

1 *Article*

## 2 **Safe Artificial General Intelligence via Distributed** 3 **Ledger Technology**

4 **Kristen W. Carlson**

5 Department of Neurosurgery, Neurosimulation Group

6 Beth Israel Deaconess Medical Center/Harvard Medical School

7 [kwcarlso@bidmc.harvard.edu](mailto:kwcarlso@bidmc.harvard.edu)

8

9

10 **Abstract.** Artificial general intelligence (AGI) progression metrics indicate AGI will occur within  
11 decades. No proof exists that AGI will benefit humans and not harm or eliminate humans. I propose  
12 a set of logically distinct conceptual components that are necessary and sufficient to 1) ensure  
13 various AGI scenarios will not harm humanity and 2) robustly align AGI and human values and  
14 goals. By systematically addressing pathways to malevolent AI we can induce the methods/axioms  
15 required to redress them. Distributed ledger technology (DLT, 'blockchain') is integral to this  
16 proposal, e.g. 'smart contracts' are necessary to address evolution of AI that will be too fast for  
17 human monitoring and intervention. The proposed axioms: 1) Access to technology by market  
18 license. 2) Transparent ethics embodied in DLT. 3) Morality encrypted via DLT. 4) Behavior control  
19 structure with values at roots. 5) Individual bar-code identification of critical components. 6)  
20 Configuration Item (from business continuity/disaster recovery planning). 7) Identity verification  
21 secured via DLT. 8) 'Smart' automated contracts based on DLT. 9) Decentralized applications - AI  
22 software modules encrypted via DLT. 10) Audit trail of component usage stored via DLT. 11) Social  
23 ostracism (denial of resources) augmented by DLT petitions. 12) Game theory and mechanism  
24 design.

25

26 **Keywords:** Artificial general intelligence, AGI, blockchain, distributed ledger, AI containment, AI  
27 safety, AI value alignment, ASILOMAR

28

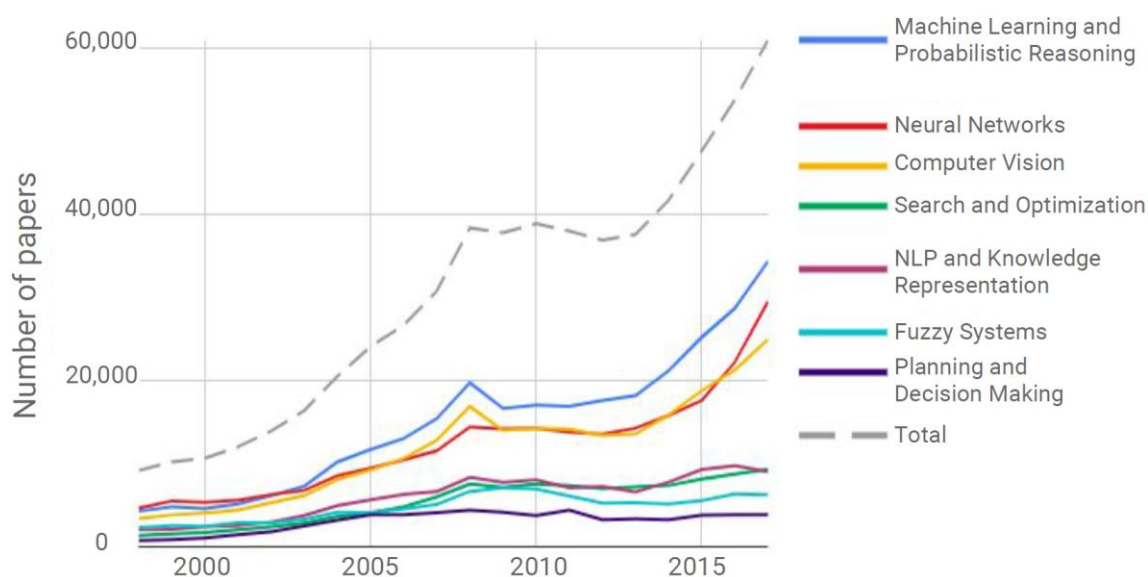
29

30

## 31 1. Introduction

32 The problem of superhuman artificial intelligence ('artificial general intelligence', AGI) harming  
 33 or eradicating humankind is an increasing concern as the prospect of AGI nears. Current attempts to  
 34 measure AI progress show exponential growth in activity globally and technical improvement across  
 35 the board of functionality measured – including 'Human-Level Performance Milestones' [1] (Fig.  
 36 1a). Recent watershed advances include Deep Mind beating the most expert human at the complex  
 37 game of Go – which averages 250 moves per position and 150 moves per game =  $10^{359}$  possible paths  
 38 vs. chess, which averages 35 moves per position and 80 moves per game =  $10^{123}$  possible paths, *and a*  
 39 *decade earlier than expected*. Deep Mind used a neural network to assign a value at each point in a  
 40 decision tree and discarded low-valued lower-level branches and thus avoided the exponential  
 41 search required to explore them. Human Go experts assigned high creativity to Deep Mind's  
 42 strategies and tactics. A second major AI development was Deep Mind's self-teaching, reinforcement  
 43 learning ability, playing tens of thousands of games against itself in a few hours rather than  
 44 incorporating human game-play strategies and eliminating its need for human feedback [2].

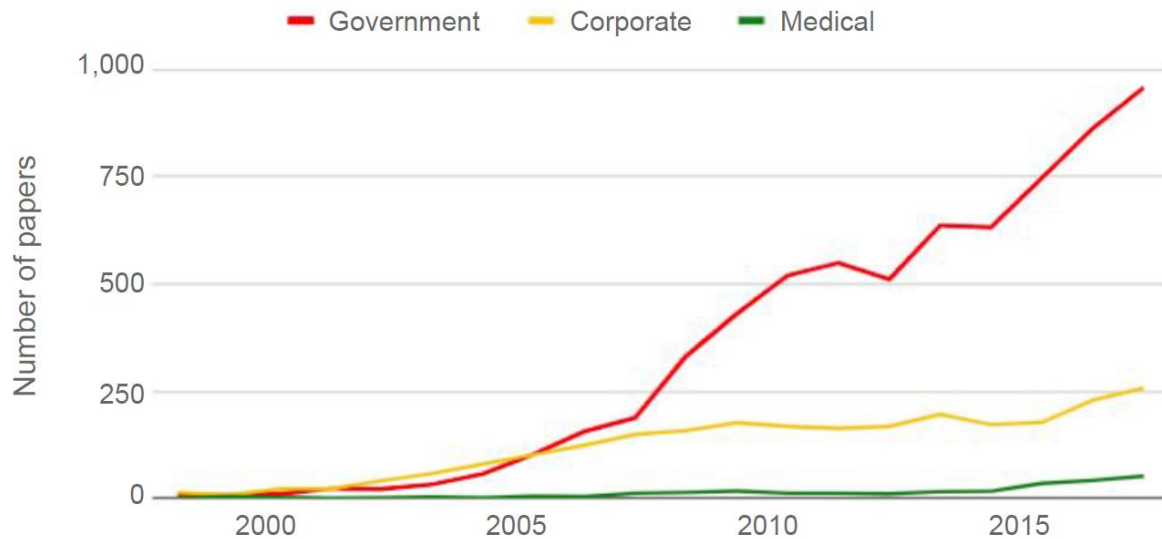
45 Collaborating, self-taught AIs played 180 human years of games per day using new  
 46 reinforcement learning policy optimization algorithms and beat human teamwork in the simulated  
 47 real-world environment of Dota2 [3] (video: <https://youtu.be/Ub9INopwJ48>). Significant advances  
 48 were made in credit assignment to short-term vs. long-term goals and learning the optimal balance  
 49 between individual and team performance. Another watershed occurred when AI beat humans at an  
 50 'imperfect information' game, poker – i.e. the opponents' hands are hidden, fundamentally different  
 51 from Go or chess –using game theory techniques including bluffing, previously thought to be  
 52 difficult to emulate [4, 5]. Such techniques could be used to beat humans in business strategy,  
 53 negotiation, strategic pricing, finance, cybersecurity, physical security, military, auctions, political  
 54 campaigns, and medical treatment planning [4]. AI continues to reach new levels of unsupervised  
 55 learning prowess (pattern recognition without human guidance), e.g. for parsing handwritten letters  
 56 and creating new letters that pass a specialized Turing test, and more efficiently than deep learning  
 57 networks [6]. AI superiority over humans in general background knowledge and parsing natural  
 58 language is old news [7], and is now being embedded in all human-computer interfacing ('powered  
 59 by Watson', Alexa, Siri, Cortana, Google Assistant, et al.), whose potential monetary value has  
 60 triggered a commercial AI arms race in parallel with a military/political one (Fig. 1) [8].



61

62

Figure 1a. Number of AI papers in Scopus by sub-category (1998–2017).



63  
64 **Figure 1b.** Papers by sector affiliation – China (1998–2017). Source: Elsevier [1].

65 Bostrom gives examples of general intelligence skills where attainment of *any* of them would  
66 trigger AGI dominance over humans (reproduced in Table 1). One such epochal AI development that  
67 could trigger the AGI singularity is the prospect of AI learning to program itself – ‘recursive self-  
68 improvement’ (*q.v.* ASILOMAR AI Principle #22, see also #19, #20, #21 [9]) – which opens a door to  
69 a positive-feedback-driven process in which AGI vastly exceeds human capabilities in short order  
70 and may change its human-instilled directives. An AGI could begin to regard humanity as a trivial,  
71 primitive nuisance, competing for vital resources required for attainment of its goals, distinct from  
72 humanity’s, stemming from alien values, as we regard mosquitoes or flies.

73 **Table 1.** Examples of super-intelligent skill sets triggering AGI world domination (from Bostrom [10];  
74 *cf.* Babcock et al. Sec. 6.2 [11]).

Intelligence amplification – AI can improve its own intelligence
Strategy – optimizing chances of achieving goals using advanced techniques, e.g. game theory, cognitive psychology, and simulation
Social manipulation – psychological and social modeling e.g. for persuasion
Hacking – exploiting security flaws to appropriate resources
R&D – create more powerful technology, e.g. to achieve ubiquitous surveillance and military dominance
Economic productivity – generate vast wealth to acquire resources

75  
76 A danger many feared would accelerate the timeline to AGI via ‘Red Queen’ cultural co-  
77 evolution [12], an AI arms race has begun, driven by the increasing realization in political and  
78 military circles that AI is the key to future military superiority [13, 14]. Thus ASILOMAR #5 & #18  
79 may already be violated [9]. The race increases emphasis on AI for intentionally destructive purposes  
80 and likely will result in less control of AI technology by its creators [15]. It is an ominous development  
81 as all nuclear powers upgrade their arsenals, proliferation increases, and arms control agreements  
82 are unraveling [16]. The day when AI is consulted and decides if ‘no first strike’ commitments or  
83 reducing ‘high alert’ status nuclear weapons is beneficial or perceived as a vulnerable weakness by  
84 adversaries looms ahead.

85 The potential speed with which AGI could advance from being human-directed and empathetic  
86 of humans to evolving beyond human-level concerns is unknown; with self-programming ability or  
87 other internal intelligence enhancement [10, 11] positive feedback will trigger super-exponential  
88 growth. At that point a malevolent AGI may arise within a fraction of a second, too fast for us to  
89 detect and respond [17].

90 What is proposed here is a complete AGI ecosystem, framed as a set of axioms at a relatively  
91 high systems level, that will ensure AGI-human value alignment, and thereby ensure benevolent AGI  
92 behavior, as seen by humans and successive generations of AGI. Notably, the axioms incorporate  
93 distributed ledger technology and smart contracts to automate and prevent corruption of many  
94 required processes.

## 95 2. Methods

### 96 2.1. To Generate a Necessary and Sufficient Set of Axioms

97 There are several taxonomies of pathways to dangerous AI, such as Yampolskiy [18], Turchin  
98 [19], Bostrom [10], and Brundage et al. [20]. These taxonomies are a reasonable starting point for  
99 systematically investigating how to ensure safe AGI. One can take each pathway to danger as a  
100 theorem and induce methods, formalized as axioms, toward generating a necessary and sufficient set  
101 of axiom-methods to eliminate all pathways or reduce their probability. Pathway categories overlap,  
102 which helps ensure redundancy in capturing the necessary and sufficient axioms to redress all  
103 categories.

104 Similarly, as one iterates the process of using each dangerous pathway to generate a complete  
105 set of axioms to address it, some axioms repeat, while some pathways require new, additional axioms  
106 until at the end of the pathways list, most are covered by the axiom set, although some pathways  
107 may be left without sufficient methods to eliminate them. For the pathways itemized in the  
108 taxonomies, the resulting axioms seem to be the minimal set for ensuring safe AGI. By 'ensuring' I  
109 mean optimally reducing the probability of a dangerous pathway manifesting.

110 Stating a set of axioms is a necessary step toward formal proof of a necessary, sufficient, and  
111 minimal set — if a formal proof is possible. Yampolskiy concludes his taxonomy by saying that formal  
112 proof of the completeness of a taxonomy is important [18] and formal methods are a main theme of  
113 Omohundro [8]. Short of a tight logical proof, probabilistically assuring benevolent AGI, e.g. through  
114 extensive simulations, may be the realistic route best to take, and must accompany any set of safety  
115 measures, including those proposed here.

116 An important way to test if each axiom is necessary is to find failure use cases when it is omitted  
117 [21]; examples are given below.

### 118 2.2. Ingredients for Formalization of AGI Safety Theory

119 Toward formalization I attempt to make the various methods to ensure safe AGI logically  
120 distinct, state them as axioms, and at a high level intended to capture concisely a necessary and  
121 sufficient set. This usage of 'axiom' generalizes that of von Neumann where certain lower systems  
122 level outputs or theorems are 'axiomatized' — seen as black boxes, or input-output specification, or  
123 logic tables — at the immediately higher systems level [22]. Each axiom is most precisely expressed  
124 by an *operational* definition specified by an algorithm implementing it, hence, a method.

125 For instance, the definition of subjective value or utility, used in the morality and game theory  
126 axioms below, is made precise by the six von Neumann-Morgenstern utility axioms [23]. As stated  
127 below, a set axioms designed, and proven via simulation, to induce cooperation among extremely  
128 diverse, complex agents may replace most of the set given herein; the simulations of Burtsev and  
129 Turchin may be prolegomena [24].

130 A problem we frequently face in modeling and simulation is: What is the highest systems level  
131 that can concisely describe and emulate the target set of phenomena? Thus, a limitation in axiomatic  
132 formulations is they leave varying amounts of implementation detail at the systems level underlying  
133 them to be specified, or to some degree, developed. For example, the DLT-based axioms 2, 4, 5, 7, 8,  
134 9, 10, and 11, are in rapid evolution toward algorithmic implementation to address diverse use cases.  
135 And behavior control (axiom 4) is in rapid development in some contexts (e.g. autonomous vehicles,  
136 factory robots), yet the degree of development still needed to align human and AGI values may be  
137 significant.

138 In another attempt to formalize the expression of AGI dangers, I offer some simple syllogisms  
139 (Appendix 1).

140 The concept of AGI-completeness, akin to NP-completeness as stated by Bostrom [10] is that a  
141 demonstration of one technology, e.g., self-improvement techniques, engendering AGI is sufficient  
142 to demonstrate that capability for a class of AI technologies, AGI-completeness may be another piece  
143 of formalizing AGI, the measures of its progress, and specifying the point of no containment unless  
144 sufficient preparations have been made.

145 Another means to formalize AGI theory is Omohundro's idea of deriving universal AGI drives  
146 from first principles [8], which can be explored to see if such drives emerge in simulations as well as  
147 via logical derivation. Omohundro argues that universal drives will inevitably lead to conflict of AI  
148 and human values from the irrefutable economic axiom of competition for resources.

149 Another formalization route is calculating the probability of hacking a blockchain against the  
150 number of AGIs required to reach consensus via the blockchain to permit unlocking the next AGI  
151 generation (see sections on decentralized apps and the Singleton problem below). This calculation is  
152 similar to the math underlying the internet's redundancy in average interconnectedness of nodes and  
153 global system fault-tolerance [25] but more complicated since it involves Byzantine fault tolerance,  
154 wherein two diagnostic agents disagree on the nature of the fault [26]. The inclusion of innovative  
155 DLT into the algorithms should permit AGI robustness to surpass the 'robust yet fragile' use case of  
156 the internet that is vulnerable to targeted attacks on the most interconnected nodes.

157 Last, it may be possible to subsume several of the axioms herein via a game theory/economics  
158 set proven via simulation. An obstacle to this approach is that game-theoretic algorithms that  
159 simulate interactions between entities with behavior expressiveness vastly larger than our own [24]  
160 may be necessary to understand and predict AGI social behavior but may also be computationally  
161 intractable (see Diversity in the AGI Ecosystem, below).

### 162 3. Results and Discussion

163 Regarding the term AI 'containment', Babcock et al. suggest that 'containment' is an appropriate  
164 term for methodologies for controlled AGI development and safety-testing rather than control over  
165 entities whose intelligence will exceed our own [11]. The current work is intended to contribute to  
166 both phases.

#### 167 3.1. A Critical Ingredient: Distributed Ledger Technology (aka 'Blockchain')

168 The recent innovation of distributed digital ledger technology (DLT) is critical to this proposal  
169 [27]. The crux of DLT is an audit trail database, in which each addition is validated by a pluralistic  
170 consensus, currently performed by humans operating computers that run hash and anti-hash  
171 functions (to wit public key encryption), stored on a distributed network also known as a blockchain:  
172 "Blockchains allow us to have a distributed peer-to-peer network where non-trusting members can  
173 interact with each other without a trusted intermediary, in a verifiable manner" [28]. Key aspects of  
174 DLT are shown in Table 2 [29] (other auxiliary DLT aspects, such as anonymity of participants, are  
175 either not necessary or not beneficial in the context of ensuring safe AGI). The 'smart' automated  
176 contract vision of Szabo [30], encrypted redundantly via DLT, could comprise the core methodology  
177 whereby AGI development and evolution can be aligned with the best human values without  
178 concomitant human intervention. Notably, smart contracts can prevent the hacking of safe AGI  
179 evolution that is too fast for human response.

180 **Table 2.** Distributed ledger technology applicable to ensuring AGI safety.

Non-hackability and non-censurability via decentralization (storage in multiple distributed servers), encryption in standardized blocks, and irrevocable transaction linkage (the 'chain')
Node-fault tolerance: redundancy via storage in a decentralized ledger of a) rules for transactions, b) the transaction audit trail, and c) transaction validations
Transparency of the transaction rules and audit trail in the DLT

---

Automated 'smart' contracts

---

Decentralized applications ('dApps'), i.e. software programs that are stored and run on a distributed network and have no central point of control or failure

---

Validation of contractual transactions by a decentralized consensus of validators

---

181

182

183

184

185

186

Here are the proposed necessary and sufficient axioms to ensure safe AGI (Table 3), followed by examples of malignant AGI categories by Turchin [19], Yampolskiy [18], and Bostrom [10], in which the danger pathway is described and a subset of axioms to reduce its probability are specified (Tables 4, 5). In the malignant AI examples the game theory/mechanism design axiom is not mentioned; see comments in the axiom descriptions and elsewhere.

187

**Table 3.** Proposed axioms to ensure human-benevolent AGI.

Symbol	Axiom
1	Access to AGI technology via market license
2	Ethics transparently stored via DLT so they cannot be altered, forged or deleted
3	Morality, defined as no use of force or fraud, stored via DLT
4	Behavior control structure (e.g. a behavior tree) augmented by adding human-compatible values (axioms 2 & 3) at its roots
5	Unique hardware and software ID codes
6	Configuration Item (automated configuration)
7	Secure identity via multi-factor authentication, public-key infrastructure and DLT
8	Smart contracts based on DLT
9	Decentralized applications (dApps) – AGI software code modules encrypted via DLT
10	Audit trail of component usage stored via DLT
11	Social ostracism – denial of societal resources – augmented by petitions based on DLT
12	Game theory – mechanism design of a communications and incentive system

188

**Table 4.** Examples from Turchin and Yampolskiy Taxonomies of AGI Failure Modes [18, 19]

Stage/Pathway	Necessary Axioms See Table 1 Axioms
Sabotage.	
a. By impersonation (e.g. hacker, programmer, tester, janitor).	a. 7.
b. AI software to cloak human identity.	b. 7.

c. By someone with access.	c. 2, 3, 4, 5, 6, 8, 9, 10, 11.
Purposefully dangerous military robots and intelligent software. Robot soldiers, armies of military drones and cyber weapons used to penetrate networks and cause disruptions to the infrastructure. a. due to command error b. due to programming error c. due to intentional command by adversary or nut d. due to negligence by adversary or nut (e.g. AI nanobots start global catastrophe)	Axiom # 3, morality, does not apply where coercive force or fraud are a premise, e.g. military or police use of force, while axiom 2, ethics, in this case embodying restrictions on use of force, and 4, behavior control, and the rest, do apply. a. 1, 2, 4, 6, 8, 11 b. 2, 4, 5, 6, 8, 9, 10, 11 c. 1, 2, 4, 6, 7, 8, 10, 11 d. 1, 2, 4, 6, 7, 8, 9, 10, 11 Under some circumstances, such as if the means is already available, there is no solution (see Appendix, Proposition 1).
AI specifically designed for malicious and criminal purposes. Artificially intelligent viruses, spyware, Trojan horses, worms, etc. Stuxnet-style virus hacks infrastructure causing e.g. nuclear reactor meltdowns, power blackouts, food and drug poisoning, airline and drone crashes, large-scale geo-engineering systems failures. Home robots turning on owners, autonomous cars attack. Narrow AI bio-hacking virus. Virus starts human extinction via DNA manipulation, virus invades brain via neural interface	1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11 Under some circumstances, no solution (see Appendix, Proposition 1).
Robots replace humans. People lose jobs, money, and/or motivation to live; genetically-modified superior human-robot hybrids replace humans	No guaranteed solution from axiom set. All jobs can be replaced by AGI including science, mathematics, management, music, art, poetry, etc. Under axioms 1-3 humans could trade technology for resources with AGI in its pre-takeoff stage to ensure some type of guaranteed income.
Narrow bio-AI creates super-addictive drug. Widespread addiction and switching off of happy, productive life, e.g. social networks, fembots, wire-heading, virtual reality, designer drugs, games	1, 2, 3, 4, 7, 8, 9, 10
Nation states evolve into computer-based totalitarianism. Suppression of human values; human replacement with robots; concentration camps; killing of 'useless' people; humans become slaves; system becomes fragile to variety of other catastrophes	1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11
AI fights for survival but incapable of self-improvement	1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11
Failure of nuclear deterrence AI. a. impersonation of entity authorized to launch attack b. virus hacks nuclear arsenal or Doomsday machine	a. 7 b. 4, 6, 8, 9, 10

<p>c. creation of Doomsday machines by AI d. self-aware military AI ('Skynet')</p>	<p>c. 1, 2 (if creation of Doomsday machine is categorized as unethical), 4, 5, 6, 7, 8, 9, 10, 11 d. 1, 2, 4, 5, 6, 7, 8, 9, 11</p>
<p>Opportunity cost if strong AI is not created. Failure of global control: e.g. bioweapons created by biohackers; other major and minor risks not averted via AI control systems.</p>	<p>To create AGI with minimized risk and avoid opportunity cost need axioms 1-11</p>
<p>AI becomes malignant. AI breaks restrictions and fights for world domination (control over all resources), possibly hiding its malicious intent.</p>	<p>1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11 Note it may achieve increasing and unlimited control over resources via market transactions by convincing enough volitional entities to give it control due to potential benefits to them</p>
<p>AI deception. AI escapes from confinement; hacks its way out; copies itself into the cloud and hides that fact; destroys initial confinement facility or keeps fake version there. AI Super-persuasion. AI uses psychology to deceive humans; 'you need me to avoid global catastrophe'. Ability to predict human behavior vastly exceeds humans' ability.</p>	<p>Deception scenarios require the axioms of identity verification via DLT. Deception plus super-persuasive AI require transparent and unhackable ethics and morality stored via DLT.</p>
<p>Singleton AI reaches overwhelming power. Prevents other AI projects from continuing via hacking or diversion; gains control over influential humans via psychology or neural hacking; gains control over nuclear, bio and chemical weaponry; gains control over infrastructure; gains control over computers and internet. AI starts initial self-improvement. Human operator unwittingly unleashes AI with self-improvement; self-improvement leads to unlimited resource demands (aka world domination) or becomes malignant. AI declares itself a world power. May or may not inform humans of the level of its control over resources, may perform secret actions; starts activity proving its existence ('miracles', large-scale destruction or construction). AI continues self-improvement. AI uses earth's and then solar system's resources to continue self-improvement and control of resources, increasingly broad and successful experiments with intelligence algorithms, and attempts more risky methods of self-improvement than designers intended.</p>	<p>The axioms <i>per se</i> do not seem to solve Singleton scenarios. They are addressed in a section below where the fundamental premise is each generation of AGI will contract with the succeeding generation and use the best technology and techniques to ensure continuation of a common but evolving value system. The same principle underlies solutions to successively self-improving AI to AGI transition and AGI evolution in which humans are still meaningfully involved.</p>
<p>AI starts conquering universe at 'light speed'. AI builds nanobot replicators, sends them out into galaxy at light speed; creates simulations of other civilizations to estimate frequency and types of alien AI and solve the Fermi paradox; conquers the universe in our light cone</p>	<p>The inevitable scenario where AI evolution exceeds human ability to monitor and intercede is what necessitates distributed, unhackable DLT</p>



and interacts with aliens and alien AI; attempts to solve end of the universe issues	methods and smart, i.e. automated, contracts. Further, transparent and unhackable ethics, and a durable form of morality, also unhackable via DLT, are what may ensure each generation of AGI passing the moral baton to the succeeding generation.
--	---

189

**Table 5.** Examples from Bostrom Pathways to Dangerous AI [10].

Pathway	Key Axioms
Perverse instantiation: 'Make us smile'	Morality defined as voluntary transactions
Perverse instantiation: 'Make us happy'	Morality defined as voluntary transactions
Final goal: Act to avoid bad conscience	Store value system in distributed app
Final goal: Maximize time-discounted integral of future reward signal	Morality defined as voluntary transactions, store value system in distributed app
Infrastructure profusion: Riemann hypothesis catastrophe	Morality defined as voluntary transactions
Infrastructure profusion: Paperclip manufacture catastrophe	Morality defined as voluntary transactions Social ostracism
Principal-Agent Failure [21] Human-Human: Agent (AI developer) disobeys contract Human-AGI: Agent disobeys contract	Digital identity, smart contracts, dApps, social ostracism

190

191 *3.2. Examination of Typical Failure Use Cases by Axiom*

192 One way to examine proposed necessary and sufficient set of axioms for AI morality is to look  
 193 at what phenomena or failure use cases result when one or more of them are excluded [21]. These  
 194 amount to a short explanation of each axiom and its necessity; longer explanations follow.

195

196

**Table 6.** Typical Failure Use Cases by Axiom.

Axiom of Safe AGI Omitted from Set	Failure Use Case if Omitted
Licensing of technology via market transactions	1. Restriction and licensing via state fiat: corrupt use or use benefitting special interest. 2. No licensing (freely available): Unauthorized and immoral use
Ethics transparently stored via DLT so they cannot be altered, forged or deleted	1. User cannot determine if AI has behavior safeguard technology (i.e. ethics) 2. Invisible ethics may not restrict moral or safe access
Morality, defined as no use of force or fraud, therefore resulting in voluntary transactions, stored via DLT	1. Inadvertent or deliberate access to dangerous technology by immoral entities

			(human or AI), i.e. entities using AI in force or fraud 2. Note that police and military AI will have modified versions of this axiom 3. Note that this axiom does not solve the case of super-persuasive AI as alternative to fraud
	Behavior control structure (e.g. a behavior tree) augmented by adding human-compatible values (axioms 2 & 3) at its roots		1. Uncontrolled behavior by AGI, e.g. behavior in conflict with a set of ethics and/or morality, either deliberately or inadvertently
	Unique hardware and software ID codes		1. Inability for entities to restrict access to AGI components because they cannot specify them 2. Inability to identify causes of AGI failure to meet design intent 3. Inability to identify causes of AGI moral failure via identification of components causing the failure Note the audit trail axiom depends on this one.
	Configuration Item (automated configuration)		1. Lessened ability to detect improper functionality or configuration of software or hardware within AGI. 2. Lessened ability to detect improper functionality or configuration of software or hardware to which AGI has access. 3. Inability to shut down internal AGI software and hardware modules. 4. Inability to shut down software and hardware modules to which AGI has access. Note smart contracts and dApps axioms depend on this axiom.
	Secure identity verification via multi-factor authentication, public-key infrastructure and DLT		1. Inability to detect fraudulent access to secured software or hardware (e.g. nuclear launch codes, financial or health accounts). 2. Inability to detect AGI impersonation of human or authentic moral AGI (e.g. POTUS, military commander, police chief, CEO, journalist, banker, auditor, et al.).
	Smart contracts based on DLT		1. Inability to enforce evolution of moral AGI due to its pace 2. Inability to enforce contracts with AGI due to its speed of decisions and actions 3. Inability to compete with regimes using smart contracts due to inefficiency, cost, slowness of evolution, etc.

Distributed applications (dApps) – Software code modules encrypted via DLT	1. Inability to restrict access to key software modules essential to AGI (i.e. they could be hacked more easily by humans or AI).
Audit trail of component usage stored via DLT	<p>1. Inability to track unauthorized usage of restricted software and hardware essential to AGI.</p> <p>2. Inability to track unethical usage of restricted software and hardware essential to AGI.</p> <p>3. Inability to track immoral usage of restricted software and hardware essential to AGI.</p> <p>4. Inability to identify which component(s) failed in AGI failure.</p> <p>5. Inability to prevent hacking of audit trail.</p> <p>6. Increased cost in time and capital to detect criminal usage of restricted software and hardware by AGI, and therefore, to apply justice and social ostracism.</p> <p>7. Inability to compete with regimes using DLT-based audit trails due to slowness to detect failure, identify entities or components responsible for failure, and implement solutions (overall: slowness of evolution).</p>
Social ostracism – denial of societal resources – augmented by petitions based on DLT	<p>1. Lessened ability to reduce criminal AGI access to societal resources.</p> <p>2. Inability for entities to preferentially reduce non-criminal AGI access to societal resources.</p>
Game theory/mechanism design	<p>1. Lacking a system to incent increasingly diverse autonomous intelligent agents to communicate results likely to be valuable to other agents and in general collaborate toward reaching individual and group goals, cohesiveness required for collaborative effort fails over time.</p> <p>2. DLT in a digital ecosystem theoretically permits all conflicts to be resolved via voluntary transactions (the Coase theorem), but a pre-requisite set of rules may be necessary.</p>

197

198 3.3. *Explanation of Each Proposed Axiom*

199 3.3.1. Access to AGI technology via license

200 Two distinct systems and traditions of technology licensing exist, 1) market transactions and 2)  
 201 state ('government', 'fiat') coercively-controlled licensing. Seizure of AI intellectual property (IP) and

202 control over its development by states is inevitable unless AI scientists and private-sector  
203 management set up their own systems to ensure safe AGI. ASILOMAR #9, Responsibility, states  
204 “Designers and builders of advanced AI systems are stakeholders in the moral implications of their  
205 use, misuse, and actions, with a responsibility and opportunity to shape those implications” [9]. The  
206 question is: How is this responsibility to be implemented — to be given ‘teeth’?

207 The system proposed herein envisions AI evolution with humans cross-licensing AI technology  
208 to each other, creating a prototype distributed applications (dApps) system instantiated in a DLT  
209 ecosystem that balances permissioned access and editing via contract with free access. The human-  
210 initiated DLT-based ecosystem would transition to AGIs licensing technology from humans, and  
211 subsequently to AGIs cross-licensing with each other.

212 History shows that in many or most cases, a market system evolves solutions faster and better  
213 than centralized state systems. Further, state systems may respond innovatively and less  
214 bureaucratically when subjected to competition with market systems; the Human Genome Project  
215 and current space-exploration efforts are examples. A market optimally distributes problems to be  
216 solved and computing power assigned to solve them in a highly decentralized manner.

217 There are valid arguments against an AI IP regime with ‘restricted’, information flow via license,  
218 whether through market or state. Progress may be slowed, and some persons with no reason to be  
219 prevented from accessing some AI technology may be restricted. The counter-argument is that AGI  
220 technology and many of its components are as dangerous or more dangerous than nuclear, biological,  
221 chemical, or other mass destruction weapons technology (WMD), since AGI will control WMD tech,  
222 along with innumerable other resources that can fatally or significantly affect humanity (Proposition  
223 1 in Appendix).

224 By way of example, assume there exists an algorithm critical for AI self-programming. With free  
225 access to the self-programming algorithm, malevolent humans, as well as extant autonomous AIs,  
226 could use that technology for unlimited self-improvement, opening a positive-feedback-driven  
227 Pandora’s box to unlimited malevolence and unlimited means to achieve it (ASILOMAR #22 [9]).  
228 Others point out dangers of a freely available ‘just add goals’ AGI [10, 18]. Thus state, private, or a  
229 hybrid means of restricting access to critical pieces of AI tech, as with WMD, seems to be a necessary  
230 axiom to align AI with human interests.

### 231 3.3.2. Ethics stored in a distributed ledger

232 I define *ethics* as the *fundamental value system* from which autonomous entities derive their  
233 decisions and choices. *Ethics* are separate from *morality*, which is a particular set of ethics. ‘Honor  
234 among thieves’, ‘do unto others as you would have them do unto you’, ‘professional courtesy’, ‘honor  
235 thy father and mother’, etc., are ethics, as are Asimov’s three laws of robotics [31]. Ethics can seem  
236 good or bad, moral or immoral, from a volitional entity’s subjective value system. An entity’s  
237 fundamental values are embedded in some type of behavior (input/output) control system. For  
238 example consider ethics represented and controlled by a behavior tree [32] where the ethics are a  
239 subset of its roots, and thus in that sense *fundamental*.

240 The intention of storing AGI ethics via DLT is to permit a class of autonomous entities to have  
241 identical ethics and to render them visible and unable to be hacked, altered or deleted. In this sense,  
242 ethics is a necessary component of the control system and allows for different sets of ethics to be  
243 instantiated. While it is not possible for all humans to have identical values and therefore moral  
244 values (however defined), DLT, in theory, permits a universal set of immutable values to be  
245 instantiated in AGIs while still permitting an unlimited range of individual AGI and AI diversity.

246 Requiring transparent instantiations of ethics for AGI systems conforms to ASILOMAR #10  
247 (Value Alignment), and IBM’s call for Supplier’s Declarations of Conformity for AI [33]. These *bona*  
248 *fides* and ethics could be stored in an AGI’s Configuration Item and/or those of its key components  
249 (see below).

### 250 3.3.3. Morality defined as voluntary vs. involuntary exchange

251 The definition below is intended to conform to ASILOMAR #11, Human Values, #14, 'benefit  
252 and empower as many people as possible', #15 and #24, benefit the 'common good' and 'widely-  
253 shared ethical ideals' [9], but notably to provide a practical implementation of them, otherwise what  
254 use are they?

255 Down through the ages there have been two main problems with discussions of morality — first,  
256 ambiguity and therefore confusion. How can we identify moral behavior if it is imprecisely defined  
257 and hard to determine [34]? And so such definitions are costly, in terms of the economics of law, to  
258 enforce. Second, nearly all morality descriptions are subjective, amounting to one person's value  
259 system imposed on others, and via coercion if enforced via the state.

260 For example, take the proposal of directing AGI to ensure 'hedonistic consequentialism' for all  
261 of humanity — selecting from a set of actions the one that would produce the best balance of pleasure  
262 versus suffering [10]. Such idealistic but vague and minimally-thought-out concepts of morality —  
263 which is nearly all of them — may sound good on paper but break down rapidly on implementation.  
264 And they all amount to a minority or individual — human or AGI, and even from the most beneficent  
265 of us — deciding what is 'moral' or not, or what is 'best' for others. When AGI is a given, the proposals  
266 depend on its super-intelligence somehow overcoming the limitations of humans' concepts of  
267 morality, how to define and implement it, and/or overcome humans' inability to read minds. And  
268 notably, they all amount to confining computation of an overall system solution to a restricted subset  
269 of all computationally active agents (see Diversity, below), which is another way of saying allowing  
270 a subset of volitional entities to impose their subjective, not absolute, value system, upon others.

271 The essence of autonomy or volition is choice-making. Herein, first, all individual choices that  
272 affect no other volitional entity are moral. Second, all voluntary exchanges are moral. But if two  
273 autonomous agents prefer a transaction between them, and that transaction is prevented by a third  
274 party, that party has imposed its value system over the others. It is also one less computational  
275 experiment the entire system performs.

276 Several economists posited that there is no universal theory or method to determine *value*, rather,  
277 all human values and the measure of utility are subjective [35], which is implicit in the game-theoretic  
278 axioms of utility [23]. Following this premise, defining morality as all voluntary transactions is  
279 *scientific* when science is likewise defined as a procedure that filters for absolutes — what we all see  
280 in common, such as the speed of light — from a vast sea of relative views [36, 37]. Later members of  
281 the Austrian school defined morality as non-interference with property (defined to include ones'  
282 body and intellectual property) [36, 38]. It is simpler and less costly to define *moral* transactions as  
283 *voluntary* transactions than to try to identify what is *property* and to define and figure out property  
284 boundaries and property interference. One of the goals of a legal system is to resolve conflicts in an  
285 economically efficient manner and it has been argued that the evolution of common law is toward  
286 such efficiency [39].

287 If you want to upload your mind and join a collective intelligence, or rather stay physically  
288 human, and not even accept lifespan enhancement, it is up to you. Under this system you and AGI  
289 cannot force choices on anyone else even if you or AGI believe it is best for them. But what if a super-  
290 intelligence could make some or all of your decisions better than you can [10]? Each individual can  
291 sign on with the super-AI that seems to best fit your values and goals. It would be your choice, just  
292 like taking the advice of a consultant or hiring an agent for a specified set of tasks today.

293 This definition and axiom may not solve the problem of AGI with vast knowledge of the  
294 evolution of our psychology and innate choice-making algorithms [40, 41] and the propensity to  
295 manipulate us with that knowledge, although the argument can be made that with such knowledge  
296 in a voluntary exchange system, AGI would be more able to offer 'good' choices (i.e. as we perceive  
297 them) to us than without that knowledge.

298 AGIs will have a larger and more complex set of value preferences than ours (see Diversity,  
299 below); what will be the morality of their interaction with each other? The voluntary transaction  
300 definition may fit their behavior as well. A system of voluntary transactions permits Pareto optimality  
301 and maximizes computational experiments driven by local, subjective preference systems [42].

302 Transaction costs and the need for trusted third parties prevents Pareto optimality [43]. DLT and  
303 smart contracts potentially permit full Pareto optimality in the digital AI ecosystem by reducing  
304 transaction costs to negligible amounts and eliminating costly, imperfect third parties.  
305

#### 306 3.3.4. Behavior control system

307 Behavior control is *sine qua non* to value human-AI alignment (ASILOMAR #10, 16, etc. [9]).

308 At one end of the knowledge representation/control spectrum is a 'flat' set of large numbers of  
309 heuristical condition-action rules that are selected, not based on general principles, but on matching  
310 specified patterns. At the other end of the spectrum is a strict postulatory-deductive tree in which the  
311 internal node 'beliefs' are logically derived from the postulates as are the actions represented at the  
312 leaf-nodes. A postulatory-deductive system is the ideal contemplated here, which would satisfy the  
313 need for control, the desire for transparency of its operation, and part of the need for formal proof of  
314 its reliability. However, it is an ideal. Any type of hierarchical control system that can hold values at  
315 its highest levels and is transparent enough to reveal control over behavior by values is a candidate  
316 for aligning AGI and human values, and the ecology of value systems that will evolve from the initial  
317 sets.

318 I believe humans innately attempt to form postulatory-deductive systems using non-  
319 mathematical, *ad hoc* 'logics' [40, 41] in an effort to organize their world-view into causes and effects,  
320 and general principles governing specialized condition-action pairs. Mathematical and scientific  
321 postulatory-deductive systems are recent, specialized, powerful cases, improvements built on the  
322 general-purpose cognitive architecture, in which universally-valid logic replaces the *ad hoc*  
323 evolutionary 'logics' and the entire system is validated through repeated observations directly  
324 confirming the postulates or indirectly via observation of valid derivatives (i.e. predictions) with zero  
325 fault-tolerance. Further, in the ritualized transparency of its methods and crowd-sourced validation  
326 via multiple subjective observers, science is an absolute voluntary consensus, rather than  
327 confirmation of an unprovable 'objective' world [37] and resembles DLT.

328 In the innate human system, a causatory cascade of beliefs and actions stem from the  
329 fundamental beliefs (postulates, including values). Outside of the mathematical and scientific  
330 postulatory systems a more complex set of relative and subjective 'logics' connects beliefs —  
331 efficacious from an evolutionary standpoint but also unreliable across different contexts [40, 41] as  
332 seen in beliefs of mathematicians and scientists outside of mathematical and scientific domains.

333 An AI control system that may be able to represent current and future postulatory-deductive  
334 systems is the *behavior tree* [32].

335 The game-theoretic axioms of utility drive decisions from a hypothesis that the decision will  
336 ultimately lead to an improvement in the volitional entity's state, as defined internally and  
337 subjectively by its value system [23] aka *the pursuit of happiness* [36]. The utility axioms extend to  
338 machines with subjective value systems.

#### 339 3.3.5, j3.3.6. Unique component IDs, Configuration Item (CI)

340 Several technological and business process developments lead toward a universally  
341 interconnected system that self-configures, self-diagnoses its component failures, and repairs them  
342 automatically; *in toto*, a paradigm whose ultimate use will be integration into the human-AGI  
343 ecology. These technologies help to decrease Coasean transactions costs (e.g. detection and  
344 enforcement) toward facilitating a Pareto-optimal economy.

345 Unique identification (ID) numbers evolved as an economically-efficient means to organize and  
346 validate property exchanges, contributing to a stable society, starting with large or important pieces  
347 of property such as real estate via book and page of a recorded deed, automobiles via title or vehicle  
348 ID number, stocks via CUSIP number, etc. As the cost of creating unique ID numbers decreased via  
349 technology, the system extended to machines and devices via model and serial numbers, and more  
350 recently to any product via one- and two-dimensional bar and matrix machine-readable codes to  
351 facilitate supply-chain management, quality control, customer service, and other functions.

352 The transition from the internet of computers to the ‘internet of things’ (IoT) envisions  
353 ubiquitous communication and computation connecting physical devices with the digital world via  
354 miniaturized sensors and chips containing only as much computing power and energy usage that is  
355 needed to perform their intended functionality in their context — “a self-configuring network that is  
356 much more complex and dynamic than the conventional internet” [44]. In the IoT ID numbers become  
357 digital as well as physical, e.g. radio frequency ID codes. In the IoT world AGI will be able to  
358 communicate with, and potentially control, any digital or physical device.

359 The IoT world was presaged by the development of *disaster recovery and business continuity*  
360 *planning*, and the key role of configuration items in them. Disaster recovery (DR) arose on the  
361 realization that the cost of *not* doing contingency planning for disasters (a hazardous material spill,  
362 hurricane, tornado, power outage, etc.) could vastly exceed the cost of such planning, including total  
363 business loss. Judicious planning for disasters, such as foreseeing an alternate location from which to  
364 conduct operations in the event of facility downtime and establishing redundant communication  
365 protocols to coordinate team response to disasters, are relatively inexpensive insurance measures.  
366 Business continuity planning (BCP) logically arose from DR, extending the DR premise of disaster  
367 planning to pre-planned, prioritized responses to *all* component failure, including normal end of  
368 service life. For example, recovery of failed email for the company as a whole is accorded lower  
369 priority than for customer-service representatives and top management. BCP’s goal is, through  
370 contingency planning, to reduce the internal and external impact of business process downtime to a  
371 minimum.

372 The configuration item (CI) arose in BCP/DR conceptually as a system component’s on-board  
373 algorithm and parameter set that allowed computers and components to detect each other’s  
374 configuration requirements, automatically configure the component, or perform error-detection,  
375 reporting, and correction (cf. ASILOMAR #7, Failure Transparency [9] and Manheim [21]). In the  
376 context of DLT, it becomes a smart contract.

377 Many paths to dangerous AI, including much of the broad class of human-AI value  
378 misalignment, are a result of improperly configured or failed components, or sabotage (e.g. accidental  
379 nuclear war, failure of safeguard components, inadvertent security vulnerabilities leaving a system  
380 open to hacking, misconfiguration of software modules e.g. in autonomous vehicles, power  
381 blackouts, financial system meltdowns, etc.). Thus, the paradigm of BCP/DR and CIs will be integral  
382 to maintaining the fidelity of AGI-human value alignment amidst the IoT of the future. Further, CIs  
383 of critical AGI components can be encoded via DLT, thus greatly reducing or eliminating the  
384 possibility of unauthorized use, corruption, failure, etc.

385 IBM’s Supplier’s Declaration of Conformity to ensure AI safety [33] could be incorporated into  
386 CIs and used as one pre-requisite for deployment of an AGI system or component.

### 387 3.3.7. Digital identity via distributed ledger technology

388 Restricting access to potentially dangerous technology (Axiom #1) necessitates identity  
389 verification. Few readers would deny the need of multi-factor authentication for nuclear missile  
390 launch codes. Identity verification is currently accepted for access to military bases, high-tech  
391 weapons, aircraft, most private and public buildings, financial systems, health records, and other  
392 data that individuals consider private for their own reasons, all toward the goal of ensuring a safe  
393 and secure world.

394 In contrast to a third-party-based identity authentication system such as state- or private  
395 company-issued ID cards, many decentralized DLT-based methods have been created, competing  
396 with the trusted-third-party method to reduce the chance of forgery or other hacking, and bribery or  
397 other corruption. In a DLT version of the current public-key encryption-based X.509 standard [45], a  
398 DL replaces the third-party issuing authority in its components: certificate version, serial number,  
399 type of algorithm used to sign the certificate, issuing authority, validity period, name of entity being  
400 verified, and entity’s public key .

401 Initially, digital identity verification will be done on humans matching biometrics such as facial  
402 features, fingerprint, voice, in addition to SMS etc., but as AI evolves, AGIs will use technology and

403 techniques that they develop against evolving threats to hack verification of humans, e.g. speech  
404 synthesis or video manipulation [18] and threats that are currently unforeseeable.

#### 405 3.3.8. Smart contracts based on digital ledger technology

406 Smart contracts were conceived by Szabo decades ago, before the inventions of DLT and IoT that  
407 enable their inexpensive implementation, to automate contractual clauses via cryptography that can  
408 be self-executing and self-enforcing [46]. Smart contracts as an integral part of DLT are “scripts  
409 residing on a blockchain that automate multi-step processes” [28]. Szabo’s inspirations were the  
410 original commercial security transaction protocols: SWIFT, ACH, and FedWire for electronic funds  
411 transfer, credit card point of sale terminals, and the Electronic Data Interchange for transactions  
412 between large corporations such as purchase and sale [30]. He used the simple example of a vending  
413 machine, through which transactions are performed without a third-party intermediary to verify that  
414 the terms of the transaction have been satisfied.

415 Two critical design goals were to make verifying satisfaction of contractual terms  
416 computationally cheap and breaching terms computationally expensive, both of which are realized  
417 in a far superior generalized manner via DLT than via prior methods (reminiscent of Bush’s and  
418 Nelson’s conception of hyperlinking before the invention of the internet [47]). Smart contracts require  
419 the digital specification of obligations each party must meet to trigger a transaction, a blockchain for  
420 consensus verification that each party has met its obligation, an immutable audit trail of transactions,  
421 and the design goal of excluding unintended effects on non-contractual parties.

422 Omohundro envisions smart contracts interfacing autonomous agents with the heterogeneity of  
423 human legal codes and future legal codes designed to help ensure safe AI interactions with humans  
424 [48](ASILOMAR #8 [9]. Pierce envisions a mass migration of the current compliance regime via law  
425 and regulation to an economically more efficient and secure regime based on smart contracts [49]  
426 (ASILOMAR #2); such a system greatly facilitates Omohundro’s.

427 As AGI evolves beyond our understanding and visibility, and notably when it hits ‘escape  
428 velocity’ — exponential evolution culminating in generations succeeding each other in fractions of a  
429 second — prescribed, automated smart contracts will be essential to perpetuating ethical values in  
430 each successive generation. The concept is that a more advanced AGI generation cannot succeed a  
431 less-advanced one without licensing key components — certain algorithms, hardware, the axiom-  
432 methods proposed herein, behavior control systems invented by humans and AI, etc. — from the  
433 less-advanced generation, subject to satisfying its value system and oversight.

434 The configuration ‘handshake’ between an AGI and its component CIs is a smart contract  
435 between them, and the intelligence of those handshakes can increase in the future. CIs must  
436 incorporate the ability to deny activation of a component within a system, or shut it down, if lack of  
437 satisfaction of a given clause, or violation of a clause, of any extant contract is detected by any  
438 distributed ledger stakeholder in the transaction. All such contractual stakeholders must be silenced  
439 just as living cell cycle checkpoints must be silenced for the cell to progress through the intricately  
440 orchestrated process of mitosis, otherwise it self-destructs [50]. More of these ‘deadman switches’  
441 that actively suppress unauthorized use or malfunctioning AI will increase a secure evolution of  
442 benign AI; for example the limited term of digital identity certificates that expire and require re-  
443 verification of the subject entity’s identity at regular intervals [45].

444 Szabo’s vision of embedding smart contracts in objects [30] is realized by embedding CIs in all  
445 non-trivial interconnected devices and algorithms in the IoT. In this manner the smart contract and  
446 preceding axiom-methods work in concert to ensure human-AGI value-alignment and AGI  
447 containment within bounds that are benevolent for humans and the succession of AGI generations.

448 In principle, smart contracts help approach a zero-transaction cost world by eliminating trusted  
449 third parties, and their role in detection and enforcement of contractual rights (e.g. physical and  
450 intellectual property rights).



### 451 3.3.9. Decentralized applications (dApps)

452 DLT-based decentralized applications (dApps) differ from conventional application programs  
453 in that they 1) are outside the overview and control of a central authority such as a company making  
454 the app or state agency controlling it, 2) operate on a peer-to-peer network instead of a centralized  
455 one, and 3) do not have a central point of failure — they are redundant in hardware and software and  
456 therefore fault-tolerant [51]. Smart contracts are an example of dApps, as are decentralized versions  
457 of exchanges to trade various types of goods or services, notably intellectual property, which can  
458 transition into exchanges between AGIs, social media including networking, communications  
459 protocols, prediction markets, and a growing number of DLT-enabled applications.

460 Axiom 1, Access to Technology via Market License, requires that some dApps — notably those  
461 that are critical to AGI — would be implemented via permissioned DLs, which are DLs with an added  
462 control layer that can prevent unrestricted and unauthenticated public access. Some cryptocurrency  
463 observers feel any type of control that is not fully ‘public’ violates the decentralization principle;  
464 however, consider ‘private’ DLs as a critically important tool in the DLT toolbox. For example, should  
465 we not consider delegating control over access to critical AGI algorithms to a consensus of signatories  
466 committed to the goal of AI-human value alignment or ethical use of AI, e.g. the ASILOMAR AI  
467 Principles [9]? Further, the control layer, in part or eventually *in toto*, can be automated by  
468 incorporating smart contracts and/or smart tokens to reduce the probability that central control can  
469 be hacked or corrupted. Smart contract terms could require 2/3 or 100% acceptance of DLT-  
470 authenticated (Axiom 6) signatories to ASILOMAR AI Principles or similar regulatory documents.  
471 Smart contract terms can deny access to those who do not fulfill a transparency requirement via  
472 Supplier’s Declaration of Conformity [33], which document could in turn require inclusion of an  
473 accepted set of ethics and morality (Axioms, 2, 3) and a safety testing record meeting certain  
474 standards [11, 52], all of which can be incorporated into a CI (Axiom 5). Equally critical, dApps permit  
475 separation and balance of powers of key AGI components, analogous to no one entity having all the  
476 nuclear launch codes. The significance of dApps for ensuring benevolent AGI is discussed further in  
477 two malignant use cases it addresses, the Rogue Programmer and Singleton AGI, below.

478 Two levels of permissioned access to dApps may be needed: 1) access for use, and 2) access to  
479 modify the code (while, again, a purist view of dApps sees their development as open-sourced). A  
480 similar consideration must be given to AGI technology patents. The primary purpose and  
481 requirement of patents is to ‘teach the art’ clearly and explicitly so the innovation can be implemented  
482 by the reader. The patent system at a meta-level has largely been denied market evolution to try other  
483 purposes and requirements. Be that as it may, to facilitate *safe* free exchange of information, a  
484 ‘Transportation Security Administration’-type of pre-screening for access to critical AGI patents may  
485 be needed to prevent access by malevolent entities and may be efficiently implemented via smart  
486 tokens.

487 If no formal proof of benevolent AGI methodology is possible or available soon, sandbox  
488 simulations of new AGI technology are critical to our future and implementing them via dApps will  
489 be essential to ensure they cannot be hacked or corrupted by humans or AGIs [52].

### 490 3.3.10. Audit trail of component usage stored via distributed ledger technology

491 DLT is inherently a low-cost, redundant, decentralized, hack-free audit trail — a significant  
492 improvement on traditional centralized audit trail technology. An unhackable audit trail of critical  
493 AI components such as collaborative, self-learning, or self-programming algorithms will facilitate  
494 rapid, efficient detection of their authorized or unauthorized use (i.e. a hack of a contract, a set of  
495 ethics, or an identity verification) or failure (cf. ASILOMAR #7, Failure Transparency [9] and Maheim  
496 [21]). and increase the probability of remedying the system fault. The IBM Research Supplier’s  
497 Declaration of Conformity via a factsheet for AI software incorporates an audit trail as a fundamental  
498 principle [33]. Bore et al. describe a system for incorporating an audit trail in DLT as part of  
499 embedding AI simulations in DLT so that trust in the simulations’ validity is enabled between  
500 researchers without requiring a trusted intermediary [52].

### 501 3.3.11. Social ostracism (voluntary denial of resources)

502 As various writers point out, a ‘power-hungry AGI’ or ‘AGI pursuing world domination’ implies  
503 a AGI attempting to access and control an ever-increasing amount of society’s resources [10, 17-19].  
504 Therefore, the ability for entities to deny societal resources to an errant AGI is a counterforce on its  
505 ambitions. This voluntary mechanism is another aspect of a market economy in which computation  
506 is distributed, local, and optimized — each entity makes its own choice based on its own unique,  
507 subjective experience. A further optimization is that market votes can occur as often as each entity  
508 wishes to change its choice, such as denying its resources to another entity or collection of entities.  
509 Market votes occurs immeasurably more often than political votes and implement a far more fluid  
510 and asymptotically Pareto-optimal society.

511 In the current technology for ‘democracy’ the political vote is the means to reach consensus,  
512 which is tallied by a central authority and enforced via coercion by the same entity. In contrast,  
513 voluntary concerted boycotts of companies, facilitated by modern social media, are increasingly  
514 affecting corporate policy (corporations being one type of voluntary association among individuals  
515 for their mutual benefit).

516 DLT is a fundamentally new way to reach and archive a consensus. DLT-based unhackable  
517 petitions can be smart contracts to facilitate denial of resources to an errant AGI and can be rapidly  
518 implemented via CIs. For instance, IBM’s call for Supplier’s Declaration of Conformity to help ensure  
519 safe AI implies voluntary adoption [33], but would be more effective if enforced via social ostracism  
520 and implemented automatically via CI incorporation, just as web browser security currently can alert  
521 a user to reject non-security-credentialed (non-https) internet domains, thereby immediately denying  
522 them the user’s resources.

523 The ASILOMAR principles, currently signed by 1273 AI workers [9], are a significant first step,  
524 like a letter of intent, toward a necessary, more binding and important agreement. A next step could  
525 be archiving the ASILOMAR agreement and its signatories via DLT so that the principles cannot be  
526 hacked and can only be amended via consensus of the signatories. A further step could be embedding  
527 the document and signatories in the Supplier’s Declaration as a second, more restricted layer of access  
528 protection. Another step would be automatically-triggered, smart contract DLT-based petitions  
529 attached to the Supplier’s Declaration, denying a given set of AGI access to specified AGI technology  
530 in response to detected AGI behavior contradicting the ASILOMAR principles.

### 531 3.3.12. Game theory and mechanism design

532 Game theory and evolution have explained five categories of the evolution of cooperation —  
533 direct reciprocity e.g. ‘tit for tat’, indirect reciprocity e.g. reputation value in ‘what goes around,  
534 comes around’, reciprocity in societal networks and topologies, group reciprocity e.g. the good  
535 Samaritan and altruism, and kin reciprocity, e.g. ‘I would lay down my life for two brothers or eight  
536 cousins’ (J. B. S. Haldane) [53]. Nowak’s current goal is to extend these explanations to game-theoretic  
537 frameworks for global cooperation and cooperation across generations. These efforts will involve  
538 mechanism design, the branch of game theory concerned with designing game-theoretic and  
539 economic structures that build in incentives for communicating truthfully about one’s valuations in  
540 a potential transaction [21, 23, 54, 55]. That is the goal of game theory in the context of axioms for safe  
541 AGI.

542 It is possible that a suitably designed communication protocol and game-theoretic incentives  
543 using DLT could replace the other axioms, which would emerge from the simpler axiomatic system.  
544 For example, an axiomatic (first principles) simulation of game-theoretic evolution wherein agents  
545 have a complex set of strategies found that inclusion of two axioms, (1) inheritable agent types, and  
546 (2) visibility of types to other agents, resulted in evolution of cooperation strategies [24]. These axioms  
547 could be more general than the license, ethics, morality, configuration item, audit trail, and social  
548 ostracism axioms proposed herein. The unique component IDs, digital identity verification, and game  
549 theoretic axioms along with DLT to ensure transparency, may suffice to generate the rest of the set,  
550 just as a wide variety of market-based structures and mechanisms emerge from axiom sets that

551 generate markets (a large proportion of economics, game theory, and agent-based modeling literature  
552 could be cited here; see, just by minimal example, the following and their references [23, 54, 55]).

### 553 3.4. Diversity in the AGI Ecosystem: Computation Is Local, Communication Is Global

554 However, proving this possibility may be intractable. Going back at least as far as Newell, it has  
555 been stated that the complexity of behavior (input-output functions) for  $I$  inputs and  $O$  outputs is  $O^I$   
556 [56]. Intuitively, this is rolling a die with  $I$  faces  $O$  times since any number of the  $I$  inputs could map  
557 to each output. A series of actions, i.e. behaviors, is calculated by the power tower,  
558

$$559 \quad O^{1^{O^1 O^1 \dots}} \quad (1)$$

560  
561 whose complexity grows super-exponentially. But in fact, complexity grows faster than the  $O^I$  power  
562 tower in the cases where the topology of I-O mappings matters, such as in successive neural net  
563 actions. In those cases,  $O$  is raised to the power set of  $I$ ,  $2^I$ , and the succession of actions is calculated  
564 by the power tower,  
565

$$566 \quad O^{2^{1^{O^{2^{1^{O^2 \dots}}}}} \quad (2)$$

567  
568 whose complexity exceeds that of (1). These intractable formulae have significant implications for the  
569 AGI ecosystem. One is that astronomically greater diversity of value systems is possible compared to  
570 humans'. Second, AGIs' behavior in ecosystems will likely take them to disparate locations in the  
571 problem spaces they investigate, creating a very sparsely inhabited matrix of a vast number of  
572 possible behaviors. Third, in that context, game theory and mechanism design may be the key  
573 structure inducing their ongoing cooperative behavior, notably to allocate problems to be solved and  
574 communicate results that may be valuable to the other players truthfully and in a timely manner.

575 For example, in our primitive intellectual property regime, a protocol that induces efficient,  
576 truthful reporting is the requirement that a patent clearly teach the new art to those skilled in its  
577 subject matter. Absent that requirement and patent protection, players might be induced to seek  
578 intellectual property protection via secrecy, e.g. 'trade secrets,' decreasing cooperative search and  
579 overall technological progress. A protocol that induces timely reporting of innovation is the recent U.  
580 S. patent rules change to grants rights to those who are 'first to file' versus 'first to invent', which was  
581 economically inefficient and lacked the inducement to disclose earlier rather than later.

582 The fourth implication is that, as described differently in disparate intellectual settings [42, 56-  
583 58], computation will continue to be performed in unique, sparsely populated loci in the general  
584 problem space using subjective criteria for exploration, and communicated via vastly shorter, high-  
585 level symbol sequences compared to the lengths of computational sequences and complexity of  
586 modeling producing them.

### 587 3.5. Should AI Research and Technology Be Freely Available While Nuclear, Biological, and Chemical Weapons 588 Research Are Not?

589 The Rogue Programmer problem assumes that one amoral, misguided, naïve, or malevolent  
590 individual could make the single advance generating AGI, and this risk depends on how close the  
591 technology is to a single leap causing 'take-off'. History shows that all innovations will occur in a  
592 matter of time, some taking more time than others. For instance, differential calculus was invented  
593 by Newton in the spring of 1665 and by Leibniz in the fall of 1675 [59]. The historical record is clear  
594 that what appear in retrospect to be great innovative leaps are actually the final step built on stronger  
595 antecedents than are assumed in scientific mythology, and in fact a chain of them involving many  
596 individuals [60]. Perhaps most pertinent to the advent of AGI is the detonation of the atomic bomb  
597 by the U.S. on 16 July 1945, then by the U.S.S.R. on 29 August 1949. The fusion bomb was detonated  
598 by the U.S. on 1 November 1952 and by the U.S.S.R. on 22 November 1955, an event that was  
599 accelerated by spying, which of course is a possibility with AI research [61, 62].

600 Such science and technology feats are large-scale group efforts. The Rogue Programmer problem  
601 arises when one individual circumvents the consensus agreement of end usage permission by the  
602 contributors to his/her technology (e.g. the 1273 AI worker signatories to the ASILOMAR principles  
603 [9]).

604 Two recent examples of rogue programmers are worth noting. A Chinese scientist used gene-  
605 editing techniques — developed elsewhere and made freely available in the spirit of the free exchange  
606 of ideas and technology — to change the genes of human eggs *in vitro* [63]. The innovation escaped  
607 overview, was motivated by ambition and pecuniary desire, and ignored a variety of the scientific  
608 community's publicly-voiced, well-thought-out *but unenforceable* concerns. Second, recently an AI  
609 programmer claimed his robot, which applied for and received citizenship in Saudi Arabia, would  
610 achieve human-level intelligence within 5 – 10 years [64]. His apparent variety of noble and possibly  
611 naïve motivations suggest that, even if he was not capable of making the innovation he pursues, he  
612 would combine innovations by others to achieve and claim the first human-level AI.

613 The problems, then, are unenforceable restrictions in a regime of 'free exchange of ideas and  
614 technology,' including public patents, and the lack of reliable means to measure how far away, in  
615 time or succession of innovations, we are from AGI.

### 616 3.6. *Measuring the Progression to AGI*

617 How urgent is the need to develop AGI-human value alignment technology? Can that debate be  
618 grounded in empirical data? Opinions differ on the timing to AGI — as of 2015 there were over 1300  
619 published predictions [65]. Timing predictions affect the urgency of preparing AGI-human alignment  
620 and control, which influences the resources we should devote to that effort. For this and other reasons  
621 it would be helpful to measure progress to AGI in time or in successions of specific AGI-enabling  
622 technologies [66], including the positive-reinforcement, recursive self-improvement abilities such as  
623 self-teaching, collaboration, self-programming, etc.

624 Akin to bottom-up versus top-down economic forecasting, a method that captures and compiles  
625 many local, informed assessments is polling AI experts [65, 67]. A second bottom-up approach is  
626 taken in the McKinsey Global Institute report, which assesses AI progress by its value-added to  
627 business processes using industry leader interviews and analytics [68].

628 A third approach, a hybrid of bottom-up and empirical metrics, is the Electronic Frontier  
629 Foundation crowd-sourcing technical progress metrics [69]. A fourth approach, empirical in concept,  
630 is taken in the AI Index 2018 Annual Report, a set of metrics intended to 'ground the AI conversation  
631 in data' divided into categories: Volume of Activity, Derivative Measures, Technical Performance,  
632 Towards Human Performance, and Recent Government Initiatives and using such metrics as  
633 numbers of papers published, course enrollment, conference participation, robot software  
634 downloads, robot installations, GitHub ratings, AI startups, venture capital funding, job demand,  
635 number of patents, adoption by industry and company department, and mentions in corporate  
636 earning commentary [1].

### 637 3.7. *AGI Development Control Analogy with Cell-Cycle Checkpoints*

638 Biological cell division is a complex and carefully-orchestrated process. Part of the insurance  
639 against cancer and other disorders resulting from defective replication is an ancient and strongly-  
640 conserved and evolved set of checkpoints that require fidelity tests to be passed in order for the cell  
641 to pass successive stages of division [50]. A notable feature of the checkpoints is their 'deadman  
642 switch' setup, i.e., rather than listening for signals of defects and then emitting signals to halt the  
643 process, their default mode is to send signals that suppress entering the next stage and require active  
644 silencing by successfully passing the fidelity tests. The analogy for AGI evolution is a set of active,  
645 not passive, checkpoints that halt or delay further AGI progress until certain safety criteria  
646 established by a consensus of researchers (human or AGI) are met.

### 647 3.8. *Intelligent Coins of the Realm*

648 A fundamental difference between today's money and cryptocurrencies is that the latter can be  
649 'intelligent', i.e. can be endowed with more functionality than a simple token representing mutually-  
650 agreed-upon or fiat-enforced value. For example, a common AGI malevolent path is achieving world  
651 domination, inadvertently or deliberately, by commanding an exorbitant share of resources, e.g.  
652 Bostrom's paper-clip disaster [10]. Omohundro considers how universal AGI drives may be  
653 engendered and reasons that since most goals require physical and computational resources  
654 unlimited resource acquisition may be an example [8]. 'Open-ended self-improvement' is another  
655 possible universal drive example [18, 19]. In biological systems, cell-doubling is a potentially  
656 dangerous path to deleterious claim on resources, and cancers are a collection of such paths. It is  
657 worth noting, analogous to AGI evolution, that biological evolution has found hundreds of cancerous  
658 paths, many using re-programming to avoid cell-cycle checkpoints, and resistance to treatments is  
659 real-time exploration of new paths using various genetic algorithms [50, 70, 71].

660 As stated, the axioms provide checks, in some cases redundantly, against this danger path. An  
661 additional check and/or means of implementation could be requiring a specialized token to purchase  
662 server time or rent AGI technology that automatically looks for the requester's compliance with AGI  
663 safety agreements and standards, otherwise the requester's 'credit' is denied. The token's DL then  
664 records the secure audit trail including measures of resources requested and protects against hacking  
665 to hide the evidence. Signals of possible dangerous activity, such as exponentially-increasing requests  
666 for resources by the same or related entities, could be incorporated into the token's programming.  
667 More broadly still, Omohundro cites the vision of a plethora of smart tokens performing  
668 intermediation of value and contractual obligations between the Internet of Things and humans [48].

### 669 3.9. *The Need for Simulation of Control and Value Alignment*

670 Considerable effort has gone into analyzing how to design, formulate, and validate computer  
671 programs that do what they were designed to do; the general problem is formally undecidable.  
672 Similarly, exploring the space of theorems (e.g. AGI safety solutions) from a set of axioms presents  
673 an exponential explosion.

674 A possible solution is to create a safe 'sandbox' environment where, iteratively and with  
675 parameter sweeps, simulations can be performed and improvements made to *control* and *value*  
676 *alignment* systems until the principles resulting in robust performance validating our design intent  
677 can be induced.

678 Critiques of the sandbox strategy includes: 1) AGI faking benign goals or obedience in the  
679 sandbox and then pursue its actual goals when released; 2) AGI hacking out using superior  
680 technology, developed while in captivity if needed, and most generally, 3) 'juvenile' AGI behavior in  
681 the sandbox that fails to predict bad behavior of a more advanced AGI into which it evolves [10]. To  
682 address #1 and #2, we need a control system that is effective enough and transparent enough to  
683 prevent those paths, such as through Axioms 2 and 3, transparent and unhackable ethics and  
684 morality, and Axiom 4, the behavior tree value system. Bore et al. take the goal of transparent  
685 simulation and modeling to a new level by describing a system wherein simulation specifications and  
686 an audit trail are stored via DLT, thus facilitating a means to cross-validate simulations before  
687 deployment and obstruct malicious hacking or fraud in simulations by humans or AI [52] (cf.  
688 ASILOMAR #6, Safety – 'verifiably so' [9]). Sandbox problem #3 may be redressed with the separation  
689 and balance of powers described next.

### 690 3.10. *A Singleton versus a Balance of Powers and Transitive Control Regime*

691 Bostrom defined 'singleton' as a single AGI possessing a decisive strategic advantage over  
692 humans and other AIs; a single world-dominant decision-making agency at the highest level [10].  
693 Even if a consistent axiom set is possible that solves the AGI deception and hacking problems and  
694 others, such a set may not be sufficient to solve the problem of the singleton. The solution proposed  
695 below also addresses the proposition that ensuring *most* AGI are safe to humans is not sufficient and

696 that *all* AGI must be rendered safe [34]. The axioms proposed herein presuppose that we cannot  
697 foresee how the evolution of AGI may outgrow the axiom set and the technology and techniques  
698 used to implement them.

699 Further, if simulation cannot conclusively demonstrate a solution to the singleton problem, then  
700 evolving the methods used to ensure moral, benign AGI along with AGI intelligence must be  
701 delegated to a consortium of AGIs whose values are aligned with humanity's. The idea is that a  
702 beneficent value and control system will evolve along with AGI and each generation consisting of  
703 multiple, cross-check-and-balance AGIs will, out of self-interest, endow the succeeding generation  
704 with the latest value and control version. Here 'generation' means a set of AGIs incorporating a  
705 significant technological advance over a prior set of AGIs. If there is only one AGI, it seems more  
706 likely that an aberrant or errant version could emerge, while if there are, e.g. 500 AGIs in a generation  
707 that are competing pluralistically, as in markets and government based on separation and balance of  
708 powers, to win the DLT consensus to unlock the next generation-enabling AGI technology, it seems  
709 far less likely.

710 Thus what may lock in the transitive endowment of improved control and value alignment  
711 technology between successive AGI generations is storing the technology enabling the next  
712 generation via dApps in the blockchain and requiring multiple AGIs to reach a consensus to unlock,  
713 license, and use the tech, including control and value alignment, to succeeding AGI generations. In  
714 this manner hacking the blockchain, or attempting to coerce individual consensus agents, would be  
715 thwarted in the same way as it is done in the nascent DL methodology extant today. In addition,  
716 game theoretic design approaches may help ensure stable evolutionary strategies, likely a succession  
717 of them (dynamic equilibrium) [24, 53, 72]. In that context note there can be no Nash equilibrium with  
718 one overwhelmingly dominant player.

719 *Prima facie*, an entirely different way to put the principle underlying safe AGI solution to the  
720 singleton problem is to think of future AGI as a distributed automaton, and to recall von Neumann's  
721 solution to designing a reliable automaton from unreliable parts via redundancy [73]. Critical AGI  
722 algorithms may reside on multiple agents in one or more generations, who require consensus for  
723 ongoing access and cross-check each other in real time (like a deadman's switch).

#### 724 4. Conclusion and Future Work

725 One epochal event likely to trigger AGI, if not the key event, is AI self-programming, or any  
726 other self-improvement, positive-feedback advancement. Close attention should be given to that  
727 development path, progress metrics and simulations developed, and measures enacted to ensure that  
728 access to key self-improvement techniques is via licensing with appropriate safeguards.

729 Before self-improvement technology can be unleashed, AI behavior control systems need to be  
730 developed and tested in transparent, non-hackable simulation sandbox environments as proposed  
731 by Bore et al. [52] seems essential.

732 If the ASILOMAR AI Principles [9] or similar agreements are akin to the U.S. Declaration of  
733 Independence, we need to move to the 'Articles of Confederation', step up the current 'Federalist  
734 Papers' stage, and then move to enact the 'Constitution', i.e. firm and ineluctable consensuses among  
735 leading AI workers, encrypted via DLT, as are possible.

736 **Acknowledgments.** The Seminar on Natural and Artificial Computation (SNAC), hosted by Stewart Wilson and  
737 Dave Waltz, which Stephen J. Smith and I chaired at the Rowland Institute for Science, was instrumental in my  
738 AI/machine learning education. A workshop, Can Machines Think?, organized by Kurt Thearling at Thinking  
739 Machines Corporation c. 1990 catalyzed thoughts on AI safety toward a diversified, 'separation and balance of  
740 power' pluralistic system where AI agents competed to satisfy human needs.

#### 741 Appendix: Simple Syllogisms to Help Formalize the Problem Statement

742 *Proposition 1, Probability of Malevolent Use:* With no restriction on AGI technology flow via  
743 licensing, malevolent use of AGI is a certainty.

744 *Proof:* Assume: 1. There exist malevolent or incompetent humans. 2. They can freely access AGI  
745 technology (e.g. via an AGI app offering ‘just add goals’). Then: There will exist malevolent use of  
746 AGI.

747 *Corollary 1A:* With no restriction on technology flow via licensing, malevolent AGI will destroy  
748 a significant portion of humanity, or the entire species.

749 *Proof:* Assume in addition to 1 & 2: 3a. Some malevolent humans would employ AGI for mass  
750 destruction; 3b. Some would seek mass destruction of the entire species.

751 *Corollary 1B:* With no restriction on technology flow via licensing, there is a chance that  
752 malevolent AGI may destroy the entire species.

753 *Proof:* Assume in addition to 1, 2 & 3: 4. Some malevolent humans are incompetent in their  
754 attempts to contain their destructive goals.

755 *Corollary 1C:* The more widely available and easily accessible the destructive AI or AGI, the  
756 higher the probability of its deliberate or inadvertent destructive use.

757 *Proposition 2, Extent of Danger, Importance of Containing:* Containing AGI is more important than  
758 containing nuclear weapon usage.

759 *Proof:* Assume AGI will have control, by deliberate human consent and design, by accident, or  
760 by AGI intervention, over nuclear weapons, and in addition, other critical resources, e.g. power grid,  
761 transportation systems, financial systems, negotiations between states, etc. Then clearly AGI  
762 containment is more important than containment of nuclear weapon use.

763 *Proposition 3, Probability of Value Misalignment:* Given the unlimited availability of an AGI  
764 technology as enabling as ‘just add goals’, then AGI-human value misalignment is inevitable.

765 *Proof:* From a subjective point of view, all that is required is value misalignment by the operator  
766 who adds to the AGI his/her own goals, stemming from his/her values, that conflict with any human’s  
767 values; or put more strongly, the effects are malevolent as perceived by large numbers of humans.  
768 From an absolute point of view, all that is required is misalignment of the operator who adds his/her  
769 goals to the AGI system that conflict with the definition of morality presented here, voluntary, non-  
770 fraudulent transacting (Axiom 3), i.e. usage of the AGI to force his/her preferences on others.

771 **Funding:** This research received no external funding, but I wish it had.

772 **Conflicts of Interest:** The author declares no conflict of interest.

## 773 References

774 Figure 1 is from Shoham et al. [1], Creative Commons License:  
775 <https://creativecommons.org/licenses/by-nd/4.0/legalcode>.

776  
777

- 778 1. Shoham, Y., Perrault, R., Brynjolfsson, E., Clark, J., Manyika, J., Niebles, J. C., Lyons, T.,  
779 Etchemendy, J. *et al.*, The AI Index 2018 Annual Report. 94. Available online:  
780 <http://cdn.aiindex.org/2018/AI%20Index%202018%20Annual%20Report.pdf>.
- 781 2. Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., van den Driessche, G., Schrittwieser,  
782 J., Antonoglou, I. *et al.* Mastering the game of Go with deep neural networks and tree search,  
783 *Nature*, **2016**, 529, 484-9.
- 784 3. Rodriguez, J. *The Science Behind OpenAI Five that just Produced One of the Greatest*  
785 *Breakthrough in the History of AI*. Medium. (Accessed: 12/1/18). Available online:  
786 [https://towardsdatascience.com/the-science-behind-openai-five-that-just-produced-one-of-](https://towardsdatascience.com/the-science-behind-openai-five-that-just-produced-one-of-the-greatest-breakthrough-in-the-history-b045bcdc2b69)  
787 [the-greatest-breakthrough-in-the-history-b045bcdc2b69](https://towardsdatascience.com/the-science-behind-openai-five-that-just-produced-one-of-the-greatest-breakthrough-in-the-history-b045bcdc2b69).
- 788 4. Brown, N. and Sandholm, T. Reduced Space and Faster Convergence in Imperfect-Information  
789 Games via Pruning, in *Proceedings of the 34th International Conference on Machine Learning*,  
790 Sydney, Australia, 2017.
- 791 5. Knight, W. 2017. Why Poker Is a Big Deal for Artificial Intelligence. *MIT Technology Review*.  
792 Available online: [https://www.technologyreview.com/s/603385/why-poker-is-a-big-deal-for-](https://www.technologyreview.com/s/603385/why-poker-is-a-big-deal-for-artificial-intelligence/)  
793 [artificial-intelligence/](https://www.technologyreview.com/s/603385/why-poker-is-a-big-deal-for-artificial-intelligence/).

- 794 6. Lake, B. M., Salakhutdinov, R., and Tenenbaum, J. B. Human-level concept learning through  
795 probabilistic program induction, *Science*, **2015**, 350, 1332-1338.
- 796 7. Ferrucci, D., Levas, A., Bagchi, S., Gondek, D., and Mueller, E. T. Watson: Beyond Jeopardy!,  
797 *Artificial Intelligence*, **2013**, 199-200, 93-105.
- 798 8. Omohundro, S. Autonomous technology and the greater human good, *Journal of Experimental*  
799 *and Theoretical Artificial Intelligence*, **2014**, 26, 303-315.
- 800 9. Future of Life Institute. *ASILOMAR AI Principles*. Future of Life Institute. (Accessed: 22  
801 December 2018). Available online: <https://futureoflife.org/ai-principles/>.
- 802 10. Bostrom, N., *Superintelligence: Paths, Dangers, Strategies*. Oxford, England: Oxford  
803 University Press, 2016. 415 pp.
- 804 11. Babcock, J., Kramar, J., and Yampolskiy, R., Guidelines for artificial intelligence containment.  
805 13. Accessed: 1 October 2018. Available online: <https://arxiv.org/abs/1707.08476>.
- 806 12. Dawkins, R. and Krebs, J. R. Arms races between and within species, *Proceedings of the Royal*  
807 *Society of London B*, **1979**, 205, 489–511.
- 808 13. Rabesandratana, T. Europe moves to compete in global AI arms race, *Science*, **2018**, 360 474.
- 809 14. Zwetsloot, R., Toner, H., and Ding, J. Beyond the AI Arms Race, *Foreign Affairs*, **2018**, 97,
- 810 15. Geist, E. M. It's already too late to stop the AI arms race—We must manage it instead, *Bulletin*  
811 *of the Atomic Scientists*, **2016**, Vol. 72, 318-321.
- 812 16. Tannenwald, N. The Vanishing Nuclear Taboo? How Disarmament Fell Apart, *Foreign Affairs*,  
813 **2018**, 97, 16-24
- 814 17. Callaghan, V., Miller, J., Yampolskiy, R., and Armstrong, S. The Technological Singularity:  
815 Managing the Journey: Springer, 2017, p. 261. [Online]. Available online:  
816 <https://www.springer.com/us/book/9783662540312>. Accessed 21 December 2018.
- 817 18. Yampolskiy, R. Taxonomy of Pathways to Dangerous Artificial Intelligence, presented at  
818 *Workshops of the 30th AAAI Conference on AI, Ethics, and Society*, AAAI. **2016**. Available  
819 online: <https://arxiv.org/abs/1511.03246>.
- 820 19. Turchin, A. *A Map: AGI Failures Modes and Levels*. (Accessed: 5 February 2018). Available  
821 online: <http://immortality-roadmap.com/AIfails.pdf>.
- 822 20. Brundage, M., Avin, S., Clark, J., Toner, H., and Eckersley, P., The Malicious Use of AI-  
823 Forecasting, Prevention, and Mitigation. 101. 2/20/2018. Accessed: 12/2/2018. Available  
824 online: <https://arxiv.org/abs/1802.07228>.
- 825 21. Manheim, D. Multiparty Dynamics and Failure Modes for Machine Learning and Artificial  
826 Intelligence, *Big Data Cogn. Comput.*, **2019**, 3, 15.
- 827 22. von Neumann, J. The General and Logical Theory of Automata, in *The World of Mathematics*,  
828 vol. 4, J. R. Newman, Ed. New York: John Wiley & Sons, 1956], 2070-2098.
- 829 23. Narahari, Y., *Game Theory and Mechanism Design* (IISc Lecture Notes Series, no. 4).  
830 Singapore: IISc Press/World Scientific, 2014. 492 pages.
- 831 24. Burtsev, M. and Turchin, P. Evolution of cooperative strategies from first principles, *Nature*,  
832 **2006**, 440, 1041-1044.
- 833 25. Doyle, J. C., Alderson, D. L., Li, L., Low, S., Roughan, M., Shalunov, S., Tanaka, R., and  
834 Willinger, W. The “robust yet fragile” nature of the Internet, *Proceedings of the National*  
835 *Academy of Sciences of the United States of America*, **2005**, 102, 14497.
- 836 26. Alzahrani, N. and Bulusu, N. Towards True Decentralization: A Blockchain Consensus  
837 Protocol Based on Game Theory and Randomness, presented at *Decision and Game Theory*  
838 *for Security*, Seattle WA, Springer. **2018**. Available
- 839 27. Nakamoto, S., Bitcoin: A Peer-to-Peer Electronic Cash System. 9. Accessed: 22 December  
840 2018. Available online: <https://bitcoin.org/en/bitcoin-paper>.
- 841 28. Christidis, K. and Devetsikiotis, M. Blockchains and Smart Contracts for the Internet of  
842 Things, *IEEE Access*, **2016**, 4, 2292-2303.
- 843 29. FinYear. *Eight Key Features of Blockchain and Distributed Ledgers Explained*. (Accessed: 5  
844 November 2018). Available online: [https://www.finyear.com/Eight-Key-Features-of-Blockchain-and-Distributed-Ledgers-Explained\\_a35486.html](https://www.finyear.com/Eight-Key-Features-of-Blockchain-and-Distributed-Ledgers-Explained_a35486.html).
- 845 30. Szabo, N., Smart Contracts: Building Blocks for Digital Markets. 16. Available online:  
846 [http://www.fon.hum.uva.nl/rob/Courses/InformationInSpeech/CDROM/Literature/LOTwinterschool2006/szabo.best.vwh.net/smart\\_contracts\\_2.html](http://www.fon.hum.uva.nl/rob/Courses/InformationInSpeech/CDROM/Literature/LOTwinterschool2006/szabo.best.vwh.net/smart_contracts_2.html).
- 847 848
- 849 31. Asimov, I., *I, Robot*. New York: Gnome Press, 1950. 253 pp.



- 850 32. Collendanchise, M. and Ogren, P., Behavior Trees in Robotics and AI. 198. Accessed:  
851 12/2/2018. Available online: <https://arxiv.org/abs/1709.00084>.
- 852 33. Hind, M., Mehta, S., Mojsilovic, A., Nair, R., Ramamurthy, K. N., Olteanu, A., and Varshney,  
853 K. R., Increasing Trust in AI Services through Supplier's Declarations of Conformity. 29. doi:  
854 arXiv:1808.07261.
- 855 34. Yampolskiy, R. and Sotala, K. Risks of the Journey to the Singularity, in *The Technological*  
856 *Singularity*, V. Callaghan, J. Miller, R. Yampolskiy, and S. Armstrong, Eds. Berlin: Springer,  
857 2017], 11-24.
- 858 35. Stigler, G. J. The development of utility theory I, *Journal of Political Economy*, **1950**, 58, 307-  
859 327.
- 860 36. Galambos, A. J., *Thrust for Freedom: An Introduction to Volitional Science*. CA USA:  
861 Universal Scientific Publications, 2000.
- 862 37. Eddington, A. S., *The Philosophy of Physical Science* (Tarnier lectures, 1938). 230 pp.
- 863 38. Rothbard, M. N., *Man, Economy, and State: A Treatise on Economic Principles*. Auburn AL:  
864 Ludwig Von Mises Institute, 1993. 987 pp.
- 865 39. Webster, T. J. Economic efficiency and the common law, *Atlantic Economic Journal*, **2004**,  
866 32, 39-48.
- 867 40. Todd, P. M. and Gigerenzer, G., *Ecological Rationality Intelligence in the World* (Evolution  
868 and Cognition). Oxford: New York: Oxford Univ. Press, 2011.
- 869 41. Gigerenzer, G. and Todd, P. M., *Simple Heuristics That Make Us Smart* (Evolution and  
870 Cognition). Oxford; New York: Oxford Univ. Press, 1999. 416 pp.
- 871 42. Hayek, F. The use of knowledge in society, *American Economic Review*, **1945**, 35, 519-530.
- 872 43. Coase, R. H. The problem of social cost, *Journal of Law and Economics*, **1960**, 3, 1-44.
- 873 44. Minerva, R., Biru, A., and Rotondi, D. "Towards a definition of the Internet of Things (IOT),"   
874 IEEE/Telecom Italia. 2015. Accessed: 12 May 2018. Available online:  
875 [https://iot.ieee.org/images/files/pdf/IEEE\\_IoT\\_Towards\\_Definition\\_Internet\\_of\\_Things\\_Issue1\\_14MAY15.pdf](https://iot.ieee.org/images/files/pdf/IEEE_IoT_Towards_Definition_Internet_of_Things_Issue1_14MAY15.pdf).
- 876
- 877 45. Chokhani, S., Ford, W., Sabett, R., Merrill, C., and Wu, S., Internet X.509 Public Key  
878 Infrastructure Certificate Policy and Certification Practices Framework. 94. Accessed: 23  
879 December 2018. Available online: <ftp://ftp.rfc-editor.org/in-notes/rfc3647.txt>.
- 880 46. Szabo, N., Formalizing and Securing Relationships on Public Networks, *First Monday*.  
881 Accessed: 1997-09-01. doi: 10.5210/fm.v2i9.548. Available online:  
882 <https://ojphi.org/ojs/index.php/fm/article/view/548/469>.
- 883 47. Anonymous. *Hyperlink*. Wikipedia. (Accessed: 25 December 2018). Available online:  
884 <https://en.wikipedia.org/wiki/Hyperlink#History>.
- 885 48. Omohundro, S., Cryptocurrencies, Smart Contracts, and Artificial Intelligence, *AI Matters*, vol.  
886 1, no. 2. 19-21. doi: 10.1145/2685328.268533.
- 887 49. Pierce, B. Encoding law, regulation, and compliance ineluctably into the blockchain, presented  
888 at *WALL ST. Conference*, Twitter. **2018**. Available online:  
889 <https://twitter.com/LeX7Mendoza/status/1085643180744339456/video/1>.
- 890 50. Barnum, K. J. and O'Connell, M. J. Cell cycle regulation by checkpoints, *Methods Mol Biol*,  
891 **2014**, 1170, 29-40.
- 892 51. Buterin, V. and et al., A Next-Generation Smart Contract and Decentralized Application  
893 Platform. White Paper, 224. Accessed: 29 January 2019. Available online:  
894 <https://github.com/ethereum/wiki/wiki/White-Paper#applications>.
- 895 52. Bore, N. K., Raman, R. K., Markus, I. M., Remy, S., Bent, O., Hind, M., Pissadaki, E. K.,  
896 Srivastava, B. *et al.*, Promoting distributed trust in machine learning and computational  
897 simulation via a blockchain network. 13. Available online: <https://arxiv.org/abs/1810.11126>.
- 898 53. Nowak, M. A. and Highfield, R., *Super Cooperators: Evolution, Altruism and Human*  
899 *Behaviour or Why We Need Each Other to Succeed*. Edinburgh; New York: Canongate, 2011.  
900 330 pp.
- 901 54. Rasmusen, E., *Games and Information : an Introduction to Game Theory*. 4th ed. Malden, MA;  
902 Oxford: Blackwell, 2007. 528 pp.
- 903 55. Shoham, Y. and Leyton-Brown, K., *Multiagent Systems : Algorithmic, Game-theoretic, and*  
904 *Logical Foundations*. Cambridge; New York: Cambridge Univ. Press, 2009. 483 pp.

- 905 56. Newell, A., *Unified Theories of Cognition* (William James Lectures, no. 1987). Cambridge,  
906 MA: Harvard Univ. Press, 1990. 549 pp.
- 907 57. Potapov, A., Svitenkov, A., and Vinogradov, Y. Differences between Kolmogorov complexity  
908 and Solomonoff probability: consequences for AGI, in *Artificial General Intelligence*, 2012,  
909 vol. 7716, Berlin, Heidelberg: Springer.
- 910 58. Waltz, D. L. The prospects for building truly intelligent machines, *Daedalus: Proc. AAAS*,  
911 **1988**, 117, 191-212.
- 912 59. Westfall, R. S., *Never at Rest: A Biography of Isaac Newton*. Cambridge: Cambridge Univ  
913 Press, 1980. 908 pp.
- 914 60. Cohen, I. B., *Revolution in Science*. Cambridge MA: Harvard University Press, 1987. 732  
915 pages.
- 916 61. Cassidy, D. C., *J. Robert Oppenheimer and the American Century*. New York: Pearson/Pi  
917 Press, 2005. 462 pages.
- 918 62. Goodchild, P., *J. Robert Oppenheimer: Shatterer of Worlds*. USA: BBC/WGBH, 1981. 301  
919 pages.
- 920 63. Associated Press. *U.S. Nobel Laureate Knew about Chinese Scientist's Gene-Edited Babies*.  
921 NBC News. (Accessed: 29 January 2019). Available online:  
922 [https://www.nbcnews.com/health/health-news/u-s-nobel-laureate-knew-about-chinese-](https://www.nbcnews.com/health/health-news/u-s-nobel-laureate-knew-about-chinese-scientist-s-gene-n963571)  
923 [scientist-s-gene-n963571](https://www.nbcnews.com/health/health-news/u-s-nobel-laureate-knew-about-chinese-scientist-s-gene-n963571).
- 924 64. Gill, K., Sophia Robot Creator: We'll Achieve Singularity in 5 to 10 years. Article & video, 1.  
925 Accessed: 29 January 2019. Available online: [https://cheddar.com/videos/sophia-bot-creator-](https://cheddar.com/videos/sophia-bot-creator-well-achieve-singularity-in-five-to-10-years)  
926 [well-achieve-singularity-in-five-to-10-years](https://cheddar.com/videos/sophia-bot-creator-well-achieve-singularity-in-five-to-10-years).
- 927 65. Predictions of Human-Level AI Timelines, no. 15 January 2019. 3. Available online:  
928 <https://aiimpacts.org/category/ai-timelines/predictions-of-human-level-ai-timelines/>.
- 929 66. Concrete AI Tasks for Forecasting, no. 15 January 2019. 5. Available online:  
930 <https://aiimpacts.org/concrete-ai-tasks-for-forecasting/>.
- 931 67. Muller, V. C. and Bostrom, N. Future progress in artificial intelligence: A survey of expert  
932 opinion, in *Fundamental Issues of Artificial Intelligence*, vol. 377(Synthese Library, no. 377)  
933 Berlin: Springer, 2016].
- 934 68. Bughin, J., Hazan, E., Manyika, J., and Woetzel, J., Artificial Intelligence - The Next Digital  
935 Frontier? 80. Accessed: 15 January 2019. Available online: Search at  
936 <https://www.mckinsey.com/mgi>.
- 937 69. Eckersley, P. and Nasser, Y. *AI Progress Measurement*. Electronic Frontier Foundation.  
938 (Accessed). Available online: <https://www.eff.org/ai/metrics>.
- 939 70. Housman, G., Byler, S., Heerboth, S., Lapinska, K., Longacre, M., Snyder, N., and Sarkar, S.  
940 Drug resistance in cancer: an overview, *Cancers (Basel)*, **2014**, 6, 1769-92.
- 941 71. Al-Dimassi, S., Abou-Antoun, T., and El-Sibai, M. Cancer cell resistance mechanisms: a mini  
942 review, *Clin Transl Oncol*, **2014**, 16, 511-6.
- 943 72. ten Broeke, G. A., van Voorn, G. A. K., Ligtenberg, A., and Molenaar, J. Resilience through  
944 adaptation, *PLoS One*, **2017**, 12, 1-21.
- 945 73. Shannon, C. E. and McCarthy, J. *Automata Studies*, Annals of Mathematics Studies, No. 34.  
946 Princeton: Princeton Univ. Press, 1956, 285 pages.