

Global Catastrophic Risk INSTITUTE

March 4, 2022

RFI Response: National Artificial Intelligence Research and Development Strategic Plan— White House Office of Science and Technology Policy
87 FR 5876; Document Number 2022-02161

Dr. Alondra Nelson, Deputy Director of Science and Society of the Office of Science and Technology Policy (OSTP) and Performing the Duties of OSTP Director, the National Science and Technology Council's (NSTC) Select Committee on Artificial Intelligence (Select Committee), the NSTC Machine Learning and AI Subcommittee (MLAI-SC), the National AI Initiative Office (NAIIO), and the Networking and Information Technology Research and Development (NITRD) National Coordination Office (NCO):

Thank you for the invitation to submit comments in response to the Request For Information (RFI) to the National Artificial Intelligence Research and Development Strategic Plan. We, the Global Catastrophic Risk Institute (GCRI), are researchers with expertise on AI ethics and AI governance. We offer the following submission for your consideration.

We support the eight strategic aims described in the 2019 Update. As detailed below, we encourage specific changes to seven of the aims that support a more robust and inclusive approach to AI ethics and governance. These changes seek to ensure that the National AI R&D Strategic Plan supports broad and sustainable success and meets high social and ethical standards.

Our recommendations are as follows:

Strategy 1: Make long-term investments in AI research.

We recommend an emphasis on interdisciplinary research in which technological progress is oriented according to social and ethical values and in which technology governance is informed by a sound understanding of the nature of AI technology.

Investment in basic research is essential for improving the capability and reliability of AI systems. However, there has been a tendency for some research on AI systems to be focused on improving capabilities without substantial consideration of social and ethical

dimensions.¹ This risks the development of AI technology that can have inappropriate impacts on society. For example, machine learning research has often pursued larger neural network models to improve model accuracy without regard for the adverse energy resource and climate change consequences of larger models.² To align research on AI systems with social and ethical values, it is essential to have these values built into the core of the research, including the selection of which research directions to pursue.³

The ability of society to successfully govern AI technology additionally depends on basic research. AI technology has only recently risen to prominence as a societal issue, and the study of AI governance is likewise at an early stage. Research on AI governance has often focused on general concepts such as ethical principles, with less regard for how to operationalize them.⁴ Given rapid ongoing changes in AI technology, an important challenge is to ensure that governance concepts are informed by a state-of-the-art understanding of the technology.⁵ To formulate practical and technologically sound AI governance concepts, it is essential to invest in AI governance research in which the AI technology is not treated as a black box but instead is considered in detail.

Strategy 2: Develop effective methods for human-AI collaboration.

We recommend an emphasis on methods to ensure that human-AI collaboration is in the common public good and not just in the interests of the select few with control of advanced AI tools.

AI technology has enabled major advances in workplace productivity, bringing major economic benefits. However, this often comes at the expense of disadvantaged populations. For example, low-wage workers are forced to work unpredictable schedules that are optimized according to AI processing of last-minute data⁶ and social

¹ Seth D. Baum, "On the promotion of safe and socially beneficial artificial intelligence", *AI & Society*, vol. 32, no. 4 (2017), pp. 543-551, <https://doi.org/10.1007/s00146-016-0677-0>.

² Lynn Kaack, Priya Donti, Emma Strubell, George Kamiya, Felix Creutzig, and David Rolnick, "Aligning artificial intelligence with climate change mitigation", 2021, <https://hal.archives-ouvertes.fr/hal-03368037>.

³ Steven Umbrello and Ibo Van de Poel, "Mapping value sensitive design onto AI for social good principles", *AI and Ethics*, vol. 1, no. 3 (2021), pp. 283-296, <https://doi.org/10.1007/s43681-021-00038-3>.

⁴ Jessica Morley, Libby Kinsey, Anat Elhalal, Francesca Garcia, Marta Ziosi, and Luciano Floridi, "Operationalising AI ethics: barriers, enablers and next steps", *AI & Society* (2021), <https://doi.org/10.1007/s00146-021-01308-8>.

⁵ Wendell Wallach and Gary Marchant, "Toward the Agile and Comprehensive International Governance of AI and Robotics", *Proceedings of the IEEE*, vol. 107, no. 3 (2019), pp. 505-508, <https://doi.org/10.1109/JPROC.2019.2899422>.

⁶ Cathy O'Neil, *Weapons of Math Destruction: How Big Data Increases Inequality and Threatens Democracy* (Crown, 2016).

media websites use AI recommender algorithms that increase user engagement by promoting politically extremist content.⁷

The 2016 National AI R&D Strategic Plan recognizes that “the challenge of understanding and designing human-AI ethics and value alignment into systems remains an open research area.” We agree that this is an important open research area. However, it is important to address: to whom specifically are AI systems aligned?⁸ AI systems that are designed to support a single individual require fundamentally different designs than those that are designed to support society as a whole.⁹ Only by including the common good of society in AI system design can human-AI collaboration advance the interests of all members of society instead of the select few.

Strategy 3: Understand and address the ethical, legal, and societal implications of AI.

We recommend a holistic and pluralistic approach to the ethical, legal, and societal (ELS) implications of AI, in particular, to address the full range of implications and to evaluate them in terms of a diversity of social and ethical perspectives.

It is vital that the ELS implications of AI are central to all work on AI and not tacked on as an afterthought. There are ELS implications in, among other things, (1) the selection of algorithms, determining which ethical concepts can be implemented in an AI system¹⁰; (2) the selection of the scale at which to implement the algorithm, determining the energy consumption of the AI system¹¹; (3) the selection of training data, determining which problems and issues an AI system can be applied to¹² and the

⁷ Joe Whittaker, Seán Looney, Alastair Reed, and Fabio Votta, "Recommender systems and the amplification of extremist content", *Internet Policy Review*, vol. 10, no. 2 (2021), <https://doi.org/10.14763/2021.2.1565>.

⁸ Seth D. Baum, "Social choice ethics in artificial intelligence", *AI & Society*, vol. 35, no. 1 (2020), pp. 165-176, <https://doi.org/10.1007/s00146-017-0760-1>.

⁹ Roger Lera-Leri, Filippo Bistaffa, Marc Serramia, Maite Lopez-Sanchez, and Juan Rodriguez-Aguilar, "Towards Pluralistic Value Alignment: Aggregating Value Systems through ℓ_p -Regression", Proceedings of the 21st International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2022, May 9–13), <https://filippobistaffa.github.io/papers/2022aamas.pdf>.

¹⁰ Suzanne Tolmeijer, Markus Kneer, Cristina Sarasua, Markus Christen, and Abraham Bernstein, "Implementations in machine ethics: A survey", *ACM Computing Surveys*, vol. 53, no. 6 (2020), article 132, <https://doi.org/10.1145/3419633>.

¹¹ Lynn Kaack, Priya Donti, Emma Strubell, George Kamiya, Felix Creutzig, and David Rolnick, "Aligning artificial intelligence with climate change mitigation", 2021, <https://hal.archives-ouvertes.fr/hal-03368037>.

¹² For example, an AI language processing system trained in a dominant language such as English will not function in other languages.

potential for biases in AI system outputs¹³; and (4) the deployment of AI systems, determining the specific societal and environmental impacts¹⁴. The breadth of ELS implications underscores the importance of embedding ELS at all points on the AI system lifecycle.

Work on ELS should additionally welcome a range of perspectives and support constructive and open debate. Recent work on AI ethics has emphasized the development of consensus-driven sets of ethics principles or guidelines. However, these guidelines come overwhelmingly from North American and European organizations, many of which had little or no public participation.¹⁵ Marginalized populations, such as Indigenous peoples, may have differing perspectives on ELS issues.¹⁶ Inclusion of diverse perspectives can help to overcome apparent gaps in existing ELS work, such as pertaining to the moral status of nonhumans.¹⁷ Furthermore, active AI ELS debates remain unresolved, such as on the relative importance of near-term and long-term dimensions of AI.¹⁸ Given the lack of universal consensus on AI ELS issues, it is important to support an inclusive and open-minded conversation about the issues.

Strategy 4: Ensure the safety and security of AI systems.

We recommend an emphasis on a dynamic research and development program to ensure that AI systems remain safe and secure as the technology and its usage evolve over time.

An essential aspect of AI safety and security is that as AI systems become more capable and become used more widely, the safety and security challenges become

¹³ Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan, "A survey on bias and fairness in machine learning", *ACM Computing Surveys*, vol. 54, no. 6 (2021), article 115, <https://doi.org/10.1145/3457607>.

¹⁴ For example, decisions of when to deploy autonomous weapons; Hendrik Huelss, "Deciding on Appropriate Use of Force: Human-machine Interaction in Weapons Systems and Emerging Norms", *Global Policy* vol. 10, no. 3 (2019), pp. 354-358, <https://doi.org/10.1111/1758-5899.12692>.

¹⁵ Daniel Schiff, Jason Borenstein, Justin Biddle, and Kelly Laas, "AI ethics in the public, private, and NGO sectors: a review of a global document collection", *IEEE Transactions on Technology and Society*, vol. 2, no. 1 (2021), pp. 31-42, <https://doi.org/10.1109/TTS.2021.3052127>.

¹⁶ Jason Edward Lewis, Angie Abdilla, Noelani Arista, Kaipulaumakaniolono Baker, Scott Benesiinaabandan, Michelle Brown, Melanie Cheung et al., "Indigenous protocol and artificial intelligence position paper", Indigenous Protocol and Artificial Intelligence Working Group and the Canadian Institute for Advanced Research, Honolulu, HI (2020), <https://spectrum.library.concordia.ca/id/eprint/986506>.

¹⁷ Andrea Owe and Seth D. Baum, "Moral consideration of nonhumans in the ethics of artificial intelligence", *AI and Ethics*, vol. 1, no. 4 (2021), pp. 517-528, <https://doi.org/10.1007/s43681-021-00065-0>.

¹⁸ Charlotte Stix and Matthijs M. Maas, "Bridging the gap: the case for an 'Incompletely Theorized Agreement on AI policy'", *AI and Ethics*, vol. 1, no. 3 (2021), pp. 261-271, <https://doi.org/10.1007/s43681-020-00037-w>.

more important. Prior to the deep learning revolution, AI technology had limited usage¹⁹ and likewise little need for safety and security. By now, AI technology is used widely across economic sectors and other areas of human society. Barring a new AI winter, i.e. another years-long drop in AI research if AI progress fails to live up to expectations, AI technology will only grow in its importance. As it does, the potential for AI systems to cause harm is likely to increase. In risk terms, harm from AI systems could become more frequent, due to their wider usage, and more severe, due to their usage in more high-stakes applications. Prospects for catastrophic harm may further increase via the use of AI in crucial sectors such as agriculture²⁰ and via increasingly general-purpose AI systems that are becoming widely used across sectors.²¹

Given the growing stakes, it is vital for AI safety and security to keep up with the changing technology. This is not a trivial challenge. Some aspects of AI safety and security may remain viable even as the technology becomes more advanced,²² whereas other aspects may need to be customized for more advanced systems.²³ Research and development on AI safety and security must be forward-looking in order to keep society safe in the face of more capable and more widely deployed AI systems. Doing so will further be of economic and strategic benefit because safer and more secure AI technology would permit the technology to be deployed more widely, especially in more sensitive settings.

Strategy 5: Develop shared public datasets and environments for AI training and testing.

We have no comments on Strategy 5.

Strategy 6: Measure and evaluate AI technologies through standards and benchmarks.

¹⁹ Terrence J. Sejnowski, *The Deep Learning Revolution* (MIT Press, 2018).

²⁰ Victor, Galaz, Miguel A. Centeno, Peter W. Callahan, Amar Causevic, Thayer Patterson, Irina Brass, Seth Baum et al., "Artificial intelligence, systemic risks, and sustainability", *Technology in Society*, vol. 67 (2021), article 101741, <https://doi.org/10.1016/j.techsoc.2021.101741>.

²¹ Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein et al., "On the opportunities and risks of foundation models", Stanford Institute for Human-Centered Artificial Intelligence (2021), <https://arxiv.org/abs/2108.07258>.

²² Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané, "Concrete problems in AI safety", (2016), <https://arxiv.org/abs/1606.06565>.

²³ Stuart Russell, *Human Compatible: Artificial Intelligence and the Problem of Control* (Viking, 2019).

We recommend support for programs that can facilitate the ongoing development and adoption of standards across the AI industry, including via the development of regimes for the certification of AI systems.

A variety of AI standards and frameworks are in the process of being developed, including the National Institute of Standards and Technology (NIST) AI Risk Management Framework²⁴ and the Institute of Electrical and Electronics Engineers (IEEE) Standards Association P2863 Recommended Practice for Organizational Governance of Artificial Intelligence.²⁵ These standards and frameworks will provide voluntary guidance for developers and deployers of AI systems. As these and other frameworks are completed, two challenges will be faced. One is to ensure that the standards remain relevant and appropriate in the face of ongoing changes in AI technology and its various applications. The other is to facilitate the adoption of the standards by AI developers and deployers.

One valuable approach to facilitating adoption of standards is via certification regimes.²⁶ Certification regimes serve to address information asymmetries between insiders within an organization and outsiders who wish to know if the organization is complying with relevant standards. Certification is already widely used in many sectors, such as in the US EnergyStar program for consumer appliances, the ISO 9001 program for supply chains, and the LEED program for building design. AI certification regimes have been developed or proposed by, among others, the European Commission,²⁷ IEEE,²⁸ and the government of Malta.²⁹ When implemented effectively, certification regimes can incentivize adoption of standards. Certification regimes are also flexible in terms of being public or private (or both), voluntary or mandatory, and geographically local or global. These attributes make certification regimes an important element of AI standards adoption.

Strategy 7: Better understand the national AI R&D workforce needs.

²⁴ <https://www.nist.gov/itl/ai-risk-management-framework>

²⁵ <https://sagroups.ieee.org/2863>

²⁶ Peter Cihon, Moritz J. Kleinaltenkamp, Jonas Schuett, and Seth D. Baum, "AI certification: Advancing ethical practice by reducing information asymmetries", *IEEE Transactions on Technology and Society*, vol. 2, no. 4 (2021), pp. 200-209, <https://doi.org/10.1109/TTS.2021.3077595>.

²⁷ European Commission, "White paper on artificial intelligence: A European approach to excellence and trust," (2020), https://ec.europa.eu/info/publications/white-paper-artificial-intelligence-european-approach-excellence-and-trust_en.

²⁸ IEEE, "The ethics certification program for autonomous and intelligent systems (ECPAIS)", <https://standards.ieee.org/industry-connections/ecpais.html>.

²⁹ Malta Digital Innovation Authority, "AI-ITA blueprint guidelines", (2019), <https://mdia.gov.mt/wp-content/uploads/2019/10/AI-ITA-Blueprint-Guidelines-03OCT19.pdf>.

We recommend an emphasis on a multidisciplinary and multitalented workforce that is capable of developing and applying AI technology for the common good.

As AI technology grows in its ethical, legal, and societal (ELS) significance, it is vital for the computer scientists and engineers who design and deploy AI systems to be conversant in ELS topics. Historically, AI technology had little societal impact; AI computer scientists and engineers were likewise focused narrowly on how to increase the capabilities of AI systems with little regard for ELS implications.³⁰ The emergence of AI as a class of technology with significant ELS implications has created a need for cultivating an understanding of ELS among AI computer scientists and engineers. Some progress on this front has been made,³¹ but this remains a major weakness of the AI workforce.

Concurrently, there is a need for ELS experts who are conversant in AI technology. Computer scientists and engineers play a vital role in AI technology design and in relating design details to AI governance. However, it is inappropriate to expect computer scientists and engineers to have a comparable depth of knowledge about ELS as people who are professionally trained in ELS fields such as moral philosophy, law, and the social sciences. In order for ELS experts to successfully contribute to AI governance, it is vital for them to have some knowledge (the more the better) of how AI technology works. This is not typically taught in ELS university programs. It is therefore important to support dedicated programs. The Technology, Management, and Policy Consortium³² provides a valuable set of benchmarks for how such programs can be designed and executed.

Strategy 8: Expand public-private partnerships to accelerate advances in AI.

We recommend that public-private partnerships be pursued both domestically and internationally to ensure that the AI industry as a whole is oriented toward the common good.

In AI, as is the case in many industries, there is often tension between the private interest and the public good. For example, as noted above, companies running social media websites sometimes use AI recommender algorithms that increase user

³⁰ John Bohannon, "Fears of an AI pioneer", *Science*, vol. 349, no. 6245 (2015), p.252, <https://doi.org/10.1126/science.349.6245.252>.

³¹ For example, the AAAI/ACM Conference on AI, Ethics, and Society, <https://www.aies-conference.com>.

³² <https://tmpconsortium.org>

engagement by promoting politically extremist content³³ and AI developers sometimes build larger neural network models to improve model accuracy without regard for the adverse energy resource and climate change consequences of larger models.³⁴ Additionally, in the international context, competition between military adversaries can result in both sides pursuing unsafe AI technology.³⁵ These situations are known as collective action problems; they constitute an important class of challenges in AI governance.³⁶

Public-private partnerships can play a valuable role in addressing AI collective action problems. Domestically, these partnerships can help support private firms as they orient their activities toward the common good. Internationally, the partnerships can facilitate the cooperation needed to solve AI collective action problems at the global scale. The US already participates in international public-private partnerships for mutual benefit with other countries.³⁷ Additionally, some international organizations already work to bring together public and private AI stakeholders; these include the Global Partnership on Artificial Intelligence (GPAI)³⁸ and the OECD Artificial Intelligence Policy Observatory.³⁹ The US should participate in these organizations as part of its program on public-private partnerships. The US should further seek to include diverse participants from the international community, including to support global justice in AI technology.⁴⁰ Additionally, where appropriate, including rival powers such as China in these forums would further facilitate the resolution of AI collective action problems. If successful, these partnerships could support a “race to the top” dynamic in which competition between companies and countries advances benefits for the national and global common good.⁴¹

³³ Joe Whittaker, Seán Looney, Alastair Reed, and Fabio Votta, "Recommender systems and the amplification of extremist content", *Internet Policy Review*, vol. 10, no. 2 (2021), <https://doi.org/10.14763/2021.2.1565>.

³⁴ Lynn Kaack, Priya Donti, Emma Strubell, George Kamiya, Felix Creutzig, and David Rolnick, "Aligning artificial intelligence with climate change mitigation", 2021, <https://hal.archives-ouvertes.fr/hal-03368037>.

³⁵ Richard Danzig, "Managing Loss of Control as Many Militaries Pursue Technological Superiority", Center for New American Security (2018), <https://www.cnas.org/publications/reports/technology-roulette>.

³⁶ Robert de Neufville and Seth D. Baum, "Collective action on artificial intelligence: A primer and review", *Technology in Society*, vol. 66 (2021), article 101649, <https://doi.org/10.1016/j.techsoc.2021.101649>.

³⁷ For example, the US-Israel Binational Industrial Research and Development Foundation (BIRD), established in 1977, aims to provide funding to joint projects of mutual benefit to both the US and Israel; see <https://www.birdf.com>.

³⁸ <https://gpai.ai>

³⁹ <https://oecd.ai/en>

⁴⁰ Eugenio V. Garcia, "The International Governance of AI: Where is the Global South?" (2021), <https://www.researchgate.net/publication/348848134>.

⁴¹ Will Hunt, "The Flight to Safety-Critical AI: Lessons in AI Safety from the Aviation Industry", Center for Long-Term Cybersecurity, University of California, Berkeley (2020), <https://cltc.berkeley.edu/2020/08/11/new-report-the-flight-to-safety-critical-ai-lessons-in-ai-safety-from-the-aviation-industry>.