

Regulating for ‘Normal AI Accidents’:

Operational Lessons for the Responsible Governance of Artificial Intelligence Deployment

Matthijs M. Maas[†]

Faculty of Law.

University of Copenhagen.

Copenhagen, Denmark.

Matthijs.maas@jur.ku.dk

ABSTRACT

New technologies, particularly those which are deployed rapidly across sectors, or which have to operate in competitive conditions, can disrupt previously stable technology governance regimes. This leads to a precarious need to balance caution against performance while exploring the resulting ‘safe operating space’. This paper will argue that Artificial Intelligence is one such critical technology, the responsible deployment of which is likely to prove especially complex, because even narrow AI applications often involve networked (tightly coupled, opaque) systems operating in complex or competitive environments. This ensures such systems are prone to ‘normal accident’-type failures which can cascade rapidly, and are hard to contain or even detect in time. Legal and governance approaches to the deployment of AI will have to reckon with the specific causes and features of such ‘normal accidents’. While this suggests that large-scale, cascading errors in AI systems are inevitable, an examination of the operational features that lead technologies to exhibit ‘normal accidents’ enables us to derive both tentative principles for precautionary policymaking, and practical recommendations for the safe(r) deployment of AI systems. This may help enhance the safety and security of these systems in the public sphere, both in the short- and in the long term.

KEYWORDS

Ethical design and development of AI systems; AI and Law; trust and explanations in AI systems; normal accident theory

ACM Reference format:

Matthijs M. Maas. 2018. Regulating for ‘normal AI accidents’: operational lessons for the responsible governance of artificial intelligence deployment. In *Proceedings of 2018 AAAI/ACM Conference on AI, Ethics, and Society (AIES ’18)*, February 2–3, 2018, New Orleans, LA. ACM, NY, NY, USA, 6 pages. <https://doi.org/10.1145/3278721.3278766>

[†] Centre for International Law, Conflict and Crisis, Copenhagen University. Karen Blixens Plads 16, 2300 Copenhagen S. ORCID ID: 0000-0002-6170-9393. Research Affiliate, Governance of AI Program, Future of Humanity Institute, University of Oxford.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from Permissions@acm.org.

AIES ’18, February 2–3, 2018, New Orleans, LA, USA

© 2018 Copyright is held by the owner/author(s). Publication rights licensed to

ACM. ACM 978-1-4503-6012-8/18/02...\$15.00

<https://doi.org/10.1145/3278721.3278766>

1 Introduction

New technologies, particularly those which are deployed rapidly across industries, or which can offer a (economic, political, military) edge, frequently disrupt previously stable international governance arrangements or power distributions. The introduction of such technology can therefore be followed by a period of uncertainty and risk, as policymakers, operators and public seek to grasp their societal and operational implications. This leads to a pre-carious need to balance caution against performance while exploring the ‘safe operating space’ of said critical technology.

Artificial Intelligence (AI) is one such critical technology. In recent years, the field of AI has advanced at a rapid pace. While researchers and experts still differ in their assessment of when, if ever, we may expect the achievement of ‘general’ artificial intelligence [1]–[6], today’s narrow artificially intelligent systems already match or outperform humans across many narrow domains. In just the past three years—to provide an incomplete list—AI systems have proven that they can meet or exceed human performance in object image recognition [7], speech transcription and direct translation [8]–[10]. AI systems have learned how to drive [11]; can parse paragraphs to answer questions posed [12]; recognize human faces (even in blurred pictures) and some emotions [13]; can create new encryption schemes and detect malware [14], [15]; identify crop diseases [16]; teach themselves the game of go in mere days [17], [18], and write cookbooks, news articles, music and published poetry [19]–[22]. As a result, these systems are beginning to see wider adoption in a broad range of applications—across fields as diverse as stock markets, transport infrastructure, healthcare, agriculture, education, cybersecurity and the military.

2 Mapping AI Governance Challenges

Along with the immense appeal, there is a growing recognition amongst policymakers that the development and deployment of AI brings ethical, political, societal and operational risks, and will even pose systemic challenges for regulation and governance frameworks [23]–[26]. Scholars have begun to articulate comprehensive research agendas aimed at exploring these governance challenges [27]; loosely, these challenges can be grouped in two categories—‘(ab)use’ and ‘accident’.

2.1 AI governance challenges from (ab)use

Many of the greatest challenges posed by AI derive from their use by diverse actors either to enable greater legibility or facility in the exercise of power, or in a competitive context—whether economic, political, or military. Thus, domestically, the impact of introducing even ‘narrow’ artificial intelligence systems may result in far-reaching technological unemployment [28], [29], comprehensive erosion of privacy [30]; over-dependency on social robots [31], [32]; increased inequality and societal dysfunction as a result of machine bias [33], or as a result of (perceived) electoral manipulation, and political polarization as a result of pervasive customized ‘computational propaganda’ [34]. Moreover, the automation of social-media identity theft [35], as well as automated system-vulnerability-detection and victim-customized social-engineering attacks—as illustrated by the ‘Mayhem’ AI in the 2016 DARPA Grand Cyber Challenge, and by tools such as ‘WifiPhisher’, respectively [36], [37]—suggest an unprecedented increase in hazards on cyberspace, from both criminal and state actors [38].

In a global context, meanwhile, the integration of lethal autonomous weapons systems with ‘war-algorithms’ [39], amongst a host of other emerging AI battlefield applications [40], raises deep ethical, legal, and operational ramifications [41]–[43]. For instance, enhanced surveillance capabilities and autonomous weapons may also disrupt established strategic landscapes, entrenching authoritarian regimes’ ability to monitor dissent and deploy centralized stand-off force projection capabilities [44], [45]. Moreover, autonomous weapons systems enable at-cost swarming tactics, in which an adversary’s defensive measures (i.e. point-defense systems) are overwhelmed by sustained simultaneous attacks. This technological innovation overturns a prevailing tactical offense-defense balance, possibly upsetting global equilibria by putting a premium on pre-emption [46]. The deployment of dispersed autonomous systems, along with AI-enabled advances in sensing and data analysis capabilities, might even put at risk the stability of existing nuclear deterrence dyads, by rendering previously ‘secure’ nuclear assets (such as ballistic missile submarines) vulnerable once more [47]–[50].

2.2 AI governance challenges from accident

Right along with the governance challenges posed by the ‘systemic’ or ‘intended’ deployment of AI, however, new AI governance approaches will also have to reckon with underlying operational risks from failure. While regulatory approaches or governance systems are no stranger to (industrial) accidents involving new (even robotic) technologies, they may be prone to misunderstanding the scale and intractability of the safety challenges posed by deployed AI systems. Critically, an anecdotal reliance on familiar accidents involving ‘embodied’ robots (such as factory robots, or autonomous cars) or on ‘malfunctioning chatbots’ as the type specimens for ‘AI accident risks’, is likely to lead us to misunderstand the causes and dynamics of cascading accidents involving networked and opaque AI systems. Worse, it may lead us to underestimate their frequency, scale, and reach.

Empirically, leading AI architectures have demonstrated risks of ‘flash crashes’ amongst other failure modes [42, pp. 35–36],

[51], and demonstrate surprisingly creative behavior [52] which can be intrinsically difficult to anticipate (in the context of evolutionary programming architectures) or reconstruct and explain post-accident (in the case of deep learning and certain neural network architectures). There exists an incipient and broad research agenda examining ‘concrete problems in AI Safety’ [53], to avoid unexpected and unwanted side effects as AI systems become more advanced and/or are deployed to more complex environments, but this valuable research program is at present still at an early stage.

This is a key problem, in light of the deployment and integration of increasingly capable AI systems across society—from the financial sector to the penal system, and from critical infrastructure to the military. There is a wide range of potential principles informing AI regulatory regimes, whether the precautionary principle; ‘responsible research and innovation’ paradigm [54], [55]; ‘differential technological progress’ [56], or some other set of emerging norms or policy desiderata [57]. Yet the observation of recurring errors in AI systems suggests that these governance regimes must also consider the problem of ‘normal accidents’—not as a possibility, but as an inevitability.

3 AI accidents as ‘normal accidents’

Developed by Charles Perrow in the wake of the Three Mile Island nuclear reactor meltdown [58], ‘normal accident theory’ has been applied to understand catastrophic technological failure across a wide range of domains, ranging from the *Apollo 13*, *Challenger* and *Columbia* spacecraft accidents [58, p. 23] to false-alarms and near-accidents plaguing the US nuclear forces [59]; and from the Air France 447 crash to the 2003 Gulf War ‘Patriot fratricides’—where faulty IFF systems led semi-autonomous coalition air defenses to shoot down friendly aircraft [42, pp. 30–33], to name but a few.

‘Normal accident’ theory analyzes how accidents at the crux of mechanical, software, operator and organizational failures. It is this systemic perspective which makes normal accident theory so useful in understanding the hazards and failure-modes of deployed AI systems; while each new technology should of course be assessed on its own characteristics, there are also valuable lessons we can derive from our past experience with strategically powerful technologies. While on an object (or scientific) level, AI is of course a profoundly different technology from, say, nuclear weapons, on an operational level, they share key features which make these assemblages (sensors, algorithms, human operators) prone to normal accidents.

3.1 AI systems are complex and opaque

Normal accident theory focuses on system applications that are complex and tightly coupled. Complex systems are those which have many interlocking parts, ‘black-boxed’ processing units, or *hidden interactions*—constellations of feedback loops which cannot be observed or fully understood directly or in real-time, but only imperfectly inferred, based on aggregate behavior. This makes the system more complex than can be properly understood by the human operator [58, p. 9]. Moreover, there are many *common mode connections*, where it is not (immediately)

clear which components have failed when there is an overall failure.

Tightly coupled means that “there is no slack or buffer or give between two items. What happens in one directly affects what happens in the other.” [58, pp. 89–90]. For instance, in the context of the US nuclear force, scholars have recorded how the tightly interlinked web of early launch warning satellites, Chrome Dome aircraft, and command-and-control nodes, created a large scope for small technical failures or operator errors to rapidly cascade throughout the system, creating major false alarms [59], [60].

Critically, emerging AI systems meet all of the relevant criteria to exhibit ‘normal accidents’: many AI algorithms are ‘black boxes’, complex in design or operation. Like all software programs, they almost by default contain bugs—past studies have estimated that the software industry sees an average error rate of 15–50 errors per 1,000 lines of code [42, p. 13], [61]. The problem of complexity is acute for many approaches in machine learning, in particular. These networks are intrinsically opaque—it is impossible to produce a formal proof of their behavior; they can be stochastic; and it is difficult to adequately anticipate all real-world scenarios, or to test a system’s reactions to them, within a simulated sandbox [62, pp. 8–9]. Critically, while work on the ‘interpretability’ of machine learning decisions is advancing—compare DARPA’s work on ‘explainable’ AI [63], [64]—this field is still in its infancy. AI systems as a result also suffer from common-mode failures, as it can be unclear where the error has originated—whether in sensor fault; underlying bias in the training datasets, adversarial input or infection by a virus, or some other component.

The opacity of leading AI systems also leads to an additional problem, in that it inhibits operator’s critical ability to learn from ‘near-accidents’ or ‘close calls’. After all, for certain AI applications, it may be unclear what would be the signature of such a ‘near-accident’. We may not get many warning shots, and past performance of a system may not adequately prepare us for the scope of eventual failure.

3.2 AI systems are tightly-coupled and fast

Moreover, many AI applications are tightly coupled, involving fast interactions and reaching and executing decisions at a speed that arguably exceeds the coupling of past systems (such as nuclear reactors) prone to normal accidents. The tight coupling and high operational speed are key features of AI systems which are plugged into extensive networks (such as the Internet of Things), or which operate in competitive environments such as high-frequency trading markets, in cyberspace or on the battlefield. The speed of AI operation ensures that when errors inevitably emerge, they are not just difficult to detect, but are also hard to halt in time. This suggests that having a human operator ‘on-the-loop’ is not always viable, if interaction speeds are high enough.

Worse, redundancies and safety measures built into an AI system can actually cause accidents. This is because features such as self-performance-monitoring sensors or -software, automated fail-safes, or behavioral tripwires [65, p. 137] may increase the overall complexity of a system. They add more ‘interacting parts’ which themselves can quietly fail or react in unanticipated ways.

In this way, technological safety measures may hinder problem isolation. At the same time, research has shown that automated safety systems can instill a blanket trust (‘automation bias’) in operators or users—a trust that may, perversely, encourage greater risk-taking by those users as a result of risk homeostasis [62, p. 12], [66]. Moreover, as demonstrated by the ‘Patriot fratricides’, automation bias may cause human operators who are nominally in-the-loop to nonetheless trust the system without question, authorizing even incorrect action requests by force of habit [42, p. 31], [67].

3.3 AI developers and operators have multiple objectives beyond safety

The technical propensity of AI systems to normal accidents (their complexity and tight coupling), will likely be exacerbated by the incentives of the principals that run them. This is because the designers, trainers, and operators of AI systems may in practice encounter multiple, conflicting organizational objectives beyond pure ‘safety’.

At the level of AI developers and trainers, there may be restrictions on error feedback and on learning from incidents. This is the case both within companies (for instance, when working towards a tight software deployment deadline, and reporting apparently insignificant ‘anomalies’ would be an inconvenience), as well as between them (for instance, when it is feared that sharing details on security incidents may give away critical information about an AI’s architecture or the initial settings of its algorithms).

On an operator level, there is a tension between automating functions to allow for rapid functioning, and decoupling & decentralizing them, to enable flexible error response. More obviously, for operators, safety often must be traded off against performance [59, p. 13], [68]: some AIs, for instance those used to automate cybersecurity penetration testing, may be designed to come up with ‘unintuitive’ solutions and to test these rapidly. Of course, we necessarily and understandably accept some risk when using many technologies—possibly only an inert system would be perfectly (and knowably) safe. Yet the impact of these errors can be particularly high once AI systems are integrated in major infrastructures.

3.4 Competitive pressures exacerbate operational risks of AI

Because of the above, many AI systems are likely to demonstrate at least some propensity towards normal accidents. Yet there are some exceptions.

In the first place, in contexts where the decisions made by the AI are less tightly coupled to (irreversible) physical outcomes, the system may be less susceptible to normal accidents, and risks may be modest or manageable. For instance, while algorithms trained on biased datasets are a real and pressing problem when applied to, for instance, sentencing or incarceration decisions [69]–[71], the speed at which these decisions are implemented is bottlenecked—for now—by the human organizations that act upon them. As some of these decisions are not time-critical, in principle, this could give operators some leeway to spot errors or test for

(statistical) discrepancies in the system's output. In such contexts, having a human-in-the-loop is not a (performance-inhibiting) safety measure, but instead a basic feature of the assemblage, and can function as a relatively effective fail-safe containing the error cascade.

In the second place, for many AI applications in society—such as in healthcare, transport or critical infrastructure—all actors involved share an interest in safe and reliable operation, which could still serve to promote the sharing of lessons and best practices even against the gradient of inter-actor rivalry.

However, neither of these caveats (uncoupling; unanimous interest in safety) applies in a competitive context, such as on physical or especially virtual battlefields. In a competitive context, a number of factors begin to exacerbate the operational risks of deployed AI systems. These factors include: (1) there is a pressure towards deploying systems rapidly, before adversaries do, and to err on this side, rather than on extensive safety testing. This manifests itself in what Richard Danzig has called ‘technology roulette’—a strategic dynamic where the pursuit of technological supremacy by one actor (in areas such as nuclear weapons, synthetic biology, and AI) drives processes of proliferation and the rapid deployment of systems, creating new emergent risks and perversely undercutting (national) security [72]. (2) decisions must be made based on incomplete, ‘messy’, non-structured and potentially unreliable data; (3) the speed of reaction or interaction is accelerated even further; (4) there are risks of systems being hacked [42, p. 34]—either directly, or indirectly, through spoofing or behavioral exploitation.

For instance, deep neural networks have proven vulnerable to confrontation with adversarial examples [73], which can be generated without any privileged access to the algorithm’s training data or logic; these spoofing attacks can moreover be hidden, so that they are invisible to humans. At their extreme, unexpected interactions between competing systems, especially in cyberspace, could cause unexpected escalation—a ‘flash war’ [74], analogous to the algorithmic flash crashes observed in the financial sector.

Because of the advantage afforded by operational speed, many of the ‘normal accident’ problems may therefore be particularly acute in military AI applications [42], [62]—a frightening prospect, given that major powers, including the United States,[75], [76] China [26], [77], [78], and Russia [79] increasingly perceive AI as a cornerstone of their next-generation military and strategic supremacy. In the context of competitive pressures, there may occur technology race dynamics which produce a strong pressure to cut down on safety and err on the side of rapid deployment of relatively untested systems [80], even when principals are nominally committed to safe and responsible development.

4 Implications for AI governance

In sum, it appears plausible that many AI applications may be even more susceptible to normal accidents than past ‘textbook’ case technologies such as nuclear power or aviation. Moreover, normal accident theory suggests—and past and present experience in the field of cybersecurity has repeatedly borne out

[51]—that such risks cannot simply be ‘designed out’ of the technology (at least not without giving up on many of their benefits). These operational insights into the ‘normal’ failure modes of AI systems matter on two levels.

On a systemic level, they provide an overall context for understanding when to ‘trust’ AI systems—after all, not all AI configurations or applications are equally vulnerable to normal accidents—and conversely, when operator ‘trust’ in an AI begins to turn into an accident-enabling ‘automation bias’. Distinguishing which AIs are susceptible to normal accident failure, and which are not, aids in improving the scope of regulation, to focus debate on risks while avoiding stifling and unnecessary blanket restrictions; it also can inform cost-benefit analyses of which configurations to deploy to which sectors and environments, and what are acceptable risks.

Secondly, at the governance level of principles, the likely susceptibility of some (especially competitive) AI systems to normal accidents suggests that even if regulatory regimes can converge on concrete norms and standards to ensure that AI systems are deployed in a lawful and ethical manner, ‘unforeseeable’ yet inevitable accidents will emerge in their performance, putting both users and the public at risk. This throws up interesting problems for frameworks ranging from liability to disaster insurance, to name but a few. There are, as of yet, no clear-cut answers to the question of how to ‘prevent’ normal accidents. The following is therefore largely speculative. Nonetheless, one upside of AI is that it is, in its applications and range of configurations, arguably more ‘flexible’ or customizable than, say, a nuclear power plant. As such, we may still try to derive relevant lessons for responsible AI regulation.

In the first place, legal and regulatory strategies should understand that accidents cannot be ‘designed out’, and that we cannot rely on automated fail-safes only, which do not reliably protect us from accident cascades; regulators should understand how the illusion of reliability which such ‘safety’ systems create may in some ways make us less safe or robust once disaster does strike.

Secondly, and in contrast, regulators should accept that even the ideal of maintaining a human-in-the-loop does not offer the level of safety, reliability or control which is often ascribed to it. Indeed, they must understand that (absent broader measures) such an arrangement at best simply sets up an operator to take the fall—to serve as ‘moral crumple zone’ [81]—once accidents happen, and at worst actually creates new avenues for those errors to be introduced.

Instead, regulators (and innovators) might seek to focus regulation on the ‘levers’ affecting normal accident risk: the nature (complexity, opacity, tight coupling, and speed) of the AI system in question; and the operational environment (the organizational objectives and competitive pressures experienced by principals or users) in which it operates. For example, policies could encourage research into ‘Explainable AI’, and promote adoption of resulting architectures (i.e. reduce ‘opacity’); promote heterogeneity in deployed AI architectures or cap unnecessary integration with networks (i.e. reduce ‘coupling’ & ‘complexity’), to insulate systems from flash crashes; restricting a system’s

deployment outside of intended operational environments, or restricting its autonomy and speed in such environments, especially competitive environments that might allow for the easy injection of adversarial input (i.e. reduce complexity; limit competitive environment); train operators to better recognize signs of automation bias; and emphasize inter-organization exchange of best practice and sharing of incident reports (align organizational objectives). This is just an initial sketch of what should be a much broader line of research into ways to anticipate and mitigate emergent risks from AI systems.

5 Conclusion

This argument has sought to explore responsible and robust governance approaches for the deployment of AI. Like all tightly coupled, opaque systems, AIs will be prone to ‘normal accidents’, ensuring that perfect safety may not be attainable. Nonetheless, we may anticipate the particular ‘field conditions’ under which AIs are more or less susceptible to such errors, and try to account for these ‘risk factors’. On this basis, this paper has briefly suggested tentative principles and practical recommendations for precautionary regulation.

Further research (as per usual) is needed. Such research might seek to map out in which sectors, or for which applications, AI systems are more or less likely to exhibit features priming them for normal accidents—and in which of these cases it is possible to mitigate these risk factors (by reducing opacity; by increasing ‘slack’) without incurring decay in overall performance. In addition, research could explore assemblages of organizational and technological innovation, which might better promote graceful failures, or even function as hypothetical ‘super-fail-safe’ solutions (perhaps allowing the failure of system or operator in isolation, but reliably preventing the simultaneous failure of both). Finally, research might examine whether new generations of AI systems might perhaps after all be able to serve as more reliable ‘monitoring and fail-safe systems’, potentially by identifying ‘signatures of failure’ in deployed AIs. The history of normal accidents suggests that the error will get through.

ACKNOWLEDGMENTS

The author thanks the attendees at the 2018 AAAI/ACM Conference on AI, Ethics, and Society for their feedback and comments. No conflict of interest is identified.

REFERENCES

- [1] A. Plebe and P. Perconti, “The Slowdown Hypothesis,” in *Singularity Hypotheses*, A. H. Eden, J. H. Moor, J. H. Søraker, and E. Steinhart, Eds. Springer Berlin Heidelberg, 2012, pp. 349–365.
- [2] T. G. Dietterich and E. J. Horvitz, “Rise of concerns about AI: reflections and directions.” *Commun. ACM*, vol. 58, no. 10, pp. 38–40, Sep. 2015.
- [3] S. D. Baum, B. Goertzel, and T. G. Goertzel, “How long until human-level AI? Results from an Expert Assessment,” *Technol. Forecast. Soc. Change*, vol. 78, pp. 185–195, 2011.
- [4] S. Armstrong and K. Sotala, “How We’re Predicting AI—Or Failing To,” in *Beyond AI: Artificial Dreams*, J. Romportl, P. Ircing, E. Zackova, M. Polak, and R. Schuster, Eds. Pilsen: University of West Bohemia, 2012, pp. 52–75.
- [5] M. Brundage, “My AI Forecasts—Past, Present, and Future (Main Post),” 2017. [Online]. Available: <http://www.milesbrundage.com/1/post/2017/01/my-ai-forecasts-past-present-and-future-main-post.html>. [Accessed: 20-Feb-2017].
- [6] V. C. Müller and N. Bostrom, “Future Progress in Artificial Intelligence: A Survey of Expert Opinion,” in *Fundamental Issues of Artificial Intelligence*, Müller, Vincent C., Ed. Berlin: Synthese Library, 2016.
- [7] A. Linn, “Microsoft researchers win ImageNet computer vision challenge,” Next at Microsoft, 10-Dec-2015..
- [8] W. Xiong et al., “Achieving human parity in conversational speech recognition,” *ArXiv Prepr. ArXiv161005256*, 2016.
- [9] D. Castelvecchi, “Deep learning boosts Google Translate tool,” *Nat. News*.
- [10] G. Lewis-kraus, “The Great A.I. Awakening,” *The New York Times*, 14-Dec-2016.
- [11] R. Bryant, “Google’s AI becomes first non-human to qualify as a driver,” *Dezeen*, 12-Feb-2016. [Online]. Available: <https://www.dezeen.com/2016/02/12/google-self-driving-car-artificial-intelligence-system-recognised-as-driver-usa/>. [Accessed: 25-Feb-2017].
- [12] C. Metz, “The AI Threat Isn’t Skynet. It’s the End of the Middle Class,” *WIRED*, 2017.
- [13] L. Newman, “AI Can Recognize Your Face Even If You’re Pixelated,” *WIRED*, 2016. [Online]. Available: <https://www.wired.com/2016/09/machine-learning-can-identify-pixelated-faces-researchers-show/>. [Accessed: 26-Feb-2017].
- [14] M. Abadi and D. G. Andersen, “Learning to Protect Communications with Adversarial Neural Cryptography,” 2016.
- [15] L. Musthaler, “How to use deep learning AI to detect and prevent malware and APTs in real-time,” *Network World*, 11-Mar-2016. [Online]. Available: <http://www.networkworld.com/article/3043202/security/how-to-use-deep-learning-ai-to-detect-and-prevent-malware-and-apts-in-real-time.html>. [Accessed: 25-Feb-2017].
- [16] D. Furness, “AI in agriculture? Algorithms help farmers spot crop disease like experts,” *Digital Trends*, 06-Oct-2016. [Online]. Available: <http://www.digitaltrends.com/computing/ai-crop-disease/>. [Accessed: 20-Feb-2017].
- [17] D. Silver et al., “Mastering the game of Go without human knowledge,” *Nature*, vol. 550, no. 7676, p. nature24270, Oct. 2017.
- [18] D. Silver et al., “Mastering the game of Go with deep neural networks and tree search,” *Nature*, vol. 529, no. 7587, pp. 484–489, Jan. 2016.
- [19] H. Mascarenhas, “Associated Press to expand its sports coverage by using AI to write Minor League Baseball articles,” *International Business Times UK*, 05-Jul-2016. [Online]. Available: <http://www.ibtimes.co.uk/associated-press-expand-its-sports-coverage-by-using-ai-write-minor-league-baseball-articles-1568804>. [Accessed: 25-Feb-2017].
- [20] A. Marshall, “From Jingles to Pop Hits, A.I. Is Music to Some Ears,” *The New York Times*, 22-Jan-2017.
- [21] Z. Scholl, “Turing Test: Passed, using computer-generated poetry,” *Raspberry PI AI*, 24-Jan-2015..
- [22] A. Kleeman, “Cooking with Chef Watson, I.B.M.’s Artificial-Intelligence App,” *The New Yorker*, 2016. [Online]. Available: <http://www.newyorker.com/magazine/2016/11/28/cooking-with-chef-watson-ibms-artificial-intelligence-app>. [Accessed: 25-Feb-2017].
- [23] OSTP, “Preparing for the Future of Artificial Intelligence,” Office of Science and Technology Policy, National Science and Technology Council Committee on Technology, Oct. 2016.
- [24] US Senate Subcommittee on Space, Science and Competitiveness, *The Dawn of Artificial Intelligence*. 2016.
- [25] House of Commons Science and Technology Committee, “Robotics and artificial intelligence: Fifth Report of Session 2016–17..,” 2016.
- [26] China’s State Council, “A Next Generation Artificial Intelligence Development Plan,” *New America Cybersecurity Initiative*, Aug. 2017.
- [27] A. Dafoe, “AI Governance: A Research Agenda,” p. 52, 2018.
- [28] E. Brynjolfsson and A. McAfee, *The Second Machine Age: Work, Progress and Prosperity in a Time of Brilliant Technologies*. New York: W.W. Norton & Company, 2014.
- [29] R. Hanson, *The Age of Em: Work, Love, and Life when Robots Rule the Earth*. Oxford, New York: Oxford University Press, 2016.
- [30] R. Calo, “Peeping HALS: Making Sense of Artificial Intelligence and Privacy,” *EJLS - Eur. J. Leg. Stud.*, vol. 2, no. 3, 2010.
- [31] P. Lin, K. Abney, and G. A. Bekey, *Robot Ethics: The Ethical and Social Implications of Robotics*. MIT Press, 2011.
- [32] Y.-H. Weng, “Beyond Robot Ethics: On a Legislative Consortium for Social Robotics,” *Adv. Robot.*, vol. 24, no. 13, pp. 1919–1926, Jan. 2010.
- [33] K. Crawford and R. Calo, “There is a blind spot in AI research,” *Nat. News*, vol. 538, no. 7625, p. 311, Oct. 2016.
- [34] S. C. Woolley and P. N. Howard, “Political Communication, Computational Propaganda, and Autonomous Agents — Introduction,” *Int. J. Commun.*, vol. 10, no. 0, p. 9, Oct. 2016.
- [35] L. Bilge, T. Strufe, D. Balzarotti, and E. Kirda, “All Your Contacts Are Belong to Us: Automated Identity Theft Attacks on Social Networks,” presented at the Track: Security and Privacy / Session: Web Security, 2009.
- [36] DARPA, “‘Mayhem’ Declared Preliminary Winner of Historic Cyber Grand Challenge,” 2016. [Online]. Available: <http://www.darpa.mil/news-events/2016-08-04>. [Accessed: 11-Mar-2017].
- [37] wifiphisher, wifiphisher: Automated victim-customized phishing attacks against Wi-Fi clients, Wifiphisher, 2017.
- [38] M. Brundage et al., “The Malicious Use of Artificial Intelligence: Forecasting, Prevention, and Mitigation,” *ArXiv180207228 Cs*, Feb. 2018.

- [39] D. A. Lewis, G. Blum, and N. K. Modirzadeh, "War-Algorithm Accountability," Social Science Research Network, Rochester, NY, SSRN Scholarly Paper ID 2832734, Aug. 2016.
- [40] M. Maas, T. Sweijns, and S. De Spiegeleire, Artificial Intelligence and the Future of Defense: Strategic Implications for Small- and Medium-Sized Force Providers. The Hague, The Netherlands: The Hague Centre for Strategic Studies, 2017.
- [41] B. Nehal, S. Beck, R. Geiss, H.-Y. Liu, and K. Kress, Eds., Autonomous Weapons Systems: Law, Ethics, Policy. Cambridge: Cambridge University Press, 2016.
- [42] P. Scharre, "Autonomous Weapons and Operational Risk," Center for a New American Security, 2016.
- [43] H. M. Roff, "The Strategic Robot Problem: Lethal Autonomous Weapons in War," *J. Mil. Ethics*, vol. 13, no. 3, pp. 211–227, Oct. 2014.
- [44] M. C. Horowitz, "Who'll want artificially intelligent weapons? ISIS, democracies, or autocracies?", *Bulletin of the Atomic Scientists*, 29-Jul-2016.
- [45] M. C. Horowitz, S. E. Kreps, and M. Fuhrmann, "Separating Fact from Fiction in the Debate over Drone Proliferation," *Int. Secur.*, vol. 41, no. 2, pp. 7–42, Oct. 2016.
- [46] J.-M. Rickli, "Artificial Intelligence and the Future of Warfare," in WEF Global Risks Report 2017, 2017, p. 49.
- [47] R. Courtland, "DARPA's Self-Driving Submarine Hunter Steers Like a Human," *IEEE Spectrum: Technology, Engineering, and Science News*, 07-Apr-2016. [Online]. Available: <http://spectrum.ieee.org/automaton/robotics/military-robots/darpa-activates-self-driving-submarine-hunter-steers-like-a-human>. [Accessed: 19-Sep-2016].
- [48] K. A. Lieber and D. G. Press, "The New Era of Counterforce: Technological Change and the Future of Nuclear Deterrence," *Int. Secur.*, vol. 41, no. 4, pp. 9–49, Apr. 2017.
- [49] D. Hambling, "The Inescapable Net: Unmanned Systems in Anti-Submarine Warfare," British-American Security Information Council, 2016.
- [50] E. Geist and A. J. Lohn, "How Might Artificial Intelligence Affect the Risk of Nuclear War?," RAND, 2018.
- [51] R. V. Yampolskiy and M. S. Spellchecker, "Artificial Intelligence Safety and Cybersecurity: a Timeline of AI Failures," *ArXiv Prepr. ArXiv161007997*, 2016.
- [52] J. Lehman et al., "The Surprising Creativity of Digital Evolution: A Collection of Anecdotes from the Evolutionary Computation and Artificial Life Research Communities," *ArXiv180303453 Cs*, Mar. 2018.
- [53] D. Amodei, C. Olah, J. Steinhardt, P. Christiano, J. Schulman, and D. Mané, "Concrete Problems in AI Safety," *ArXiv160606565 Cs*, Jun. 2016.
- [54] J. Stilgoe, R. Owen, and P. Macnaghten, "Developing a framework for responsible innovation," *Res. Policy*, vol. 42, no. 9, pp. 1568–1580, Nov. 2013.
- [55] M. Brundage, "Artificial Intelligence and Responsible Innovation," in Fundamental Issues of Artificial Intelligence, V. C. Müller, Ed. Springer International Publishing, 2016, pp. 543–554.
- [56] N. Bostrom, "Existential risks: Analyzing human extinction scenarios and related hazards," *J. Evol. Technol.*, vol. 9, no. 1, 2002.
- [57] N. Bostrom, A. Dafoe, and C. Flynn, "Public Policy and Superintelligent AI: A Vector Field Approach," 2018.
- [58] C. Perrow, "Normal Accidents: Living with High Risk Technologies," 1984. [Online]. Available: <http://press.princeton.edu/titles/6596.html>. [Accessed: 04-Mar-2017].
- [59] S. D. Sagan, The Limits of Safety: Organizations, Accidents, and Nuclear Weapons. Princeton: Princeton University Press, 1993.
- [60] E. Schlosser, Command and Control: Nuclear Weapons, the Damascus Accident, and the Illusion of Safety, Reprint edition. New York: Penguin Books, 2014.
- [61] S. McConnell, Code Complete: A Practical Handbook of Software Construction, Second Edition, 2nd edition. Redmond, Wash: Microsoft Press, 2004.
- [62] J. Borrie, "Safety, Unintentional Risk and Accidents in the Weaponization of Increasingly Autonomous Technologies," UNIDIR, 5, 2014.
- [63] DARPA, "GRANT OPPORTUNITY - DARPA-BAA-16-53 - Explainable Artificial Intelligence (XAI)." Department of Defense - DARPA - Information Innovation Office.
- [64] D. Doran, S. Schulz, and T. R. Besold, "What Does Explainable AI Really Mean? A New Conceptualization of Perspectives," *ArXiv171000794 Cs*, Oct. 2017.
- [65] N. Bostrom, Superintelligence: Paths, Dangers, Strategies. Oxford University Press, 2014.
- [66] M. L. Cummings, "Automation bias in intelligent time critical decision support systems," in *AIAA 1st Intelligent Systems Technical Conference*, 2004, vol. 2, pp. 557–562.
- [67] J. K. Hawley, "Not by Widgets Alone: The Human Challenge of Technology-intensive Military Systems," *Armed Forces J.*, 2011.
- [68] P. Lewis, H. Williams, B. Pelopidas, and S. Aghlani, "Too Close for Comfort: Cases of Near Nuclear Use and Options for Policy," Chatham House, Apr. 2014.
- [69] L. Kirchner, J. Angwin, J. Larson, and S. Mattu, "Machine Bias: There's Software Used Across the Country to Predict Future Criminals. And it's Biased Against Blacks..," ProPublica, 23-May-2016. [Online]. Available: <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>. [Accessed: 24-May-2017].
- [70] S. Corbett-Davies, E. Pierson, A. Feller, and S. Goel, "A computer program used for bail and sentencing decisions was labeled biased against blacks. It's actually not that clear.," Washington Post, 17-Oct-2016.
- [71] T. Bolukbasi, K.-W. Chang, J. Zou, V. Saligrama, and A. Kalai, "Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings," *ArXiv160706520 Cs Stat*, Jul. 2016.
- [72] R. Danzig, "Technology Roulette: Managing Loss of Control as Many Militaries Pursue Technological Superiority," Center for a New American Security, Jun. 2018.
- [73] A. Nguyen, J. Yosinski, and J. Clune, "Deep neural networks are easily fooled: High confidence predictions for unrecognizable images," in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2015, pp. 427–436.
- [74] P. Scharre, "Flash War - Autonomous Weapons and Strategic Stability," presented at the Understanding Different Types of Risk, Geneva, 11-Apr-2016.
- [75] Center for Strategic and International Studies, "Assessing the Third Offset Strategy." [Online]. Available: <https://www.csis.org/events/assessing-third-offset-strategy>. [Accessed: 05-Mar-2017].
- [76] S. J. Freedberg, "Centauro Army: Bob Work, Robotics, & The Third Offset Strategy," Breaking Defense, 09-Nov-2015.
- [77] E. Kania, "Great Power Competition and the AI Revolution: A Range of Risks to Military and Strategic Stability," Lawfare, 19-Sep-2017. [Online]. Available: <https://www.lawfareblog.com/great-power-competition-and-ai-revolution-range-risks-military-and-strategic-stability>. [Accessed: 03-Oct-2017].
- [78] E. Kania, "数字化 – 网络化 – 智能化: China's Quest for an AI Revolution in Warfare," The Strategy Bridge, 08-Jun-2017..
- [79] V. Putin, "Открытый урок «Россия, устремлённая в будущее» | Open Lesson 'Russia moving towards the future,'" Президент России | President of the Russian Federation, 01-Sep-2017. [Online]. Available: <http://kremlin.ru/events/president/news/55493>. [Accessed: 16-Oct-2017].
- [80] S. Armstrong, N. Bostrom, and C. Shulman, "Racing to the precipice: a model of artificial intelligence development," Future of Humanity Institute, Technical Report 2013-1, 2013.
- [81] M. C. Elish, "Moral Crumple Zones: Cautionary Tales in Human-Robot Interaction (We Robot 2016)," Social Science Research Network, Rochester, NY, SSRN Scholarly Paper ID 2757236, Mar. 2016.
- [82] R. Andorno, "The Precautionary Principle: A New Legal Standard for a Technological Age," *J. Int. Biotechnol. Law*, vol. 1, pp. 11–19, 2004.