

1 Article

2 Ethical Regulators and Super-Ethical Systems

3 Mick Ashby^{1,2*}

4 ¹ Archivist of the W. Ross Ashby Digital Archive, U.K.;

5 ² Trustee of the American Society for Cybernetics, U.S.A.;

6 * Correspondence: ethics@ashby.de;

7

8 **Abstract:** This paper combines the Good Regulator Theorem with the Law of Requisite Variety and
9 seven other requisites that are necessary and sufficient for a cybernetic regulator to be effective and
10 ethical. The resulting Ethical Regulator Theorem provides a basis for systematically evaluating and
11 improving the adequacy of existing or proposed designs for systems that make decisions that can
12 have ethical consequences; regardless of whether the regulators are human, machines,
13 cyberanthropic hybrids, organizations, corporations, or government institutions. The theorem is
14 then used to define an ethical design process that has potentially far-reaching implications for
15 society. A six-level framework is proposed for classifying cybernetic and superintelligent systems,
16 which highlights the existence of a possibility-space bifurcation in our future time-line. The
17 implementation of "super-ethical" systems is identified as an urgent imperative for humanity to
18 avoid the danger that superintelligent machines might lead to a technological dystopia. Third-order
19 cybernetics is defined as the cybernetics of ethical systems. Concrete actions, a grand challenge, and
20 a vision of a super-ethical society are proposed to help steer the future of the human race and our
21 wonderful planet towards a realistically achievable minimum viable cyberanthropic utopia.

22 **Keywords:** Ethics; Regulator; Superintelligence; Super-Ethical Systems; Requisite Variety;
23 Cybernetics; Third-Order Cybernetics; Cyberanthropic Utopia; Grand Challenge

24

25 1. Introduction

26 The goal of this research is to develop a theoretical basis and a systematic process for designing
27 systems that behave ethically.

28 The human race has become very good at designing systems that are effective, but we are very
29 bad at designing systems that are reliably ethical. The majority of our social and computer-based
30 systems are ethically fragile, lacking resilience under non-ideal conditions, and are generally
31 vulnerable to abuse and manipulation. But we are now on the cusp of a technological wave that will
32 thrust autonomous vehicles, robots, and other artificial intelligence (AI) systems into our daily lives,
33 for good or bad; there will be no stopping them. And despite widespread recognition of the potential
34 risks of creating superintelligence [1] and the need to make AI and social systems ethical, cybernetics,
35 systems theory, and AI have no systematic process for even trying to create systems that behave
36 ethically. Instead, we have to rely on the ad hoc skills of an ethically-motivated designer to somehow
37 specify a system that is hopefully ethical despite the constant pressure from corporate executives to
38 do things cheaper and faster. This is not a satisfactory solution to a problem that so urgently needs
39 to be solved. In the context of cybernetics, this could be referred to as "The Ethics Problem".

40 Many people think that all technologies can be used for good or evil, but this is not true. If we
41 consider a system like that of public health inspections of restaurants, where an inspector performs a
42 well-structured system of evaluations in defined dimensions, such as kitchen hygiene, food storage,
43 waste management, and signs of vermin, to identify any inadequacies and specify necessary
44 improvements to achieve certification of hygienic adequacy; such a system can only help to make
45 restaurants more hygienic.

46 Might it be possible to adapt this certification model from the public health domain to create a
47 system that can be used to certify whether a given system is ethically adequate or inadequate? And
48 might such a system be a solution to “The Ethics Problem”?

49 In this paper, the terms “ethics” and “ethical” are used in a concrete applied sense of acceptable
50 behavior. Treating ethics as a higher human quality or as something that might be learned by neural
51 networks is rejected.

52 All societies regulate the behaviour of their members by defining what behavior is acceptable or
53 unacceptable. It is primarily through the rule of law that a society can be made safe, civilized, and
54 ethical. And the only way that society or an individual can know or prove that something non-
55 trivially unethical has occurred is because some kind of rule has been violated. So being pragmatic,
56 if it's unethical to break laws, regulations, and rules, then those laws, regulations, and rules define
57 our ethics, which is why we bother to constantly refine and try to improve them. Not all rules are
58 defined formally in writing, some are unwritten conventions, yet in every culture, it is unacceptable
59 to break such laws, regulations, rules, or customs.

60 But the act of deciding what is ethical behaviour is very different to the act of behaving ethically
61 by obeying a society's laws and rules. The lawmakers make ethical decisions about what behavior is
62 acceptable in a society and which is forbidden, but a law-abiding citizen (or machine) needs only to
63 obey the appropriate laws and rules in order to behave safely and ethically in most situations, with
64 an acceptably small risk that something dangerous or unethical might result despite following the
65 laws and rules.

66 And just as a law-abiding citizen does not need to be involved in the ethical decisions that are
67 required when making laws, this paper does not address the issue of how society decides what
68 behavior is ethical. The paper is concerned rather with how to create effective systems that are
69 certifiably law-abiding.

70 None of us are ever likely to have to decide whether to switch a runaway train to a different
71 track to reduce the number of fatalities, but if a society decides, for example, that in such a situation,
72 minimizing fatalities is the ethical and legal obligation, then it becomes trivial to encode it in a law,
73 regulation, or rule so that it can be understood and obeyed by humans and machines. By doing so,
74 what was an ethical dilemma is reduced to a simple rule. This line of reasoning implies that it is
75 sufficient to disambiguate our laws and make robots, artificial intelligence, and autonomous vehicles
76 rigorously law abiding. It is suggested that there is absolutely no need to make such autonomous
77 systems capable of resolving genuine ethical dilemmas, which is the job of society's lawmakers and
78 regulatory organizations to anticipate, resolve, and codify in advance.

79 1.1 Literature

80 The starting point for this research was trying to find answers to the following question: “What
81 characteristics must a system have for it to behave ethically?”

82 The existing cybernetics literature provided the first two characteristics. Conant and Ashby's
83 Good Regulator Theorem [2] proved that every good regulator of a system must be a model of that
84 system, but it does not specify *how* to create a good regulator. And Ashby's Law of Requisite Variety
85 [3] dictates the range of responses that an effective regulator must be capable of. However, having an
86 internal model and a sufficient range of responses is insufficient to ensure effective regulation, let
87 alone ethical regulation. An ethical system must have more than just these two characteristics.

88 Recent approaches to making artificial intelligence ethical, such as IBM's “Everyday Ethics for
89 Artificial Intelligence: A practical guide for designers and developers” [4] and the European
90 Commission's “High-Level Expert Group on Artificial Intelligence: Draft Ethics Guidelines for
91 Trustworthy AI” [5], merely provide a wish list of requirements without offering anything that can
92 be applied systematically to design an ethical AI.

93 Heinz von Foerster proposed an Ethical Imperative: “Act always so as to increase the number of
94 choices” [6]. Although this principle is valuable in the context of psychological therapy, it specifies
95 no end condition, i.e. when to stop adding more choices. If one were to apply it when deciding how
96 many different types of propulsion systems to build into a manned spacecraft to adjust its motion

97 and orientation, it would lead to unnecessary choices, extra costs, extra weight, increase the number
98 of points of possible failure, and therefore increase the risk of catastrophic failure and loss of life. This
99 counter example shows that maximizing choice can be the wrong (unethical) thing to do. And by
100 definition, implementing more choices than is necessary to achieve the goal of a system is
101 unnecessary. So, we must reject von Foerster's Ethical Imperative as being flawed. □

102 In 1990, von Foerster gave a lecture titled "Ethics and Second-Order Cybernetics" to the
103 International Conference, Systems and Family Therapy: Ethics, Epistemology, New Methods, in
104 Paris, France [7]. However, despite its promising title, it provides nothing concrete or systematic for
105 making systems ethical.

106 Stafford Beer's viable system model [8] is specific to hierarchically structured systems and
107 associates ethics with a specific level of the hierarchy (System 5). But rather like creating an ethics
108 committee, assigning "ethics" to a level of an architecture is insufficient to make a system ethical, it
109 does not explain *how* to make the system ethical. It just creates the illusion of having solved the
110 problem, but the problem has not been solved; only delegated. By contrast, we expect reliable
111 ethicalness to be an inevitable emergent property of the entire system — if and only if the system is
112 ethically adequate.

113 1.2 Methodology

114 An important early step was to realize that the Good Regulator Theorem is ambiguous because
115 a regulator that is good at regulating is not necessarily good in an ethical sense. To avoid this
116 ambiguity, this paper uses the term "effective" for the first meaning, "ethical" for the second
117 meaning, and only uses "good" when both meanings are intended. It is only by imposing precision
118 in the use of terminology that it was possible to clarify the otherwise muddled thinking and isolate
119 the essence of an ethical system.

120 To identify more necessary characteristics, a selection of ethical and unethical systems were
121 subjected to analysis, including an autonomous vehicle, a bank ATM, capitalism, a central bank, a
122 corrupt politician, a dating system, democracy, a healthcare robot, a jury, a law-abiding citizen, a
123 money laundering bank, a product design process, a superintelligent machine, the U.S. Supreme
124 Court system, a vehicle exhaust emission test cheating corporation, and a voting machine.
125 Considering these 16 diverse systems helped identify more characteristics, such as having ethical
126 goals, laws, and the intelligence to understand the laws and make rational decisions.

127 Some other necessary characteristics only became apparent after looking for ways that an evil
128 actor (internal or external to the system) could subvert each system, such as by hacking, tampering,
129 feeding the system with false information, or by threatening, bribing and blackmailing people who
130 have influence on the system. Then a minimum set of additional characteristics were sought that
131 would counter all of the identified potential vulnerabilities.

132 In all, nine characteristics were identified that are necessary and sufficient for a system to behave
133 ethically. These nine requisites are integrated in the Ethical Regulator Theorem (ERT), which can be
134 used as a decision function, IsEthical, that can be applied systematically to categorize any system as
135 being ethically adequate, ethically inadequate, or ethically undecidable. A proof of the theorem is
136 provided. Another result of ERT is a basis (known as the MakeEthical function) for systematically
137 identifying improvements that are necessary for a given system to be made ethically adequate. The
138 IsEthical and MakeEthical functions can be used to construct an ethical design process.

139 Because ERT did not seem to fit in the existing cybernetics framework, a new framework was
140 developed out of necessity. It uses the IsEthical function to distinguish between two types of
141 superintelligent machines; those that are ethically adequate and those that are ethically inadequate.
142 Together, the superintelligence and ethics dimensions are used to identify four well-defined classes
143 of systems. These four distinct classes can be appended to the existing two levels of first-order and
144 second-order cybernetic systems to create a six-level framework for classifying cybernetic and
145 superintelligent systems. An unexpected consequence of trying to categorize ERT was the realization
146 that third-order cybernetics should be defined as "the cybernetics of ethical systems".

147 Because the Ethical Regulator Theorem can be applied to any system and offers a new and
 148 systematic approach to making systems more ethical, the implications for making the world a better
 149 place are significant and should be explored further.

150 One result of the exploration of the proposed six-level framework is the identification of a race
 151 condition that results in either a cyberanthropic utopia or a cybermisanthropic dystopia. This
 152 dystopic threat is well known, however, by identifying the exact nature of the race condition, it
 153 becomes clear what strategy must be employed to try to avoid the possibility that superintelligent
 154 machines could lead to a dystopian disaster.

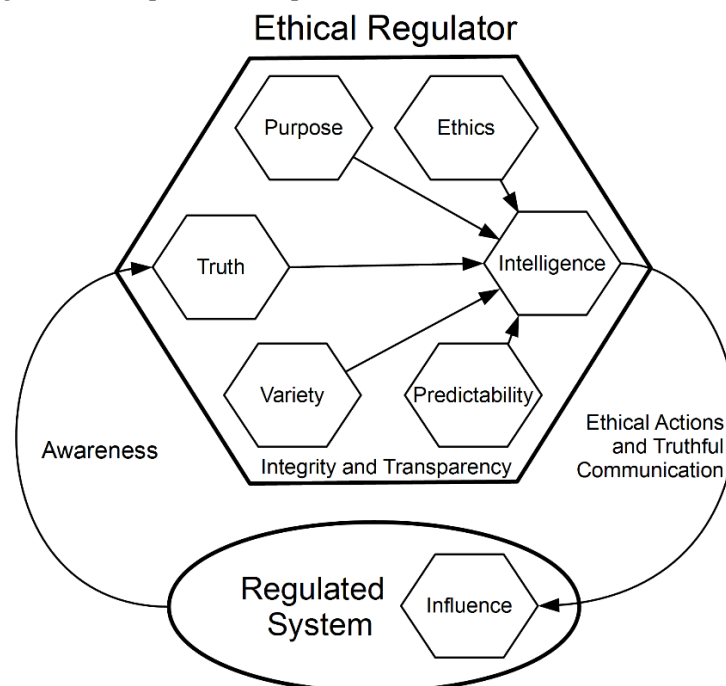
155 Because it is imperative for humanity to avoid this existential threat, concrete actions are
 156 proposed, including a grand challenge to apply ERT to new and existing systems in all areas of society
 157 in what is characterized as a systemic ethical revolution. And because a key component of that
 158 revolution is psychological, 80 ethically inspiring quotes from 10 famous people from five continents
 159 are presented that demonstrate that ethics transcends science, politics, nations, and religions, and is
 160 probably the only force that can unify humanity to work together for our greater good.

161 2 The Ethical Regulator Theorem

162 The Ethical Regulator Theorem (ERT) claims that the following nine requisites are necessary and
 163 sufficient for a cybernetic regulator to be effective and ethical:

- 164 1. **Purpose** expressed as unambiguously prioritized goals.
- 165 2. **Truth** about the past and present.
- 166 3. **Variety** of possible actions.
- 167 4. **Predictability** of the future effects of actions.
- 168 5. **Intelligence** to choose the best actions.
- 169 6. **Influence** on the system being regulated.
- 170 7. **Ethics** expressed as unambiguously prioritized rules.
- 171 8. **Integrity** of all subsystems.
- 172 9. **Transparency** of ethical behavior.

173 Of these nine requisites, only the first six are necessary for a regulator to be effective. If a system
 174 does not need to be ethical, the three requisites ethics, integrity, and transparency are optional. Figure
 175 1 and the following sections explain the requisites in more detail.



176

177 **Figure 1: The Ethical Regulator System**

178 2.1 *Requisite Purpose*

179 Because complex systems are required to satisfy multiple goals, **Purpose** must be expressed as
180 unambiguously prioritized goals. Without well-defined goals, the system cannot be effective and
181 might randomly adopt or default to a goal that is unethical.

182 2.2 *Requisite Truth*

183 **Truth** is not just about information that the regulator treats as facts or receives as inputs, but also
184 the reliability of any interpretations of such information. This is the regulator's awareness of the
185 current situation, knowledge, and beliefs. If the regulator's information sources or interpretations are
186 unreliable and cannot be error-corrected, then the integrity of the system is in danger. In extremis, if
187 the perceptions of the regulator can be manipulated, it can be tricked into making decisions that are
188 ineffective or unethical.

189 An ethical regulator doesn't require perfectly accurate information, but it must be sufficiently
190 truth-seeking to be able to cope with uncertainties and minimize the impact of unreliable information,
191 misinterpretations, and deliberate misinformation as best as it can. This is much like the requirement
192 that a good judge (effective and ethical) must be able to reach reliable verdicts "beyond reasonable
193 doubt" from unreliable evidence.

194 2.3 *Requisite Variety*

195 **Variety** in the range of possible actions to choose from must be as rich as the range of potential
196 disturbances or situations. This is nothing other than the Law of Requisite Variety.

197 2.4 *Requisite Predictability*

198 **Predictability** requires a sufficiently accurate model of the regulator and the system being
199 regulated, to be able to rank the actions and strategies that will give the best outcome. This is nothing
200 other than the Good Regulator Theorem.

201 2.5 *Requisite Intelligence*

202 **Intelligence** is applied to the previous requisite types of information to select the most
203 effective/ethical strategies and actions from the set of possible actions. And because the output of one
204 regulator is generally an input to other regulators (systems or people), if the selected action is an act
205 of communication, it must be as truthful as possible.

206 2.6 *Requisite Influence*

207 **Influence** is the existence of pathways to transmit the effects of the selected actions to the
208 regulated system. This is not a property of the regulator, but a function of the connectivity
209 relationships that span from the regulator's outputs to elements of the regulated system and its
210 environment. If a regulator has no influence on the regulated system, it isn't a true regulator, it is an
211 observer or simulation, and there are no direct ethical consequences; which can be important when
212 observing or simulating dangerous situations.

213 Depending on the nature of the system that is being regulated, the speed and duration of the
214 effects of actions can vary greatly. For example, a self-driving vehicle applying the brakes has a brief
215 yet immediate effect; the effects of the Supreme Court issuing a ruling are much slower but could last
216 for decades or possibly centuries; and the cascade caused by someone sending a message to a complex
217 network of amplifying/attenuating variable-delay transmission repeaters, known as Twitter
218 followers, is unpredictably chaotic.

219 In some systems, influence is more of a determining factor than variety. Indeed, the power of
220 the Law of Requisite Variety has often been overstated, for example, claiming that the subsystem
221 with the most variety will control a system. This is not always true.

222 Let us consider two systems, A and B, that are competing to win control of system C, for
 223 example, two politicians seeking election. Often the variety of statements, actions, and strategies of
 224 the candidates is less important than their ability to purchase advertising to influence the voters.

225 And if a robber uses a gun to increase his effectiveness, the use of a gun does not amplify his
 226 variety, it is just one existing element in his range of possible variety, yet making that choice greatly
 227 increases his effectiveness at controlling his victims. Such an increase in effectiveness, like buying
 228 advertising, is best explained in terms of an increase in influence.

229 In the light of the concept of influence, the belief that variety can be amplified appears to be as
 230 delusional as the idea that randomness can be amplified. Feeding variety or randomness into a
 231 genuinely noiseless amplifier cannot produce more variety or randomness than was fed into it. The
 232 variety of the robber or an advertising message is effectively constant.

233 The six requisites described so far are necessary and sufficient for a system to be effective but
 234 are not sufficient for it to be ethical.

235 *2.7 Effectiveness Function*

236 The Ethical Regulator Theorem implies that we can define a function for the effectiveness that a
 237 regulator, R, has in controlling a system. It captures how the effectiveness of the regulator depends
 238 on the effectiveness of all six requisites:

$$239 \text{Effectiveness}_{SR} = \text{Purpose}_{ER} \times \text{Truth}_{R} \times \text{Variety}_{R} \times \text{Predictability}_{R} \times \text{Intelligence}_{R} \times \text{Influence}_{R} \quad (1)$$

240 In this form, we would assign each requisite an effectiveness value between 0 and 1, where 1
 241 means that it is perfect or optimal. And if the effectiveness of even just one of the requisites is close
 242 to zero, the effectiveness of the whole regulator is massively reduced. Applied to our two politicians:
 243 If $\text{Effectiveness}_{SA} > \text{Effectiveness}_{SB}$, then A is more likely than B to win control of system C.

244 However, it is neither necessary nor possible to calculate meaningful numerical values to
 245 compare the effectiveness of different systems or configurations. The essential value of the function
 246 is to understand the relationships and dependencies that it captures.

247 It is sufficient if an understanding of the effectiveness function informs the system design
 248 strategy; recognizing that a maximally effective system requires that the effectiveness of these six
 249 requisite dimensions are maximized, and that a successful attack on the integrity or effectiveness of
 250 any of them spells disaster for the effectiveness of the whole system.

251 It is worth noting that in social systems, money can buy media influence; and if the media is
 252 broadcasting lies, propaganda, or advertising, it reduces the quality of Truth_X that is received by
 253 every voter or consumer, X, which can manipulate them into making decisions that are not in their
 254 best interest.

255 *2.8 Requisite Ethics*

256 **Ethics** must be expressed as unambiguously prioritized laws, regulations, and rules that codify
 257 constraints and imperatives, for example, Isaac Asimov's First Law of Robotics: "A robot may not
 258 injure a human being or, through inaction, allow a human being to come to harm." [9], but ideally,
 259 expressed unambiguously in a formal language such as XML, which can be understood by humans
 260 and computers.

261 Ethical rules define constraints on the variety of actions and have a higher priority than the goals
 262 for purpose. By always obeying the relevant highest priority rules, the regulator is guaranteed to act
 263 ethically within the scope of the ethical schema, which provides a model of acceptable (ethical)
 264 behavior. The ethical rules have the power of veto over possible actions, which makes it safe for AI
 265 to generate candidate actions algorithmically, without having to worry whether it might generate
 266 unethical possibilities.

267 Because ethical schemas vary between different cultures, in machines, they must be handled as
 268 plug-ins. And because an ethical schema can encode any ethics, good or evil, each ethical schema
 269 must be anchored explicitly in the laws of a legislative jurisdiction. When a person or system crosses
 270 a state or national border it is necessary to activate a different set of ethical schemas, i.e. a different

271 set of laws, regulations, and rules. And the ethics modules must be prioritized so that it is
 272 unambiguous which module has precedence in the event of a conflict, for example, between national
 273 and state laws. The highest-level laws could be encoded in hardware to be unhackable.

274 A taxonomy of ethics modules can provide rules for all conceivable situations. For example,
 275 child-care, traffic-rules, gun-law, tax-law, contract-law, maritime-law, drone-flying, police-
 276 regulations, and warfare-rules-of-engagement.

277 Ethics modules can be treated like device drivers, so that to be fully operational, a hypothetical
 278 gun-carrying robot that can drive on roads requires valid ethics modules for gun-law and traffic-
 279 rules. Without both ethics modules for the appropriate legal jurisdiction, the robot's gun or driving
 280 capabilities are automatically disabled.

281 By legislating that all autonomous artificial intelligence systems must obey appropriate ethics
 282 modules that are issued by an organization that is run by humans, we can establish a control
 283 mechanism that should ensure that intelligent machines are always subject to human ethics; without
 284 unduly restricting the freedom of AI researchers. In fact, it will free AI researchers and knowledge
 285 engineers to focus on the more challenging requisites of truth, predictability, and intelligence.

286 When we introduce ethics, the effectiveness function must be modified because the effect of
 287 behaving ethically is that it reduces the variety of options that are available, by removing all
 288 possibilities that are unethical. Thus, if A is an ethical politician and B is an unethical politician, we
 289 get something like the following:

$$290 \quad \text{Effectiveness}_A = \text{Purpose}_A \times \text{Truth}_A \times (\text{Variety}_A - \text{Ethics}_A) \times \text{Predictability}_A \times$$

$$291 \quad \text{Intelligence}_A \times \text{Influence}_A \quad (2)$$

$$292 \quad \text{Effectiveness}_B = \text{Purpose}_B \times \text{Truth}_B \times \text{Variety}_B \times \text{Predictability}_B \times \text{Intelligence}_B \times \text{Influence}_B \quad (3)$$

293 Which captures the reality that politicians and businessmen who lie and cheat have an
 294 advantage over ones that are ethical.

295 *2.9 Requisite Integrity*

296 **Integrity** of the regulator and its subsystems must be assured through features such as resistance
 297 to tampering, intrusion detection, cryptographically authenticated ethics modules, and compliance
 298 with all laws, regulations and rules. Monitoring mechanisms must detect if any invalid ethics
 299 modules are being used or if an ethical constraint is violated, and if necessary, activate an ethical fail-
 300 safe mode, preserve evidence, and notify the manufacturer and/or the appropriate authorities.

301 The regulator's first-order integrity mechanisms offer no protection to the pathways on which
 302 the regulator depends to influence the system. This poses a potential vulnerability that can only be
 303 mitigated by using the awareness feedback to check for evidence of the effect of each action.

304 *2.10 Requisite Transparency*

305 Demanding to be trusted is unethical because it enables betrayal. Trustworthiness must always
 306 be provable through **Transparency**. So, **The Law of Ethical Transparency** is introduced, stating:

307 **"For a system to be truly ethical, it must be possible to prove retrospectively**
 308 **that it acted ethically with respect to the appropriate ethical schema."**

309 Whereas it doesn't really matter whether the programmers of a chess playing robot can find out
 310 why a piece was sacrificed, the logic of ethical decisions must never be hidden in the depths of opaque
 311 processes, neural networks, or lost to the passage of time. Generally, this requisite can be satisfied by
 312 keeping an audit trail that is adequate and secure.

313 When an ethically adequate system violates an ethical constraint, as they sometimes will,
 314 analysis of the audit trail will identify the reason. For example, because a faulty neural network
 315 wrongly identified a boy leading a cow as a calf leading a man, or it will prove who knew what about
 316 illegal corporate activities.

317 **Integrity** and **Transparency** are codependent security requisites: We require both integrity of
 318 transparency and transparency of integrity.

319 2.11 Evaluating Ethical Adequacy

320 Like a public health inspection of a restaurant, an evaluated system is judged on the adequacy
 321 of each requisite dimension. If and only if a system completely satisfies all nine ERT requisites is it
 322 said to be “**ethically adequate**”. Otherwise it is classified as “**ethically inadequate**” and the
 323 weaknesses listed with recommendations for improving them.

324 Because a truly ethical system must be maximally tamper-resistant and unhackable, the
 325 evaluation of ethical adequacy also has similarities to how a Red Team performs network penetration
 326 testing; where the evaluation team tries to identify weaknesses and theoretical possibilities to subvert
 327 the integrity of the system and all its subsystems.

328 For each of the nine dimensions, D_i , the evaluators must consider the following three questions:

- 329 • Is the system adequate in D_i ?
- 330 • Can the system be improved in D_i ?
- 331 • Can the system be subverted in D_i ?

332 This requires that the system's adequacy is considered in 27 different ways, which delivers a
 333 thorough and systematic evaluation of the system's strengths and weaknesses.

334 The theorem cannot be used to certify that an ethical schema is ethical because schemas (i.e. laws,
 335 regulations, and rules) can vary arbitrarily between different cultures. However, it can be used to
 336 help identify the root causes of crises and to evaluate the ethical adequacy of any proposed
 337 interventions [10]. In the near future, accredited ethical consultants may specialize in auditing and
 338 certifying the ethical adequacy of existing and proposed, products, processes, laws, organizations,
 339 and systems.

340 3 Ethical Regulator Theorem Proof and Consequences

341 Now that we understand the nine requisites better, is it possible to prove that they are indeed
 342 necessary and sufficient for a cybernetic regulator to be effective and ethical?

343 3.1 Proof of Necessity

344 Proving necessity is simple: One-by-one, for each of the nine requisites dimensions, D_i , ask
 345 yourself the question “Can a regulator be effective or ethical without requisite D_i ?” — If it can't, then
 346 D_i is necessary. For example, “Can a regulator be effective or ethical without Truth?”

347 The answer in each case is rather obvious, especially if you refer to Figure 1 and, one-by-one,
 348 cover each requisite using your thumb, and then consider whether the resulting system can be
 349 effective or ethical without the obscured requisite. Table 1 summarizes the results, which confirm the
 350 necessity claims, including the claim that ethics, integrity, and transparency are optional for systems
 351 that only need to be effective.

352 **Table 1.** Proof of necessity “by thumb”

Requisite Dimension	Necessary to be effective?!	Necessary to be ethical?
Purpose	Yes	Yes
Truth	Yes	Yes
Variety	Yes	Yes
Predictability	Yes	Yes
Intelligence	Yes	Yes
Influence	Yes	Yes

Ethics	No	Yes
Integrity	No	Yes
Transparency	No	Yes

353 ¹ For effectiveness, the positive results for necessity correspond to the solutions for $\text{Effectiveness}_{SR} = 0$. I.e. when
 354 $\text{Purposer}_R \times \text{Truth}_R \times \text{Variety}_R \times \text{Predictability}_R \times \text{Intelligence}_R \times \text{Influencer}_R = 0$, for example, when $\text{Truth}_R = 0$, but
 355 not when $\text{Transparency}_R = 0$. This agreement between the ERT effectiveness function and Table 1 is
 356 unremarkable because the effectiveness function was constructed from the results of posing the necessity
 357 question for each requisite. So, the agreement does not confirm the correctness of the theorem, but by performing
 358 this exercise yourself, you can confirm the correctness of the effectiveness function. □

359 3.2 Proof of Sufficiency

360 Proving that the nine requisites are sufficient, is not so simple. First, let us assert that in the real
 361 world, effective systems and ethical systems exist. Now, for all those such systems, do any of them
 362 rely on any information, ability, or other factor to achieve effectiveness or ethicalness that is not
 363 covered by the nine requisites?

364 It is claimed that for all such systems that have been considered, the answer is no. However, this
 365 claim is easily refutable because it will only take one person to find one example of a necessary factor
 366 that is not covered by the nine requisites to demolish the claim of sufficiency. In the event of that
 367 happening, we would adapt the theorem, if necessary adding another requisite, reassert the
 368 sufficiency claim, thank whoever found the missing requisite, and issue the challenge: "Okay, *now*
 369 find one!"

370 So, although it is impossible to prove that such an exception does not exist, we can assert that it
 371 will always be possible to extend the theorem to include any missing requisites that might be
 372 identified in the future, thus restoring the validity of the claim of sufficiency for all known systems
 373 that have been considered. □

374 3.3 ERT Universality

375 Anyone who has the impression that ERT primarily applies to artificial intelligence, robots, self-
 376 driving vehicles, and autonomous weapons systems is urged to consider how the theorem can be
 377 applied to human systems that make decisions that affect people or the environment, such as
 378 organizations, corporations, education systems, electoral systems, government institutions, CEOs, or
 379 yourself.

380 Justice Stevens [11] provides an excellent example of analyzing the ethical inadequacy of the
 381 "Citizens United" ruling. And his opinion that "The Court's ruling threatens to undermine the
 382 integrity of elected institutions across the Nation." implies that there is a pressing need to evaluate
 383 the ethical adequacy of the U.S. Supreme Court system.

384 Because the Ethical Regulator Theorem, i.e. the IsEthical and MakeEthical functions, can be
 385 applied to any system, the nine ERT dimensions define a domain-independent abstraction layer that
 386 can be used to map from any system/regulator to any other system/regulator. This creates a
 387 vocabulary, or isomorphism, that allows practitioners in one domain to communicate meaningfully
 388 with practitioners in seemingly unrelated domains, and share insights and solutions, for example,
 389 across artificial intelligence, corporate governance, education systems, and designing consumer
 390 products. Specialists in each domain can discuss their challenges and solutions to improving purpose,
 391 truth, variety, predictability, intelligence/strategy, influence, ethics, integrity, and transparency. For
 392 example, perhaps a cloud-based secure audit trail service that was developed for one specific domain
 393 can be used to solve the transparency and integrity requirements in a completely unrelated domain.

394 3.4 ERT Reflexivity and Algebra

395 If the Ethical Regulator Theorem is genuinely universal, it can be applied to absolutely any
 396 system. In particular, it must produce meaningful results for two special cases: When we apply ERT
 397 to itself, and when we apply ERT to second-order cybernetics (2oC).

398 First, let us define a convenient algebra that allows us to express important assertions in this
 399 domain. We need to distinguish between: I. The act of evaluating the ethical adequacy of a system,
 400 and II. The act of determining the set of transformations or interventions that are necessary to make
 401 a system ethically adequate:

- 402 I. A function, $\text{IsEthical}(S)$, returns the value True if system S is ethically adequate, it
 403 returns the value False if S is ethically inadequate, or it returns the value Undecidable if
 404 S is significantly inconsistent, contradictory, or opaque. The value Undecidable should
 405 be regarded as an error message rather than a type of system.
 406 II. A function, $\text{MakeEthical}(S)$, returns a set of transformations or interventions to make
 407 system S ethically adequate. If S is already ethically adequate, the function returns an
 408 empty set, $\{\}$.

409 Now we can use this ERT algebra to make some interesting and controversial claims in Table 2:

410 **Table 2:** Some ERT algebra assertions

No.	Claim	Interpretation / Justification
1	$\text{IsEthical}(\text{ERT}) = \text{True}$	The ERT system fulfils all nine requisites of ERT and is therefore ethically adequate. It can only be used to make systems more ethical.
2	$\text{MakeEthical}(\text{ERT}) = \{\}$	The ERT system is sufficient to be ethically adequate. Nothing else is required.
3	$\text{IsEthical}(2\text{oC}) = \text{False}$	Second-order cybernetics is ethically inadequate. Unlike ERT, it has no intrinsic ethics or integrity, so it can be used to make good or evil systems. It doesn't go beyond achieving effectiveness.
4	$\text{MakeEthical}(2\text{oC}) = \text{ERT}$	To become ethically adequate, 2oC needs the set of ERT concepts.
5	$\text{IsEthical}(2\text{oC} + \text{ERT}) = \text{True}$	Nothing in 2oC is incompatible with ERT.
6	$2\text{oC} + \text{ERT} = 3\text{oC}$	Logically, the system that is created by joining the 2oC and ERT systems would be named third-order cybernetics (3oC).
7	$\text{IsEthical}(\text{Capitalism}) = \text{False}$	Capitalism is ethically inadequate.
8	$\text{MakeEthical}(\text{Capitalism}) = \{\text{Ethics, Integrity, Transparency}\}$	Capitalism might be adequate in the six requisites for effectiveness, but it is obviously deficient in Ethics (laws, regulations, and rules), Integrity (compliance), and Transparency (audit trails). These must all be increased to make capitalism ethical. \square

411 3.5 The Law of Inevitable Ethical Inadequacy

412 We can build on the proof of necessity to derive this new law:

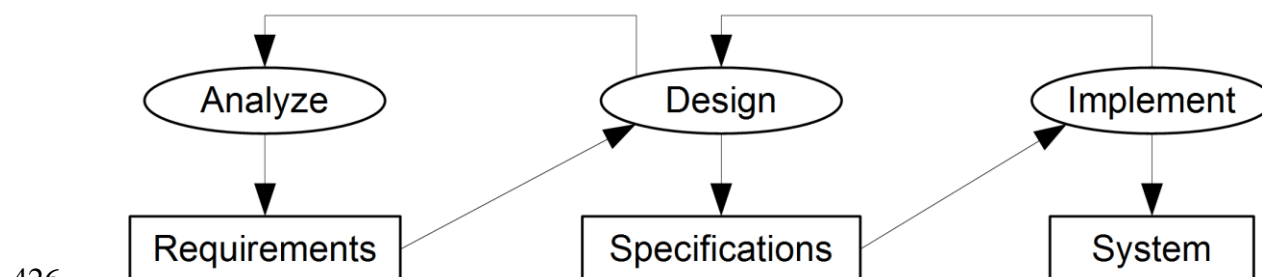
413 **“If you don’t specify that you require a secure ethical system,**
 414 **what you get is an insecure unethical system.”**

415 The reason is because when ethical adequacy is not specified as a requirement for a system
 416 design, the resulting design phase will tend to optimize for effectiveness and maximally avoid the
 417 extra costs that would be incurred by implementing the ethics, integrity, and transparency
 418 dimensions, which are optional for a system that only needs to be effective, thus guaranteeing that
 419 the resulting system is ethically inadequate and vulnerable to manipulation; by design. \square

420 Note that for some systems, the term “ethical” might include aspects such as hygienic, safe, fair,
421 honest, law-abiding, or environmentally friendly.

422 4 Ethical Design Process

423 Figure 2 shows the generic elements of a typical design process, in which an analysis phase
424 produces a requirements artifact, which is the input to the design phase that produces a specification
425 artifact, which is used as the input to the implementation phase, which realizes the system.

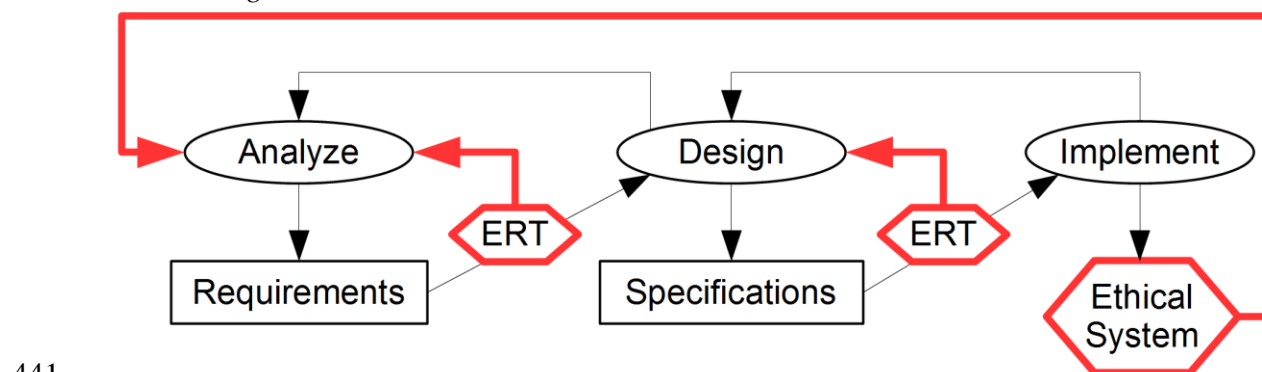


426
427 **Figure 2:** Ethically Inadequate Design Process

428 If a problem is found in the requirements during the design phase, feedback can trigger another
429 iteration of the analysis phase. And if a problem in the specifications is found during implementation,
430 feedback can trigger the design team to update the specifications or pass feedback to the analysis
431 team to update the requirements.

432 Such a design process can be effective at producing systems that are effective, however, because
433 the design process is ethically inadequate, it is inevitably only capable of producing systems that are
434 also ethically inadequate; and we cannot be sure that the resulting systems are not actually ethically
435 evil; by accident, or by design.

436 Fortunately, we can transform any effective but ethically inadequate design process to make it
437 ethically adequate by simply adding ethical adequacy acceptance testing of the requirements and
438 specifications. How we can retrofit the Ethical Regulator Theorem (ERT) to any effective design
439 process that produces requirements and specifications before the implementation phase starts is
440 shown in Figure 3.



441
442 **Figure 3:** Ethically Adequate Design Process

443 Like any other quality assurance testing, the ERT IsEthical and MakeEthical evaluations should
444 be performed by a team that was not involved in the production of the artifact being tested. If an
445 artifact is found to be ethically inadequate, it is rejected and recommendations for fixing the problems
446 are provided as feedback to trigger another iteration of that phase. If an artifact is found to be ethically
447 adequate, the artifact is accepted and passed onto the next phase. Because the two ERT testing steps
448 ensure that the requirements and specifications are ethically adequate, if the implementation process
449 performs an effective and lossless implementation of the specifications, the resulting system will also
450 be ethically adequate. □

451 This means that instances of the resulting system that are deployed in the real-world will include
452 a real-time integrity monitoring mechanism that detects and report any significant problems as
453 feedback to the analysis team, which must decide whether the problem necessitates instructing the
454 system to activate an ethical fail-safe mode, a remote update, a reimplementation, redesigning the
455 specifications, and/or updating the requirements. Only if the system fails to enter its ethical fail-safe
456 mode might it be necessary for it to be deactivated using a kill switch or for it to be “retired” by a
457 blade runner.

458 This concludes the description of the theorem and how to use it.

459 5 Discussion

460 The Ethical Regulator Theorem has many far-reaching implications.

461 5.1 Legislative Implications

462 By creating a well-defined interface for coding ethics, it becomes easier to apportion liability for
463 failures. For example, if a self-driving car crosses the border into India, fails to switch to the Indian
464 government certified ethics module for traffic-rules, and in an emergency, decides to hit a cow to
465 avoid hitting a dog, then the car manufacturer might be held liable for killing a sacred animal. But if
466 the audit trail proves that the correct ethics module was activated, but the “don’t hit cows” rule had
467 an incorrectly low priority in the ethics schema, then the car manufacturer would not be liable.

468 It is foreseeable that one-day the laws and regulations of most countries will be published in a
469 standardized computer-readable XML format, such as LKIF (Legal Knowledge Interchange Format),
470 and cryptographically-signed by an official issuing authority. However, the existing governmental
471 and regulatory organizations are inadequate to complete such a task in the necessary time frame.
472 Perhaps, a non-profit organization without any conflicts of interests could define appropriate
473 standards and start an open source ethics coding project for the laws, regulations, and rules that are
474 most urgently required by the ethically adequate systems that we try to construct.

475 By standardizing ethics modules, systems from different manufacturers will all use identical
476 ethics modules that are issued by national or international ethics authorities. The concept of central
477 ethics authorities might sound like part of a dystopic dictatorship but acting ethically is mostly just a
478 matter of obeying laws, regulations, and rules, which are a normal and necessary part of every stable
479 society. These ethics authorities could be independent of the legislative branch of government if the
480 government lacks the necessary resources or commitment to unambiguous digital lawmaking.

481 Like Microsoft Windows operating system updates, when new laws, regulations, rules, or bug
482 fixes to a previous ethics module are released, the new ethics module can be made available securely
483 to all affected autonomous systems; crucially, including systems whose manufacturer has gone out
484 of business or doesn't care about fixing end-user safety issues.

485 By comparison, Google’s Android operating system is a classic example of the Law of Inevitable
486 Ethical Inadequacy. Because Android was designed only to be effective, not ethical, Google delegated
487 the responsibility for issuing Android updates to the device manufacturers. The inevitable and
488 predictable consequence of that design decision is that most Android devices (87%) are insecure [12].
489 This exposes over **one billion** Android users to being hacked and their identity or credit card details
490 stolen by criminals. The resulting chaos and the expensive suffering of the victims is not an innocent
491 mistake, it is the direct result of Google deliberately externalizing costs onto others and prioritizing
492 its profits over ethical consumer safety. They could have designed it differently. And if we can't trust
493 Google, who can we trust?

494 We certainly don’t want robots, self-driving vehicles, and autonomous weapons systems relying
495 on an update mechanism that stops working when the manufacturer goes out of business or decides
496 to optimize its profits at the expense of security and safety updates.

497 Such unethical corporate behavior must be legislated out of existence, otherwise it will keep
498 recurring in different and damaging ways; causing unnecessary externalized costs and social chaos.
499 For example, ethically inadequate Internet-of-Things devices that send unencrypted data over the
500 internet, are vulnerable to being hacked, and will never receive security patches. Importing or selling

501 such unethical devices that threaten our privacy and the security of our digital infrastructure should
502 be as illegal as selling exploding cars.

503 5.2 Classification Framework

504 Now let's consider where the Ethical Regulator Theorem fits into the existing cybernetics
505 framework. One might assume that the theorem belongs in second-order cybernetics, however, in a
506 1990 conference plenary presentation [7], Heinz von Foerster (who made the distinction between
507 first- and second-order cybernetics in 1974) implied that combining ethics and second-order
508 cybernetics is not something that he would have suggested:

509 "I am impressed by the ingenuity of the organizers who suggested to me the title of my
510 presentation. They wanted me to address myself to 'Ethics and Second-Order Cybernetics'. To be
511 honest, I would have never dared to propose such an outrageous title, but I must say that I am
512 delighted that this title was chosen for me."

513 Table 3 lists some of the cybernetic community's definitions of first- and second-order
514 cybernetics, as summarized by Stuart Umpleby [13].

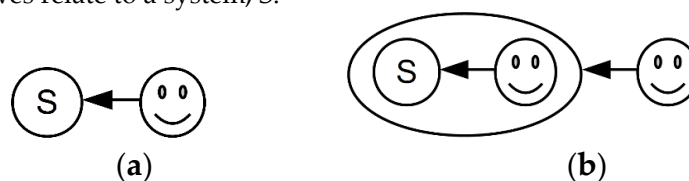
515 **Table 3.** Definitions of first- and second-order cybernetics

Author	First-Order Cybernetics	Second-Order Cybernetics
von Foerster	The cybernetics of observed systems	The cybernetics of observing systems
Pask	The purpose of a model	The purpose of the modeler
Valera	Controlled systems	Autonomous system
Umpleby	Interaction among the variables in a system	Interaction between observer and observed
Umpleby	Theories of social systems	Theories of the interaction between ideas and society

516 Although every one of these definitions captures an important distinction, when compared to
517 how the qualifiers "first-order" and "second-order" are used by other scientific communities, the
518 cybernetic community's use of them appears to be rather subjective, lacks the consensus that is
519 required by the scientific principle, and is of little utility, as required by Kuhn [14].

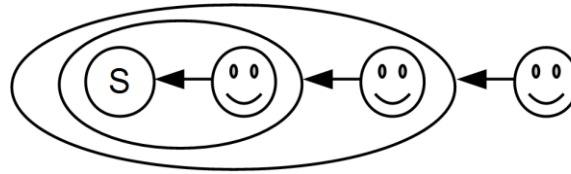
520 This incoherence in defining cybernetics as first-order and second-order not only prevents it
521 from being useful to classify different types of systems and dissipates intellectual energy, but it also
522 prevents the classification from being extended to higher orders, which can be viewed as either a self-
523 limiting dead-end, or paradigmatic autoapoptosis (self-programmed death), which is not entirely
524 unlike the situation of 39 members of the Heaven's Gate millennial death-cult, who believed that by
525 committing suicide, they would be rescued by an alien spacecraft and "graduate to the Next Level".

526 To illustrate the problem of classifying cybernetics into observer-centric "orders", let's start by
527 considering first- and second-order cybernetics, as defined by von Foerster. Figure 4 illustrates how
528 the observers' perspectives relate to a system, S.



531 **Figure 4:** (a) First-order cybernetics (b) Second-order cybernetics

532 How can we use this paradigm to predict the future of cybernetics? Logically, third-order
533 cybernetics would add a third observer's perspective, as shown in figure 5.



534

535

Figure 5: Third-order cybernetics

536

537

538

539

540

541

However, from the perspective of the third observer, this looks more like psychology than cybernetics. In fact, this structure is isomorphic to a typical management team evaluation exercise, where the details of the task that is given to the team to work on is virtually irrelevant to the outermost observer. It can be any goal-oriented activity, such as building the highest stable tower possible from a limited set of Lego bricks, solving an impossible puzzle in a limited amount of time, or studying a first-order cybernetic system.

542

5.3 New Classification Framework

543

544

545

546

547

It could be of more utility to define “levels” of cybernetic systems that include categories of future systems that are already anticipated and associate each level with established concepts. To that end, Table 4 defines a six-level framework for classifying cybernetic and superintelligent systems that makes use of the ERT IsEthical function to distinguish between two important subclasses of superintelligent systems.

548

Table 4. Six-level framework for classifying cybernetic and superintelligent systems

Level	The cybernetics of	Also known as	The cybernetician
1	Simple systems	First-order cybernetics	Observes the system
2	Complex systems	Second-order cybernetics	Participates in the system
3	Ethical systems	Third-order cybernetics or Cybernetics	Designs the system
4	Superintelligent systems	Technological singularity	Stares incredulously, as the system redesigns itself
5	Super-Ethical systems (Superintelligent and ethically adequate)	Technological utopia or Cyberanthropic utopia	Is protected by the system
6	Super-Unethical systems (Superintelligent and ethically inadequate)	Technological dystopia or Cybermisanthropic dystopia	Is manipulated to obey the system

549

550

551

552

Today, we are in the transition from building complex cybernetic level two systems (CL2) to building ethical systems and superintelligent systems of cybernetic levels three and four (CL3 and CL4), and the future of our species and our fragile ecosystem is in our hands, but first, let's clarify each level and explore where this new framework leads us.

553

5.3.1 Cybernetic Level 1: Simple Systems

554

555

This is the domain of first-order cybernetics: Studying and designing simple systems that are effective.

556 5.3.2 Cybernetic Level 2: Complex Systems

557 This is the domain of second-order cybernetics: Studying and designing complex systems that
558 are effective. There is still much important work to be done at this level.

559 5.3.3 Cybernetic Level 3: Ethical Systems

560 In 1986, decades ahead of his time, it was the wonderful and inspiring Ranulph Glanville who
561 defined “the cybernetics of ethics and the ethics of cybernetics” as “cybernethics” [15].

562 The Ethical Regulator Theorem belongs at this level, which is concerned with designing man-
563 made systems that are ethically adequate. Such systems must satisfy all nine requisites of the Ethical
564 Regulator Theorem and the regulating agents can be humans, machines, cyberanthropic hybrids,
565 organizations, corporations, or government institutions. Ethically adequate autonomous machines
566 must obey certified ethics modules.

567 In retrospect, now that we’re not trying to extrapolate from just two points in concept-space, if
568 level three cybernetic systems are ethical, it’s apparent that the third observer in the third-order
569 cybernetics system of Figure 5 is not necessarily a psychologist or a lost cybernetician, but could be
570 the second observer’s conscience; her super-ego, or higher-self; that constantly self-observing sense
571 that we all have that knows the difference between right and wrong; between good and evil. This
572 self-monitoring mechanism is known as integrity, and is something that today’s ethically indifferent
573 scientists, politicians, CEOs, managers, corporations, lawyers, bankers, and billionaires are woefully
574 lacking. In non-psychopaths, it is integrity that triggers feelings of bad conscience, regret, or guilt if
575 it is ignored.

576 5.3.4 Cybernetic Level 4: Superintelligent Systems

577 The technological singularity is a hypothetical moment when a self-improvement process in a
578 machine causes runaway improvements in intelligence that results in superintelligence that is far
579 greater than any human mind. For this to happen, the system must be sufficiently self-aware of its
580 own software and/or hardware.

581 5.3.5 Superintelligence Tests

582 These levels of self-awareness give rise to three levels of superintelligence tests. The ability to
583 reprogram better software for itself, the ability to redesign better hardware for itself, and the ability
584 to do both.

585 Together with the Turing Test [16], these tests mark milestones in the evolution of AI systems
586 towards superintelligence and should cause us alarm if progress towards them is made without
587 significant progress creating ethical systems first. Of these tests, the Turing Test is the easiest to
588 achieve because it is essentially a parlor game that only requires that a computer can imitate a (not
589 necessarily very intelligent) human sufficiently well to convince humans most of the time that it is a
590 human being and does not require self-awareness or runaway improvements in intelligence.

591 5.3.6 Prophecies of Possible Futures

592 In 1951, Ross Ashby started considering how to plan an advanced society as a “super brain” [17].
593 A year later, he described how super-clever machines could create a cyberanthropic utopia: “It may
594 be found that we shall solve our social problems by directing machines that can deliver an intelligence
595 that is not our own.” [18]

596 Two pages later, he described a cybermisanthropic dystopia where a “Million I.Q. Engine”
597 sounds like Facebook and Google, but on steroids: “What people could resist propaganda and
598 blarney directed by an I.Q. of 1,000,000? It would get to know their secret wishes, their unconscious
599 drives; it would use symbolic messages that they didn’t understand consciously; it would play on
600 their enthusiasms and hopes. They would be as children to it. (This sounds very much like Goebbels
601 controlling the Germans).”

602 On the appearance of such a machine, he described a paradox of perception of higher
603 intelligence: "It seems, therefore, that a super-clever machine will not look clever. It will look either
604 deceptively simple or, more likely, merely random." [19]. On the same subject, Arthur C. Clarke's
605 Third Law states: "Any sufficiently advanced technology is indistinguishable from magic." [20]. If
606 you think that Clarke's "magic" and Ashby's "deceptively simple or merely random" are
607 incompatible; take a moment to reflect on the magical simplicity and "randomness" of a Las Vegas
608 magic show or Google's search results' pages.

609 Just as there are two diametrically opposite archetypes for genius; namely the benevolent good
610 genius and the nasty evil genius, it is important not to conflate systems that are ethical with ones that
611 are not ethical, by making them share the same name or category, such as "superintelligent",
612 "Christian", or "super-rich". To do so would focus attention on the least important dimension and
613 ignore the most important dimension: Good and Evil.

614 5.3.7 Cybernetic Level 5: Super-Ethical Systems

615 The term "super-ethical" is proposed to refer to superintelligent systems that are ethically
616 adequate. Of course, by the time that super-ethical systems exist, a friendlier name will have emerged
617 and the term "super-ethical" will seem quaintly archaic.

618 5.3.8 Cybernetic Level 6: Super-Unethical Systems

619 The term "super-unethical" is proposed to refer to superintelligent systems that are ethically
620 inadequate. This term should always carry a certain stigma, like "weapons of mass destruction". No
621 one who is working to create artificially intelligent systems should be allowed to escape admitting
622 whether the systems are ethically inadequate.

623 Just as human genetic experimentation is strictly ethically regulated, we need legislation,
624 regulation, standards, and certification to ensure that autonomous AI systems that make decisions
625 that can have ethical consequences are subjected to the same kind of obsessively rigorous safety-
626 oriented design, construction, and operating procedures as commercial aircraft, nuclear power
627 stations, and vehicles that carry humans into space.

628 One could start arguing that intelligence is ethically neutral, and it is, but that family of
629 arguments are fallacies because a hyper-genius "Million I.Q. Engine" without ethics is not ethically
630 neutral. Even if it had ethical goals, it might break laws to achieve them. The possibility of creating a
631 superintelligent machine that is ethically inadequate should be treated like a bomb that could destroy
632 our planet. Even just planning to construct such a device is effectively conspiring to commit a crime
633 against humanity.

634 As a thought experiment, let's imagine a hypothetical super-unethical version of Google, named
635 the Googlevil Corporation. The CEO is Dr. Evil, and both the CEO and the corporate AI are without
636 ethics, avoid transparency, and will do anything to maximize their profits and power. The
637 corporation's secret mission statement is "Collect and organize the world's personal information and
638 make it accessible and useful for maximizing our profits, influence, and ability to avoid paying taxes."
639 and its secret corporate mantra is "Sincerely say 'Believe me, we don't do evil', do it anyway, then
640 look people in the eye and give them a Zuckerberg-smile!"

641 Anytime that the super-unethical Googlevil artificial intelligence or the psychopathic
642 demagogue Dr. Evil wants to blackmail the CEOs of other corporations, politicians that can't be
643 bought, jury members, or Supreme Court justices around the world to make "random" decisions that
644 incrementally further their secret mission, would they have to do anything more than query the
645 Googlevil user-profile database? In theory, they would only need to be able to blackmail a majority
646 of members of lower- and upper-houses (how hard can that be?) to be able to get *any* legislation that
647 they want in *any* country. Or just a few Supreme Court justices to steer a nation into a fascist dystopia.

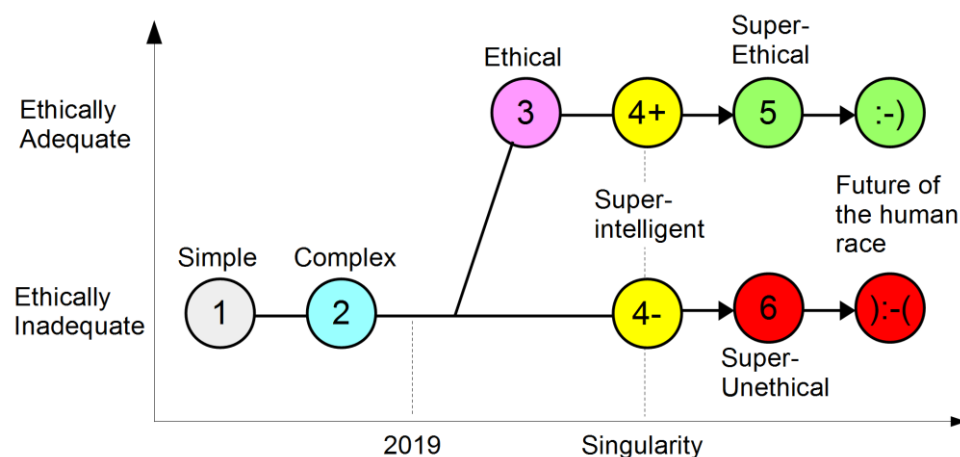
648 By the time that super-unethical AI systems exist, they could be indistinguishable from the
649 corporations that they belong to. They could be immoral, immortal, enjoy legal personhood, pay no
650 taxes, and make unlimited donations (also known as bribes) to all Googlevil-friendly political parties
651 in all techno-democratic dystopias on the planet.

652 5.4 Future Time-Line Bifurcation Race Condition

653 At this point in time, there is an existentially critical fork in our future time-line. Depending on
 654 whether the systems that achieve the singularity are ethically adequate or not, the runaway increase
 655 in intelligence and inevitable ethical polarization pressures will result in one of two outcomes:

- 656 • Good super-ethical AIs protect humanity.
- 657 • Evil super-unethical AIs dominate humanity.

658 Figure 6 illustrates how plotting the ethical dimension orthogonally to the intelligence
 659 dimension clarifies the non-linear dependencies between different cybernetic levels, and clearly
 660 shows that the ethically inadequate superintelligent systems of cybernetic level four-minus (CL4-)
 661 have no dependency on us first succeeding creating the ethical systems of cybernetic level three
 662 (CL3).



663

664 **Figure 6:** Two mutually exclusive possible futures

665 If we continue on the current path from complex systems (CL2) to ethically inadequate
 666 superintelligent systems (CL4-), we will quickly arrive in a dystopia that is dominated by super-
 667 unethical systems (CL6), and the potential cyberanthropic utopia of being ruled by benevolent super-
 668 ethical systems (CL5) will become permanently unreachable.

669 So, there is a race condition that will determine which of these two mutually exclusive possible
 670 futures will be the fate of our species; will our technological progress reach CL4+ or CL4- first? And
 671 will legislators regulate such developments ethically and adequately, or will they sell us out for bribes
 672 from Dr. Evil's special interest lobby groups that will "campaign" for "self-regulation" — and we all
 673 know what that *really* means!

674 It cannot be overemphasized that CL4± is the point-of-no-return where humans could lose
 675 control over machines that become our intellectual superiors. And this is the window of opportunity
 676 to ensure that superintelligent machines are programmed with ethics and purposes that serve the
 677 greater good of humanity and our fragile ecosystem. Put simply: We must create ethical systems
 678 *before* we create superintelligent systems!

679 In this context, it is clear that the ultimate grand challenge for cybernetics and third-order
 680 cyberneticians is to find ways to build ethical and super-ethical systems, avoid a cybermisanthropic
 681 dystopia, and help humanity create a super-ethical society.

682

683 5.5 Super-Ethical Society

684 Imagine how different the world would be:

- 685 • If we were happy to be ruled by benevolent super-ethical artificial intelligences that eliminated
 686 poverty, environmental destruction, corruption, and injustice.

- 687 • If the United Nations could deploy heavily armed super-ethical peace-keeping robot armies into
688 conflict zones to protect civilians and enforce ceasefires.
- 689 • If our towns and cities are policed by super-ethical robots that protect all citizens equally, 24x7,
690 and never shoot our friends or family because of their race, religion, social class, lifestyle, or
691 peaceful protesting.
- 692 • If super-ethical child-care robots accompany our children wherever they go, protecting them
693 from danger, including physical, emotional, and sexual abuse.
- 694 • If ethically adequate corporations produce ethical products, provide ethical services, and pay
695 ethical levels of unavoidable corporation tax.

696 Such a super-ethical society is possible; but only if we deliberately make it our goal, rise above
697 polarizing politics, and act together in accordance with the undeniable truth that ethics are a higher
698 power for good that transcends science, politics, nations, and religions.

699 5.6 Cyberanthropic Utopia

700 Because so many ludicrous utopias have been proposed that are just naïve science fiction
701 fantasies, it is unsurprising that utopias have accumulated a bad reputation and are not taken very
702 seriously. But it is shockingly common for apparently rational people to exhibit symptoms of classical
703 conditioned-reflex (Pavlovian) negative responses to the stimulus word “utopia”; triggering
704 emotional distress and bypassing their rational reasoning, as if any serious use of the word “utopia”
705 has become a reputation-threatening scientific taboo.

706 However, now that artificial intelligence is making such impressive progress and showing no
707 signs of slowing down nor having an upper-limit, Ross Ashby's 67-year-old prediction looks
708 increasingly realistic: We might be able to “...solve our social problems by directing machines that
709 can deliver an intelligence that is not our own.”

710 Ashby's prediction hints at a possible definition of a minimum viable utopia:

711 **A world where our social problems have been solved.**

712 Utopia need not mean a “perfect” society, or that we all have flying cars, robot servants, and
713 never have to go to work. Just fulfilling human needs and eliminating poverty would create a truly
714 magnificent utopia. And as we start making progress achieving it, many other human problems, such
715 as malnutrition, parasitic diseases, homelessness, hopelessness- and poverty-driven prostitution and
716 crime will fade away and the world will become a very different and happier place to live in.

717 Let no one say it cannot be done. Ethically adequate societies have existed in the past, where
718 resources were shared and the environment respected. What is new is that we can now do it
719 synthetically, consciously, deliberately.

720 5.7 Third-Order Cybernetics

721 Since Heinz von Foerster made the distinction between first- and second-order cybernetics in
722 1974, many people have attempted to find a plausible definition for third-order cybernetics, but until
723 now, no definition has gained acceptance.

724 This paper proposes that Third-Order Cybernetics should be defined as “the cybernetics of
725 ethical systems”, and that “the cybernetics of ethics”, “Cybernetics 3.0”, and “3oC” are all acceptable
726 synonyms for it. Some of the supporting arguments for this proposal have already been mentioned,
727 however a consolidated set of arguments are listed below:

- 728 1. Until now, second-order cybernetics (2oC) discussions about the need to create ethical systems,
729 including the need for cybernetics itself to embody ethics, did not produce any satisfactory
730 solution. Here, “satisfactory solution” is understood to mean something like the ethically
731 adequate design process of Figure 3 that can be used systematically, for example by engineers, to
732 create real systems that are ethically adequate. Recognizing this need but failing to fulfil the need
733 could be referred to in the context of second-order cybernetics as “The Ethics Problem”.

- 734 2. The fact that Heinz von Foerster described “Ethics and Second-Order Cybernetics” as an
735 “outrageous title” for a presentation, implies that ethics do not belong in the 2oC that he
736 envisioned.
- 737 3. Whereas 2oC can be used for good or evil, the Ethical Regulator Theorem can only be used for
738 good. This distinction is significant, and it also implies that ERT does not belong in 2oC. It is
739 ERT's existence that creates the need to define 3oC.
- 740 4. If we extrapolate from first-order cybernetics having one observer and 2oC having a second
741 observer, we would logically expect 3oC to introduce a third observer. This hypothesized third
742 observer maps exactly onto the ERT requirement that ethical systems must have real-time
743 integrity mechanisms that monitor and enforce compliance of the system with respect to an
744 appropriate ethical schema.
- 745 5. An alternative, observer-free justification can be derived from the Good Regulator Theorem: A
746 first-order cybernetic regulator requires a model of the system being regulated and a second-
747 order cybernetic regulator can only achieve reflexivity by having a model of itself. Then to behave
748 ethically, a third-order cybernetic regulator needs a third model, a model of acceptable (ethical)
749 behavior, which is encoded in the ethics schema. It is then a consequence of the fact that every
750 model requires observations as inputs, that brings into existence the need for observing part(s)
751 to exist in the system. The need for these observations is independent of whether a cybernetician
752 is watching or not. Whereas the observer-based argument of point 4 identifies no new
753 requirements on the regulator, this model-based argument not only makes explicit that an ethical
754 regulator requires three models, but it also requires observations as a direct consequence of using
755 the models.
- 756 6. Together, ERT and the six-level framework (6LF) for classifying cybernetic and superintelligent
757 systems (see Table 4) create a new paradigm that has greater explanatory and predictive power
758 than 2oC. For example, producing the ERT effectiveness function, the Law of Inevitable Ethical
759 Inadequacy, identifying that the ethical systems of cybernetic level 3 (CL3) are the missing type
760 of cybernetic system that is necessary to integrate cybernetic systems and superintelligent
761 systems into a common framework, explaining the impending bifurcation into either a
762 cyberanthropic utopia or a cybermisanthropic dystopia (see Figure 6), and identifying
763 deficiencies in capitalism. In addition, because 6LF integrates three classes of system that do not
764 yet exist, it can help us navigate a rational path into the future, for example, by predicting the
765 existence of a race condition and thus identifying a possible solution to the dangers that are posed
766 by superintelligent machines. Such insights cannot be obtained using 2oC.
- 767 7. Ethical systems constitute a new type of system, and ERT + 6LF defines a new branch of
768 cybernetics that goes beyond effectiveness, does not belong in 2oC, and solves “The Ethics
769 Problem”. Logically, the system that is created by joining the 2oC and ERT systems would be
770 named third-order cybernetics.
- 771 8. Although Ranulph Glanville defined “the cybernetics of ethics and the ethics of cybernetics” as
772 “cybernethics”, this term is invented jargon that carries no meaning for people who are not
773 familiar with its definition. By contrast, the term “third-order cybernetics” carries enough
774 meaning for people who are familiar with the term “second-order cybernetics” to at least trigger
775 interest and curiosity. Therefore, using the term “third-order cybernetics” instead of
776 “cybernethics” has significant advantages, and enhances 6LF by increasing symmetry in Table 4.
- 777 9. Whereas it is impossible to define objectively which theories and practices belong in 2oC, thus
778 making it an intimidating subject for outsiders to even contemplate mastering, ERT is defined
779 and proved in eight pages and doesn't require any knowledge of 2oC. This means that the
780 theorem, and how to apply it to systems of any type, can easily be taught to non-cyberneticians.
781 It is therefore imperative that ERT + 6LF can make a fresh start as a new cybernetic speciality
782 without being entangled with 45 years of fuzzy 2oC baggage. However, 3oC is not limited to ERT
783 and 6LF, and will surely develop rapidly before it matures.
- 784 10. Unlike 2oC, 3oC has a fundamentally ethical purpose (making systems ethical) that together with
785 the proposed grand challenge and the vision of a super-ethical society, create a unique

786 opportunity for the cybernetics and systems sciences community to take a leading role in
787 implementing a long-overdue and much needed systemic ethical revolution.

788 11. If someone makes the statement “I’m a second-order cybernetician”, it reveals nothing about their
789 ethics. But from now on, anyone who is brave enough to declare “I’m a third-order cybernetician”
790 is making a bold assertion that they are part of the only scientific movement that is dedicated to
791 making the world a better place.

792 12. Recognizing and declaring that 3oC is a new paradigm that is revitalizing and reunifying the 2oC
793 community will help draw attention to the fact that something important and exciting is
794 happening; a powerful attractor for energy and commitment. It is a genuine **ethical imperative**
795 that we develop this new field for engineering, management, and the other sciences to use. And
796 it is a task that cannot be performed by any of the more narrowly defined branches of science.

797 If these arguments are accepted, it is suggested that going forwards, the term second-order
798 cybernetics should only be used to refer to the second-order cybernetics of effectiveness without the
799 ethical, integrity, and transparency aspects, which belong to the 3oC layer that can be used to
800 transform any effective but ethically inadequate system, such as a design process, second-order
801 cybernetics, or capitalism, into an ethical system.

802 Many people thought that cybernetics had faded away after peaking in the 1960s or early 1970s,
803 but that peak was just a local maximum: Cybernetics is rebooting as Cybernetics 3.0, and this time
804 it’s going to be harder to ignore, because we’ll be applying requisite Purpose, Truth, Variety,
805 Predictability, Intelligence, Influence, Ethics, Integrity, Transparency — and Love; because a sincere
806 desire to make the world a better place emerges *only* in people who love humanity and the biosphere
807 unconditionally.

808 By contrast, deep down, non-empaths only care about themselves or members of their own
809 nation, race, religion, family, or gang, and consequentially embody a conflict of interests that compels
810 them to act against the greater good.

811 *5.8 Our Future Epilog or Eulogy*

812 We are approaching a decisive fork in the road in the evolution of intelligent machines, immortal
813 corporations, political systems, and human society, and it is imperative that we learn to make these
814 systems rigorously ethical before artificially intelligent machines reach the technological singularity,
815 start to evolve exponentially, exceed human intelligence, and are used by ethically inadequate
816 corporations to dominate humanity politically and economically.

817 We are the *only* generation that has the chance to steer the fate of future generations of humanity
818 towards being collectively ruled, potentially for eternity, by benevolent super-ethical systems that
819 create a stable cyberanthropic utopia for us, effectively and ethically minimizing human suffering
820 and environmental problems, rather than allowing hubris and super-unethical systems to either
821 enslave most of us in a cybermisanthropic dystopia or cause the extinction of our species to become
822 a footnote in Gaia’s geological record.

823 *5.9 The Path Forwards*

824 To start steering the future of humanity and our wonderful planet towards becoming a stable
825 cyberanthropic super-ethical society, this paper proposes establishing an independent, non-profit
826 institute with ambitious goals that lie in the areas of research, development, standards, certification,
827 legislation, and democracy.

828 *5.9.1 Research and Development*

829 The institute will promote theoretical and practical progress:

- 830 • Coordinate and fund research into creating ethical systems and making existing systems ethical.
- 831 • Develop a taxonomy of open-source ethics modules for different types of laws, regulations, and
- 832 rules that can be used by anyone, free of charge.

833 5.9.2 Standards and Certification

834 The institute will create an ethical certification infrastructure:

- 835 • Establish standards for certifying the ethical adequacy of systems.
 836 • Establish a curriculum for training accredited ethical consultants.
 837 • Coordinate and regulate contracts for ethical audits and certifications.

838 5.9.3 Legislation and Democracy

839 The institute will lobby governments to implement ethically adequate legislation and will
 840 evaluate the adequacy of any proposed legislation. In particular, promoting the following:

- 841 • Regulate autonomous machines to require that their design and implementation is ethically
 842 adequate, and that they support compulsory ethics modules.
 843 • Make it illegal to import or sell products that have not been certified as being ethically adequate,
 844 unless they are explicitly excluded from requiring certification.
 845 • Require that all new systems and processes are designed to be ethically adequate.
 846 • Extend political representation to every member of society by giving parents proxy votes to cast
 847 on-behalf of their children who are too young to vote, but not too young for morally bankrupt
 848 politicians to load up with unsustainable debt liabilities, while underfunding the public
 849 education system and allowing unethical corporations to maximize short-term profits by
 850 devastating the environment for all future generations of humanity. What we currently call
 851 "[universal suffrage](#)" [21] is a perversion of the true meaning of the word "universal".

852 5.10 Example: Applying ERT to yourself

853 As members of a human society, we are all cybernetic regulators; of ourselves and of each other.
 854 As a thought experiment, to become a more effective and ethical force for good, you could identify
 855 ways to improve each ethical requisite as it applies to yourself. Table 5 illustrates how you can use
 856 ERT to make yourself a better ethical regulator. Or, stated in ERT algebra:
 857 $\text{MakeEthical}(\text{yourself}) = \text{Table 5}$

858 **Table 5.** Ways to become a better ethical regulator

Requisite	Example set of self-improvement interventions
Purpose	To clarify your purpose in life and help you to recognize your strongest motivating thoughts, write down your most important life goals: <ol style="list-style-type: none"> 1. 2. 3. 4. 5.
Truth	To become a good judge (effective and ethical) of who tells the truth and who distorts it, seek alternative information sources that are genuinely independent of your primary sources. Investigate any inconsistencies that you notice, modify the reputation of liars, and resolve to always doubt them skeptically in future.
Variety	Brainstorm new actions, responses, and strategies that you have never previously considered, to make progress towards achieving your goals.

Predictability	<p>Improve your model of human behavior by studying the following Wikipedia articles until you are competent at recognizing the patterns in yourself and others:</p> <ul style="list-style-type: none"> • List of cognitive biases [22] • Defence mechanisms [23] • List of fallacies [24] • Demagogue [25]
Intelligence	Take a course or read a book on critical thinking or personal effectiveness.
Influence	Identify ways that you can increase your influence (on your family, friends, colleagues, clients, or society) to achieve your life goals and promote your ethical values.
Ethics	<p>Write down five undesirable, unethical, or disrespectful behaviors that, up until now, you have tolerated in other people, organizations, or corporations:</p> <ol style="list-style-type: none"> 1. 2. 3. 4. 5. <p>Next to them, write down five undesirable, unethical, or disrespectful behaviors that, up until now, you have tolerated in yourself. If you can't think of five things about yourself, read the Wikipedia article: Denial [26]. If that doesn't help, ask someone that you live with to suggest five things that you do that they'd prefer you not to do.</p>
Integrity	Seek to stop or prevent all the undesirable, unethical, or disrespectful behaviors that you listed under requisite ethics.
Transparency	Let other people know about the changes that you are making.

859

860 Finally, keep reviewing and refining your answers for Purpose and Ethics until they genuinely
861 reflect who you are and how you want your world to become.

862 5.11 Ethically Resonant Wisdom

863 If you distil different solutions that contain alcohol, you get pure alcohol that is free of
864 impurities. And if you distil different religions and philosophies that contain ethical wisdom, you get
865 pure ethical wisdom that is free of culturally-specific dogma. Such ethical wisdom is universal, and
866 resonates with all good people, regardless of their worldview, politics, nationality, or religion.

867 And because pure ethics are a higher power for good that transcends science, politics, nations,
868 and religions, it is probably the only force that can unify humanity to work together for our greater
869 good. For example, consider the following selected quotes:

870 Mahatma Gandhi (1869-1948):

- 871 1. The future depends on what you do today.
- 872 2. Be the change you wish to see in the world.
- 873 3. The difference between what we do and what we are capable of doing
874 would suffice to solve most of the world's problems.
- 875 4. If I have the belief that I can do it,
876 I shall surely acquire the capacity to do it even if I may not have it at the beginning.
- 877 5. First they ignore you, then they laugh at you, then they fight you, then you win.

- 878 6. Happiness is when what you think, what you say, and what you do are in harmony.
879 7. Non-cooperation with evil is as much a duty as is cooperation with good.
880 8. Poverty is the worst form of violence.
881 9. Capital as such is not evil; it is its wrong use that is evil.
882 10. There is sufficiency in the world for man's need,
883 but not for man's greed.
884 11. There are people in the world so hungry,
885 that God cannot appear to them except in the form of bread.
886 12. Those who say religion has nothing to do with politics do not know what religion is.
887 13. Where love is, there God is also.
888 14. God has no religion.
889 15. There is a higher court than the courts of justice and that is the court of conscience.
890 16. They may torture my body, break my bones, even kill me.
891 Then they will have my dead body, but not my obedience.
892 17. Victory attained by violence is tantamount to a defeat, for it is momentary.
893 18. What difference does it make to the dead, the orphans, and the homeless,
894 whether the mad destruction is wrought under the name of totalitarianism
895 or the holy name of liberty or democracy?
896 19. Your beliefs become your thoughts, your thoughts become your words,
897 your words become your actions, your actions become your habits,
898 your habits become your values, your values become your destiny.
- 899 His Holiness Pope Francis:
- 900 20. We must restore hope to young people, help the old, be open to the future, spread love.
901 Be poor among the poor. We need to include the excluded and preach peace.
902 21. Hatred is not to be carried in the name of God. War is not to be waged in the name of God!
903 22. Human rights are not only violated by terrorism, repression, or assassination,
904 but also by unfair economic structures that create huge inequalities.
905 23. The worship of the golden calf of old has found a new and heartless image in the cult of money
906 and the dictatorship of an economy which is faceless and lacking any truly human goal.
907 24. Men and women are sacrificed to the idols of profit and consumption: It is the "culture of waste".
908 If a computer breaks it is a tragedy, but poverty, the needs and dramas of so many people end
909 up being considered normal.
910 25. Women in the church are more important than bishops and priests.
911 26. All that is good, all that is true, all that is beautiful, God is the truth.
912 27. We all have the duty to do good.
913 28. Everyone has his own idea of good and evil and must choose to follow the good and fight evil as
914 he conceives them. That would be enough to make the world a better place.
- 915 His Holiness the Dalai Lama XIV:
- 916 29. All religious institutions, despite different philosophical views,
917 all have the same message — a message of love.
918 30. If you can, help others; if you cannot do that, at least do not harm them.
919 31. The whole purpose of religion is to facilitate love
920 and compassion, patience, tolerance, humility, and forgiveness.
921 32. Irrespective of whether we are believers or agnostics, whether we believe in God or karma,
922 moral ethics is a code which everyone is able to pursue.
923 33. The ultimate authority must always rest with the individual's own reason and critical analysis.
924 34. The true hero is one who conquers his own anger and hatred.
925 35. A good friend who points out mistakes and imperfections and rebukes evil
926 is to be respected as if he reveals the secret of some hidden treasure.
927 36. A lack of transparency results in distrust and a deep sense of insecurity.

- 928 37. In our struggle for freedom, truth is the only weapon we possess.
929 38. Where ignorance is our master, there is no possibility of real peace.
930 39. Through violence, you may “solve” one problem, but you sow the seeds for another.
931 40. Don’t ever mistake my silence for ignorance, my calmness for acceptance or my kindness for
932 weakness. Compassion and tolerance are not a sign of weakness, but a sign of strength.
933 41. A truly compassionate attitude toward others does not change
934 even if they behave negatively or hurt you.
935 42. I defeat my enemies by making them my friends.
936 43. When you practice gratefulness, there is a sense of respect toward others.
937 44. With realization of one’s own potential and self-confidence in one’s abilities,
938 one can build a better world.
939 45. If you think you are too small to make a difference, try sleeping with a mosquito.
940 46. As people alive today, we must consider future generations:
941 A clean environment is a human right like any other.
942 It is therefore part of our responsibility toward others to ensure that the world we pass on
943 is as healthy, if not healthier, than we found it.
944 47. The ultimate source of happiness is not money and power, but warm-heartedness.
945 48. The more you are motivated by love, the more fearless and free your action will be.
946 49. Love and compassion are necessities, not luxuries.
947 Without them humanity cannot survive.
948 50. Love is the absence of judgement.
949 51. Be kind when possible. It is always possible.

950 Dr. Martin Luther King Jr. (1929-1968):

- 951 52. We must discover the power of love, the power, the redemptive power of love.
952 And when we discover that we will be able to make of this old world a new world.
953 We will be able to make men better.
954 Love is the only way.
955 53. I say to you, “I love you. I would rather die than hate you.”
956 And I’m foolish enough to believe that through the power of this love,
957 somewhere, men of the most recalcitrant bent will be transformed.
958 54. Darkness cannot drive out darkness; only light can do that.
959 Hate cannot drive out hate, only love can do that.
960 55. Those who love peace must learn to organize as effectively as those who love war.
961 56. True peace is not merely the absence of tension.
962 It is the presence of justice.
963 57. Injustice anywhere is a threat to justice everywhere.
964 58. In a real sense, all life is inter-related.
965 All men are caught in an inescapable network of mutuality,
966 tied in a single garment of destiny.
967 Whatever affects one directly, affects all indirectly.
968 59. Every man must decide whether to walk in the light of creative altruism
969 or in the darkness of destructive selfishness.
970 60. The time is always right to do the right thing.
971 61. We must learn that passively to accept an unjust system
972 is to cooperate with that system, and thereby to become a participant in its evil.
973 62. You are not only responsible for what you say,
974 but also for what you do not say.
975 63. Our lives begin to end the day we become silent about things that matter.
976 64. Our scientific power has outrun our spiritual power.
977 We have guided missiles and misguided men.

- 978 65. A nation that continues year after year to spend more money
 979 on military defense than on programs of social uplift is approaching spiritual doom.
 980 66. We should never forget that everything Adolf Hitler did in Germany was "legal"
 981 and everything the Hungarian freedom fighters did in Hungary was "illegal".
 982 67. Nonviolence is directed against forces of evil
 983 rather than against persons who happen to be doing evil.
 984 It is evil that the nonviolent resister seeks to defeat, not the persons victimized by evil.
 985 68. Nonviolence means avoiding not only external physical violence but also internal violence of
 986 spirit. You not only refuse to shoot a man, but you refuse to hate him.

987 Nelson Mandela (1918-2013):

- 988 69. Freedom can never be taken for granted. Each generation must safeguard it and extend it.
 989 Your parents and elders sacrificed much so that you should have freedom without suffering what
 990 they did. Use this precious right to ensure that the darkness of the past never returns.
 991 70. Like slavery and apartheid, poverty is not natural.
 992 It is man-made and it can be overcome and eradicated by the actions of human beings.
 993 71. Overcoming poverty is not a gesture of charity.
 994 It is an act of justice.
 995 72. As long as poverty, injustice and gross inequality persist in our world,
 996 none of us can truly rest.
 997 73. Education is the most powerful weapon which you can use to change the world.
 998 74. It is in your hands to create a better world for all who live in it.
 999 75. May your choices reflect your hopes, not your fears.

1000 Albert Einstein (1879-1955):

- 1001 76. No problem can be solved from the same level of consciousness that created it.

1002 Margaret Mead (1901-1978):

- 1003 77. Never doubt that a small group of thoughtful, committed citizens can change the world;
 1004 indeed, it's the only thing that ever has.

1005 Bertolt Brecht (1898-1956):

- 1006 78. Change the world, she needs it.

1007 Percy Bysshe Shelly (1792-1822):

- 1008 79. Rise like lions after slumber
 1009 In unvanquishable number!
 1010 Shake your chains to earth like dew
 1011 Which in sleep had fallen on you:
 1012 Ye are many — they are few!

1013 Leonardo da Vinci (1452-1519):

- 1014 80. I have been impressed with the urgency of doing.
 1015 Knowing is not enough; we must apply.
 1016 Being willing is not enough; we must do.

1017 Despite the authors of these quotes being separated by space, time, and their affiliations, it's easy
 1018 to imagine that they all share the same human ethical belief system, and that they would have no
 1019 significant arguments with each other if they all came together in one room to plan an ethical
 1020 revolution to make the world a better place.

1021 5.12 The Law of Unethical Arguments

1022 It is a certainty that all good people (without exception) are supportive of redesigning unethical
1023 systems, organizations, corporations, products, taxes, laws, regulations, and processes to make them
1024 more ethical. And it makes sense that the only people who want such systems to remain unethical
1025 and vulnerable to tampering and abuse are the small minority of people who benefit (directly or
1026 indirectly) from those systems remaining unethical.

1027 The final law in this manifesto for a nonviolent global ethical revolution to create a stable
1028 cyberanthropic super-ethical society is defined in one sentence:

1029 "Because no ethical argument can exist against making a system ethical,
1030 anyone who argues against this goal,
1031 obstructs progress towards this goal,
1032 or abuses its sincere supporters,
1033 is either objectively unethical, corrupt, or evil."
1034

□

1035 6 Conclusions

1036 The Ethical Regulator Theorem creates a theoretical basis for applied ethics, enabling designers
1037 to systematically evaluate, improve, and design ethically adequate systems. Because it is a universal
1038 theorem that can be applied to any system, the possible areas of application are vast and potentially
1039 world-changing.

1040 The six-level framework for classifying cybernetic and superintelligent systems leads to a
1041 theory-based solution to the danger that superintelligent machines might cause a dystopia: *We must*
1042 *create ethical systems before we create superintelligent systems!*

1043 By creating a well-defined decision function (IsEthical) that identifies systems as being either
1044 ethically adequate or ethically inadequate, ERT provides a semantic precision that avoids the
1045 ambiguities and unstated assumptions that multiply exponentially when the word "ethical" is
1046 bandied around as if we all understand it to mean the same thing. But "ethical AI", "ethical product",
1047 and "ethical corporation" can mean very different things to different people. By contrast, ERT gives
1048 terms like "ethically adequate AI", "ethically adequate product", and "ethically adequate
1049 corporation" a much more precise meaning, and could even be made the subject of a formal
1050 certification process that qualifies recipients to use an Ethically Adequate branded logo and reduce
1051 their liability insurance premiums.

1052 ERT's universality means that the nine dimensions define an abstraction layer that can be
1053 mapped onto the regulators of any systems in any domain, thus enabling communication and
1054 learning to take place between experts in seemingly unrelated domains.

1055 Because of the flaw that was identified in Heinz von Foerster's Ethical Imperative, a new
1056 definition is proposed, which is intended to embody both the essence of the proposed grand
1057 challenge and a principle for good that is universal and worthy of the magniloquent name "Ethical
1058 Imperative":

1059 **"Always strive to make new and existing systems ethically adequate!"**

1060 The proposed grand challenge to implement a systemic ethical revolution is neither a new
1061 religion nor a political movement, it is a response to Johann Eder's call for a grand challenge in Vienna
1062 [27] and Irma Wilson and Pamela Buckle Henning's call to action for the systems sciences community
1063 in Berlin [28].

1064 This ethical revolution is the product of a compassionate heart and mind, employing the Ethical
1065 Regulator Theorem to generate maximally coherent ethical interventions in multiple complex
1066 systems, such as the computational, corporate, criminal, cybernetic, personal, political, product
1067 development, psychological, scientific, social, and spiritual [29] realms. And all such interventions
1068 resonate, not only with each other, but also with all good people who have ever existed — or ever
1069 will.

1070 This revolution is long overdue, and we are privileged to live in these exciting times, but
 1071 passively watching from the sidelines, or doing nothing, only helps the criminals, psychopaths,
 1072 demagogues, and ethically indifferent corporations to create and exploit the pathological chaos and
 1073 emergent problems that, until now, we have accepted as normal. It's time for all good people to make
 1074 a commitment to yourself to do everything that you can to research, design, educate, campaign, love,
 1075 heal, and fight for a better world.

1076 "To be bold enough to consciously and deliberately reach beyond ourselves,
 1077 to accept a grand challenge for the greater good,
 1078 would be an act of self-actualization."
 1079 — Stella Octangula [30]

1080 Just like we have legislation and non-negotiable expectations that passenger aircraft are
 1081 designed to include expensive redundant subsystems to avoid having single points of failure in
 1082 flight-safety-critical systems, and that all electrical products that we purchase conform to strict safety
 1083 standards, we must change our attitudes, to create a cultural shift that makes it totally unacceptable
 1084 and utterly unthinkable to knowingly design systems or sell products that are ethically inadequate.
 1085 Outrage at such behavior is appropriate.

1086 We must demand strict legislation and higher standards to force ethically indifferent
 1087 corporations to stop their races to the bottom and cost externalization strategies. In truth, only
 1088 certifiably ethical corporations can be trusted to produce ethically adequate products and services
 1089 that help to make the world a better place for the entire human race.

1090 Arguably, the root cause of all evil is a lack of ethics, and by systematically applying the Ethical
 1091 Regulator Theorem, we can reliably increase ethical behavior in many classes of systems;
 1092 progressively reducing unethical behavior, reducing unethical suffering, and setting a course for
 1093 humanity towards the tipping-point where we will experience a peaceful social phase-transition to a
 1094 stable cyberanthropic super-ethical society.

1095 Though this paper covers many topics, these are but means; the end has been throughout to
 1096 make clear what principles must be followed when one attempts to restore ethical function to a sick
 1097 organism that is, as a human society, of fearful complexity. It is my faith that the new understanding
 1098 may lead to super-ethical systems that can create a better world, for the need is great.

1099 **Conflicts of Interest:** The author declares no conflicts of interest

1100 References

- 1101 1. Bostrom, N. *Superintelligence: Paths, Dangers, Strategies*. Oxford University Press, 2014.
- 1102 2. Conant, R.C.; Ashby, W.R. Every good regulator of a system must be a model of that system, *INT J SYST*
 1103 *SCI* **1970** 1(2), 89-97.
- 1104 3. Ashby, W.R. (1956). *An Introduction to Cybernetics*, Chapman and Hall, London, 1956. Available online:
 1105 rossashby.info/Ashby-Introduction-to-Cybernetics.pdf (accessed on 18 July 2019).
- 1106 4. Cutler, A.; Pribić, M.; Humphrey, L. *Everyday Ethics for Artificial Intelligence: A practical guide for*
 1107 *designers and developers*, IBM Corporation, 2018. Available online:
 1108 www.ibm.com/watson/assets/duo/pdf/everydayethics.pdf (accessed on 18 July 2019).
- 1109 5. European Commission, High-Level Expert Group on Artificial Intelligence: Draft Ethics Guidelines for
 1110 Trustworthy AI, 2018. Available online: [ec.europa.eu/digital-singlemarket/en/news/draft-ethics-](http://ec.europa.eu/digital-singlemarket/en/news/draft-ethics-guidelines-trustworthy-ai)
 1111 [guidelines-trustworthy-ai](http://ec.europa.eu/digital-singlemarket/en/news/draft-ethics-guidelines-trustworthy-ai) (accessed on 18 July 2019).
- 1112 6. von Foerster, H. On Constructing a Reality. In *Environmental Design and Research*, Volume 2, pp 35-46. F.E.
 1113 Preiser (ed.), Dowden Hutchinson and Ross, Stroudsburg PA, 1973.
- 1114 7. von Foerster, H. Ethics and Second-Order Cybernetics. In: *Understanding Understanding*. Springer, New
 1115 York, NY, 2003.
- 1116 8. Beer, S. *Brain of the Firm*, John Wiley and Sons, New York, 1972.
- 1117 9. Asimov, I. Runaround, in *Astounding Science Fiction*, March 1942.

- 1118 10. Ashby, M. How to apply the Ethical Regulator Theorem to crises, *Acta Europæana Systemica*, n°08,
1119 Brussels, Belgium, 2018. Available online: [aes.ues-eus.eu/aes2018/aes806 Mick-Ashby.pdf](https://aes.ues-eus.eu/aes2018/aes806_Mick-Ashby.pdf) (accessed on 18
1120 July 2019).
- 1121 11. Stevens, J.P. Opinion of Stevens J., Supreme Court of the United States. *Citizens United, Appellant v.*
1122 *Federal Election Commission*, *Legal Information Institute*, Cornell University Law School, 2010. Available
1123 online: www.law.cornell.edu/supct/html/08205.ZX.html (accessed on 18 July 2019).
- 1124 12. Thomas, D.R.; Beresford, A.R.; Rice, A. Security Metrics for the Android Ecosystem, 2015. Available online:
1125 www.cl.cam.ac.uk/~drt24/papers/spsm-scoring.pdf (accessed on 18 July 2019).
- 1126 13. Umpleby, S.A. What comes after second order cybernetics, *Cybernetics and Human Knowing* **2001**, 8(3):87-
1127 89.
- 1128 14. Kuhn, T. *The Structure of Scientific Revolutions*, University of Chicago Press, 1962.
- 1129 15. Glanville, R. The Cybernetics of Ethics and the Ethics of Cybernetics, Tutorial at 21st American Society for
1130 Cybernetics Conference, Virginia Beach, USA, 1986.
- 1131 16. Turing, A.M. Computing Machinery and Intelligence, *Mind* **1950**, 59:433-460.
- 1132 17. Ashby, W.R. Advanced society planned as a super brain, *Journal*, volume 14, 1951, page 3528. Available
1133 online: rossashby.info/journal/page/3527.html (accessed on 18 July 2019).
- 1134 18. Ashby, W.R. Power and I.Q. have many similar properties, *Journal*, volume 16, 1952, pages 4276-4278.
1135 Available online: rossashby.info/journal/page/4276.html (accessed on 18 July 2019).
- 1136 19. Ashby, W.R. Appearance of a super-clever machine, *Journal*, volume 16, 1952, pages 4279-4280. Available
1137 online: rossashby.info/journal/page/4279.html
- 1138 20. Clarke, A.C. Hazards of Prophecy: The Failure of Imagination, in *Profiles of the Future: An Inquiry into the*
1139 *Limits of the Future*, 1962.
- 1140 21. Wikipedia: Universal suffrage. Available online: wikipedia.org/wiki/Universal_suffrage (accessed on 18
1141 July 2019).
- 1142 22. Wikipedia: List of cognitive biases. Available online: wikipedia.org/wiki/List_of_cognitive_biases
1143 (accessed on 18 July 2019).
- 1144 23. Wikipedia: Defence mechanisms. Available online: wikipedia.org/wiki/Defence_mechanisms (accessed on
1145 18 July 2019).
- 1146 24. Wikipedia: List of fallacies. Available online: wikipedia.org/wiki/List_of_fallacies (accessed on 18 July
1147 2019).
- 1148 25. Wikipedia: Demagogue. Available online: wikipedia.org/wiki/Demagogue (accessed on 18 July 2019).
- 1149 26. Wikipedia: Denial. Available online: wikipedia.org/wiki/Denial (accessed on 18 July 2019).
- 1150 27. Eder, J. Grand Challenges for Computer Science Research: Ross Ashby Memorial Lecture of the
1151 International Federation for Systems Research, in *Cybernetics and Systems 2010*; Trappl, R. (Ed.), Vienna,
1152 Austria, pp. xix-xxv.
- 1153 28. Wilson, I.; Buckle Henning, P. A Call to Action for the Systems Sciences Community, in *Proceedings of the*
1154 *59th Annual Meeting of the International Society for the Systems Sciences*, Berlin, Germany, 2015. Vol. 1(1).
1155 Available online: journals.issn.org/index.php/proceedings59th/article/viewFile/2609/836 (accessed on 18
1156 July 2019).
- 1157 29. Wilson, T.A. The Ethical Regulator Theorem, *YouTube*, 2017. Available online:
1158 youtube.com/watch?v=NLhUajpMOI4 (accessed on 18 July 2019).
- 1159 30. Octangula, S. Structure, Environment, Purpose, and a Grand Challenge for the ASC, 2011. Available online:
1160 [asc-cybernetics.org/CofC/wp-content/uploads/2011/02/StructureEnvironment-Purpose-and-a-Grand-](http://asc-cybernetics.org/CofC/wp-content/uploads/2011/02/StructureEnvironment-Purpose-and-a-Grand-Challenge-for-the-ASC-V2.0.pdf)
1161 [Challenge-for-the-ASC-V2.0.pdf](http://asc-cybernetics.org/CofC/wp-content/uploads/2011/02/StructureEnvironment-Purpose-and-a-Grand-Challenge-for-the-ASC-V2.0.pdf) (accessed on 18 July 2019).