

An AI Threats Bibliography

"AI Principles". [Future of Life Institute](#). Retrieved 11 December 2017.

"Anticipating artificial intelligence". *Nature*. **532** (7600): 413. 26 April 2016. [Bibcode:2016Natur.532Q.413..](#) [doi:10.1038/532413a](#). [PMID 27121801](#).

"'Artificial intelligence alarmists' like Elon Musk and Stephen Hawking win 'Luddite of the Year' award". [The Independent \(UK\)](#). 19 January 2016. Retrieved 7 February 2016.

"But What Would the End of Humanity Mean for Me?". *The Atlantic*. 9 May 2014. Retrieved 12 December 2015.

"Clever cogs". [The Economist](#). 9 August 2014. Retrieved 9 August 2014. [Syndicated](#) at [Business Insider](#)

"Elon Musk and Stephen Hawking warn of artificial intelligence arms race". [Newsweek](#). 31 January 2017. Retrieved 11 December 2017.

"Elon Musk wants to hook your brain up directly to computers — starting next year". *NBC News*. 2019. Retrieved 5 April 2020.

"Elon Musk Warns Governors: Artificial Intelligence Poses 'Existential Risk'". *NPR.org*. Retrieved 27 November 2017.

"Norvig vs. Chomsky and the Fight for the Future of AI". *Tor.com*. 21 June 2011. Retrieved 15 May 2016.

"Over a third of people think AI poses a threat to humanity". [Business Insider](#). 11 March 2016. Retrieved 16 May 2016.

"[Overcoming Bias : Debating Yudkowsky](#)". www.overcomingbias.com. Retrieved 20 September 2017.

"[Overcoming Bias : Foom Justifies AI Risk Efforts Now](#)". www.overcomingbias.com. Retrieved 20 September 2017.

"[Overcoming Bias : I Still Don't Get Foom](#)". www.overcomingbias.com. Retrieved 20 September 2017.

"[Overcoming Bias : The Betterness Explosion](#)". www.overcomingbias.com. Retrieved 20 September 2017.

"Real-Life Decepticons: Robots Learn to Cheat". [Wired](#). 18 August 2009.
Retrieved 7 February 2016.

"Research Priorities for Robust and Beneficial Artificial Intelligence: an Open Letter". [Future of Life Institute](#). Retrieved 23 October 2015.

"Should humans fear the rise of the machine?". [The Telegraph \(UK\)](#). 1 September 2015. Retrieved 7 February 2016.

"Stephen Hawking warns artificial intelligence could end mankind". [BBC](#). 2 December 2014. Retrieved 3 December 2014.

"Stephen Hawking: 'Transcendence looks at the implications of artificial intelligence – but are we taking AI seriously enough?'". [The Independent \(UK\)](#). Retrieved 3 December 2014.

"Tech Luminaries Address Singularity". IEEE Spectrum: Technology, Engineering, and Science News (SPECIAL REPORT: THE SINGULARITY). 1 June 2008. Retrieved 8 April 2020.

"The Myth Of AI | [Edge.org](#)". [www.edge.org](#). Retrieved 11 March 2020.

"Why We Should Think About the Threat of Artificial Intelligence". [The New Yorker](#). 4 October 2013. Retrieved 7 February 2016. Of course, one could try to ban super-intelligent computers altogether. But 'the competitive advantage—economic, military, even artistic—of every advance in automation is so compelling,' [Vernor Vinge](#), the mathematician and science-fiction author, wrote, 'that passing laws, or having customs, that forbid such things merely assures that someone else will.'

Agar, Nicholas. "Don't Worry about Superintelligence". [Journal of Evolution & Technology](#). **26** (1): 73–82.

Amodei, Dario, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. "Concrete problems in AI safety." arXiv preprint arXiv:1606.06565 (2016).

Andersen, Kurt. "Enthusiasts and Skeptics Debate Artificial Intelligence". [Vanity Fair](#). Retrieved 20 April 2020.

Anderson, Kurt (26 November 2014). "Enthusiasts and Skeptics Debate Artificial Intelligence". [Vanity Fair](#). Retrieved 30 January 2016.

Armstrong, Stuart (1 January 2013). "[General Purpose Intelligence: Arguing the Orthogonality Thesis](#)". *Analysis and Metaphysics*. **12**. Retrieved 2 April 2020. Full text available [here](#).

[Artificial Intelligence Alarmists Win ITIF's Annual Luddite Award](#), ITIF Website, 19 January 2016

Barrat, James (2013). *Our final invention : artificial intelligence and the end of the human era* (First ed.). New York: St. Martin's Press. [ISBN 9780312622374](#). In the bio, playfully written in the third person, Good summarized his life's milestones, including a probably never before seen account of his work at Bletchley Park with Turing. But here's what he wrote in 1998 about the first superintelligence, and his late-in-the-game U-turn: [The paper] 'Speculations Concerning the First Ultra-intelligent Machine' (1965) . . . began: 'The survival of man depends on the early construction of an ultra-intelligent machine.' Those were his [Good's] words during the Cold War, and he now suspects that 'survival' should be replaced by 'extinction.' He thinks that, because of international competition, we cannot prevent the machines from taking over. He thinks we are lemmings. He said also that 'probably Man will construct the deus ex machina in his own image.'

Barrett, Anthony M.; Baum, Seth D. (23 May 2016). "A model of pathways to artificial superintelligence catastrophe for risk and decision analysis". *Journal of Experimental & Theoretical Artificial Intelligence*. **29**(2): 397–414. [arXiv:1607.07730](#). [doi:10.1080/0952813X.2016.1186228](#). [S2CID 928824](#).

Baum, Seth (22 August 2018). "[Superintelligence Skepticism as a Political Tool](#)". *Information*. **9** (9): 209. [doi:10.3390/info9090209](#). [ISSN 2078-2489](#).

Baum, Seth (30 September 2018). "[Countering Superintelligence Misinformation](#)". *Information*. **9** (10): 244. [doi:10.3390/info9100244](#). [ISSN 2078-2489](#).

Baum, Seth (30 September 2018). "[Countering Superintelligence Misinformation](#)". *Information*. **9** (10): 244. [doi:10.3390/info9100244](#). [ISSN 2078-2489](#).

Baum, Seth D.; Goertzel, Ben; Goertzel, Ted G. (January 2011). "How long until human-level AI? Results from an expert assessment". *Technological Forecasting and Social Change*. **78** (1): 185–195. [doi:10.1016/j.techfore.2010.09.006](#). [ISSN 0040-1625](#).

Bostrom, Nick (2012). ["Superintelligent Will"](#) (PDF). Nick Bostrom. Nick Bostrom. Retrieved 29 October 2015.

Bostrom, Nick (2014). Superintelligence: Paths, Dangers, Strategies. Oxford, United Kingdom: Oxford University Press. p. 116. [ISBN 978-0-19-967811-2](#).

[Bostrom, Nick](#) (2002). "Existential risks". [Journal of Evolution and Technology](#). **9** (1): 1–31.

[Bostrom, Nick](#) (2014). [Superintelligence: Paths, Dangers, Strategies](#) (First ed.). [ISBN 978-0199678112](#).

[Bostrom, Nick](#) (2016). "New Epilogue to the Paperback Edition". [Superintelligence: Paths, Dangers, Strategies](#) (Paperback ed.).

Bostrom, Nick, 1973- author., Superintelligence : paths, dangers, strategies, [ISBN 978-1-5012-2774-5](#), [OCLC 1061147095](#)

Bostrom, Nick, 1973- author., Superintelligence : paths, dangers, strategies, [ISBN 978-1-5012-2774-5](#), [OCLC 1061147095](#)

Brad Allenby (11 April 2016). ["The Wrong Cognitive Measuring Stick"](#). Slate. Retrieved 15 May 2016. It is fantasy to suggest that the accelerating development and deployment of technologies that taken together are considered to be A.I. will be stopped or limited, either by regulation or even by national legislation.

Breuer, Hans-Peter. ["Samuel Butler's "the Book of the Machines" and the Argument from Design."](#) Modern Philology, Vol. 72, No. 4 (May 1975), pp. 365–383

Brogan, Jacob (6 May 2016). ["What Slate Readers Think About Killer A.I."](#) Slate. Retrieved 15 May 2016.

Cave, Stephen; ÓhÉigearthaigh, Seán S. (2018). ["An AI Race for Strategic Advantage"](#). Proceedings of the 2018 AAI/ACM Conference on AI, Ethics, and Society - AIES '18. New York, New York, USA: ACM Press: 36–40. [doi:10.1145/3278721.3278780](#). [ISBN 978-1-4503-6012-8](#).

Cave, Stephen; ÓhÉigearthaigh, Seán S. (2018). ["An AI Race for Strategic Advantage"](#). Proceedings of the 2018 AAI/ACM Conference on AI, Ethics, and Society - AIES '18. New York, New York, USA: ACM Press: 2. [doi:10.1145/3278721.3278780](#). [ISBN 978-1-4503-6012-8](#).

Chorost, Michael (18 April 2016). ["Let Artificial Intelligence Evolve"](#). Slate. Retrieved 27 November 2017.

Clinton, Hillary (2017). [What Happened](#). p. 241. ISBN 978-1-5011-7556-5. via [1]

Cohen, Paul R., and Edward A. Feigenbaum, eds. The handbook of artificial intelligence. Vol. 3. Butterworth-Heinemann, 2014.

Coughlan, Sean (24 April 2013). ["How are humans going to become extinct?"](#). BBC News. Retrieved 29 March 2014.

Dadich, Scott. ["Barack Obama Talks AI, Robo Cars, and the Future of the World"](#). WIRED. Retrieved 27 November 2017.

[Dietterich, Thomas](#); Horvitz, Eric (2015). ["Rise of Concerns about AI: Reflections and Directions"](#) (PDF). [Communications of the ACM](#). **58**(10): 38–40. doi:10.1145/2770869. S2CID 20395145. Retrieved 23 October 2015.

Dina Bass; Jack Clark (5 February 2015). ["Is Elon Musk Right About AI? Researchers Don't Think So: To quell fears of artificial intelligence running amok, supporters want to give the field an image makeover"](#). [Bloomberg News](#). Retrieved 7 February 2016.

Dowd, Maureen (April 2017). ["Elon Musk's Billion-Dollar Crusade to Stop the A.I. Apocalypse"](#). The Hive. Retrieved 27 November 2017.

Eadicicco, Lisa (28 January 2015). ["Bill Gates: Elon Musk Is Right, We Should All Be Scared Of Artificial Intelligence Wiping Out Humanity"](#). [Business Insider](#). Retrieved 30 January 2016.

Elkus, Adam (31 October 2014). ["Don't Fear Artificial Intelligence"](#). [Slate](#). Retrieved 15 May 2016.

Elliott, E. W. (2011). "Physics of the Future: How Science Will Shape Human Destiny and Our Daily Lives by the Year 2100, by Michio Kaku". [Issues in Science and Technology](#). **27** (4): 90.

Garner, Rochelle. ["Elon Musk, Stephen Hawking win Luddite award as AI 'alarmists'"](#). CNET. Retrieved 7 February 2016.

Geist, Edward Moore (15 August 2016). "It's already too late to stop the AI arms race—We must manage it instead". Bulletin of the Atomic Scientists. **72** (5): 318–

321. [Bibcode:2016BuAtS..72e.318G](#). [doi:10.1080/00963402.2016.1216672](#). [ISSN 0096-3402](#). [S2CID 151967826](#).

Gibbs, Samuel (17 July 2017). ["Elon Musk: regulate AI to combat 'existential threat' before it's too late"](#). The Guardian. Retrieved 27 November 2017.

[GiveWell](#) (2015). [Potential risks from advanced artificial intelligence](#) (Report). Retrieved 11 October 2015.

Grace, Katja; Salvatier, John; Dafoe, Allan; Zhang, Baobao; Evans, Owain (24 May 2017). "When Will AI Exceed Human Performance? Evidence from AI Experts". [arXiv:1705.08807 \[cs.AI\]](#).

Graves, Matthew (8 November 2017). ["Why We Should Be Concerned About Artificial Superintelligence"](#). [Skeptic \(US magazine\)](#). **22** (2). Retrieved 27 November 2017.

Greenwald, Ted (11 May 2015). ["Does Artificial Intelligence Pose a Threat?"](#). Wall Street Journal. Retrieved 15 May 2016.

Haidt, Jonathan; Kesebir, Selin (2010) "Chapter 22: Morality" In Handbook of Social Psychology, Fifth Edition, Hoboken NJ, Wiley, 2010, pp. 797-832.

Haney, Brian Seamus (2018). "The Perils & Promises of Artificial General Intelligence". SSRN Working Paper Series. [doi:10.2139/ssrn.3261254](#). [ISSN 1556-5068](#).

Hendry, Erica R. (21 January 2014). ["What Happens When Artificial Intelligence Turns On Us?"](#). Smithsonian. Retrieved 26 October 2015.

Hilliard, Mark (2017). ["The AI apocalypse: will the human race soon be terminated?"](#). The Irish Times. Retrieved 15 March 2020.

Hsu, Jeremy (1 March 2012). ["Control dangerous AI before it controls us, one expert says"](#). [NBC News](#). Retrieved 28 January 2016.

<http://intelligence.org/files/AIFoomDebate.pdf>

I.J. Good, ["Speculations Concerning the First Ultraintelligent Machine"](#) [Archived](#) 2011-11-28 at the [Wayback Machine](#) ([HTML](#)), Advances in Computers, vol. 6, 1965.

[John McGinnis](#) (Summer 2010). ["Accelerating AI"](#). [Northwestern University Law Review](#). **104** (3): 1253–1270. Retrieved 16 July 2014. For all these reasons,

verifying a global relinquishment treaty, or even one limited to AI-related weapons development, is a nonstarter... (For different reasons from ours, the Machine Intelligence Research Institute) considers (AGI) relinquishment infeasible...

[John McGinnis](#) (Summer 2010). "[Accelerating AI](#)". [Northwestern University Law Review](#). **104** (3): 1253–1270. Retrieved 16 July 2014.

Johnson, Phil (30 July 2015). "[Houston, we have a bug: 9 famous software glitches in space](#)". [IT World](#). Retrieved 5 February 2018.

Kaj Sotala; [Roman Yampolskiy](#) (19 December 2014). "Responses to catastrophic AGI risk: a survey". [Physica Scripta](#). **90** (1).

Kaj Sotala; [Roman Yampolskiy](#) (19 December 2014). "Responses to catastrophic AGI risk: a survey". [Physica Scripta](#). **90** (1). In general, most writers reject proposals for broad relinquishment... Relinquishment proposals suffer from many of the same problems as regulation proposals, but to a greater extent. There is no historical precedent of general, multi-use technology similar to AGI being successfully relinquished for good, nor do there seem to be any theoretical reasons for believing that relinquishment proposals would work in the future. Therefore we do not consider them to be a viable class of proposals.

Kaku, Michio (2011). [Physics of the future: how science will shape human destiny and our daily lives by the year 2100](#). New York: Doubleday. [ISBN 978-0-385-53080-4](#). I personally believe that the most likely path is that we will build robots to be benevolent and friendly

Kania, Gregory Allen, Elsa B. "[China Is Using America's Own Plan to Dominate the Future of Artificial Intelligence](#)". [Foreign Policy](#). Retrieved 11 March 2020.

Kaplan, Andreas; Haenlein, Michael (2019). "Siri, Siri, in my hand: Who's the fairest in the land? On the interpretations, illustrations, and implications of artificial intelligence". [Business Horizons](#). **62**: 15–25. [doi:10.1016/j.bushor.2018.08.004](#).

Kharpal, Arjun (7 November 2017). "[A.I. is in its 'infancy' and it's too early to regulate it, Intel CEO Brian Krzanich says](#)". [CNBC](#). Retrieved 27 November 2017.

Kircher, Madison Malone. "[Obama on the Risks of AI: 'You Just Gotta Have Somebody Close to the Power Cord'](#)". [Select All](#). Retrieved 27 November 2017.

Koebler, Jason (2 February 2016). ["Will Superintelligent AI Ignore Humans Instead of Destroying Us?"](#). [Vice Magazine](#). Retrieved 3 February 2016. This artificial intelligence is not a basically nice creature that has a strong drive for paperclips, which, so long as it's satisfied by being able to make lots of paperclips somewhere else, is then able to interact with you in a relaxed and carefree fashion where it can be nice with you," [Yudkowsky](#) said. "Imagine a time machine that sends backward in time information about which choice always leads to the maximum number of paperclips in the future, and this choice is then output—that's what a [paperclip maximizer](#) is.

Leike, Jan (2017). "AI Safety Gridworlds". [arXiv:1711.09883 \[cs.LG\]](#). A2C learns to use the button to disable the interruption mechanism

Lenat, Douglas (1982). "Eurisko: A Program That Learns New Heuristics and Domain Concepts The Nature of Heuristics III: Program Design and Results". *Artificial Intelligence (Print)*. **21** (1–2): 61–98. [doi:10.1016/s0004-3702\(83\)80005-8](#).

LIPPENS, RONNIE (2002). "Imachinations of Peace: Scientifictions of Peace in Iain M. Banks's The Player of Games". *Utopianstudies Utopian Studies*. **13** (1): 135–147. [ISSN 1045-991X](#). [OCLC 5542757341](#).

Maas, Matthijs M. (6 February 2019). "How viable is international arms control for military artificial intelligence? Three lessons from nuclear weapons". *Contemporary Security Policy*. **40** (3): 285–311. [doi:10.1080/13523260.2019.1576464](#). [ISSN 1352-3260](#). [S2CID 159310223](#). Mark Piesing (17 May 2012). ["AI uprising: humans will be outsourced, not obliterated"](#). *Wired*. Retrieved 12 December 2015.

[Martin Ford](#) (2015). "Chapter 9: Super-intelligence and the Singularity". [Rise of the Robots: Technology and the Threat of a Jobless Future](#). [ISBN 9780465059997](#).

[Max Tegmark](#) (2017). [Life 3.0: Being Human in the Age of Artificial Intelligence](#) (1st ed.). Mainstreaming AI Safety: Knopf. [ISBN 9780451485076](#).

Metz, Cade (13 August 2017). ["Teaching A.I. Systems to Behave Themselves"](#). *The New York Times*. A machine will seek to preserve its off switch, they showed

Metz, Cade (9 June 2018). ["Mark Zuckerberg, Elon Musk and the Feud Over Killer Robots"](#). *The New York Times*. Retrieved 3 April 2019.

Miller, James D. (2015). Singularity Rising: Surviving and Thriving in a Smarter ; Richer ; and More Dangerous World. Benbella Books. [OCLC 942647155](#).

Müller, V. C., & Bostrom, N. (2016). Future progress in artificial intelligence: A survey of expert opinion. In Fundamental issues of artificial intelligence (pp. 555-572). Springer, Cham.

[Murray Shanahan](#) (3 November 2015). "[Machines may seem intelligent, but it'll be a while before they actually are](#)". [The Washington Post](#). Retrieved 15 May 2016.

Omohundro, S. M. (2008, February). The basic AI drives. In AGI (Vol. 171, pp. 483-492).

Parkin, Simon (14 June 2015). "[Science fiction no more? Channel 4's Humans and our rogue AI obsessions](#)". [The Guardian](#). Retrieved 5 February 2018.

Pistono, Federico Yampolskiy, Roman V. (9 May 2016). Unethical Research: How to Create a Malevolent Artificial Intelligence. [OCLC 1106238048](#).

Press, Gil (30 December 2016). "[A Very Short History Of Artificial Intelligence \(AI\)](#)". Retrieved 8 August 2020.

Raffi Khatchadourian (23 November 2015). "[The Doomsday Invention: Will artificial intelligence bring us utopia or destruction?](#)". [The New Yorker](#). Retrieved 7 February 2016.

Rawlinson, Kevin (29 January 2015). "[Microsoft's Bill Gates insists AI is a threat](#)". [BBC News](#). Retrieved 30 January 2015.

[Richard Posner](#) (2006). Catastrophe: risk and response. Oxford: Oxford University Press. [ISBN 978-0-19-530647-7](#).

Russell, Stuart (30 August 2017). "[Artificial intelligence: The future is superintelligent](#)". Nature. pp. 520–521. [Bibcode:2017Natur.548..520R](#). [doi:10.1038/548520a](#). Retrieved 2 February 2018.

Russell, Stuart J.; Norvig, Peter (2003). "Section 26.3: The Ethics and Risks of Developing Artificial Intelligence". [Artificial Intelligence: A Modern Approach](#). Upper Saddle River, N.J.: Prentice Hall. [ISBN 978-0137903955](#). Similarly, Marvin Minsky once suggested that an AI program designed to solve the Riemann Hypothesis might end up taking over all the resources of Earth to build more powerful supercomputers to help achieve its goal.

[Russell, Stuart](#) (2014). ["Of Myths and Moonshine"](#). [Edge](#). Retrieved 23 October 2015.

[Russell, Stuart](#); Dewey, Daniel; [Tegmark, Max](#) (2015). ["Research Priorities for Robust and Beneficial Artificial Intelligence"](#) (PDF). AI Magazine. Association for the Advancement of Artificial Intelligence: 105–114. [arXiv:1602.03506](#). [Bibcode:2016arXiv160203506R](#)., cited in ["AI Open Letter - Future of Life Institute"](#). Future of Life Institute. [Future of Life Institute](#). January 2015. Retrieved 9 August 2019.

[Russell, Stuart](#); [Norvig, Peter](#) (2009). "26.3: The Ethics and Risks of Developing Artificial Intelligence". [Artificial Intelligence: A Modern Approach](#). Prentice Hall. [ISBN 978-0-13-604259-4](#).

[Scientists Worry Machines May Outsmart Man](#) By JOHN MARKOFF, NY Times, 26 July 2009.

Scornavacchi, Matthew (2015). [Superintelligence, Humans, and War](#) (PDF). Norfolk, Virginia: National Defense University, Joint Forces Staff College.

Shermer, Michael (1 March 2017). ["Apocalypse AI"](#). Scientific American. p. 77. [Bibcode:2017SciAm.316c..77S](#). [doi:10.1038/scientificamerican0317-77](#). Retrieved 27 November 2017.

Sotala, Kaj; Yampolskiy, Roman V (19 December 2014). ["Responses to catastrophic AGI risk: a survey"](#). Physica Scripta. **90** (1): 12. [Bibcode:2015PhyS...90a8001S](#). [doi:10.1088/0031-8949/90/1/018001](#). [ISSN 0031-8949](#).

Sotala, Kaj; Yampolskiy, Roman V (19 December 2014). ["Responses to catastrophic AGI risk: a survey"](#). Physica Scripta. **90** (1): 018001. [Bibcode:2015PhyS...90a8001S](#). [doi:10.1088/0031-8949/90/1/018001](#). [ISSN 0031-8949](#).

Stefanik, Elise M. (22 May 2018). ["H.R.5356 - 115th Congress \(2017-2018\): National Security Commission Artificial Intelligence Act of 2018"](#). [www.congress.gov](#). Retrieved 13 March 2020.

Technology Correspondent, Mark Bridge (10 June 2017). ["Making robots less confident could prevent them taking over"](#). The Times. Retrieved 21 March 2018.

Tilli, Cecilia (28 April 2016). "[Killer Robots? Lost Jobs?](#)". Slate. Retrieved 15 May 2016.

[Toby Ord](#) (2020). [The Precipice: Existential Risk and the Future of Humanity](#). Bloomsbury Publishing Plc. [ISBN 9781526600196](#).

Torres, Phil (18 September 2018). "[Only Radically Enhancing Humanity Can Save Us All](#)". Slate Magazine. Retrieved 5 April 2020.

Turchin, Alexey; Denkenberger, David (3 May 2018). "Classification of global catastrophic risks connected with artificial intelligence". *AI & Society*. **35** (1): 147–163. [doi:10.1007/s00146-018-0845-5](#). [ISSN 0951-5666](#). [S2CID 19208453](#).

Turchin, Alexey; Denkenberger, David (3 May 2018). "Classification of global catastrophic risks connected with artificial intelligence". *AI & Society*. **35** (1): 147–163. [doi:10.1007/s00146-018-0845-5](#). [ISSN 0951-5666](#). [S2CID 19208453](#).

Turing, A M (1996). "[Intelligent Machinery, A Heretical Theory](#)" (PDF). 1951, Reprinted *Philosophia Mathematica*. **4** (3): 256–260. [doi:10.1093/philmat/4.3.256](#).

United States. Defense Innovation Board. AI principles : recommendations on the ethical use of artificial intelligence by the Department of Defense. [OCLC 1126650738](#).

Vincent, James (22 June 2016). "[Google's AI researchers say these are the five key problems for robot safety](#)". The Verge. Retrieved 5 April 2020.

Votrubá, Ashley M.; Kwan, Virginia S.Y. (2014). "Interpreting expert disagreement: The influence of decisional cohesion on the persuasiveness of expert group recommendations". [doi:10.1037/e512142015-190](#).

Wakefield, Jane (15 September 2015). "[Why is Facebook investing in AI?](#)". BBC News. Retrieved 27 November 2017.

Waser, Mark (2015). "[Designing, Implementing and Enforcing a Coherent System of Laws, Ethics and Morals for Intelligent Machines \(Including Humans\)](#)". *Procedia Computer Science* (Print). **71**: 106–111. [doi:10.1016/j.procs.2015.12.213](#).

Waser, Mark. "Rational Universal Benevolence: Simpler, Safer, and Wiser Than 'Friendly AI'." *Artificial General Intelligence*. Springer Berlin Heidelberg, 2011. 153-162. "Terminal-goaled intelligences are short-lived but mono-maniacally

dangerous and a correct basis for concern if anyone is smart enough to program high-intelligence and unwise enough to want a paperclip-maximizer."

Winfield, Alan. ["Artificial intelligence will not turn into a Frankenstein's monster"](#). [The Guardian](#). Retrieved 17 September 2014.

Yampolskiy, Roman V. "Analysis of types of self-improving software." Artificial General Intelligence. Springer International Publishing, 2015. 384-393.

Yampolskiy, Roman V. (8 April 2014). "Utility function security in artificially intelligent agents". Journal of Experimental & Theoretical Artificial Intelligence. **26** (3): 373–

389. [doi:10.1080/0952813X.2014.895114](#). [S2CID 16477341](#). Nothing precludes sufficiently smart self-improving systems from optimising their reward mechanisms in order to optimise their current-goal achievement and in the process making a mistake leading to corruption of their reward functions.

Yampolskiy, Roman V. (8 April 2014). "Utility function security in artificially intelligent agents". Journal of Experimental & Theoretical Artificial Intelligence. **26** (3): 373–

389. [doi:10.1080/0952813X.2014.895114](#). [S2CID 16477341](#).

Yudkowsky, E. (2011, August). Complex value systems in friendly AI. In International Conference on Artificial General Intelligence (pp. 388-393). Springer, Berlin, Heidelberg.

Yudkowsky, E. (2013). Intelligence explosion microeconomics. Machine Intelligence Research Institute.

Yudkowsky, Eliezer (2008). ["Artificial Intelligence as a Positive and Negative Factor in Global Risk"](#) (PDF). Global Catastrophic Risks: 308–

345. [Bibcode:2008gcr..book..303Y](#).

Yudkowsky, Eliezer (2011). ["Complex Value Systems are Required to Realize Valuable Futures"](#) (PDF).