create an undue risk to human survival.

Wilson [301] makes possibly the most detailed AGI regulation proposal so far, recommending a new international treaty where a body of experts would determine whether there was a 'reasonable level of concern' about AGI or some other possibly dangerous research. States would be required to regulate research or even temporarily prohibit it once experts agreed upon there being such a level of concern. He also suggests a number of other safeguards built into the treaty, such as the creation of ethical oversight organizations for researchers, mechanisms for monitoring abuses of dangerous technologies and an oversight mechanism for scientific publications.

**3.3.2. Encourage research into safe AGI.** In contrast, McGinnis [183] argues that the government should not attempt to regulate AGI development. Rather, it should concentrate on providing funding for research projects intended to create safe AGI.

Goertzel [115] argue for an open-source approach to safe AGI development instead of regulation. Hibbard [149] has likewise suggested developing AGI via open-source methods, but not as an alternative to regulation.

Legg [172] proposes funding safe AGI research via an organization that takes a venture capitalist approach to funding research teams, backing promising groups and cutting funding to any teams that fail to make significant progress. The focus of the funding would be to make AGI as safe as possible.

**3.3.3. Differential technological progress.** Both review boards and government funding could be used to implement 'differential intellectual progress':

Differential intellectual progress consists in prioritizing risk-reducing intellectual progress over risk-increasing intellectual progress. As applied to AI risks in particular, a plan of differential intellectual progress would recommend that our progress on the scientific, philosophical and technological problems of AI safety outpace our progress on the problems of AI capability such that we develop safe superhuman AIs before we develop (arbitrary) superhuman AIs [200].

Examples of research questions that could constitute philosophical or scientific progress in safety can be found in later sections of this paper—for instance, the usefulness of different internal constraints on ensuring safe behavior, or ways of making AGIs reliably adopt human values as they learn what those values are like.

Earlier, Bostrom [47] used the term ‘differential technological progress’ to refer to differential intellectual progress in technological development. Bostrom defined differential technological progress as ‘trying to retard the implementation of dangerous technologies and accelerate implementation of beneficial technologies, especially those that ameliorate the hazards posed by other technologies’.

One issue with differential technological progress is that we do not know what kind of progress should be accelerated and what should be retarded. For example, a more advanced communication infrastructure could make AGIs more dangerous, as there would be more networked machines that could be accessed via the Internet. Alternatively, it could be that the world will already be so networked that AGIs will be a major threat anyway and further advances will make the networks more resilient to attack. Similarly, it can be argued that AGI development is dangerous for as long as we have yet to solve the philosophical problems related to safe AGI design and do not know which AGI architectures are safe to pursue [309]. But it can also be argued that we should invest in AGI development now, when the related tools and hardware are still primitive enough that progress will be slow and gradual [115].

**3.3.4. International mass surveillance.** For AGI regulation to work, it needs to be enacted on a global scale. This requires solving both the problem of effectively enforcing regulation within a country and the problem of getting many different nations to all agree on the need for regulation.

Shulman [241] discusses various factors influencing the difficulty of AGI arms control. He notes that AGI technology itself might make international cooperation more feasible. If narrow AIs and early-stage AGIs were used to analyze the information obtained from wide-scale mass surveillance and wiretapping, this might make it easier to ensure that nobody was developing more advanced AGI designs.

Shulman [242] similarly notes that machine intelligences could be used to enforce treaties between nations. They could also act as trustworthy inspectors which would be restricted to communicating only information about treaty violations, thus not endangering state secrets even if they were allowed unlimited access to them. This could help establish a ‘singleton’ [48] regulatory regimen capable of effectively enforcing international regulation, including AGI-related treaties. Goertzel [115] also discuss the possibility of having a network of AGIs monitoring the world in order to police other AGIs and to prevent any of them from suddenly obtaining excessive power.

Another proposal for international mass surveillance is to build an ‘AGI Nanny’ [111, 115], a proposal discussed in section 5.4.

Large-scale surveillance efforts are ethically problematic and face major political resistance, and it seems unlikely that current political opinion would support the creation of a far-reaching surveillance network for the sake of AGI risk alone. The extent to which such extremes would be necessary depends on exactly how easy it would be to develop AGI in secret. Although several authors make the point that AGI is much easier to develop unnoticed than something like nuclear weapons [183, 189], cutting-edge high-tech research does tend to require major investments which might plausibly be detected even by less elaborate surveillance efforts.

To the extent that surveillance does turn out to be necessary, there is already a strong trend toward a ‘surveillance society’ with increasing amounts of information about people being collected and recorded in various databases [9]. As a reaction to the increased surveillance, Mann *et al* [179] propose to counter it with *sousveillance*—giving private individuals the ability to document their life and subject the authorities to surveillance in order to protect civil liberties. This is similar to the proposals of Brin [57], who argues that technological progress might eventually lead to a ‘transparent society’, where we will need to redesign our societal institutions in a way that allows us to maintain some of our privacy despite omnipresent surveillance. Miller [189] notes that intelligence agencies are already making major investments in AI-assisted analysis of surveillance data.

If social and technological developments independently create an environment where large-scale surveillance or *sousveillance* is commonplace, it might be possible to take advantage of those developments in order to police AGI risk<sup>25</sup>. Walker [279] argues that in order for mass surveillance to become effective, it must be designed in such a way that it will not excessively violate people’s privacy, for otherwise the system will face widespread sabotage<sup>26</sup>. Even under such conditions, there is no clear way to define what counts as dangerous AGI. Goertzel [115] point out that there is no clear division between narrow AI and AGI and attempts to establish such criteria have failed. They argue that since AGI has a nebulous definition, obvious wide-ranging economic benefits and potentially significant penetration into multiple industry sectors, it is unlikely to be regulated due to speculative long-term risks.

AGI regulation requires global cooperation, as the non-cooperation of even a single nation might lead to catastrophe. Historically, achieving global cooperation on tasks such as

<sup>25</sup> An added benefit would be that this could also help avoid other kinds of existential risks, such as the intentional creation of dangerous new diseases.

<sup>26</sup> Walker also suggests that surveillance systems could be designed to automatically edit out privacy-endangering details (such as pictures of people) from the data that they transmit, while leaving in details which might help in revealing dangerous ploys (such as pictures of bombs). However, this seems impossible to implement effectively, as research has found ways to extract personally identifying information and details from a wide variety of supposedly anonymous datasets [66, 95, 116, 203, 204, 206, 263]. Narayanan [205] even go as far as to state that ‘the false dichotomy between personally identifiable and non-personally identifiable information should disappear from privacy policies, laws, etc. Any aspect of an individual’s ... personality can be used for de-anonymization, and this reality should be recognized by the relevant legislation and corporate privacy policies’.

nuclear disarmament and climate change has been very difficult. As with nuclear weapons, AGI could give an immense economic and military advantage to the country that develops it first, in which case limiting AGI research might even give other countries an incentive to develop AGI faster [65, 82, 183, 189].

To be effective, regulation also needs to enjoy support among those being regulated. If developers working in AGI-related fields only follow the letter of the law, while privately viewing all regulations as annoying hindrances, and fears about AGI as overblown, the regulations may prove ineffective. Thus, it might not be enough to convince governments of the need for regulation; the much larger group of people working in the appropriate fields may also need to be convinced.

While Shulman [241] argues that the unprecedentedly destabilizing effect of AGI could be a cause for world leaders to cooperate more than usual, the opposite argument can be made as well. Gubrud [124] argues that increased automation could make countries more self-reliant and international cooperation considerably more difficult. AGI technology is also much harder to detect than, for example, nuclear technology is—nuclear weapons require a substantial infrastructure to develop, while AGI needs much less [183, 189].

Miller [189] even suggests that the mere possibility of a rival being close to developing AGI might, if taken seriously, trigger a nuclear war. The nation that was losing the AGI race might think that being the first to develop AGI was sufficiently valuable that it was worth launching a first strike for, even if it would lose most of its own population in the retaliatory attack. He further argues that, although it would be in the interest of every nation to try to avoid such an outcome, the ease of secretly pursuing an AGI development program undetected, in violation of treaty, could cause most nations to violate the treaty.

Miller also points out that the potential for an AGI arms race exists not only between nations, but between corporations as well. He notes that the more AGI developers there are, the more likely it is that they will all take more risks, with each AGI developer reasoning that if they do not take this risk, somebody else might take that risk first.

Goertzel [115] suggest that for regulation to be enacted, there might need to be an ‘AGI Sputnik’—a technological achievement that makes the possibility of AGI evident to the public and policy makers. They note that after such a moment, it might not take very long for full human-level AGI to be developed, while the negotiations required to enact new kinds of arms control treaties would take considerably longer.

So far, the discussion has assumed that regulation would be carried out effectively and in the pursuit of humanity’s common interests, but actual legislation is strongly affected by lobbying and the desires of interest groups [210] Mueller:2003. Many established interest groups would have an economic interest in either furthering or retarding AGI development, rendering the success of regulation uncertain.

**3.3.5. ‘Regulate research’ proposals: our view.** Although there seem to be great difficulties involved with regulation, there also remains the fact that many technologies have been successfully subjected to international regulation. Even if one were skeptical about the chances of effective regulation, an AGI arms race seems to be one of the worst possible scenarios, one which should be avoided if at all possible. We are therefore generally supportive of regulation, though the most effective regulatory approach remains unclear.

### 3.4. Enhance human capabilities

While regulation approaches attempt to limit the kinds of AGIs that will be created, enhancement approaches attempt to give humanity and AGIs a level playing field. In principle, gains in AGI capability would not be a problem if humans could improve themselves to the same level.

Alternatively, human capabilities could be improved in order to obtain a more general capability to deal with difficult problems. Verdoux [275, 276] suggests that cognitive enhancement could help in transforming previously incomprehensible mysteries into tractable problems and Verdoux [275] in particular highlights the possibility of cognitive enhancement helping to deal with the problems posed by existential risks. One problem with such approaches is that increasing humanity’s capability for solving problems will also make it easier to develop dangerous technologies. It is possible that cognitive enhancement should be combined with moral enhancement, in order to help foster the kind of cooperation that would help avoid the risks of technology [215, 216].

Moravec [193, 196] proposes that humans could keep up with AGIs via ‘mind uploading’, a process of transferring the information in human brains to computer systems so that human minds could run on a computer substrate. This technology may arrive during a similar timeframe as AGI [69, 143, 165, 170, 233, 234, 254]. However, Moravec argues that mind uploading would come after AGIs, and that unless the uploaded minds (‘uploads’) would transform themselves to become radically non-human, they would be weaker and less competitive than AGIs that were native to a digital environment [194, 196]. For these reasons, Warwick [289] also expresses doubt about the usefulness of mind uploading<sup>27</sup>.

Kurzweil [170] posits an evolution that will start with brain–computer interfaces, then proceed to using brain-embedded nanobots to enhance our intelligence and finally lead to full uploading and radical intelligence enhancement. Koene [166] criticizes plans to create safe AGIs and considers uploading both a more feasible and a more reliable approach.

<sup>27</sup> Some uploading approaches also raise questions of personal identity, whether the upload would still be the same person as the original [38, 45, 72, 112, 142, 155, 193, 262, 280] and whether they would be conscious in the first place [11, 71, 142, 169, 239]. However, these concerns are not necessarily very relevant for AGI risk considerations, as a population of uploads working to protect against AGIs would be helpful even if they lacked consciousness or were different individuals than the originals.

Similar proposals have also been made without explicitly mentioning mind uploading. Cade [65] speculates on the option of gradually merging with machines by replacing body parts with mechanical components. Turney [270] proposes linking AGIs directly to human brains so that the two meld together into one entity and Warwick [289, 290] notes that cyborgization could be used to enhance humans.

Mind uploading might also be used to make human value systems more accessible and easy to learn for AGIs, such as by having an AGI extrapolate the upload's goals directly from its brain, with the upload providing feedback.

**3.4.1. Would we remain human?** Uploading might destroy parts of humanity that we value [82, 160]. De Garis [82] argues that a computer could have far more processing power than a human brain, making it pointless to merge computers and humans. The biological component of the resulting hybrid would be insignificant compared to the electronic component, creating a mind that was negligibly different from a 'pure' AGI. Kurzweil [168] makes the same argument, saying that although he supports intelligence enhancement by directly connecting brains and computers, this would only keep pace with AGIs for a couple of additional decades.

The truth of this claim seems to depend on exactly how human brains are augmented. In principle, it seems possible to create a prosthetic extension of a human brain that uses the same basic architecture as the original brain and gradually integrates with it [254]. A human extending their intelligence using such a method might remain roughly human-like and maintain their original values. However, it could also be possible to connect brains with computer programs that are very unlike human brains and which would substantially change the way the original brain worked. Even smaller differences could conceivably lead to the adoption of 'cyborg values' distinct from ordinary human values [290].

Bostrom [49] speculates that humans might outsource many of their skills to non-conscious external modules and would cease to experience anything as a result. The value-altering modules would provide substantial advantages to their users, to the point that they could outcompete uploaded minds who did not adopt the modules.

Uploading would also allow humans to make precise and deep modifications to their own minds, which carries with it dangers of a previously unencountered kind [259].

**3.4.2. Would evolutionary pressures change us?** A willingness to integrate value-altering modules is not the only way by which a population of uploads might come to have very different values from modern-day humans. This is not necessarily a bad, or even a very novel, development: the values of earlier generations have often been different from the values of later generations [138] and it might not be a problem if a civilization of uploads enjoyed very different things than a civilization of humans. Still, as there are possible outcomes that we would consider catastrophic, such as the loss of nearly all things that have intrinsic value for us

[49], it is worth reviewing some of the postulated changes in values.

For comprehensiveness, we will summarize all of the suggested effects that uploading might have on human values, even if they are not obviously negative. Readers may decide for themselves whether or not they consider any of these effects to be causes for concern.

Hanson [133] argues that employers will want to copy uploads who are good workers and that at least some uploads will consent to being copied in such a manner. He suggests that the resulting evolutionary dynamics would lead to an accelerated evolution of values. This would cause most of the upload population to evolve to be indifferent or favorable to the thought of being copied, to be indifferent toward being deleted as long as another copy of themselves remained and to be relatively uninterested in having children 'the traditional way' (as opposed to copying an already-existing mind). Although Hanson's analysis uses the example of a worker-employer relationship, it should be noted that nations or families, or even single individuals, could also gain a competitive advantage by copying themselves, thus contributing to the strength of the evolutionary dynamic.

Similarly, Bostrom [49] writes that much of human life's meaning depends on the enjoyment of things ranging from humor and love to literature and parenting. These capabilities were adaptive in our past, but in an upload environment they might cease to be such and gradually disappear entirely.

Shulman [242] likewise considers uploading-related evolutionary dynamics. He notes that there might be a strong pressure for uploads to make copies of themselves in such a way that individual copies would be ready to sacrifice themselves to aid the rest. This would favor a willingness to copy oneself and a view of personal identity which did not consider the loss of a single copy to be death. Beings taking this point of view could then take advantage of the economic benefits of continually creating and deleting vast numbers of minds depending on the conditions, favoring the existence of a large number of short-lived copies over a somewhat less efficient world of long-lived minds.

Finally, Sotala [254] consider the possibility of minds coalescing via artificial connections that linked several brains together in the same fashion as the two hemispheres of ordinary brains are linked together. If this were to happen, considerable benefits might accrue to those who were ready to coalesce with other minds. The ability to copy and share memories between minds might also blur distinctions between individual minds. In the end, most humans might cease to be individual, distinct people in any real sense of the word.

It has also been proposed that information security concerns could cause undesirable dynamics among uploads, with significant advantages accruing to those who could steal the computational resources of others and use them to create new copies of themselves. If one could seize control of the hardware that an upload was running on, it could be immediately replaced with a copy of a mind loyal to the attacker. It might even be possible to do this without being

detected, if it was possible to steal enough of an upload's personal information to impersonate it.

An attack targeting a critical vulnerability in some commonly used piece of software might quickly hit a very large number of victims. As previously discussed in section 2.3.1, both theoretical arguments and actual cases of malware show that large numbers of machines on the Internet could be infected in a very short time [191, 192, 257]. In a society of uploads, attacks such as these would be not only inconvenient, but potentially fatal. Eckersley [93] offer a preliminary analysis of information security in a world of uploads.

**3.4.3. Would uploading help?** Even if the potential changes of values were deemed acceptable, it is unclear whether the technology for uploading could be developed before developing AGI. Uploading might require emulating the low-level details of a human brain with a high degree of precision, requiring large amounts of computing power [69, 234]. In contrast, an AGI might be designed around high-level principles which have been chosen to be computationally cheap to implement on existing hardware architectures.

Yudkowsky [309] uses the analogy that it is much easier to figure out the principles of aerodynamic flight and then build a Boeing 747 than it is to take a living bird and 'upgrade' it into a giant bird that can carry passengers, all while ensuring that the bird remains alive and healthy throughout the process. Likewise, it may be much easier to figure out the basic principles of intelligence and build AGIs than to upload existing minds.

On the other hand, one can also construct an analogy suggesting that it is easier to copy a thing's function than it is to understand how it works. If a person does not understand architecture but wants to build a sturdy house, it may be easier to create a replica of an existing house than it is to design an entirely new one that does not collapse.

Even if uploads were created first, they might not be able to harness all the advantages of digitality, as many of these advantages depend on minds being easy to modify [253], which human minds may not be. Uploads will be able to directly edit their source code as well as introduce simulated pharmaceutical and other interventions, and they could experiment on copies that are restored to an unmodified state if the modifications turn out to be unworkable [242]. Regardless, human brains did not evolve to be easy to modify and it may be difficult to find a workable set of modifications which would drastically improve them.

In contrast, in order for an AGI programmed using traditional means to be manageable as a software project, it must be easy for the engineers to modify it<sup>28</sup>. Thus, even if uploading were developed before AGI, AGIs that were developed later might still be capable of becoming more powerful than uploads. However, existing uploads already enjoying some of the advantages of the newly created AGIs

would still make it easier for the uploads to control the AGIs, at least for a while.

Moravec [194] notes that the human mind has evolved to function in an environment which is drastically different from a purely digital environment and that the only way to remain competitive with AGIs would be to transform into something that was very different from a human. This suggests that uploading might buy time for other approaches, but would be only a short-term solution in and of itself.

If uploading technology were developed before AGI, it could be used to upload a research team or other group and run them at a vastly accelerated rate as compared to the rest of humanity. This would give them a considerable amount of extra time for developing any of the other approaches. If this group were among the first to be successfully emulated and sped up, and if the speed-up would allow enough subjective time to pass before anyone else could implement their own version, they might also be able to avoid trading safety for speed. However, such a group might be able to wield tremendous power, so they would need to be extremely reliable and trustworthy.

**3.4.4. 'Enhance human capabilities' proposals: our view.** Of the various 'enhance human capabilities' approaches, uploading proposals seem the most promising, as translating a human brain to a computer program would sidestep many of the constraints that come from modifying a physical system. For example, all relevant brain activity could be recorded for further analysis at an arbitrary level of detail and any part of the brain could be instantly modified without a need for time-consuming and possibly dangerous invasive surgery. Uploaded brains could also be more easily upgraded to take full advantage of more powerful hardware, while humans whose brains were still partially biological would be bottlenecked by the speed of the biological component.

Uploading does have several problems: uploading research might lead to AGI being created before the uploads, in the long term uploads might have unfavorable evolutionary dynamics and it seems likely that there will eventually be AGIs which are capable of outperforming uploads in every field of competence. Uploads could also be untrustworthy even without evolutionary dynamics. At the same time, however, uploading research does not *necessarily* accelerate AGI research very much, the evolutionary dynamics might not be as bad as they seem at the moment and the advantages gained from uploading might be enough to help control unsafe AGIs until safe ones could be developed. Methods could also be developed for increasing the trustworthiness of uploads [242]. Uploading might still turn out to be a useful tool for handling AGI risk.

### 3.5. Relinquish technology

Not everyone believes that the risks involved in creating AGIs are acceptable. *Relinquishment* involves the abandonment of technological development that could lead to AGI. This is possibly the earliest proposed approach, with Butler [64] writing that 'war to the death should be instantly proclaimed'

<sup>28</sup> However, this might not be true for AGIs created using some alternative means, such as artificial life [260].

upon machines, for otherwise they would end up destroying humans entirely. In a much-discussed article, Joy [160] suggests that it might be necessary to relinquish at least some aspects of AGI research, as well as nanotechnology and genetics research.

Hughes [157] criticizes AGI relinquishment, while Kurzweil [170] criticizes broad relinquishment but supports the possibility of ‘fine-grained relinquishment’, banning some dangerous aspects of technologies while allowing general work on them to proceed. In general, most writers reject proposals for broad relinquishment.

**3.5.1. Outlaw AGI.** Weng *et al* [297] write that the creation of AGIs would force society to shift from human-centric values to robot-human dual values. In order to avoid this, they consider the possibility of banning AGI. This could be done either permanently or until appropriate solutions are developed for mediating such a conflict of values.

McKibben [184], writing mainly in the context of genetic engineering, suggests that AGI research should be stopped. He brings up the historical examples of China renouncing seafaring in the 1400 s and Japan relinquishing firearms in the 1600 s, as well as the more recent decisions to abandon DDT, CFCs and genetically modified crops in Western countries. However, it should also be noted that Japan participated in World War II; that China now has a navy; that there are reasonable alternatives for DDT and CFCs, which probably do not exist for AGI; and that genetically modified crops are in wide use in the United States.

Hughes [157] argues that attempts to outlaw a technology will only make the technology move to other countries. He also considers the historical relinquishment of biological weapons to be a bad example, for no country has relinquished peaceful biotechnological research such as the development of vaccines, nor would it be desirable to do so. With AGI, there would be no clear dividing line between safe and dangerous research.

De Garis [82] believes that differences of opinion about whether to build AGIs will eventually lead to armed conflict, to the point of open warfare. Annas *et al* [24] have similarly argued that genetic engineering of humans would eventually lead to war between unmodified humans and the engineered ‘posthumans’, and that cloning and inheritable modifications should therefore be banned. To the extent that one accepts their reasoning with regard to humans, it could also be interpreted to apply to AGIs.

**3.5.2. Restrict hardware.** Berglas [44] suggests not only stopping AGI research, but also outlawing the production of more powerful hardware. Berglas holds that it will be possible to build computers as powerful as human brains in the very near future and that we should therefore reduce the power of new processors and destroy existing ones<sup>29</sup>. Branwen [56] argues that Moore’s law depends on the existence of a small

number of expensive and centralized chip factories, making them easy targets for regulation and incapable of being developed in secret.

### 3.5.3. ‘Relinquish technology’ proposals: our view.

Relinquishment proposals suffer from many of the same problems as regulation proposals, but to a greater extent. There is no historical precedent of general, multi-use technology similar to AGI being successfully relinquished for good, nor do there seem to be any theoretical reasons for believing that relinquishment proposals would work in the future. Therefore we do not consider them to be a viable class of proposals.

## 4. External AGI constraints

Societal approaches involving regulation or research into safe AGI assume that proper AGI design can produce solutions to AGI risks. One category of such solutions is that of *external constraints*. These are restrictions that are imposed on an AGI from the outside and aim to limit its ability to do damage.

Several authors have argued that external constraints are unlikely to work with AGIs that are genuinely far more intelligent than us [30, 72, 170, 278, 307, 309]. The consensus seems to be that external constraints might buy time when dealing with less advanced AGIs, but they are useless against truly superintelligent ones.

External constraints also limit the usefulness of an AGI, as a free-acting one could serve its creators more effectively. This reduces the probability of the universal implementation of external constraints on AGIs. AGIs might also be dangerous if they were confined or otherwise restricted. For further discussion of these points, see section 5.1.

### 4.1. AGI confinement

AGI confinement, or ‘AI boxing’ [30, 72, 92, 278, 303], involves confining an AGI to a specific environment and limiting its access to the external world.

Yampolskiy [303] makes an attempt to formalize the idea, drawing on previous computer security research on the so-called confinement problem [171]. Yampolskiy defines the *AI confinement problem* as the challenge of restricting an AGI to a confined environment from which it cannot communicate without authorization. A number of methods have been proposed for implementing AI confinement, many of which are extensively discussed in Armstrong, Sandberg and Bostrom’s [30] paper.

Chalmers [72] and Armstrong *et al* [30] mention numerous caveats and difficulties with AI-boxing approaches. A *truly* leakproof system that perfectly isolated the AGI from an outside environment would prevent us from even observing the AGI. If AGIs were given knowledge about human behavior or psychology, they could still launch social engineering attacks on us [30, 72]. An AGI that was unaware of the existence of humans would be less likely to launch such

<sup>29</sup> Berglas (personal communication) has since changed his mind and no longer believes that it is possible to effectively restrict hardware or otherwise prevent AGI from being created.

attacks, but also much more limited in the kinds of tasks that it could be used for.

Even if the AGI remained confined, it could achieve enough influence among humans to prevent itself from being reset or otherwise modified [118]. An AGI that people grew reliant on might also become impossible to reset or modify.

**4.1.1. Safe questions.** Yampolskiy [303] proposes a hierarchy of security levels for communication with the AGI, ranging from completely unrestricted communication to no communication whatsoever. He suggests that the right level for communicating with AGIs is an intermediate one, based on asking ‘safe questions’. These are defined as questions for which a human could find an answer without help from a superintelligence. For instance, a researcher who was considering two possible cures for cancer could ask the AGI which one to pursue first. Yampolskiy argues that the AGI is effectively not changing the final outcome, as the researcher would in principle be capable of finding the better treatment himself. The AGI, by recommending one of the approaches as the choice to pursue first, is merely speeding up the process. Armstrong *et al* [30] also briefly discuss such approaches.

**4.1.2. Virtual worlds.** While the confinement protocol supposes an AGI that can perceive our world directly, Chalmers [72] uses the term ‘leakproof singularity’ to refer to the project of creating isolated virtual worlds which would disallow their inhabitants from influencing the external world. An AGI placed in such a world could be monitored by humans to ensure that it was safe to release into our world. However, Chalmers notes that ‘for an AI system to be useful or interesting to us at all, it must have some effects on us. At a minimum, we must be able to observe it. And the moment we observe a virtual environment, some information leaks out from that environment into our environment and affects us’.

Hibbard [151] proposes a decision support system for evaluating AGI designs in simulated environments that model our real environment, and attempts to show that the AGI can be built to be safe.

Armstrong *et al* [30] discuss the notion that we could supplement this approach with several layers of ‘sub-worlds’: virtual worlds stacked within each other. In such a scenario, even if an AGI somehow manages to break out from one virtual reality, it is still restricted by another and cannot directly influence the world. Ultimately, the authors conclude that this approach is not very promising.

**4.1.3. Resetting the AGI.** For as long as an AGI remains confined or cooperative, its memory could be altered from the outside. Turney [270] mentions a proposal in which an AGI’s memory could be wiped after each job, thus limiting its potential. Armstrong *et al* [30] also mention the possibility of resetting an AGI after each run, though they note that there are many ways for the AGI to get around it and that it is only useful for certain kinds of AGI systems.

Good [118] warns that a superhuman AGI may be smart enough to find ways to prevent itself from being successfully reset or shut down, including methods of social manipulation. Warwick [289] makes a similar point.

**4.1.4. Checks and balances.** Armstrong *et al* [30] suggest a number of checks and balances that could be used to catch the early stages of an AGI that was behaving dangerously or trying to deceive us. Test questions with known answers could be mixed in with questions we genuinely did not know the answer to, in order to establish an AGI’s trustworthiness. ‘Honey-pots’, tempting resources with great value to the AGI, could also be left in physical or mental locations that it was forbidden from accessing. Finally, several different kinds of AGIs could be presented with the same problem, with their answers then being compared. Inconsistent answers would suggest a possible attempt at deception.

**4.1.5. ‘AI confinement’ proposals: our view.** Despite their limited ability to deal with AGIs more intelligent than us, AI-boxing techniques seem to have value as a first line of defense and it may be worthwhile to invest in developing off-the-shelf software packages for AI confinement that are easy and cheap to use. A research project that developed AGI unexpectedly might not have been motivated to make major investments in security, but the AGI might still be sufficiently limited in intelligence that confinement would work. Having a defense that is easy to deploy will make it more likely that these kinds of projects will implement better precautions.

However, at the same time there is a risk that this will promote a false sense of security and make research teams think that they have carried out their duty to be cautious merely because they are running elementary confinement protocols. Although some confinement procedures can be implemented on top of an AGI that was not expressly designed for confinement, they are much less reliable than with an AGI that was built with confinement considerations in mind [30]—and even then, relying solely on confinement is a risky strategy. We are therefore somewhat cautious in our recommendation to develop confinement techniques.

## 4.2. AGI enforcement

One problem with AI confinement proposals is that humans are tasked with guarding machines that may be far more intelligent than themselves [118]. One proposed solution for this problem is to give the task of watching AGIs to other AGIs.

Armstrong [27] proposes that the trustworthiness of a superintelligent system could be monitored via a chain of less powerful systems, all the way down to humans. Although humans could not verify and understand the workings of a superintelligence, they could verify and understand an AGI just slightly above their own level, which could in turn verify and understand an AGI somewhat above its own level, and so on.

Chaining multiple levels of AI systems with progressively greater capacity seems to be replacing the problem of

building a safe AI with a multi-system, and possibly more difficult, version of the same problem. Armstrong himself admits that there are several problems with the proposal. There could be numerous issues along the line, such as a break in the chain of communication or an inability of a system to accurately assess the mind of another (smarter) system. There is also the problem of creating a trusted bottom for the chain in the first place, which is not necessarily any easier than creating a trustworthy superintelligence.

Hall [128] writes that there will be a great variety of AGIs, with those that were designed to be moral or aligned with human interests keeping the non-safe ones in check. Goertzel [115] also propose that we build a community of mutually policing AGI systems of roughly equal levels of intelligence. If an AGI started to ‘go off the rails’, the other AGIs could stop it. This might not prevent a single AGI from undergoing an intelligence explosion, but a community of AGIs might be in a better position to detect and stop it than humans would.

Having AGIs police each other is only useful if the group of AGIs actually has goals and values that are compatible with human goals and values. To this end, appropriate internal constraints are needed.

The proposal of a society of mutually policing AGIs would avoid the problem of trying to control a more intelligent mind. If a global network of mildly superintelligent AGIs could be instituted in such a manner, it might detect and prevent any nascent takeoff. However, by itself such an approach is not enough to ensure safety—it helps guard against individual AGIs ‘going off the rails’, but it does not help in a scenario where the programming of *most* AGIs is flawed and leads to non-safe behavior. It thus needs to be combined with the appropriate internal constraints.

A complication is that a hard takeoff is a *relative* term—an event that happens too fast for any outside observer to stop. Even if the AGI network were a hundred times more intelligent than a network composed of only humans, there might still be a more sophisticated AGI that could overcome the network.

**4.2.1. ‘AGI enforcement’ proposals: our view.** AGI enforcement proposals are in many respects similar to social integration proposals (section 3.2), in that they depend on the AGIs being part of a society which is strong enough to stop any single AGI from misbehaving. The greatest challenge is then to make sure that most of the AGIs in the overall system are safe and do not unite against humans rather than against misbehaving AGIs. Also, there might not be a natural distinction between a distributed AGI and a collection of many different AGIs and AGI design is in any case likely to make heavy use of earlier AI/AGI techniques. AGI enforcement proposals therefore seem like implementation variants of various internal constraint proposals (section 5), rather than independent proposals.

## 5. Internal constraints

In addition to external constraints, AGIs could be designed with internal motivations designed to ensure that they would take actions in a manner beneficial to humanity. Alternatively, AGIs could be built with internal constraints that make them easier to control via external means.

With regard to internal constraints, Yudkowsky distinguishes between *technical failure* and *philosophical failure* [309]:

Technical failure is when you try to build an AI and it does not work the way you think it does—you have failed to understand the true workings of your own code. Philosophical failure is trying to build the wrong thing, so that even if you succeeded you would still fail to help anyone or benefit humanity. Needless to say, the two failures are not mutually exclusive.

In practice, it is not always easy to distinguish between the two. Most of the discussion below focuses on the philosophical problems of various proposals, but some of the issues, such as whether or not a proposal is actually possible to implement, are technical.

### 5.1. Oracle AI

An *Oracle AI* is a hypothetical AGI that executes no actions other than answering questions. This might not be as safe as it sounds, however. Correctly defining ‘take no actions’ might prove surprisingly tricky [30] and the oracle could give flawed advice even if it did correctly restrict its actions.

Some possible examples of flawed advice: as extra resources are useful for the fulfilment of nearly all goals [211, 212], the oracle may seek to obtain more resources—such as computing power—in order to answer questions more accurately. Its answers might then be biased toward furthering this goal, even if this temporarily reduces the accuracy of its answers, if it believes this to increase the accuracy of its answers in the long run. Another example is that if the oracle had the goal of answering as many questions as possible as fast as possible, it could attempt to manipulate humans into asking it questions that were maximally simple and easy to answer.

Holden Karnofsky has suggested that an Oracle AI could be safe if it was ‘just a calculator’, a system which only computed things that were asked of it, taking no goal-directed actions of its own. Such a ‘tool-Oracle AI’ would keep humans as the ultimate decision makers. Furthermore, the first team to create a tool-Oracle AI could use it to become powerful enough to prevent the creation of other AGIs [162, 163].

An example of a tool-Oracle AI approach might be Omohundro’s [213] proposal of ‘safe-AI scaffolding’: creating highly constrained AGI systems which act within limited, predetermined parameters. These could be used to develop

formal verification methods and solve problems related to the design of more intelligent, but still safe, AGI systems.

*5.1.1. Oracles are likely to be released.* As with a boxed AGI, there are many factors that would tempt the owners of an Oracle AI to transform it to an autonomously acting agent. Such an AGI would be far more effective in furthering its goals, but also far more dangerous.

Current narrow-AI technology includes HFT algorithms, which make trading decisions within fractions of a second, far too fast to keep humans in the loop. HFT seeks to make a very short-term profit, but even traders looking for a longer-term investment benefit from being faster than their competitors. Market prices are also very effective at incorporating various sources of knowledge [135]. As a consequence, a trading algorithm's performance might be improved both by making it faster and by making it more capable of integrating various sources of knowledge. Most advances toward general AGI will likely be quickly taken advantage of in the financial markets, with little opportunity for a human to vet all the decisions. Oracle AIs are unlikely to remain as pure oracles for long.

Similarly, Wallach [283] discuss the topic of autonomous robotic weaponry and note that the US military is seeking to eventually transition to a state where the human operators of robot weapons are 'on the loop' rather than 'in the loop'. In other words, whereas a human was previously required to explicitly give the order before a robot was allowed to initiate possibly lethal activity, in the future humans are meant to merely supervise the robot's actions and interfere if something goes wrong.

Human Rights Watch [90] reports on a number of military systems which are becoming increasingly autonomous, with the human oversight for automatic weapons defense systems—designed to detect and shoot down incoming missiles and rockets—already being limited to accepting or overriding the computer's plan of action in a matter of seconds. Although these systems are better described as automatic, carrying out pre-programmed sequences of actions in a structured environment, than autonomous, they are a good demonstration of a situation where rapid decisions are needed and the extent of human oversight is limited. A number of militaries are considering the future use of more autonomous weapons.

In general, any broad domain involving high stakes, adversarial decision making and a need to act rapidly is likely to become increasingly dominated by autonomous systems. The extent to which the systems will need general intelligence will depend on the domain, but domains such as corporate management, fraud detection and warfare could plausibly make use of all the intelligence they can get. If one's opponents in the domain are also using increasingly autonomous AI/AGI, there will be an arms race where one might have little choice but to give increasing amounts of control to AI/AGI systems.

Miller [189] also points out that if a person was close to death, due to natural causes, being on the losing side of a war,

or any other reason, they might turn even a potentially dangerous AGI system free. This would be a rational course of action as long as they primarily valued their own survival and thought that even a small chance of the AGI saving their life was better than a near-certain death.

Some AGI designers might also choose to create less constrained and more free-acting AGIs for aesthetic or moral reasons, preferring advanced minds to have more freedom.

*5.1.2. Oracles will become authorities.* Even if humans were technically kept in the loop, they might not have the time, opportunity, motivation, intelligence, or confidence to verify the advice given by an Oracle AI. This may be a danger even with narrower AI systems. Friedman [102] discuss APACHE, an expert system that provides medical advice to doctors. They write that as the medical community puts more and more trust into APACHE, it may become common practice to act automatically on APACHE's recommendations and it may become increasingly difficult to challenge the 'authority' of the recommendations. Eventually, APACHE may in effect begin to dictate clinical decisions.

Likewise, Bostrom [53] point out that modern bureaucrats often follow established procedures to the letter, rather than exercising their own judgment and allowing themselves to be blamed for any mistakes that follow. Dutifully following all the recommendations of an AGI system would be an even better way of avoiding blame.

Wallach [283] note the existence of robots which attempt to automatically detect the locations of hostile snipers and to point them out to soldiers. To the extent that these soldiers have come to trust the robots, they could be seen as carrying out the robots' orders. Eventually, equipping the robot with its own weapons would merely dispense with the formality of needing to have a human to pull the trigger.

Thus, even AGI systems that function purely to provide advice will need to be explicitly designed to be safe in the sense of not providing advice that would go against human values [282]. Yudkowsky [313] further notes that an Oracle AI might choose a plan that is beyond human comprehension, in which case there is still a need to design it as explicitly safe and conforming to human values.

*5.1.3. 'Oracle AI' proposals: our view.* Much like with external constraints, it seems like Oracle AIs could be a useful stepping stone on the path toward safe, freely acting AGIs. However, because any Oracle AI can be relatively easily turned into a free-acting AGI and because many people will have an incentive to do so, Oracle AIs are not by themselves a solution to AGI risk, even if they are safer than free-acting AGIs when kept as pure oracles.

## 5.2. Top-down safe AGI

AGIs built to take autonomous actions will need to be designed with safe motivations. Wallach and Allen divide approaches for ensuring safe behavior into 'top-down' and 'bottom-up' approaches [14, 15, 281, 282, 284]. They define top-down approaches as ones that take a specified ethical

theory and attempt to build a system capable of implementing that theory [282].

Wallach and Allen [14, 15, 281, 282, 284] have expressed skepticism about the feasibility of both pure top-down and bottom-up approaches, arguing for a hybrid approach<sup>30</sup>. With regard to top-down approaches, which attempt to derive an internal architecture from a given ethical theory, Wallach [281] finds three problems:

1. 'Limitations already recognized by moral philosophers: For example, in a utilitarian calculation, how can consequences be calculated when information is limited and the effects of actions cascade in never-ending interactions? Which consequences should be factored into the maximization of utility? Is there a stopping procedure?' [281].
2. The 'frame problem' refers to the challenge of discerning relevant from irrelevant information without having to consider all of it, as all information could be relevant in principle [8, 85]. Moral decision-making involves a number of problems that are related to the frame problem, such as needing to know what effects different actions have on the world and needing to estimate whether one has sufficient information to accurately predict the consequences of the actions.
3. 'The need for background information. What mechanisms will the system require in order to acquire the information it needs to make its calculations? How does one ensure that this information is up to date in real time?' [281].

To some extent, these problems may be special cases of the fact that we do not yet have AGI with good general learning capabilities: creating an AGI would also require solving the frame problem, for instance. These problems might therefore not all be as challenging as they seem at first, presuming that we manage to develop AGI in the first place.

**5.2.1. Three laws.** Probably the most widely known proposal for machine ethics is Isaac Asimov's [33] three laws of robotics:

1. A robot may not injure a human being or, through inaction, allow a human being to come to harm.
2. A robot must obey orders given to it by human beings except where such orders would conflict with the first law.
3. A robot must protect its own existence as long as such protection does not conflict with either the first or second law.

Asimov and other writers later expanded the list to include a number of additional laws, including the zeroth law:

A robot may not harm humanity, or through inaction allow humanity to come to harm.

<sup>30</sup> For a definition of bottom-up approaches, see section 5.3.

Although the three laws are widely known and have inspired numerous imitations, several of Asimov's own stories were written to illustrate the fact that the laws contained numerous problems. They have also drawn heavy critique from others [23, 75, 76, 120, 180, 201, 224, 282, 295, 296] and are not considered a viable approach for safe AI. Among their chief shortcomings is the fact that they are too ambiguous to implement and, if defined with complete accuracy, contradict each other in many situations.

**5.2.2. Categorical imperative.** The best-known universal ethical axiom is Kant's categorical imperative. Many authors have discussed using the categorical imperative as the foundation of AGI morality [14, 40, 41, 63, 222, 256, 282]. All of these authors conclude that Kantian ethics is a problematic goal for AGI, though Powers [222] still remains hopeful about its prospects.

**5.2.3. Principle of voluntary joyous growth.** Goertzel [106, 107] considers a number of possible axioms before settling on what he calls the 'principle of voluntary joyous growth', defined as 'maximize happiness, growth and choice'. He starts by considering the axiom 'maximize happiness', but then finds this to be problematic and adds 'growth', which he defines as 'increase in the amount and complexity of patterns in the universe'. Finally he adds 'choice' in order to allow sentient beings to 'choose their own destiny'.

**5.2.4. Utilitarianism.** Classic utilitarianism is an ethical theory stating that people should take actions that lead to the greatest amount of happiness and the smallest amount of suffering. The prospects for AGIs implementing a utilitarian moral theory have been discussed by several authors [14, 19, 21, 77, 104, 119, 121, 199, 282]. The consensus is that pure classical utilitarianism is problematic and does not capture all human values. For example, a purely utilitarian AGI could reprogram the brains of humans so that they did nothing but experience the maximal amount of pleasure all the time and that prospect seems unsatisfactory to many<sup>31</sup>.

**5.2.5. Value learning.** Freeman [101] describes a decision-making algorithm which observes people's behavior, infers their preferences in the form of a utility function and then attempts to carry out actions which fulfil everyone's preferences. Similarly, Dewey [88] discusses *value learners*, AGIs which are provided a probability distribution over possible utility functions that humans may have. Value learners then attempt to find the utility functions with the best match for human preferences. Hibbard [153] builds on Dewey's work to offer a similar proposal.

<sup>31</sup> Note that utilitarianism is not the same thing as having a utility function. Utilitarianism is a specific kind of ethical system, while utility functions are general-purpose mechanisms for choosing between actions and can in principle be used to implement very different kinds of ethical systems, such as egoism and possibly even rights-based theories and virtue ethics [217].

One problem with conceptualizing human desires as utility functions is that human desires change over time [272] and also violate the axioms of utility theory required to construct a coherent utility function [271]. While it is possible to treat inconsistent choices as random deviations from an underlying ‘true’ utility function that is then learned [207], this does not seem to properly describe human preferences. Rather, human decision making and preferences seem to be driven by multiple competing systems, only some of which resemble utility functions [81]. Even if a true utility function could be constructed, it does not take into account the fact that we have second-order preferences, or desires about our desires: a drug addict may desire a drug, but also desire that he not desire it [98]. Similarly, we often wish that we had stronger desires toward behaviors which we consider good but cannot make ourselves engage in. Taking second-order preferences into account leads to what philosophers call ‘ideal preference’ theories of value [55, 176, 225, 247, 249, 264, 314].

Taking this into account, it has been argued that we should aim to build AGIs which act according to humanity’s *extrapolated* values [199, 265, 308]. Yudkowsky proposes attempting to discover the ‘coherent extrapolated volition’ (CEV) of humanity, which he defines as [308]

our wish if we knew more, thought faster, were more the people we wished we were, had grown up farther together; where the extrapolation converges rather than diverges, where our wishes cohere rather than interfere; extrapolated as we wish that extrapolated, interpreted as we wish that interpreted.

CEV remains vaguely defined and has been criticized by several authors [109, 115, 148, 189, 294]. However, Tarleton [265] finds CEV a promising approach and suggests that CEV has five desirable properties, and that many different kinds of algorithms could possess these features:

*Meta-algorithm:* Most of the AGI’s goals will be obtained at runtime from human minds, rather than explicitly programmed in before runtime.

*Factually correct beliefs:* The AGI will attempt to obtain correct answers to various factual questions, in order to modify preferences or desires that are based upon false factual beliefs.

*Singleton:* Only one superintelligent AGI is to be constructed and it is to take control of the world with whatever goal function is decided upon.

*Reflection:* Individual or group preferences are reflected upon and revised.

*Preference aggregation:* The set of preferences of a whole group are to be combined somehow.

At least two CEV variants have been proposed: coherent aggregated volition [109] and coherent blended volition [115]. Goertzel [115] describe a methodology which was used to help end the apartheid in South Africa. The methodology involves people with different views exploring different

future scenarios together and in great detail. By exploring different outcomes together, the participants build emotional bonds and mutual understanding, seeking an outcome that everyone can agree to live with. The authors characterize the coherent blended volition of a diverse group as analogous to the ‘conceptual blend’ resulting from the methodology, incorporating the most essential elements of the group into a harmonious whole.

Christiano [73] attempts to sketch out a formalization of a value extrapolation approach called ‘indirect normativity’. It proposes a technique that would allow an AI to approximate the kinds of values a group of humans would settle on if they could spend an unbounded amount of time and resources considering the problem.

Other authors have begun preliminary work on simpler value learning systems, designed to automatically learn moral principles. Anderson *et al* [18, 20, 21] have built systems based around various moral duties and principles. As lists of duties do not in and of themselves specify what to do when they conflict, the systems let human experts judge how each conflict should be resolved and then attempt to learn general rules from the judgments. As put forth, however, this approach would have little ability to infer ethical rules which did not fit the framework of proposed duties. Improved computational models of ethical reasoning [18, 20, 21, 32, 185, 186], perhaps incorporating work from neuroscience and moral psychology [42, 199, 245], could help address this. Potapov [221] propose an approach by which an AGI could gradually learn the values of other agents as its understanding of the world improved.

A value extrapolation process seems difficult to specify exactly, as it requires building an AGI with programming that formally and rigorously defines human values. Even if it manages to avoid the first issue in Wallach’s [281] list (section 5.2), top-down value extrapolation may fall victim to other issues, such as computational tractability. One interpretation of CEV would seem to require modeling not only the values of everyone on Earth, but also the evolution of those values as the people in question interacted with each other, became more intelligent and more like their ideal selves, chose which of their values they wanted to preserve, etc. Even if the AGI could eventually obtain the enormous amount of computing power required to run this model, its behavior would need to be safe from the beginning, or it could end up doing vast damage to humanity before understanding what it was doing wrong.

Goertzel and Pitt’s [115] hybrid approach, in which AGIs cooperate with humans in order to discover the values humans wish to see implemented, seems more likely to avoid the issue of computational tractability. However, it will fail to work in a hard takeoff situation where AGIs take control before being taught the correct human values. Another issue with coherent blended volition is that schemes which require absolute consensus are unworkable with large groups, as anyone whose situation would be worsened by a change of events could block the consensus. A general issue with value extrapolation approaches is that there may be several valid

ways of defining a value extrapolation process, with no objective grounds for choosing one rather than another.

Goertzel [109] notes that in formal reasoning systems a set of initially inconsistent beliefs which the system attempts to resolve might arrive at something very different than the initial belief set, even if there existed a consistent belief set that was closer to the original set. He suggests that something similar might happen when attempting to make human values consistent, though whether this would be a bad thing is unclear.

**5.2.6. ‘Top-down safe AGI’ proposals: our view.** Of the various top-down proposals, value learning proposals seem to be the only ones which properly take into account the complexity of value thesis (section 2.2), as they attempt to specifically take into account considerations such as ‘would humanity have endorsed this course of action if it had known the consequences?’ Although there are many open questions concerning the computational tractability as well as the general feasibility of such approaches, they seem like some of the most important ones to work on.

### 5.3. Bottom-up and hybrid safe AGI

Wallach [281] defines bottom-up approaches as ones that favor evolving or simulating the mechanisms that give rise to our moral decisions. Another alternative is hybrid approaches, combining parts of both top-down and bottom-up approaches.

A problem with pure bottom-up approaches is that techniques such as artificial evolution or merely rewarding the AGI for the right behavior may cause it to behave correctly in tests, but would not guarantee that it would behave safely in any other situation. Even if an AGI *seems* to have adopted human values, the actual processes driving its behavior may be very different from the processes that would be driving the actions of a human who behaved similarly. It might then behave very unexpectedly in situations which are different enough [152, 229, 309, 312].

Armstrong, Sandberg and Bostrom discuss various problems related to such approaches and offer examples of concepts which seem straightforward to humans but are not as simple as they may seem on the surface. One of their examples relates to the concept of time [30]:

If the [AGI] had the reasonable-sounding moral premise that ‘painlessly killing a human being, who is going to die in a micro-second anyway, in order to gain some other good, is not a crime’, we would not want it to be able to redefine millennia as seconds.

All humans have an intuitive understanding of time and no experience with beings who could arbitrarily redefine their own clocks and might not share the same concept of time. Such differences in the conceptual grounding of an AGI’s values and of human values might not become apparent until too late.

**5.3.1. Evolutionary invariants.** Human morality is to a large extent shaped by evolution [83, 161] and evolutionary approaches attempt to replicate this process with AGIs.

Hall [128, 131] argues that self-improving AGIs are likely to exist in competition with many other kinds of self-improving AGIs. Properties that give AGIs a significant disadvantage might then be strongly selected against and disappear. We could attempt to identify *evolutionary invariants*, or evolutionarily stable strategies, which would both survive in a competitive environment and cause an AGI to treat humans well.

Hall [131] lists self-interest, long planning horizons, knowledge, an understanding of evolutionary ethics and guaranteed honesty as invariants that are likely to make an AGI more moral as well as to persist even under intense competition. He suggests that, although self-interest may sound like a bad thing in an AGI, non-self-interested creatures are difficult to punish. Thus, enlightened self-interest might be a good thing for an AGI to possess, as it will provide an outside community with both a stick and a carrot to control it with.

Similarly, Waser [292] suggests that minds which are intelligent enough will, due to game-theoretical and other considerations, become altruistic and cooperative. Waser [294] proposed the principle of rational universal benevolence (RUB), the idea that the moral course of action is cooperation while letting everyone freely pursue their own goals. Waser proposes that, instead of making human-friendly behavior an AGI’s only goal, the AGI would be allowed to have and form its own goals. However, its goals and actions would be subject to the constraint that they should respect the principle of RUB and not force others into a life those others would disagree with.

Kornai [167] cites Gewirth’s [103] work on the principle of generic consistency, which holds that respecting others’ rights to freedom and well-being is a logically necessary conclusion for any rational agents. Kornai suggests that if the principle is correct, then AGIs would respect humanity’s rights to freedom and well-being, and that AGIs which failed to respect the principle would be outcompeted by ones which did.

Something similar was also proposed by Gips [104] and Versenyi [277], who advocates the creation of ‘wise robots’ that would recognize the extent to which their own well-being depended on cooperation with humans and would act accordingly.

Although these approaches expect AGI either to evolve altruism or to find it the most rational approach, true altruism or even pure tit-for-tat [34] is not actually the best strategy in evolutionary terms. Rather, a better strategy is *Machiavellian* tit-for-tat: cultivating an appearance of altruism and cooperation when it benefits oneself and acting selfishly when one can get away with it. Humans seem strongly disposed toward such behavior [127].

Another problem is that tit-for-tat as a good strategy assumes that both players are equally powerful and both have the same options at their disposal. If the AGI became far more powerful than most humans, it might no longer be in its

interests to treat humans favorably [97]. This hypothesis can be tested by looking at human behavior: if exploiting the weak is an evolutionarily useful strategy, then humans should engage in it when given the opportunity. Humans who feel powerful do indeed devalue the worth of the less powerful and view them as objects of manipulation [164]. They also tend to ignore social norms [274] and to experience less distress and compassion toward the suffering of others [273].

Thus, even if an AGI cooperated with other similarly powerful AGIs, the group of AGIs might still decide to collectively exploit humans. Similarly, even though there might be pressure for AGIs to make themselves more transparent and easily inspected by others, this only persists for as long as the AGI needs others more than the others need the AGI.

**5.3.2. Evolved morality.** Another proposal is to create AGIs via algorithmic evolution, selecting in each generation the AGIs which are not only the most intelligent, but also the most moral. These ideas are discussed to some extent by Wallach [282].

**5.3.3. Reinforcement learning.** In machine learning, *reinforcement learning* is a model in which an agent takes various actions and is differentially rewarded for the actions, after which it learns to perform the actions with the greatest expected reward. In psychology, it refers to agents being rewarded for certain actions and thus learning behaviors which they have a hard time breaking, even if some other kind of behavior is more beneficial for them later on.

Applied to AGI, the machine learning sense of reinforcement involves teaching an AGI to behave in a safe manner by rewarding it for ethical choices and letting it learn for itself the underlying rules of what constitutes ethical behavior. In an early example of this kind of proposal, McCulloch [181] described an ‘ethical machine’ that could infer the rules of chess by playing the game and suggested that it could also learn ethical behavior this way.

Hibbard [146, 148] suggested using reinforcement learning to give AGIs positive emotions toward humans. Early AGIs would be taught to recognize happiness and unhappiness in humans, and the results of this learning would be hard-wired as emotional values in more advanced AGIs. This training process would then be continued—for example, by letting the AGIs predict how human happiness would be affected by various actions and using those predictions as emotional values.

A reinforcement learner is supplied with a reward signal and it always has the explicit goal of maximizing the sum of this reward, any way it can. In order for this goal to align with human values, humans must engineer the environment so that the reinforcement learner is prevented from receiving rewards if human goals are not fulfilled [88]. A reinforcement-learning AGI only remains safe for as long as humans are capable of enforcing this limitation and will become unpredictable if it becomes capable of overcoming it. Hibbard [152] has retracted his earlier reinforcement learning-based proposals,

as they would allow the AGI to maximize its reinforcement by modifying humans to be maximally happy, even against their will [88].

**5.3.4. Human-like AGI.** Another kind of proposal involves building AGIs that can learn human values by virtue of being similar to humans.

Connectionist systems, based on artificial neural nets, are capable of learning patterns from data without being told what the patterns are. As some connectionist models have learned to classify problems in a manner similar to humans [187, 219, 267], it has been proposed that connectionist AGI might learn moral principles that are too complex for humans to specify as explicit rules<sup>32</sup>. This idea has been explored by Guarini [123] and Wallach [282].

One specific proposal that draws upon connectionism is to make AGIs act according to virtue ethics [14, 281, 282, 284]. These authors note that previous writers discussing virtuous behavior have emphasized the importance of learning moral virtues through habit and practice. As it is impossible to exhaustively define a virtue, virtue ethics has traditionally required each individual to learn the right behaviors through ‘bottom-up processes of discovery or learning’ [282]. Models that mimicked the human learning process well enough could then potentially learn the same behaviors as humans do.

Another kind of human-inspired proposal is the suggestion that something like Stan Franklin’s LIDA architecture [99, 227, 248], or some other approach based on Bernard Baars’s [35, 36] ‘global workspace’ theory, might enable moral reasoning. In the LIDA architecture, incoming information is monitored by specialized *attention codelets*, each of which searches the input for specific features. In particular, *moral codelets* might look for morally relevant factors and ally themselves with other codelets to promote their concerns to the level of conscious attention. Ultimately, some coalitions will win enough support to accomplish a specific kind of decision [282, 285, 286].

Goertzel [115] consider human memory systems (episodic, sensorimotor, declarative, procedural, attentional and intentional) and ways by which human morality might be formed via their interaction. They briefly discuss the way that the OpenCog AGI system [113, 141] implements similar memory systems and how those systems could enable it to learn morality. Similarly, Goertzel [114] discuss the stages of moral development in humans and suggest ways by which they could be replicated in AGI systems, using the specific example of Novamente, a proprietary version of OpenCog.

Waser [293] also proposes building an AGI by studying the results of evolution and creating an implementation as close to the human model as possible.

Human-inspired AGI architectures would intuitively seem the most capable of learning human values, though what would be human-like enough remains an open question.

<sup>32</sup> But it should be noted that there are also promising non-connectionist approaches for modeling human classification behavior: see, e.g., Tenenbaum *et al* [266].

It is possible that even a relatively minor variation from the norm could cause an AGI to adopt values that most humans would consider undesirable. Getting the details right might require an extensive understanding of the human brain.

There are also humans who have drastically different ethics than the vast majority of humanity and argue for the desirability of outcomes such as the extinction of mankind [43, 89]. There remains the possibility that even AGIs which reasoned about ethics in a completely human-like manner would reach such conclusions.

Humans also have a long track record of abusing power, or of undergoing major behavioral changes due to relatively small injuries—the ‘safe *Homo sapiens*’ problem also remains unsolved. On the other hand, it seems plausible that human-like AGIs could be explicitly engineered to avoid such problems.

The easier that an AGI is to modify the more powerful it might become [253] and very close recreations of the human brain may turn out to be difficult to extensively modify and upgrade. Human-inspired safe AGIs might then end up outcompeted by AGIs which were easier to modify and which might or might not be safe.

Even if human-inspired architectures could be easily modified, the messiness of human cognitive architecture means that it might be difficult to ensure that their values remain beneficial during modification. For instance, in LIDA-like architectures, beneficial behavior will depend on the correct coalitions of morality codelets winning each time. If the system undergoes drastic changes, this can be very difficult if not impossible to ensure.

Most AGI builders will attempt to create a mind that displays considerable advantages over ordinary humans. Some such advantages might be best achieved by employing a very non-human architecture [194], so there will be reasons to build AGIs that are not as human-like. These could also end up outcompeting the human-like AGIs.

### 5.3.5. ‘Bottom-up and hybrid safe AGI’ proposals: our view.

We are generally very skeptical about pure bottom-up methods, as they only allow a very crude degree of control over an AGI’s goals, giving it a motivational system which can only be relied on to align with human values in the very specific environments that the AGI has been tested in. Evolutionary invariants seem incapable of preserving complexity of value and they might not even be capable of preserving human survival. Reinforcement learning, on the other hand, depends on the AGI being incapable of modifying the environment against the will of its human controllers. Therefore, none of these three approaches seems workable.

Human-like AGI might have some promise, depending on exactly how fragile human values were. If the AGI reasoning process could be made sufficiently human-like, there is the possibility that the AGI could remain relatively safe, though less safe than a well-executed value extrapolation-based AGI.

## 5.4. AGI Nanny

A more general proposal, which could be achieved by either top-down, bottom-up, or hybrid methods, is the proposal of an ‘AGI Nanny’ [111, 115]. This is an AGI that is somewhat more intelligent than humans and is designed to monitor Earth for various dangers, including more advanced AGI.

The AGI Nanny would be connected to powerful surveillance systems and would control a massive contingent of robots. It would help abolish problems such as disease, involuntary death and poverty, while preventing the development of technologies that could threaten humanity. The AGI Nanny would be designed not to rule humanity on a permanent basis, but to give us some breathing room and time to design more advanced AGIs. After some predetermined amount of time, it would cede control of the world to a more intelligent AGI.

Goertzel [115] briefly discuss some of the problems inherent in the AGI Nanny proposal. The AGI would have to come to power in an ethical way and might behave unpredictably despite our best efforts. It might also be easier to create a dramatically self-improving AGI than to create a more constrained AGI Nanny.

**5.4.1. ‘AGI Nanny’ proposals: our view.** Upon asserting control, the AGI Nanny would need to have precisely specified goals, so that it would stop other AGIs from taking control but would also not harm human interests. It is not clear to what extent defining these goals would be easier than defining the goals of a more free-acting AGI [200]. Overall, the AGI Nanny seems to have promise, but it is unclear whether it can be made to work.

## 5.5. Formal verification

Formal verification methods prove specific properties about various algorithms. If the complexity and fragility of value theses hold, it follows that safe AGI requires the ability to verify that proposed changes to the AGI will not alter its goals or values. If even a mild drift from an AGI’s original goals might lead to catastrophic consequences, then utmost care should be given to ensuring that the goals will not change inadvertently. This is particularly the case if there are no external feedback mechanisms which would correct the drift. Before modifying itself, an AGI could attempt to formally prove that the changes would not alter its existing goals and would therefore keep them intact even during extended self-modification [309]. Such proofs could be required before self-modification was allowed to occur and the system could also be required to prove that this verify-before-modification property itself would always be preserved during self-modification.

Formal verification is also an essential part of Omohundro’s [213] safe-AI scaffolding strategy, as noted in section 5.6.4.

Spears [255] combines machine learning and formal verification methods to ensure that AIs remain within the bounds of pre-specified constraints after having learned new

behaviors. She attempts to identify ‘safe’ machine learning operators, which are guaranteed to preserve the system’s constraints.

One AGI system built entirely around the concept of formal verification is the Gdel machine [236, 258]. It consists of a *solver*, which attempts to achieve the goals set for the machine, and a *searcher*, which has access to a set of axioms that completely describe the machine. The searcher may completely rewrite any part of the machine, provided that it can produce a formal proof showing that such a rewrite will further the system’s goals.

Goertzel [110] proposes goal-oriented learning meta-architecture (GOLEM), a meta-architecture that can be wrapped around a variety of different AGI systems. GOLEM will only implement changes that are predicted to be more effective at achieving the original goal of the system. Goertzel argues that GOLEM is likely to be both self-improving and steadfast: either it pursues the same goals it had at the start, or it stops acting altogether.

Unfortunately, formalizing the AGI’s goals in a manner that will allow formal verification methods to be used is a challenging task. Within cryptography, many communications protocols have been proven secure, only for successful attacks to be later developed against their various implementations. While the formal proofs were correct, they contained assumptions which did not accurately capture the way the protocols worked in practice [84]. Proven theorems are only as good as their assumptions, so formal verification requires good models of the AGI hardware and software.

**5.5.1. ‘Formal verification’ proposals: our view.** Compared to the relatively simple domain of cryptographic security, verifying things such as ‘does this kind of a change to the AGI’s code preserve its goal of respecting human values?’ seems like a much more open-ended and difficult task, one which might even prove impossible. Regardless, it is the only way of achieving high confidence that a system is safe, so it should at least be attempted.

## 5.6. Motivational weaknesses

Finally, there is a category of internal constraints that, while not making an AGI’s *values* safer, make it easier to control AGI via external constraints.

**5.6.1. High discount rates.** AGI systems could be given a high discount rate, making them value short-term goals and gains far more than long-term goals and gains [30, 243]. This would inhibit the AGI’s long-term planning, making it more predictable. However, an AGI can also reach long-term goals through a series of short-term goals [30]. Another possible problem is that it could cause the AGI to pursue goals which were harmful for humanity’s long-term future. Humanity may arguably be seen as already behaving in ways that imply an excessively high discount rate, such as by consuming finite natural resources without properly taking into account the well-being of future generations.

**5.6.2. Easily satiable goals.** Shulman [243] proposes designing AGIs in such a way that their goals are easy to satisfy. For example, an AGI could receive a near-maximum reward for simply continuing to receive an external reward signal, which could be cut if humans suspected misbehavior. The AGI would then prefer to cooperate with humans rather than trying to attack them, even if it was very sure of its chances of success<sup>33</sup>. Likewise, if the AGI could receive a maximal reward with a relatively small fraction of humanity’s available resources, it would have little to gain from seizing more resources.

An extreme form of this kind of a deal is Orseau and Ring’s [214] ‘simpleton gambit’, in which an AGI is promised everything that it would ever want, on the condition that it turn itself into a harmless simpleton. Orseau and Ring consider several hypothetical AGI designs, many of which seem likely to accept the gambit, given certain assumptions.

In a related paper, Ring [229] consider the consequences of an AGI being able to modify itself to receive the maximum possible reward. They show that certain types of AGIs will then come to only care about their own survival. Hypothetically, humans could promise not to threaten such AGIs in exchange for them agreeing to be subject to AI-boxing procedures. For this to work, the system would have to believe that humans will take care of its survival against external threats better than it could itself. Hibbard [150, 153] discusses the kinds of AGIs that would avoid the behaviors described by Ring and Orseau [214, 229].

**5.6.3. Calculated indifference.** Another proposal is to make an AGI indifferent to a specific event. For instance, the AGI could be made indifferent to the detonation of explosives attached to its hardware, which might enable humans to have better control over it [28, 30].

**5.6.4. Programmed restrictions.** Goertzel [115] suggest we ought to ensure that an AGI does not self-improve too fast, because AGIs will be harder to control as they become more and more cognitively superior to humans. To limit the rate of self-improvement in AGIs, perhaps AGIs could be programmed to extensively consult humans and other AGI systems while improving themselves, in order to ensure that no unwanted modifications would be implemented.

Omohundro [213] discusses a number of programmed restrictions in the form of constraints on what the AGI is allowed to do, with formal proofs being used to ensure that an AGI will not violate its safety constraints. Such limited AGI systems would be used to design more sophisticated AGIs.

Programmed restrictions are problematic, as the AGI might treat these merely as problems to solve in the process of meeting its goals and attempt to overcome them [212]. Making an AGI not want to quickly self-improve might not solve the problem by itself. If the AGI ends up with a second-order desire to rid itself of such a disinclination, the stronger

<sup>33</sup> On the other hand, this might incentivize the AGI to deceive its controllers into believing it was behaving properly and also to actively hide any information which it even suspected might be interpreted as misbehavior.

desire will eventually prevail [259]. Even if the AGI wanted to maintain its disinclination toward rapid self-improvement, it might still try to circumvent the goal in some other way, such as by creating a copy of itself which did not have that disinclination [212]. Regardless, such limits could help control less sophisticated AGIs.

**5.6.5. Legal machine language.** Weng *et al* [296, 297] propose a ‘legal machine language’ which could be used to formally specify actions which the AGI is allowed or disallowed to do. Governments could then enact laws written in legal machine language, allowing them to be programmed into robots.

**5.6.6. ‘Motivational weaknesses’ proposals: our view.** Overall, motivational weaknesses seem comparable to external constraints: possibly useful and worth studying, but not something to rely on exclusively, particularly in the case of superintelligent AGIs. As with external constraints and Oracle AIs, an arms race situation might provide a considerable incentive to loosen or remove such constraints.

## 6. Conclusion

We began this paper by noting that a number of researchers are predicting AGI in the next twenty to one hundred years. One must not put excess trust in this time frame: as Armstrong [29] show, experts have been terrible at predicting AGI. Muehlhauser [200] consider a number of methods other than expert opinion that could be used for predicting AGI, but find that they too provide suggestive evidence at best.

It would be a mistake, however, to leap from ‘AGI is very hard to predict’ to ‘AGI must be very far away’. Our brains are known to think about uncertain, abstract ideas like AGI in ‘far mode’, which also makes it feel like AGI must be temporally distant [198, 268], but something being *uncertain* is not strong evidence that it is *far away*. When we are highly ignorant about something, we should widen our error bars in both directions. Thus, we should not be highly confident that AGI will arrive this century and we should not be highly confident that it *will not*.

Next, we explained why AGIs may be an existential risk. A trend toward automatization would give AGIs increased influence in society and there might be a discontinuity in which they gained power rapidly. This could be a disaster for humanity if AGIs do not share our values and, in fact, it looks difficult to make them share our values because human values are complex and fragile, and therefore problematic to specify.

The recommendations given for dealing with the problem can be divided into proposals for societal action (section 3), external constraints (section 4) and internal constraints (section 5). Many proposals seem to suffer from serious problems, or seem to be of limited effectiveness. Others seem to have enough promise to be worth exploring. We will conclude by reviewing the proposals which we feel are worthy of further study.

As a brief summary of our views, in the medium term, we think that the proposals of AGI confinement (section 4.1), Oracle AI (section 5.1) and motivational weaknesses (section 5.6) would have promise in helping create safer AGIs. These proposals share in common the fact that, although they could help a cautious team of researchers create an AGI, they are not solutions to the problem of AGI risk, as they do not prevent others from creating unsafe AGIs, nor are they sufficient in guaranteeing the safety of sufficiently intelligent AGIs. Regulation (section 3.3) as well as human capability enhancement (section 3.4) could also help to somewhat reduce AGI risk. In the long run, we will need the ability to guarantee the safety of freely acting AGIs. For this goal, value learning (section 5.2.5) would seem like the most reliable approach if it could be made to work, with human-like architecture (section 5.3.4) a possible alternative which seems less reliable but possibly easier to build. Formal verification (section 5.5) seems like a very important tool in helping to ensure the safety of our AGIs, regardless of the exact approach that we choose.

Responses to catastrophic AGI risk			
Societal proposals			
Do nothing		AGI is distant	
		Little risk, no action needed	
		Let them kill us	
Integrate to society		Legal and economic controls	
		Foster positive values	
Regulate research		Review boards	
		Encourage safety research	
		Differential technological progress	
		International mass surveillance	
Enhance human capabilities			
Relinquish technology		Outlaw AGI	
		Restrict hardware	
AGI design proposals			
External constraints		Internal constraints	
AGI confinement	Safe questions	Oracle AI	
	Virtual worlds	Top-down approaches	Three laws
	Resetting the AGI		Categorical imperative
	Checks and balances		Principle of voluntary joyous growth
AGI enforcement		Bottom-up and hybrid approaches	Utilitarianism
			Value learning
			Evolutionary invariants
			Evolved morality
			Reinforcement learning
			Human-like AGI

(Continued.)

**AGI design proposals**

External constraints	Internal constraints
	<b>AGI Nanny</b>
	<b>Formal</b>
	<b>verification</b>
	<b>Motivational</b>
	<b>weaknesses</b>
	High discount rates
	Easily satiable goals
	Calculated indifference
	Programmed restrictions
	Legal Machine Language

Of the societal proposals, we are supportive of the calls to regulate AGI development, but we admit there are many practical hurdles which might make this infeasible. The economic and military potential of AGI, and the difficulty of verifying regulations and arms treaties restricting it, could lead to unstoppable arms races.

We find ourselves in general agreement with the authors who advocate funding additional research into safe AGI as the primary solution. Such research will also help establish the kinds of constraints which would make it possible to successfully carry out integration proposals.

Uploading approaches, in which human minds are made to run on computers and then augmented, might buy us some time to develop safe AGI. However, it is unclear whether they can be developed before AGI and large-scale uploading could create strong evolutionary trends which seem dangerous in and of themselves. As AGIs seem likely to eventually outpace uploads, uploading by itself is probably not a sufficient solution. What uploading could do is to reduce the initial advantages that AGIs enjoy over (partially uploaded) humanity, so that other responses to AGI risk can be deployed more effectively.

External constraints are likely to be useful in controlling AGI systems of limited intelligence and could possibly help us develop more intelligent AGIs while maintaining their safety. If inexpensive external constraints were readily available, this could encourage even research teams skeptical about safety issues to implement them. Yet it does not seem safe to rely on these constraints once we are dealing with a superhuman intelligence and we cannot trust everyone to be responsible enough to contain their AGI systems, especially given the economic pressures to ‘release’ AGIs. For such an approach to be a solution for AGI risk in general, it would have to be adopted by all successful AGI projects, at least until safe AGIs were developed. Much the same is true of attempting to design Oracle AIs. In the short term, such efforts may be reinforced by research into motivational weaknesses, internal constraints that make AGIs easier to control via external means.

In the long term, the internal constraints that show the most promise are value extrapolation approaches and human-like architectures. Value extrapolation attempts to learn human values and interpret them as we would wish them to be interpreted. These approaches have the advantage of potentially maximizing the preservation of human values and the disadvantage that such approaches may prove intractable or impossible to properly formalize. Human-like architectures seem easier to construct, as we can simply copy mechanisms that are used within the human brain, but it seems hard to build such an exact match as to reliably replicate human values. Slavish reproductions of the human psyche also seem likely to be outcompeted by less human, more efficient architectures.

Both approaches would benefit from better formal verification methods, so that AGIs which were editing and improving themselves could verify that the modifications did not threaten to remove the AGIs’ motivation to follow their original goals. Studies which aim to uncover the roots of human morals and preferences also seem like candidates for research that would benefit the development of safe AGI [42, 199, 245], as do studies into computational models of ethical reasoning [186].

We reiterate that when we talk about ‘human values’, we are not making the claim that human values would be static, nor that *current* human values would be ideal. Nor do we wish to imply that the values of other sentient beings would be unimportant. Rather, we are seeking to guarantee the implementation of some very basic values, such as the avoidance of unnecessary suffering, the preservation of humanity and the prohibition of forced brain reprogramming. We believe the vast majority of humans would agree with these values and be sad to see them lost.

## Acknowledgments

We extend special thanks to Luke Muehlhauser for extensive assistance throughout the writing process. We are grateful to Olle Häggström for organizing the event that led to this paper being formally published. We would also like to thank Abram Demski, Alexei Turchin, Alexey Potapov, Anders Sandberg, Andras Kornai, Anthony Berglas, Aron Vallinder, Ben Goertzel, Ben Noble, Ben Sterrett, Brian Rabkin, Bill Hibbard, Carl Shulman, Dana Scott, Daniel Dewey, David Pearce, Evelyn Mitchell, Evgenij Thorstensen, Frank White, gvern branwen, Harri Valpola, Jaan Tallinn, Jacob Steinhardt, James Babcock, James Miller, Joshua Fox, Louie Helm, Mark Gubrud, Mark Waser, Michael Anissimov, Michael Vassar, Miles Brundage, Moshe Looks, Randal Koene, Robin Hanson, Risto Saarelma, Steve Omohundro, Suzanne Lidström, Steven Kaas, Stuart Armstrong, Tim Freeman, Ted Goertzel, Toni Heinonen, Tony Barrett, Vincent Mller, Vladimir Nesov, Wei Dai and two anonymous reviewers as well as several users of <http://www.LessWrong.com> for their helpful comments.

## References

- [1] Anderson M and Anderson S L (ed) 2011 *Machine Ethics* (New York: Cambridge University Press)
- [2] Bostrom N and Ćirković M M (ed) 2008 *Global Catastrophic Risks* (New York: Oxford University Press)
- [3] Bostrom N 2014 *Superintelligence: Paths, Dangers, Strategies* (Italy: Oxford University Press)
- [4] National Defense Authorization 2001 Public Law 106-398, 114 Stat. 1654 (An act by the US Congress, [www.gpo.gov/fdsys/pkg/PLAW-106publ398/html/PLAW-106publ398.htm](http://www.gpo.gov/fdsys/pkg/PLAW-106publ398/html/PLAW-106publ398.htm))
- [5] Eden A, Søraker J, Moor J H and Steinhart E (ed) 2012 *Singularity Hypotheses: A Scientific and Philosophical Assessment (The Frontiers Collection)* (Berlin: Springer)
- [6] IEEE Spectrum 2008 Tech luminaries address singularity *The Singularity: Special Report* <http://spectrum.ieee.org/computing/hardware/tech-luminaries-address-singularity>
- [7] Kringsbach M L and Berridge K C (ed) 2009 *Pleasures of the Brain (Series in Affective Science)* (New York: Oxford University Press)
- [8] Pylyshyn Z W (ed) 1987 *The Robot's Dilemma: The Frame Problem in Artificial Intelligence* (Norwood, NJ: Ablex)
- [9] Wood D M and Kirstie B (ed) 2006 *A Report on the Surveillance Society: For the Information Commissioner, by the Surveillance Studies Network* (Wilmslow, UK: Office of the Information Commissioner) ([www.ico.org.uk/about\\_us/research/~media/documents/library/Data\\_Protection/Practical\\_application/SURVEILLANCE\\_SOCIETY\\_SUMMARY\\_06.ashx](http://www.ico.org.uk/about_us/research/~media/documents/library/Data_Protection/Practical_application/SURVEILLANCE_SOCIETY_SUMMARY_06.ashx))
- [10] Adams S S *et al* 2012 Mapping the landscape of human-level artificial general intelligence *AI Mag.* **33** 25–42
- [11] Agar N 2011 Ray Kurzweil and uploading *J. Evolution Technol.* **22** 23–36
- [12] Agliata D and Tantleff-Dunn S 2004 The impact of media exposure on males' body image *J. Social Clinical Psych.* **23** 7–22
- [13] Allen C and Wallach W 2012 Moral machines: contradiction in terms of abdication of human responsibility *Robot Ethics: The Ethical and Social Implications of Robotics* (Cambridge, MA: MIT Press) pp 55–68
- [14] Allen C, Varner G and Zinser J 2000 Prolegomena to any future artificial moral agent *J. Exp. Theor. Art. Intell.* **12** 251–61
- [15] Allen C, Smit I and Wendell W 2005 Artificial morality *Ethics Info. Technol.* **7** 149–55
- [16] Allen C, Wallach W and Iva Smit I 2006 Why machine ethics? *IEEE Intell. Systems* **21** 12–7
- [17] Amdahl G M 1967 Validity of the single processor approach to achieving large scale computing capabilities *Proc. Spring Joint Computer Conference (AFIPS '67)* (New York: ACM Press) pp 483–5
- [18] Anderson M, Anderson S L and Armen C (ed) 2005 *Machine Ethics* Technical Report FS-05-06 (Menlo Park, CA: AAAI Press)
- [19] Anderson M, Anderson S L and Armen C 2005 Towards machine ethics *Machine Ethics* Technical Report FS-05-06 (Menlo Park, CA: AAAI Press) pp 1–7
- [20] Anderson M, Anderson S L and Armen C 2005 MedE-thEx *Caring Machines* Technical Report FS-05-02, ed T Bickmore (Menlo Park, CA: AAAI Press) pp 9–16
- [21] Anderson M, Anderson S L and Armen C 2006 An approach to computing ethics *IEEE Intell. Systems* **21** 56–63
- [22] Anderson M 2010 Problem solved *H+ Magazine* [hplusmagazine.com/2010/12/15/problem-solved-unfriendly-ai](http://hplusmagazine.com/2010/12/15/problem-solved-unfriendly-ai)
- [23] Anderson S L 2011 The unacceptability of Asimov's three laws of robotics as a basis for machine ethics *Machine Ethics* (Cambridge: Cambridge University Press) pp 285–96
- [24] Annas G J, Andrews L B and Isasi R M 2002 Protecting the endangered human *Am. J. Law Med.* **28** 151–78
- [25] Anthony D and Robbins T 2004 Conversion and brainwashing in new religious movements *The Oxford Handbook of New Religious Movements* ed J R Lewis (New York: Oxford University Press) pp 243–97
- [26] Ronald C and Arkin R C 2009 *Governing Lethal Behavior in Autonomous Robots* (Boca Raton, FL: CRC)
- [27] Armstrong S 2007 Chaining god [www.neweuropeancentury.org/GodAI.pdf](http://www.neweuropeancentury.org/GodAI.pdf)
- [28] Armstrong S 2010 Utility indifference [www.fhi.ox.ac.uk/reports/2010-1.pdf](http://www.fhi.ox.ac.uk/reports/2010-1.pdf)
- [29] Armstrong S and Sotala K 2012 How we're predicting AI—or failing to *Proc. Int. Conf. Beyond AI 2012, 5–6 November 2012, Pilsen, Czech Republic* (Pilsen: University of West Bohemia) pp 52–75 ([www.kky.zcu.cz/en/publications/1/JanRomportl\\_2012\\_BeyondAIArtificial.pdf](http://www.kky.zcu.cz/en/publications/1/JanRomportl_2012_BeyondAIArtificial.pdf))
- [30] Armstrong S, Sandberg A and Bostrom N 2012 Thinking inside the box *Minds Machines* **22** 299–324
- [31] Asaro P M Robots and responsibility from a legal perspective *Proc. IEEE Conf. on Robotics and Automation, Workshop on Roboethics* p 59 ([www.peterasaro.org/writing/ASAROLegalPerspective.pdf](http://www.peterasaro.org/writing/ASAROLegalPerspective.pdf))
- [32] Ashley K D and McLaren B M 1995 Reasoning with reasons in case-based comparisons *Proc. First International Conf. on Case-Based Reasoning Research and Development* ed M M Veloso and A Aamodt (Berlin: Springer) pp 133–44 ([www.cs.cmu.edu/~bmclaren/pubs/AshleyMcLaren-ReasoningWithReasons-ICCB95.pdf](http://www.cs.cmu.edu/~bmclaren/pubs/AshleyMcLaren-ReasoningWithReasons-ICCB95.pdf))
- [33] Asimov I 1942 Runaround *Astounding Science-Fiction* pp 94–103
- [34] Axelrod R 1987 The evolution of strategies in the iterated Prisoner's Dilemma *Genetic Algorithms and Simulated Annealing* ed L Davis (Los Altos, CA: Morgan Kaufmann) pp 32–41
- [35] Baars B J 2002 The conscious access hypothesis *Trends Cogn. Sci.* **6** 47–52
- [36] Baars B J 2005 Global workspace theory of consciousness *The Boundaries of Consciousness (Progress in Brain Research no 150)* ed S Laureys (Boston: Elsevier) pp 45–53
- [37] Bach J, Goertzel B and Ikl M (ed) 2012 *Artificial General Intelligence (Lecture Notes in Artificial Intelligence no 7716)* (New York: Springer)
- [38] Bamford S 2012 A framework for approaches to transfer of a mind's sub-strate *Int. J. Machine Consciousness* **4** 23–34
- [39] Baum S D, Goertzel B and Goertzel T G 2011 How long until human-level AI? Results from an expert assessment *Technol. Forecasting Social Change* **78** 185–95
- [40] Beavers A F 2009 Between angles or animals: the question of robot ethics; or is Kantian moral agency desirable? *18th Ann. Meeting of Association for Practical and Professional Ethics, Cincinnati, OH*
- [41] Beavers A F 2012 Moral machines and the threat of ethical nihilism *Robot Ethics: The Ethical and Social Implications of Robotics* (Cambridge, MA: MIT Press) pp 333–44
- [42] Bello P and Bringsjord S 2012 On how to build a moral machine *Topoi* doi:[10.1007/s11245-012-9129-8](https://doi.org/10.1007/s11245-012-9129-8)
- [43] Benatar D 2006 *Better Never to Have Been* (New York: Oxford University Press)
- [44] Berglas A 2012 Artificial intelligence will kill our grandchildren (singularity) draft 9 <http://berglas.org/Articles/AIKillGrandchildren/AIKillGrandchildren.html>
- [45] Blackmore S 2012 She won't be me *J. Consciousness Studies* **19** 16–9
- [46] Bostrom N 1998 How long before superintelligence? *Int. J. Futures Studies* **2**
- [47] Bostrom N 2002 Existential risks *J. Evolution Technol.* **9** [www.jetpress.org/volume9/risks.html](http://www.jetpress.org/volume9/risks.html)

- [48] Bostrom N 2003 Ethical issues in advanced artificial intelligence *Cognitive, Emotive and Ethical Aspects of Decision Making in Humans and in Artificial Intelligence* vol 2 ed I Smit and G E Lasker (Windsor, ON: International Institute for Advanced Studies in Systems Research and Cybernetics) pp 12–7
- [49] Bostrom N 2004 The future of human evolution *Two Hundred Years After Kant, Fifty Years After Turing (Death and Anti-Death* vol 2) ed C Tandy (Ria University Press) pp 339–71
- [50] Bostrom N 2007 Technological revolutions *Nanoscale* ed N M de, S Cameron and M E Mitchell (Hoboken, NJ: Wiley) pp 129–52
- [51] Bostrom N 2012 The superintelligent will *Minds Machines* **22** 71–85
- [52] Bostrom N and Cirkovic M M 2011 Introduction *Global Catastrophic Risks* (Oxford: Oxford University Press) pp 1–30
- [53] Bostrom N and Yudkowsky E 2014 The ethics of artificial intelligence *Cambridge Handbook of Artificial Intelligence* ed K Frankish and W Ramsey (New York: Cambridge University Press)
- [54] Brain M 2003 Robotic nation <http://marshallbrain.com/robotic-nation.htm>
- [55] Brandt R B 1979 *A Theory of the Good and the Right* (New York: Oxford University Press)
- [56] Branwen G Slowing Moore's law [www.gwern.net/SlowingMoore'sLaw](http://www.gwern.net/SlowingMoore'sLaw)
- [57] Brin D 1998 *The Transparent Society* (Reading, MA: Perseus)
- [58] Bringsjord S and Bringsjord A 2012 Belief in the singularity is fideistic *Singularity Hypotheses* ed A H Eden, J H Moor, J H Soraker and E Steinhart (Berlin: Springer)
- [59] Brooks R A 2008 I, Rodney Brooks, am a robot *IEEE Spectrum* **45** 68–71
- [60] Brynjolfsson E and McAfee A 2011 *Race Against the Machine* (Lexington, MA: Digital Frontier)
- [61] Bryson J and Kime P 1998 Just another artefact [www.cs.bath.ac.uk/~jjb/web/aiethics98.html](http://www.cs.bath.ac.uk/~jjb/web/aiethics98.html)
- [62] Bryson J J 2010 Robots should be slaves *Close Engagements with Artificial Companions* ed Y Wilks (Philadelphia: John Benjamins) pp 107–26
- [63] Bugaj S V and Goertzel B 2007 Five ethical imperatives and their implications for human-AGI interaction *Dynamical Psychology* [http://goertzel.org/dynapsyc/2007/Five\\_Ethical\\_Imperatives\\_svbedit.htm](http://goertzel.org/dynapsyc/2007/Five_Ethical_Imperatives_svbedit.htm)
- [64] Butler S 1863 Darwin among the machines *The Press* (Christchurch, New Zealand) [www.nzetc.org/tm/scholarly/tei-ButFir-t1-g1-t1-g1-t4-body.html](http://www.nzetc.org/tm/scholarly/tei-ButFir-t1-g1-t1-g1-t4-body.html)
- [65] Cade C M 1966 *Other Worlds Than Ours* 1st edn (Museum)
- [66] Calandrino J A, Clarkson W and Felten E W 2011 Bubble trouble *Proc. 20th USENIX Security Symp* (San Francisco: USENIX) pp 267–80 [www.usenix.org/events/sec11/tech/full\\_papers/Calandrino.pdf](http://www.usenix.org/events/sec11/tech/full_papers/Calandrino.pdf)
- [67] Cassimatis N, Mueller E T and Winston P H 2006 Achieving human-level intelligence through integrated systems and research *AI Mag.* **27** 12–4 [www.aaai.org/ojs/index.php/aimagazine/article/view/1876/1774](http://www.aaai.org/ojs/index.php/aimagazine/article/view/1876/1774)
- [68] Casti J L 2012 *X-Events* (New York: William Morrow)
- [69] Cattell R and Parker A Challenges for brain emulation: why is building a brain so difficult? <http://synapticlink.org/BrainEmulationChallenges.pdf>
- [70] Commodity Futures Trading Commission and Securities and Exchange Commission 2010 Findings regarding the market events of May 6, 2010 [www.sec.gov/news/studies/2010/marketevents-report.pdf](http://www.sec.gov/news/studies/2010/marketevents-report.pdf)
- [71] Chalmers D J 1996 *The Conscious Mind (Philosophy of Mind Series)* (New York: Oxford University Press)
- [72] Chalmers D J 2010 The singularity *J. Consciousness Studies* **17** 7–65
- [73] Paul F and Christiano P F 2012 Indirect normativity write-up *Ordinary Ideas* <http://ordinaryideas.wordpress.com/2012/04/21/indirect-normativity-write-up>
- [74] Clark G 2007 *A Farewell to Alms* 1st edn (Princeton, NJ: Princeton University Press)
- [75] Clarke R 1993 Asimov's laws of robotics *Computer* **26** 53–61
- [76] Clarke R 1994 Asimov's laws of robotics *Computer* **27** 57–66
- [77] Cloos C 2005 The utilibot project *Machine Ethics* Technical Report FS-05-06 ed M Anderson, S L Anderson and C Armen (Menlo Park, CA: AAAI Press) pp 38–45
- [78] Dahm W J A 2010 Technology horizons [www.au.af.mil/au/awc/awcgate/af/tech\\_horizons\\_vol-1\\_may2010.pdf](http://www.au.af.mil/au/awc/awcgate/af/tech_horizons_vol-1_may2010.pdf)
- [79] Daley W 2011 Mitigating potential hazards to humans from the development of intelligent machines *Synthese* **2** 44–50 ([www.synesisjournal.com/vol2\\_g/2011\\_2\\_44-50\\_Daley.pdf](http://www.synesisjournal.com/vol2_g/2011_2_44-50_Daley.pdf))
- [80] Davis E 2013 The singularity and the state of the art in artificial intelligence [www.cs.nyu.edu/~davise/papers/singularity.pdf](http://www.cs.nyu.edu/~davise/papers/singularity.pdf)
- [81] Dayan P 2011 Models of value and choice *Neuroscience of Preference and Choice* ed R J Dolan and T Sharot (Waltham, MA: Academic) pp 33–52
- [82] de Garis H 2005 *The Artilect War: Cosmists vs Terrans* (Palm Springs, CA: ETC Publications)
- [83] de Waal F, Wright R, Korsgaard C M, Kitcher P and Singer P 2006 *Primates and Philosophers* 1st edn (Princeton, NJ: Princeton University Press)
- [84] Degabriele J P, Paterson K and Watson G 2011 Provable security in the real world *IEEE Security Privacy Mag.* **9** 33–41
- [85] Dennett D C 1987 Cognitive wheels *The Robot's Dilemma* ed Z W Pylyshyn (Norwood, NJ: Ablex) pp 41–64
- [86] Dennett D C 2012 The mystery of David Chalmers *J. Consciousness Studies* **19** 86–95
- [87] Deutsch D 2011 *The Beginning of Infinity* 1st edn (New York: Viking)
- [88] Dewey D 2011 Learning what to value *Artificial General Intelligence (Lecture Notes in Computer Science* no 6830) ed J Schmidhuber, K R Thirsson and M Looks (New York: Springer) pp 309–14
- [89] Dietrich E 2014 After the humans are gone *Philosophy Now* [http://philosophynow.org/issues/61/After\\_The\\_Humans\\_Are\\_Gone](http://philosophynow.org/issues/61/After_The_Humans_Are_Gone)
- [90] Docherty B and Goose S 2012 Losing humanity [www.hrw.org/sites/default/files/reports/arms1112ForUpload\\_0\\_0.pdf](http://www.hrw.org/sites/default/files/reports/arms1112ForUpload_0_0.pdf)
- [91] Douglas T 2008 Moral enhancement *J. Appl. Phil.* **25** 228–45
- [92] Drexler K E 1986 *Engines of Creation* (Garden City, NY: Anchor)
- [93] Eckersley P and Sandberg A 2013 Is brain emulation dangerous? *J. Artif. Gen. Intell.* **4** 170–94
- [94] Eisen M 2011 Amazon's \$23,698,655.93 book about flies *It is NOT Junk* [www.michaelseisen.org/blog/?p=358](http://www.michaelseisen.org/blog/?p=358)
- [95] Felten E W and Schneider M A 2000 Timing attacks on Web privacy *Proc. 7th ACM Conference on Computer and Communications Security—CCS '00* (New York: ACM Press) pp 25–32
- [96] Ferguson M J, Hassin R and Bargh J A 2007 Implicit motivation *Handbook of Motivation Science* ed J Y Shah and W L Gardner (New York: Guilford) pp 150–66
- [97] Fox J and Shulman C 2010 Superintelligence does not imply benevolence *ECAP10, VIII European Conference of Computing and Philosophy* ed K Mainzer (Munich: Dr Hut)
- [98] Frankfurt H G 1971 Freedom of the will and the concept of a person *J. Phil.* **68** 5–20
- [99] Franklin S and Patterson F G Jr 2006 The LIDA architecture *IDPT-2006 Proc.* (San Diego, CA: Society for Design and Process Science) <http://ccrg.cs.memphis.edu/assets/papers/zo-1010-lida-060403.pdf>

- [100] Freeman T 2008 Comparative advantage doesn't ensure survival [www.fungible.com/comparative-advantage.html](http://www.fungible.com/comparative-advantage.html)
- [101] Freeman T 2009 Using compassion and respect to motivate an artificial intelligence [www.fungible.com/respect/paper.html](http://www.fungible.com/respect/paper.html)
- [102] Friedman B and Kahn P H 1992 Human agency and responsible computing *J. Syst. Software* **17** 7–14
- [103] Gewirth A 1978 *Reason and Morality* (Chicago: University of Chicago Press)
- [104] Gips J 1995 Towards the ethical robot *Android Epistemology* ed K M Ford, C N Glymour and P J Hayes (Cambridge, MA: MIT Press) pp 243–52
- [105] Goertzel B 2006 Apparent limitations on the 'AI friendliness' and related concepts imposed by the complexity of the world [www.goertzel.org/papers/LimitationsOnFriendliness.pdf](http://www.goertzel.org/papers/LimitationsOnFriendliness.pdf)
- [106] Goertzel B 2010 Coherent aggregated volition *The Multiverse According to Ben* <http://multiverseaccordingtoben.blogspot.ca/2010/03/coherent-aggregated-volition-toward.html>
- [107] Goertzel B 2010 *GOLEM* <http://goertzel.org/GOLEM.pdf>
- [108] Goertzel B 2012 *CogPrime* [http://wiki.opencog.org/w/CogPrime\\_Overview](http://wiki.opencog.org/w/CogPrime_Overview)
- [109] Goertzel B 2002 Thoughts on AI morality *Dynamical Psychology* [www.goertzel.org/dynapsyc/2002/AIMorality.htm](http://www.goertzel.org/dynapsyc/2002/AIMorality.htm)
- [110] Goertzel B 2004 Encouraging a positive transcension *Dynamical Psychology* [www.goertzel.org/dynapsyc/2004/PositiveTranscension.htm](http://www.goertzel.org/dynapsyc/2004/PositiveTranscension.htm)
- [111] Goertzel B 2004 Growth, choice and joy *Dynamical Psychology* [www.goertzel.org/dynapsyc/2004/GrowthChoiceJoy.htm](http://www.goertzel.org/dynapsyc/2004/GrowthChoiceJoy.htm)
- [112] Goertzel B 2012 Should humanity build a global AI nanny to delay the singularity until it's better understood? *J. Consciousness Studies* **19** 96–111
- [113] Goertzel B 2012 When should two minds be considered versions of one another? *Int. J. Machine Consciousness* **4** 177–85
- [114] Goertzel B and Bugaj S V 2008 Stages of ethical development in artificial general intelligence systems *Artificial General Intelligence (Frontiers in Artificial Intelligence and Applications no 171)* (IOS) 448–59
- [115] Goertzel B and Pitt J 2012 Nine ways to bias open-source AGI toward friendliness *J. Evolution Technol.* **22** 116–31
- [116] Golle P and Partridge K 2009 On the anonymity of home/work location pairs *Pervasive Computing (Lecture Notes in Computer Science no 5538)* ed H Tokuda, M Beigl, A Friday, A Brush and Y Tobe (Berlin: Springer) pp 390–7
- [117] Good I J 1965 Speculations concerning the first ultraintelligent machine *Advances in Computers Volume 6* ed F L Alt and M Rubino (New York: Academic) pp 31–88
- [118] Good I J 1970 Some future social repercussions of computers *Int. J. Environ. Studies* **1** 67–79
- [119] Good I J 1982 Ethical machines *Intelligent Systems (Machine Intelligence no 10)* ed J E Hayes, D Michie and Y-H Pao (Chichester: Ellis Horwood) pp 555–60
- [120] Diana F and Gordon-Spears D F 2003 Asimov's laws *Formal Approaches to Agent-Based Systems (Lecture Notes in Computer Science no 2699)* ed M G Hinchey, J L Rash, W F Truszkowski, C Rou and D F Gordon-Spears (Berlin: Springer) pp 257–9
- [121] Christopher Grau G 2006 There is no I in Robot *IEEE Intell. Syst.* **21** 52–5
- [122] Groesz L M, Levine M P and Murnen S K. 2001 The effect of experimental presentation of thin media images on body satisfaction *Int. J. Eating Disorders* **31** 1–16
- [123] Guarini M 2006 Particularism and the classification and reclassification of moral cases *IEEE Intell. Systems* **21** 22–8
- [124] Gubrud M V 1997 Nanotechnology and international security [www.foresight.org/Conferences/MNT05/Papers/Gubrud/](http://www.foresight.org/Conferences/MNT05/Papers/Gubrud/)
- [125] Gunkel D J 2012 *The Machine Question* (Cambridge, MA: MIT Press)
- [126] Guterl F 2012 *The Fate of the Species* 1st edn (New York: Bloomsbury)
- [127] Haidt J 2006 *The Happiness Hypothesis* 1st edn (New York: Basic)
- [128] Hall J S 2007 *Beyond AI* (Amherst, NY: Prometheus)
- [129] Hall J S 2007 Ethics for artificial intellects *Nanoethics* ed F Allho, P Lin, J Moor, J Weckert and M C Roco (New York: Wiley) pp 339–52
- [130] Hall J S 2008 Engineering utopia *Artificial General Intelligence Frontiers (Artificial Intelligence and Applications no 171)* ed P Wang, B Goertzel and S Franklin (Amsterdam: IOS) pp 460–7
- [131] Hall J S Ethics for self-improving machines *Machine Ethics* ed M Anderson and S L Anderson (Cambridge: Cambridge University Press) pp 512–23
- [132] Hallevy G 2010 The criminal liability of artificial intelligence entities [http://papers.ssrn.com/sol3/papers.cfm?abstract\\_id=1564096](http://papers.ssrn.com/sol3/papers.cfm?abstract_id=1564096)
- [133] Hanson R 2007 Shall we vote on values, but bet on beliefs? <http://hanson.gmu.edu/futarchy.pdf>
- [134] Hanson R 2009 Prefer law to values *Overcoming Bias* [www.overcomingbias.com/2009/10/prefer-law-to-values.html](http://www.overcomingbias.com/2009/10/prefer-law-to-values.html)
- [135] Hanson R 1994 If uploads come first *Extropy* **6** <http://hanson.gmu.edu/uploads.html>
- [136] Hanson R 1998 Economic growth given machine intelligence <http://hanson.gmu.edu/aigrow.pdf>
- [137] Hanson R 2008 Economics of the singularity *IEEE Spectrum* **45** 45–50
- [138] Hanson R 2012 Meet the new conflict, same as the old conflict *J. Consciousness Studies* **19** 119–25
- [139] Hare R D, Clark D, Grann M and Thornton D 2000 Psychopathy and the predictive validity of the PCL-R *Behavioral Sci. Law* **18** 623–45
- [140] Harris G T and Rice M E 2006 Treatment of psychopathy *Handbook of Psychopathy* ed C J Patrick (New York: Guilford) pp 555–72
- [141] Hart D and Goertzel B 2008 OpenCog: a software framework for integrative artificial general intelligence [www.agiri.org/OpenCog\\_AGI-08.pdf](http://www.agiri.org/OpenCog_AGI-08.pdf)
- [142] Hauskeller M 2012 My brain, my mind, and I *Int. J. Machine Consciousness* **4** 187–200
- [143] Hayworth K J 2012 Electron imaging technology for whole brain neural circuit mapping *Int. J. Machine Consciousness* **4** 87–108
- [144] Heylighen F 2007 Accelerating socio-technological evolution *Globalization as Evolutionary Process (Rethinking Globalizations no 10)* ed G Modelski, T Devezas and W R Thompson (New York: Routledge) pp 284–309
- [145] Heylighen F 2012 Brain in a vat cannot break out *J. Consciousness Studies* **19** 126–42
- [146] Hibbard B 2005 The ethics and politics of super-intelligent machines [https://sites.google.com/site/whibbard/g/SI\\_ethics\\_politics.doc](https://sites.google.com/site/whibbard/g/SI_ethics_politics.doc)
- [147] Hibbard B 2005 Critique of the SIAI collective volition theory [www.ssec.wisc.edu/~billh/g/SIAI\\_CV\\_critique.html](http://www.ssec.wisc.edu/~billh/g/SIAI_CV_critique.html)
- [148] Hibbard B 2012 The error in my 2001 VisFiles column [www.ssec.wisc.edu/~billh/g/visfiles\\_error.html](http://www.ssec.wisc.edu/~billh/g/visfiles_error.html)
- [149] Hibbard B 2001 Super-intelligent machines *ACM SIGGRAPH Computer Graphics* **35** 13–5 ([www.siggraph.org/publications/newsletter/issues/v35/v35n1.pdf](http://www.siggraph.org/publications/newsletter/issues/v35/v35n1.pdf))
- [150] Hibbard B 2008 Open source AI *Artificial General Intelligence Frontiers (Artificial Intelligence and Applications no 171)* ed P Wang, B Goertzel and S Franklin (Amsterdam: IOS) pp 473–7
- [151] Hibbard B 2012 Model-based utility functions *J. Artificial Gen. Intell.* **3** 1–24
- [152] Hibbard B 2012 Decision support for safe AI design *Artificial General Intelligence (Lecture Notes in Artificial Intelligence*

- no 7716) ed J Bach, B Goertzel and M Ikl (New York: Springer) pp 117–25
- [153] Hibbard B 2012 Avoiding unintended AI behaviors *Artificial General Intelligence (Lecture Notes in Artificial Intelligence no 7716)* ed J Bach, B Goertzel and M Ikl (New York: Springer) pp 107–16
- [154] Hollerbach J M, Mason M T, Henrik I and Christensen H I 2009 A roadmap for US robotics [www.us-robotics.us/reports/CCCRReport.pdf](http://www.us-robotics.us/reports/CCCRReport.pdf)
- [155] Hopkins P D 2012 Why uploading will not work, or, the ghosts haunting transhumanism *Int. J. Machine Consciousness* **4** 229–43
- [156] Horvitz E J and Selman B 2009 Interim report from the AAAI Presidential Panel on long-term AI futures [www.aaai.org/Organization/Panel/panel-note.pdf](http://www.aaai.org/Organization/Panel/panel-note.pdf)
- [157] Hughes J 2001 Relinquishment or regulation [www.trincoll.edu/orgs/scialnce/sfr/01-02/files/Relinquishment%20or%20Regulation%203.James%20Hughes.doc](http://www.trincoll.edu/orgs/scialnce/sfr/01-02/files/Relinquishment%20or%20Regulation%203.James%20Hughes.doc)
- [158] Hutter M 2012 Can intelligence explode? *J. Consciousness Studies* **19** 143–66
- [159] Jenkins A 2003 Artificial intelligence and the real world *Futures* **35** 779–86
- [160] Joy B 2000 Why the future doesn't need us *Wired* [www.wired.com/wired/archive/8.04/joy.html](http://www.wired.com/wired/archive/8.04/joy.html)
- [161] Joyce R 2001 Evolution and morality *The Myth of Morality (Cambridge Studies in Philosophy)* (New York: Cambridge University Press)
- [162] Karnofsky H 2012 Thoughts on the singularity institute (SI) *Less Wrong* [http://lesswrong.com/lw/cbs/thoughts\\_on\\_the\\_singularity\\_institute\\_si/](http://lesswrong.com/lw/cbs/thoughts_on_the_singularity_institute_si/)
- [163] Karnofsky H and Tallinn J 2011 Karnofsky and Tallinn dialog on SIAI efficacy <http://xa.yimg.com/kq/groups/23070378/1331435883/name/Jaan+Tallinn+2011+05+-+revised.doc>
- [164] Kipnis D 1972 Does power corrupt? *J. Personality Social Psych.* **24** 33–41
- [165] Koene R A 2012 Experimental research in whole brain emulation *Int. J. Machine Consciousness* **4** 35–65
- [166] Koene R A 2012 Embracing competitive balance *Singularity Hypotheses* ed A H Eden, J H Moor, J H Soraker and E Steinhart (Berlin: Springer)
- [167] Kornai A 2014 Bounding the impact of AGI *J. Experimental Theor. Artificial Intell* **26** 417–38
- [168] Kurzweil R 2001 Response to Stephen Hawking [www.kurzweilai.net/response-to-stephen-hawking](http://www.kurzweilai.net/response-to-stephen-hawking)
- [169] Kurzweil R 2002 Locked in his Chinese room [www.kurzweilai.net/chapter-6-locked-in-his-chinese-room-response-to-john-searle](http://www.kurzweilai.net/chapter-6-locked-in-his-chinese-room-response-to-john-searle)
- [170] Kurzweil R 2005 *The Singularity Is Near* (New York: Viking)
- [171] Lampson B W 1973 A note on the confinement problem *Comm. ACM* **16** 613–5
- [172] Legg S 2009 *Funding safe AGI Vetta Project* [www.vetta.org/2009/08/funding-safe-agi/](http://www.vetta.org/2009/08/funding-safe-agi/)
- [173] Legg S and Hutter M 2007 A collection of definitions of intelligence *Advances in General Intelligence: Concepts Architectures and Algorithms (Frontiers in Artificial Intelligence and Applications no 157)* (Amsterdam: IOS) pp 17–24
- [174] Lehman-Wilzig S N 1981 Frankenstein unbound *Futures* **13** 442–57
- [175] Levy D 2009 The ethical treatment of artificially conscious robots *Int. J. Social Robotics* **1** 209–16
- [176] Lewis D 1989 Dispositional theories of value *Proc. Aristotelian Soc. Suppl. Volumes* **63** 113–37
- [177] Loosemore R and Goertzel B 2012 Why an intelligence explosion is probable *Singularity Hypotheses* ed A H Eden, J H Moor, J H Soraker and E Steinhart (Berlin: Springer)
- [178] Mainzer K (ed) 2010 *ECAP10, VIII European Conference of Computing and Philosophy* (Munich: Dr Hut)
- [179] Mann S, Nolan J and Wellman B 2003 *Sousveillance Surveillance Soc.* **1** 331–55
- [180] McCauley L 2007 AI armageddon and the three laws of robotics *Ethics Information Technol.* **9** 153–64
- [181] McCulloch W S 1956 Toward some circuitry of ethical robots; or, an observational science of the genesis of social evaluation in the mind-like behavior of artifacts *Acta Biotheoretica* **11** 147–56
- [182] McDermott D 2012 Response to the singularity by David Chalmers *J. Consciousness Studies* **19** 167–72
- [183] McGinnis J O 2010 Accelerating AI *Northwestern University Law Rev.* **104** 1253–70 [www.law.northwestern.edu/lawreview/v104/n3/1253/LR104n3McGinnis.pdf](http://www.law.northwestern.edu/lawreview/v104/n3/1253/LR104n3McGinnis.pdf)
- [184] McKibben B 2003 *Enough* (New York: Henry Holt)
- [185] McLaren B M 2003 Extensionally defining principles and cases in ethics *Artificial Intell.* **150** 145–81
- [186] McLaren B M 2006 Computational models of ethical reasoning *IEEE Intell. Syst.* **21** 29–37
- [187] McLeod P, Plunkett K and Rolls E T 1998 *Introduction to Connectionist Modelling of Cognitive Processes* (New York: Oxford University Press)
- [188] Hans Meuer H, Strohmaier E, Dongarra J and Simon H 2012 Top500 list—November 2012 [www.top500.org/list/2012/11/](http://www.top500.org/list/2012/11/)
- [189] Miller J D 2012 *Singularity Rising* (Dallas, TX: BenBella)
- [190] Minsky M, Singh P and Sloman A 2004 The St Thomas common sense symposium *AI Mag.* **25** 113–24
- [191] Moore D, Shannon C and Brown J 2002 Code-red *Proc. Second ACM SIGCOMM Workshop on Internet Measurement (IMW '02)* (New York: ACM) pp 273–84
- [192] Moore D, Paxson V, Savage S, Shannon C, Staniford S and Weaver N 2003 Inside the slammer worm *IEEE Security Privacy Mag.* **1** 33–9
- [193] Moravec H P 1988 *Mind Children* (Cambridge, MA: Harvard University Press)
- [194] Moravec H P 1992 Pigs in cyberspace [www.frc.ri.cmu.edu/~hpm/project.archive/general.articles/1992/CyberPigs.html](http://www.frc.ri.cmu.edu/~hpm/project.archive/general.articles/1992/CyberPigs.html)
- [195] Moravec H P 1998 When will computer hardware match the human brain? *J. Evolution Technol.* **1** [www.transhumanist.com/volume1/moravec.htm](http://www.transhumanist.com/volume1/moravec.htm)
- [196] Moravec H P 1999 *Robot* (New York: Oxford University Press)
- [197] Moskowitz G B, Li P and Kirk E R 2004 The implicit volition model *Adv. Experimen. Social Psychol.* **36** 317–413
- [198] Muehlhauser L 2012 *Less Wrong* [www.lesswrong.com/lw/fmf/overconfident\\_pessimism/](http://www.lesswrong.com/lw/fmf/overconfident_pessimism/)
- [199] Muehlhauser L and Helm L 2012 The singularity and machine ethics *Singularity Hypotheses* ed A H Eden, J H Moor, J H Soraker and E Steinhart (Berlin: Springer)
- [200] Muehlhauser L and Salamon A 2012 Intelligence explosion *Singularity Hypotheses* ed A H Eden, J H Moor, J H Soraker and E Steinhart (Berlin: Springer)
- [201] Murphy R and Woods D D 2009 Beyond Asimov *IEEE Intell. Syst.* **24** 14–20
- [202] Napier W 2011 Hazards from comets and asteroids *Global Catastrophic Risks* ed N Bostrom and M M Cirkovic (Oxford: Oxford University Press) pp 222–37
- [203] Narayanan A and Shmatikov V 2008 Robust de-anonymization of large sparse datasets *2008 IEEE Symposium on Security and Privacy* (Oakland, CA: IEEE Computer Society) pp 111–25
- [204] Narayanan A and Shmatikov V 2009 De-anonymizing social networks *30th IEEE Symp. on Security and Privacy* (Berkeley, CA: IEEE Computer Society) pp 173–87
- [205] Narayanan A and Shmatikov V 2009 De-anonymizing social networks [www.cs.utexas.edu/~shmat/socialnetworks-faq.html](http://www.cs.utexas.edu/~shmat/socialnetworks-faq.html)
- [206] Narayanan A, Paskov H, Gong N Z, Bethencourt J, Stefanov E, Shin E C R and Song D 2012 On the

- feasibility of internet-scale author identification 2012 *IEEE Symp. on Security and Privacy* (Oakland, CA: IEEE Computer Society) pp 300–14
- [207] Nielsen T D and Jensen F V 2004 Learning a decision maker's utility function from (possibly) inconsistent behavior *Artificial Intell.* **160** 53–78
- [208] Nordmann A 2008 Singular simplicity *IEEE Spectrum* <http://spectrum.ieee.org/robotics/robotics-software/singular-simplicity>
- [209] Nordmann A 2007 If and then *NanoEthics* **1** 31–46
- [210] Olson M 1982 *The Rise and Decline of Nations* (New Haven, CT: Yale University Press)
- [211] Omohundro S M 2007 The nature of self-improving artificial intelligence <http://selfawareness.com/2007/10/05/paper-on-the-nature-of-self-improving-artificial-intelligence/>
- [212] Omohundro S M 2008 The basic AI drives *Artificial General Intelligence Frontiers* ed P Wang, B Goertzel and S Franklin (Amsterdam: IOS) pp 483–92
- [213] Omohundro S M 2012 Rational artificial intelligence for the greater Good *Singularity Hypotheses* ed A H Eden, J H Moor, J H Soraker and E Steinhart (Berlin: Springer)
- [214] Orseau L and Ring M 2011 Self-modification and mortality in artificial agents *Artificial General Intelligence* ed J Schmidhuber, K R Thirsson and M Looks (New York: Springer) pp 1–10
- [215] Persson I and Savulescu J 2008 The perils of cognitive enhancement and the urgent imperative to enhance the moral character of humanity *J. Appl. Philosophy* **25** 162–77
- [216] Persson I and Savulescu J 2012 *Unfit for the Future* (Oxford: Oxford University Press)
- [217] Peterson N R, Pisoni D B and Miyamoto R T 2010 Cochlear implants and spoken language processing abilities *Restorative Neurology and Neuroscience* **28** 237–50
- [218] Pinker S 2002 *The Blank Slate* (New York: Viking)
- [219] Plaut D C 2003 Connectionist modeling of language *Mind, Brain, and Language* ed M T Banich and M Mack pp 143–68 (Mahwah, NJ: Lawrence Erlbaum)
- [220] Posner R A 2004 *Catastrophe* (New York: Oxford University Press)
- [221] Potapov A and Rodionov S Universal empathy and ethical bias for artificial general intelligence [http://aideus.com/research/doc/preprints/04\\_paper4\\_AGIImpacts12.pdf](http://aideus.com/research/doc/preprints/04_paper4_AGIImpacts12.pdf)
- [222] Powers T M 2006 Prospects for a Kantian machine *IEEE Intell. Syst.* **21** 46–51
- [223] Powers T M 2011 Incremental machine ethics *IEEE Robotics Automation Mag.* **18** 51–8
- [224] Pynadath D V and Tambe M 2002 Revisiting Asimov's first law *Intelligent Agents VIII* ed J-J Ch Meyer and M Tambe (Berlin: Springer) pp 307–20
- [225] Railton P 1986 Facts and values *Phil. Topics* **14** 5–31
- [226] Rajab M A, Zarfoss J, Monroe F and Terzis A 2007 My botnet is bigger than yours (maybe, better than yours) *Proc. of 1st Workshop on Hot Topics in Understanding Botnets (HotBots '07)* (Berkeley, CA: USENIX) [http://static.usenix.org/event/hotbots07/tech/full\\_papers/rajab/rajab.pdf](http://static.usenix.org/event/hotbots07/tech/full_papers/rajab/rajab.pdf)
- [227] Ramamurthy U, Baars B J, D'Mello S K and Franklin S 2006 LIDAProc. *Seventh International Conference on Cognitive Modeling* ed D Fum, F Del Missier and A Stocco (Trieste: Edizioni Goliardiche) pp 244–9 <http://ccrg.cs.memphis.edu/assets/papers/ICCM06-UR.pdf>
- [228] Reynolds C and Cassinelli A (ed) 2009 *AP-CAP 2009* <http://kant.k2.t.u-tokyo.ac.jp/ap-cap09/proceedings.pdf>
- [229] Ring M and Orseau L 2011 Delusion, survival, and intelligent agents ed J Schmidhuber, K R Thirsson and M Looks (New York: Springer) pp 11–20
- [230] Salekin R T 2010 *Treatment of child and adolescent psychopathy Handbook of Child and Adolescent Psychopathy* ed R T Salekin and D R Lynam (New York: Guilford) pp 343–73
- [231] Sandberg A 2009 An overview of models of technological singularity <http://agi-conf.org/2010/wp-content/uploads/2009/06/agi10singmodels2.pdf>
- [232] Sandberg A 2001 Friendly superintelligence [www.aleph.se/Nada/Extro5/FriendlySuperintelligence.htm](http://www.aleph.se/Nada/Extro5/FriendlySuperintelligence.htm)
- [233] Sandberg A 2012 Models of a singularity *Singularity Hypotheses* ed A H Eden, J H Moor, J H Soraker and E Steinhart (Berlin: Springer)
- [234] Sandberg A and Bostrom N 2008 Whole brain emulation *Technical Report 2008-3* (Future of Humanity Institute, University of Oxford) [www.fhi.ox.ac.uk/wp-content/uploads/brain-emulation-roadmap-report1.pdf](http://www.fhi.ox.ac.uk/wp-content/uploads/brain-emulation-roadmap-report1.pdf)
- [235] Sandberg A and Bostrom N 2011 Machine intelligence survey *Technical Report 2011-1* (Future of Humanity Institute, University of Oxford) [www.fhi.ox.ac.uk/reports/2011-1.pdf](http://www.fhi.ox.ac.uk/reports/2011-1.pdf)
- [236] Schmidhuber J 2009 Ultimate cognition à la Gödel *Cogn. Comput.* **1** 177–93
- [237] Schmidhuber J, Thirsson K R and Looks M (ed) 2011 *Artificial General Intelligence (Lecture Notes in Computer Science no 6830)* (Berlin: Springer)
- [238] Scott J C 1998 *Seeing Like a State* (New Haven, CT: Yale University Press)
- [239] Searle J R 1992 *The Rediscovery of the Mind* (Cambridge, MA: MIT Press)
- [240] Shachtman N 2007 Robot cannon kills 9, wounds 14 *Wired* [www.wired.com/dangerroom/2007/10/robot-cannon-ki/](http://www.wired.com/dangerroom/2007/10/robot-cannon-ki/)
- [241] Shulman C Arms control and intelligence explosions
- [242] Shulman C 2010 Whole brain emulation and the evolution of superorganisms <http://intelligence.org/files/WBE-Superorgs.pdf>
- [243] Shulman C 2010 Omohundro's basic AI drives and catastrophic risks <http://intelligence.org/files/BasicAIDrives.pdf>
- [244] Shulman C and Sandberg A 2010 Implications of a software-limited singularity *ECAP10, VIII European Conference of Computing and Philosophy* ed K Mainzer (Munich: Dr Hut)
- [245] Shulman C, Jonsson H and Tarleton N 2009 Which consequentialism? Machine ethics and moral divergence *AP-CAP 2009* ed C Reynolds and A Cassinelli pp 23–5 <http://kant.k2.t.u-tokyo.ac.jp/ap-cap09/proceedings.pdf>
- [246] Shulman C, Jonsson H and Tarleton N 2009 Machine ethics and superintelligence *AP-CAP 2009* ed C Reynolds and A Cassinelli pp 95–7 <http://kant.k2.t.u-tokyo.ac.jp/ap-cap09/proceedings.pdf>
- [247] Smith M 2009 Desires, values, reasons, and the dualism of practical reason *Ratio* **22** 98–125
- [248] Snider J, McCall R and Franklin S 2011 The LIDA framework as a general tool for AGI *Artificial General Intelligence (Lecture Notes in Computer Science no 6830)* ed J Schmidhuber, K R Thirsson and M Looks (New York: Springer) pp 133–42
- [249] Sobel D 1994 Full information accounts of well-being *Ethics* **104** 784–810
- [250] Sobolewski M 2012 German Cabinet to agree tougher rules on high-frequency trading (Reuters) <http://in.reuters.com/article/2012/09/25/germany-bourse-rules-idINL5E8KP8BK20120925>
- [251] Solomon R J 1985 The time scale of artificial intelligence *Human Syst. Management* **5** 149–53
- [252] Solum L B 1992 Legal personhood for artificial intelligences *North Carolina Law Rev.* **70** 1231–87 [http://papers.ssrn.com/sol3/papers.cfm?abstract\\_id=1108671](http://papers.ssrn.com/sol3/papers.cfm?abstract_id=1108671)
- [253] Sotala K 2012 Advantages of artificial intelligences, uploads, and digital minds *Int. J. Machine Consciousness* **4** 275–91
- [254] Sotala K and Valpola H 2012 Coalescing minds *Int. J. Machine Consciousness* **4** 293–312

- [255] Spears D F 2006 Assuring the behavior of adaptive agents ed C Rou, M Hinchey, J Rash, W Truszkowski and D F Gordon-Spears *Agent Technology from a Formal Perspective (NASA Monographs in Systems and Software Engineering)* (New York: Springer) pp 227–57
- [256] Stahl B C 2002 Can a computer adhere to the categorical imperative? A contemplation of the limits of transcendental ethics in IT *14th Int. Conf. on Systems Research, Informatics and Cybernetics: Symposium on Cognitive, Emotive and Ethical Aspects of Decision Making in Humans and in Artificial Intelligence* (Windsor, ON: International Institute for Advanced Studies in Systems Research/Cybernetics) pp 13–8
- [257] Staniford S, Paxson V and Weaver N 2002 How to own the internet in your spare time *Proc. 11th USENIX Security Symp.* ed D Boneh pp 149–67 [www.icir.org/vern/papers/cdc-usenix-sec02/](http://www.icir.org/vern/papers/cdc-usenix-sec02/)
- [258] Steunebrink B R and Schmidhuber J 2011 A family of Gödel machine implementations *Artificial General Intelligence (Lecture Notes in Computer Science no 6830)* ed J Schmidhuber, K R Thirsson and M Looks (New York: Springer) pp 275–80
- [259] Suber P 2002 Saving machines from themselves [www.earlham.edu/~peters/writing/selfmod.htm](http://www.earlham.edu/~peters/writing/selfmod.htm)
- [260] Sullins J P 2005 Ethics and artificial life *Ethics Inf. Technol.* **7** 139–48
- [261] Sullins J P 2006 When is a robot a moral agent? *Int. Rev. Inf. Ethics* **6** 23–30
- [262] Stillwaggon Swan L and Howard J 2012 Digital immortality *Int. J. Machine Consciousness* **4** 245–56
- [263] Sweeney L 1997 Weaving technology and policy together to maintain confidentiality *J. Law Med. Ethics* **25** 98–110
- [264] Tanyi A 2006 An essay on the desire-based reasons model [http://web.ceu.hu/polsci/dissertations/Attila\\_Tanyi.pdf](http://web.ceu.hu/polsci/dissertations/Attila_Tanyi.pdf)
- [265] Tarleton N 2010 Coherent extrapolated volition <http://intelligence.org/files/CEV-MachineEthics.pdf>
- [266] Tenenbaum J B, Griffiths T L and Kemp C 2006 Theory-based Bayesian models of inductive learning and reasoning *Trends Cogn. Sci.* **10** 309–18
- [267] Thomas M S C and McClelland J L 2008 Connectionist models of cognition *The Cambridge Handbook of Computational Psychology (Cambridge Handbooks in Psychology)* ed R Sun (Cambridge: Cambridge University Press) pp 23–58
- [268] Trope Y and Liberman N 2010 Construal-level theory of psychological distance *Psychological Rev.* **117** 440–63
- [269] Turing A M 1951 Intelligent machinery, a heretical theory
- [270] Turney P 1991 Controlling super-intelligent machines *Can. Artificial Intell.* **27** 3–35
- [271] Tversky A and Kahneman D 1981 The framing of decisions and the psychology of choice *Science* **211** 453–58
- [272] van Gelder T 1995 What might cognition be, if not computation? *J. Philosophy* **92** 345–81
- [273] van Kleef G A, Oveis C, van der Lwe I, LuoKogan A, Goetz J and Keltner D 2008 Power, distress, and compassion *Psychological Sci.* **19** 1315–22
- [274] van Kleef G A, Homan A C, Finkenauer C, Gündemir S and Stamkou E 2011 Breaking the rules to rise to power *Social Psychological Personality Sci.* **2** 500–7
- [275] Verdoux P 2010 Risk mysterianism and cognitive boosters *J. Futures Studies* **15** 1–20 ([www.jfs.tku.edu.tw/15-1/A01.pdf](http://www.jfs.tku.edu.tw/15-1/A01.pdf))
- [276] Verdoux P 2011 Emerging technologies and the future of philosophy *Metaphilosophy* **42** 682–707
- [277] Versenyi L 1974 Can robots be moral? *Ethics* **84** 248–59
- [278] Vinge V 1993 The coming technological singularity *Vision-21 (NASA Conference Publication no 10129)* (NASA Lewis Research Center) pp 11–22 ([http://ntrs.nasa.gov/archive/nasa/casi.ntrs.nasa.gov/19940022855\\_1994022855.pdf](http://ntrs.nasa.gov/archive/nasa/casi.ntrs.nasa.gov/19940022855_1994022855.pdf))
- [279] Walker M 2008 Human extinction and farsighted universal surveillance [www.nmsu.edu/~philos/documents/sept-2008-smart-dust-final.doc](http://www.nmsu.edu/~philos/documents/sept-2008-smart-dust-final.doc)
- [280] Walker M 2011 Personal identity and uploading *J. Evolution Technol.* **22** 37–51
- [281] Wallach W 2010 Robot minds and human ethics *Ethics Inf. Technol.* **12** 243–50
- [282] Wallach W and Allen C 2009 *Moral Machines* (Oxford: Oxford University Press)
- [283] Wallach W and Allen C 2012 Framing robot arms control *Ethics Inf. Technol.* **15**
- [284] Wallach W, Allen C and Smit I 2008 Machine morality *AI Society* **22** 565–82
- [285] Wallach W, Franklin S and Allen C 2010 A conceptual and computational model of moral decision making in human and artificial agents *Topics Cogn. Sci.* **2** 454–85
- [286] Wallach W, Allen C and Franklin S 2011 Consciousness and ethics *Int. J. Machine Consciousness* **3** 177–92
- [287] Wang P 2012 Motivation management in AGI systems 2012 *Artificial General Intelligence (Lecture Notes in Artificial Intelligence no 7716)* ed J Bach, B Goertzel and M Ikl (New York: Springer) pp 352–61
- [288] Wang P, Goertzel B and Franklin S (ed) 2008 *Artificial General Intelligence Frontiers (Artificial Intelligence and Applications no 171)* (Amsterdam: IOS)
- [289] Warwick K 1998 *In the Mind of the Machine* (London: Arrow)
- [290] Warwick K 2003 Cyborg morals, cyborg values, cyborg ethics *Ethics Inf. Technol.* **5** 131–7
- [291] Warwick K 2010 Implications and consequences of robots with biological brains *Ethics Inf. Technol.* **12** 223–4
- [292] Waser M R 2008 Discovering the foundations of a universal system of ethics as a road to safe artificial intelligence *Biologically Inspired Cognitive Architectures Technical Report FS-08-04* pp 195–200 ([www.aaai.org/Papers/Symposia/Fall/2008/FS-08-04/FS08-04-049.pdf](http://www.aaai.org/Papers/Symposia/Fall/2008/FS-08-04/FS08-04-049.pdf))
- [293] Waser M R 2009 A safe ethical system for intelligent machines *Biologically Inspired Cognitive Architectures* ed A V Samsonovich pp 194–9 ([www.aaai.org/ocs/index.php/FSS/FSS09/paper/view/934](http://www.aaai.org/ocs/index.php/FSS/FSS09/paper/view/934))
- [294] Waser M R 2011 Rational universal benevolence *Artificial General Intelligence* ed J Schmidhuber, K R Thirsson and M Looks (New York: Springer) pp 153–62
- [295] Weld D and Etzioni O 1994 The first law of robotics (a call to arms) *Proc. Twelfth National Conf. on Artificial Intelligence* ed B Hayes-Roth and R E Korf (Menlo Park, CA: AAAI Press) pp 1042–7 ([www.aaai.org/Papers/AAAI/1994/AAAI94-160.pdf](http://www.aaai.org/Papers/AAAI/1994/AAAI94-160.pdf))
- [296] Weng Y-H, Chen C-H and Sun C-T 2008 Safety intelligence and legal machine language *Service Robot Applications* ed Y Takahashi (InTech)
- [297] Weng Y-H, Chen C-H and Sun C-T 2009 Toward the human–robot co-existence society *Int. J. Social Robotics* **1** 267–82
- [298] Whitby B 1996 *Reflections on Artificial Intelligence* (Exeter: Intellect)
- [299] Whitby B and Oliver K 2000 How to avoid a robot takeover [www.sussex.ac.uk/Users/blayw/BlayAISB00.html](http://www.sussex.ac.uk/Users/blayw/BlayAISB00.html)
- [300] Wiener N 1960 Some moral and technical consequences of automation *Science* **131** 1355–8
- [301] Wilson G S 2014 Minimizing global catastrophic and existential risks from emerging technologies through international law *Virginia Environ Law J.* **31** 307–64
- [302] Wilson T D 2002 *Strangers to Ourselves* (Cambridge, MA: Belknap)
- [303] Yampolskiy R V 2012 Leakproofing the singularity: artificial intelligence confinement problem *J. Consciousness Studies* **1** 194–214

- [304] Yampolskiy R V 2013 Artificial intelligence safety engineering *Philosophy and Theory of Artificial Intelligence (Studies in Applied Philosophy, Epistemology and Rational Ethics* vol 5) (New York: Springer) pp 389–96
- [305] Yampolskiy R V and Fox J 2012 Safety engineering for artificial general intelligence *Topoi* doi:[10.1007/s11245-012-9128-9](https://doi.org/10.1007/s11245-012-9128-9)
- [306] Yudkowsky E 2001 Creating friendly AI 1.0 <http://intelligence.org/files/CFAI.pdf>
- [307] Yudkowsky E 2004 Coherent extrapolated volition <http://intelligence.org/files/CEV.pdf>
- [308] Yudkowsky E 2008 Hard takeoff *Less Wrong* [http://lesswrong.com/lw/wf/hard\\_takeoff/](http://lesswrong.com/lw/wf/hard_takeoff/)
- [309] Yudkowsky E 2009 Value is fragile *Less Wrong* [http://lesswrong.com/lw/y3/value\\_is\\_fragile/](http://lesswrong.com/lw/y3/value_is_fragile/)
- [310] Yudkowsky E 2012 Reply to Holden on tool AI *Less Wrong* [http://lesswrong.com/lw/cze/reply\\_to\\_holden\\_on\\_tool\\_ai/](http://lesswrong.com/lw/cze/reply_to_holden_on_tool_ai/)
- [311] Yudkowsky E 1996 Staring into the singularity <http://yudkowsky.net/obsolete/singularity.html>
- [312] Yudkowsky E 2011 Artificial intelligence as a positive and negative factor in global risk *Global Catastrophic Risks* ed N Bostrom and M M Cirkovic (Oxford: Oxford University Press) pp 308–45
- [313] Yudkowsky E 2011 Complex value systems are required to realize valuable futures <http://intelligence.org/files/ComplexValues.pdf>
- [314] Zimmerman D 2003 Why Richard Brandt does not need cognitive psychotherapy, and other glad news about idealized preference theories in meta-ethics *J. Value Inquiry* **37** 373–94