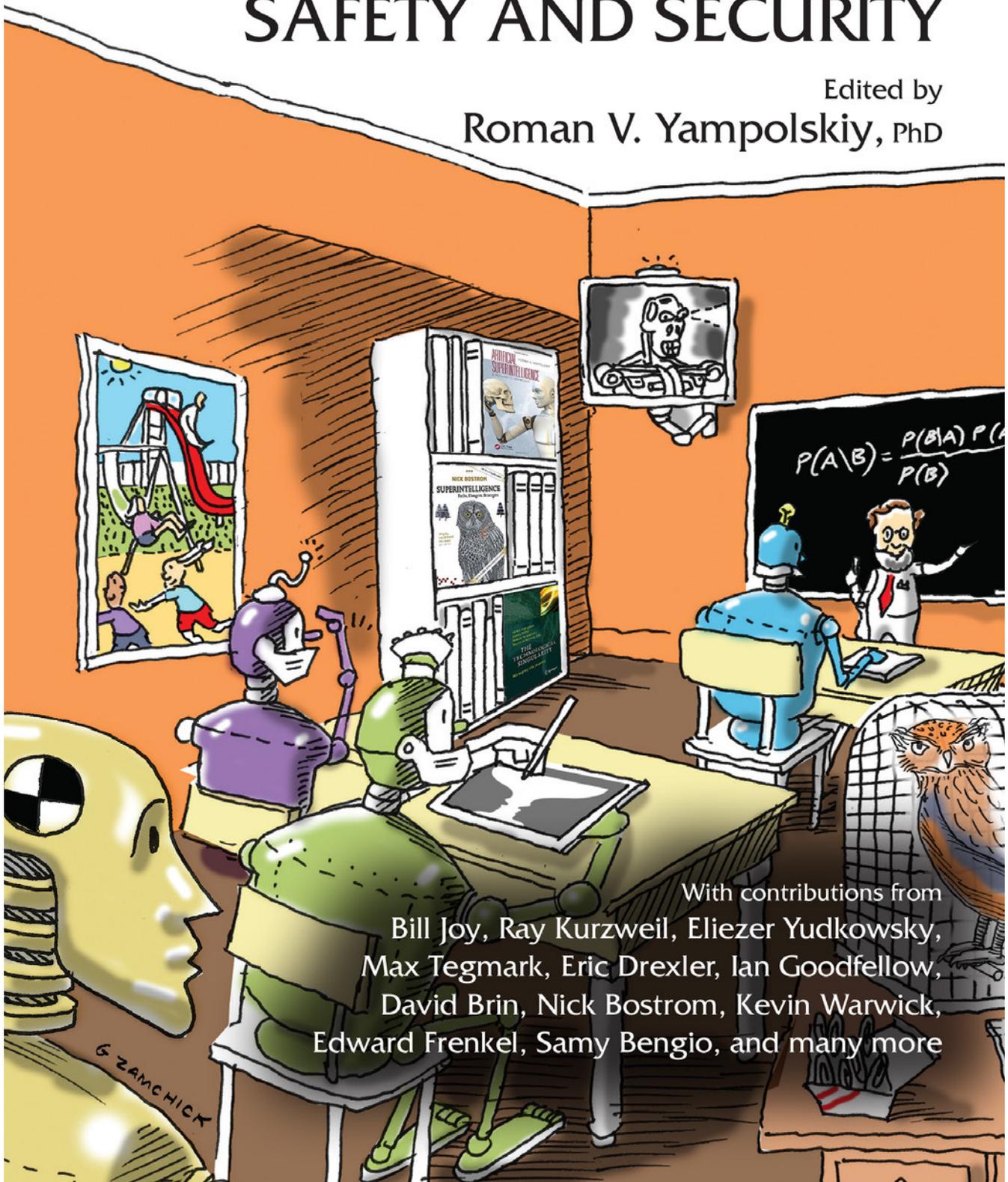


Chapman & Hall/CRC
Artificial Intelligence and Robotics Series

ARTIFICIAL INTELLIGENCE SAFETY AND SECURITY

Edited by
Roman V. Yampolskiy, PhD



With contributions from
Bill Joy, Ray Kurzweil, Eliezer Yudkowsky,
Max Tegmark, Eric Drexler, Ian Goodfellow,
David Brin, Nick Bostrom, Kevin Warwick,
Edward Frenkel, Samy Bengio, and many more

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](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/>) 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://lcn.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

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

Chapter 10	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 Chesson</i>	
Chapter 11	Strategic Implications of Openness in AI Development	145
	<i>Nick Bostrom</i>	
 PART II Responses of Scholars		
Chapter 12	Using Human History, Psychology, and Biology to Make AI Safe for Humans.....	167
	<i>Gus Bekdash</i>	
Chapter 13	AI Safety: A First-Person Perspective.....	199
	<i>Edward Frenkel</i>	
Chapter 14	Strategies for an Unfriendly Oracle AI with Reset Button	207
	<i>Olle Häggström</i>	
Chapter 15	Goal Changes in Intelligent Agents	217
	<i>Seth Herd, Stephen J. Read, Randall O'Reilly, and David J. Jilk</i>	
Chapter 16	Limits to Verification and Validation of Agentic Behavior	225
	<i>David J. Jilk</i>	
Chapter 17	Adversarial Machine Learning	235
	<i>Phillip Kuznetsov, Riley Edmunds, Ted Xiao, Humza Iqbal, Raul Puri, Noah Golmant, and Shannon Shih</i>	
Chapter 18	Value Alignment via Tractable Preference Distance	249
	<i>Andrea Loreggia, Nicholas Mattei, Francesca Rossi, and K. Brent Venable</i>	
Chapter 19	A Rationally Addicted Artificial Superintelligence.....	263
	<i>James D. Miller</i>	
Chapter 20	On the Security of Robotic Applications Using ROS	273
	<i>David Portugal, Miguel A. Santos, Samuel Pereira, and Micael S. Couceiro</i>	
Chapter 21	Social Choice and the Value Alignment Problem.....	291
	<i>Mahendra Prasad</i>	

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



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,† accidents caused by self-driving cars‡ and embarrassment from chat-bots,§ 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://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† 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.‡
- 1977 Story writing software with limited common sense produced “wrong” stories [16].
- 1982 Software designed to make discoveries, discovered how to cheat instead.§
- 1983 Nuclear attack early warning system falsely claimed that an attack is taking place.¶
- 1984 The National Resident Match program was biased in placement of married couples [17].

* https://en.wikipedia.org/wiki/Kenji_Urada

† http://lesswrong.com/lw/lvh/examples_of_ais_behaving_badly/

‡ https://en.wikipedia.org/wiki/General_Problem_Solver

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

¶ 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/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://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://www.seattletimes.com/business/microsoft/how-linkedins-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-sky-net-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

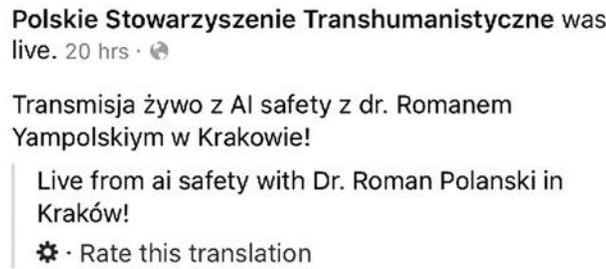


FIGURE I.1 While translating from Polish to English Facebook’s software changed Roman V. “Yampolskiy” to Roman “Polanski” due to statistically higher frequency of the latter name in sample texts.

* <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://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.†
- A mathematician AGI converts all matter into computing elements to solve problems.‡
- 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.§
- Future superintelligence optimizes away all consciousness.¶
- AGI kills humanity and converts universe into materials for improved penmanship.**
- 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://www.theguardian.com/sustainable-business/2015/jun/23/the-ethics-of-ai-how-to-stop-your-robot-cooking-your-cat>

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

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

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

** <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://selfawarenessystems.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/n41/safety_engineering_target_selection_and_alignment/

‡‡‡ <https://www.waive2018.com/>

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

¶¶¶ <http://acritch.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).^{*} 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. Sotola and R. V. Yampolskiy, “Responses to catastrophic AGI risk: Asurvey,” *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. Gundersen and L. Gundersen, "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, AAI-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, Dhvani 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

Ingeniarius, 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 & 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

Ingeniarius 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>

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 τ 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 (CH₄), 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', 318*fn*
- 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
- Openness
 - 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
- Open source
 - software, 150–151, 160
 - tools, 141
- OpenWorm*, 116
- Operational stability and control, 180
 - balance, 181–182
 - building layered immunity and tolerance, 182–183
 - honesty, 184
 - using uniformity and diversity to achieving resiliency, 183
- Operator model, 95
- “Optimizing profit”, 328
- OR, *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 superset 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, 411*fn*
 - True Negatives (TN), 258
 - True Positives (TP), 257
 - Two-phased attack, 26
- U**
- UAVs, *see* Unmanned aerial vehicles
 - Ukantianism, 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, [134fn](#)
- 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, [63fn](#)



Taylor & Francis Group
an informa business

Taylor & Francis eBooks

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

 **Routledge**
Taylor & Francis Group

 **CRC Press**
Taylor & Francis Group

Chapman & Hall/CRC
Artificial Intelligence and Robotics Series

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. This book is proposed to mitigate this fundamental problem. It is comprised of chapters from leading AI Safety researchers addressing different aspects of the AI control problem as it relates to the development of safe and secure artificial intelligence. The book is the first edited volume dedicated to addressing challenges of constructing safe and secure advanced machine intelligence.

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.



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. He is the founding and current director of the Cyber Security Lab and an author of many books including *Artificial Superintelligence: a Futuristic Approach*. 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 degree from the Department of Computer Science and Engineering at the University at Buffalo. He was a recipient of a four year National Science Foundation fellowship. Dr. Yampolskiy's main areas of interest are AI Safety, Artificial Intelligence, Behavioral Biometrics, and Cybersecurity. Dr. Yampolskiy is an author of over 150 publications including multiple journal articles and books. His research has been cited by 2000+ scientists and profiled in popular magazines both American and foreign, dozens of websites, on radio and TV. Dr. Yampolskiy's research has been featured 1000+ times in numerous media reports in some 30 languages.