



JANUARY 2026

# Artificial General Intelligence and Global Risks

*Current Status and Future Risks*

GLOBAL RISK  
FEDERATION OF AMERICAN SCIENTISTS

## ABOUT THIS REPORT

---

This report summarizes the key findings, insights, and policy options from a June 2025 roundtable event on risks at the convergence of artificial intelligence (AI) and biology. The Federation of American Scientists, in partnership with the Future of Life Institute (FLI), brought together academic, industry, and government experts spanning AI, biology, and technology policy domains for this conversation. Experts identified bottlenecks in how we understand and identify threats at the AlxBio intersection (particularly deliberate misuse), took stock of advancing AI capabilities and their impact on bioweapons and accidental releases, and considered opportunities to avert risk in even a geopolitically heating world.

This report is structured in three parts: an executive summary, a detailed analysis of the findings, and three papers authored for participants in advance of the event. Hamza Chaudhry, AI and National Security Lead at FLI, authored the first pre-read on conceptualizing and framing risks in biology, AI, and the AlxBio convergence. Dr. Oliver Stephenson, Associate Director of AI and Emerging Technology Policy at FAS, authored the second pre-read on risk mitigation approaches in biology, AI, and the AlxBio convergence. Dr. Yong-Bee Lim, Associate Director of Global Risk at FAS, authored the final pre-read on how experts frame future risk in the bio, AI, and AlxBio domains, and the opportunities some of this framing provides us for greater safety and security.

### Global Risk Program at FAS

The Global Risk Program at the Federation of American Scientists (FAS) focuses on addressing and preventing the events and threats that could permanently cripple or destroy humanity. Some key areas our team focuses on include nuclear war, the next global pandemic, biological attack, and even a collision with a massive near-earth object. Our team of policy experts, scientists, and researchers use tools including forecasting, research and analysis, and expertise in key global risk domain areas to develop modern policy solutions for a rapidly advancing and complex time in humanity's development. Find out more at our website [www.fas.org/issue/global-risk](http://www.fas.org/issue/global-risk). The project is led by Jon B. Wolfsthal the Director of the Global Risk Program at FAS.

### Funding

This report and the associated workshop were made possible through the generous support of the Future of Life Institute and are part of a wider series in our ongoing “AlxGlobal Risk Nexus” project. This project will culminate in a global summit in Spring 2026. The views expressed in this report are those of the authors and do not necessarily reflect the positions of the funders or participants.

Special thanks to Yong-Bee Lim, PhD, Associate Director of the Global Risk Program, Oliver Stephenson, PhD, Associate Director of Artificial Intelligence and Emerging Technology Policy, Elliott Gunnell, M.Sc. Project Associate, Global Risk Program, and Abhay Katoch, Visiting Fellow with the Global Risk Program, for their contribution to the event and this report. Special thanks to our colleague on FAS’s Communications team, Kate Kohn, for developing the graphic for this report.

FAS can be reached at 1150 18th St. NW, Suite 1000, Washington, DC, 20036, [fas@fas.org](mailto:fas@fas.org), or through [fas.org](http://fas.org).

COPYRIGHT © FEDERATION OF AMERICAN SCIENTISTS, 2026. ALL RIGHTS RESERVED.

## CONTENTS

---

ABOUT THIS REPORT .....	I
EXECUTIVE SUMMARY .....	1
WHAT WE HEARD .....	3
MENU OF POLICY OPTIONS .....	8
CONCLUSION .....	10
PRE-READ PAPERS FROM ROUNDTABLE .....	11

## EXECUTIVE SUMMARY

---

On September 18, 2025, the Federation of American Scientists (FAS), in partnership with the Future of Life Institute (FLI), convened a forum in the United States Capitol to examine the national security and global risk implications of increasingly capable artificial intelligence (AI) systems, including the prospect of artificial general intelligence (AGI). The discussion focused on how advanced AI capabilities should be understood for policy purposes, what warning signs might indicate transitions toward more general-purpose and autonomous systems, and what “no-regrets” actions the U.S. government could take to prepare for risks related to misuse, proliferation, and potential loss-of-control scenarios. The convening brought together policymakers—including Representative Bill Foster (D, IL-11) and Representative Ted Lieu (D, CA-36)—alongside experts from RAND, the Carnegie Endowment for International Peace, and other national security, industry, and research institutions.

## FINDINGS

Participants broadly agreed that increasingly capable, general-purpose, and autonomous AI systems could pose significant economic, national security, and global stability risks, even if timelines remain uncertain, and the concept of “AGI” is contested. Rather than hinging on a single technological milestone, these risks emerge as AI systems acquire stronger reasoning abilities, greater autonomy, and the capacity to contribute meaningfully to research, planning, and decision-making. Without adequate preparation, such systems could accelerate the proliferation of weapons of mass destruction, enable destabilizing military and cyber activities, and introduce loss-of-control risks in which systems pursue misaligned objectives in ways that are difficult to detect or interrupt.

The discussion highlighted a persistent gap between the way AI progress is typically measured and the kinds of capabilities that matter most for policy and national security. Participants emphasized that existing benchmark-driven definitions of AI are poorly suited to assessing strategic risk. Instead, they argued for a capability-focused understanding of advanced AI—particularly systems that can automate substantial portions of research and development, operate autonomously over long horizons, or meaningfully shape military or strategic outcomes. While estimates of AGI timelines varied widely, several participants noted that expert forecasts have shortened over time, reinforcing the case for early, precautionary planning under uncertainty.

Geopolitical dynamics further complicate these risks. Panelists stressed that U.S. AI strategy cannot be considered independently of China, but cautioned against framing competition as a simple or symmetric “race” toward AGI. While U.S. firms often emphasize frontier breakthroughs, China’s current approach places greater weight on the diffusion of AI across existing industries through its “AI Plus” strategy. Participants warned that these differing trajectories—combined with limited transparency, growing reliance on critical AI supply chains, and the integration of AI into military and security systems—could increase the risk of miscalculation and escalation as capabilities advance.

Across panels, there was strong convergence around the importance of “no-regrets” actions: steps that improve preparedness, situational awareness, and institutional capacity without requiring confident predictions about AI timelines or outcomes. Participants repeatedly emphasized the need to strengthen government planning, coordination, and technical understanding to manage risks as they emerge.

## SUMMARY OF POLICY OPTIONS

Based on the panel discussions and subsequent analysis, FAS identified several policy options that could help address risks at the nexus of advanced AI and global security. These options are presented for discussion purposes and do not imply endorsement by any individual participant.

- **Develop capability-focused definitions and signposts for advanced AI** to support national security planning, emphasizing operational thresholds—such as the ability to automate R&D, operate autonomously

over long horizons, or contribute meaningfully to military or strategic decision-making—rather than static academic benchmarks.

- **Build government capacity for independent evaluation, testing, and safety standards** for advanced AI systems, including the ability to assess model capabilities, reliability, and failure modes, and to apply appropriate standards for high-risk national security deployments.
- **Strengthen national security planning, wargaming, and contingency development** related to advanced AI, including scenarios involving misuse by non-state actors, loss-of-control incidents, rapid diffusion into sensitive domains, and AI-driven escalation dynamics.
- **Establish structured, trusted information-sharing mechanisms with frontier AI developers** to improve situational awareness of emerging risks, enable incident reporting, and support clearer communication between industry and government.
- **Create a government-wide coordination hub for advanced AI risk**, modeled in part on the National Counterterrorism Center, to integrate intelligence, operational planning, and policy analysis, reduce fragmentation across agencies, and support cross-government preparedness and response.
- **Deepen international coordination with allies on advanced AI development and deployment**, particularly among countries central to the AI supply chain, to promote shared situational awareness, confidence-building measures, and aligned responses to emerging risks.

## WHAT WE HEARD

---

On September 18, 2025, the Federation of American Scientists (FAS), in partnership with the Future of Life Institute (FLI), convened a policy roundtable in the United States Capitol for approximately 80 participants drawn from government, national security institutions, academia, industry, and civil society. Held under the Chatham House Rule, the discussion examined how artificial general intelligence (AGI) might emerge, the economic and national security risks it could pose, and what “no-regrets” actions the U.S. government could take to prepare. Participants also explored the implications of AGI for U.S.–China strategic competition, including risks of miscalculation, diffusion, and escalation as AI capabilities advance.

This discussion comes at a pivotal time for AI development, as increasingly capable, autonomous, and general-purpose AI systems raise questions about national security and strategic stability. Policymakers have begun to grapple with these risks, but efforts remain fragmented and largely oriented toward today’s systems rather than plausible future capabilities—such as AI-enabled R&D, long-horizon planning, or autonomous decision-making at scale. Although the U.S. government has invested heavily in AI development and adoption, it has yet to develop robust contingency plans for scenarios in which future systems introduce systemic risks. As AI capabilities advance and diffuse, opportunities for anticipatory governance and risk reduction may narrow.

This section will provide more specific insights and observations from our panelists and participants. It will also provide a curated interpretation informed by our initial analysis of the event.

## DEFINING ARTIFICIAL GENERAL INTELLIGENCE

Participants noted that discussions of artificial general intelligence (AGI) are often constrained by confusion and inconsistency in how the term is defined. As an example, OpenAI states that its mission is to ensure that “artificial general intelligence—AI systems that are generally smarter than humans—benefits all of humanity.”<sup>1</sup> At the same time, OpenAI CEO Sam Altman has argued that AGI is “not a super useful term,” reflecting the fact that the label is used to describe a wide range of distinct capabilities and assumptions.<sup>2</sup> Panelists broadly agreed that this definitional ambiguity complicates efforts by policymakers to assess risks and design appropriate responses.

Participants emphasized that current frontier AI systems, while increasingly capable and more general-purpose, largely function as tools that augment human labor rather than replace it. These systems can assist with analysis, coding, and planning, and in some cases carry out limited autonomous actions, but they remain dependent on human direction and oversight. However, as AI systems gain stronger reasoning abilities and an increased capacity to use tools autonomously, they may begin to perform tasks that require sustained human judgment, coordination, and expertise. At that point, the relevant policy question is no longer whether a system performs well on individual benchmarks, but whether it crosses capability thresholds that enable qualitatively new forms of economic, military, or strategic impact.

Accordingly, attendees repeatedly emphasized the need to define AGI in terms of underlying capabilities rather than any single benchmark or test (for example, performance on the LSAT). Several speakers offered complementary capability-focused characterizations of AGI. One defined AGI as scalable general cognitive labor—“millions of programmers limited not by their speed or ability, but solely by compute”—highlighting its potential to transform productivity, research capacity, and state power. Another characterized AGI primarily by its strategic effects, describing sufficiently advanced systems as “deception machines” capable of generating a persistent informational “fog of war” for governments and militaries. A third framing focused on recursive capability growth,

---

<sup>1</sup> OpenAI. Planning for AGI and Beyond. 24 February 2023. <https://openai.com/index/planning-for-agi-and-beyond/>

<sup>2</sup> Browne, Ryan. Sam Altman now says AGI, or human-level AI, is ‘not a super useful term’ — and he’s not alone. CNBC. 11 August 2025. <https://www.cnbc.com/2025/08/11/sam-altman-says-agi-is-a-pointless-term-experts-agree.html>

defining AGI as a system able to “automate the process of automation” by substantially accelerating AI research and development itself.

Across these definitions, speakers emphasized reasoning ability as the critical upstream capability. Strong reasoning would allow AI systems to carry out long-horizon planning, adapt to novel problems, and contribute to the improvement of future AI systems with limited human oversight—potentially enabling forms of recursive self-improvement. Speakers noted that such capabilities are likely to emerge unevenly and piecemeal rather than all at once. Nonetheless, several cautioned that even partial advances in reasoning could have outsized consequences, including, as some warned, enabling non-experts to develop weapons of mass destruction, accelerating military competition through the creation of novel weapons systems, and allowing increasingly agentic systems to pursue harmful objectives autonomously. In contrast, attendees emphasized that many existing benchmarks of AI performance are poorly suited to measuring progress toward these policy-relevant capability thresholds.

## UNCERTAINTY AND AGI TIMELINES

The timing of when AI systems may meet various definitions of AGI remains deeply uncertain. Participants noted that forecasts of so-called “AGI timelines” have historically been unreliable, shaped by shifting definitions and what several speakers described as “wild expectations” that have characterized the field since its inception. As a result, some participants cautioned against treating any single timeline estimate as a firm basis for policy.

At the same time, several attendees observed that while specific predictions have repeatedly been wrong, the time horizons forecast by many AI experts have consistently shortened as capabilities have advanced.<sup>3</sup> Some, including OpenAI CEO Sam Altman, have suggested that continued increases in computational power and model complexity could lead to AGI within this decade.<sup>4</sup> More skeptical voices project timelines of ten to twenty years,<sup>5</sup> while others argue that AGI is either too incoherent a concept or too distant to warrant sustained policymaker attention at present.

Despite these disagreements, there was broad agreement that uncertainty itself strengthens the case for early planning. Participants emphasized that once AI systems begin operating with greater autonomy and generality, risks related to loss of control, proliferation, and strategic instability could emerge rapidly. Given these dynamics, many argued that policymakers should focus less on predicting precise timelines and more on developing “no-regrets” preparations that would remain valuable across a wide range of possible futures.

## DUAL-USE RISKS FROM ADVANCED AI CAPABILITIES

Participants emphasized that many of the most consequential abilities of more capable, autonomous, and general-purpose AI systems would be inherently dual-use. As systems become able to plan, coordinate, and execute complex tasks with less human input, they may expand access to forms of scientific and technical work currently constrained by time, labor, and specialized expertise. While this could generate substantial benefits for legitimate research and economic productivity, participants warned that the same capabilities could also be exploited by harmful actors or poorly governed deployments. Several stressed that because “everything exciting is dual-use,” effective risk reduction will require focusing on capability thresholds and real-world misuse pathways rather than static academic benchmarks.

3 Toner, Helen. “Long” timelines to advanced AI have gotten crazy short. *Rising Tide* (Substack). 01 April 2025. <https://helentoner.substack.com/p/long-timelines-to-advanced-ai-have>

4 POLITICO. When Sam Altman Predicts a ‘Superintelligence’ Might Arrive. 25 September 2025. <https://www.politico.com/news/magazine/2025/09/25/sam-altman-ai-interview-axel-springer-00580997>.

5 Marcus, Gary. “Mark my words: when AGI does actually come, perhaps 10 or 20 years from now, people will laugh at the idea that o3 was close to AGI.” X. 24 December 2024. <https://x.com/GaryMarcus/status/1871605871282999760>.

## Biotechnology

Participants identified biological risk as a central national security concern arising from increasingly capable and general-purpose AI systems. Several speakers warned that advances in AI reasoning and research capabilities could lower barriers to weapons of mass destruction, including bioweapons, by enabling non-experts to carry out work that has traditionally required extensive technical training. This risk was framed as a consequence of AI-driven acceleration of research and development rather than any single domain-specific tool. Attendees emphasized the inherently dual-use nature of these capabilities, citing drug discovery as a canonical example in which the same advances that support therapeutic development could also be misused for harmful purposes. Participants further noted the importance of monitoring for signs of increased bioweapons activity as AI capabilities advance, treating such indicators as an early warning signal of emerging proliferation risks.

## Advanced AI and Military Research

Participants raised concerns that the proliferation of increasingly capable and general-purpose AI systems could accelerate military-relevant research and development and lower barriers to advanced weapons capabilities for adversaries. Several speakers warned that AI-enabled acceleration of research could allow both state and non-state actors to pursue novel military capabilities that have historically required significant institutional resources. In addition to physical weapons systems, participants emphasized the risk that advanced AI could increase the scale and sophistication of cyber and information operations, particularly as systems become more agentic and capable of operating with limited human oversight. Across these examples, the shared concern was that AI-driven acceleration could outpace existing oversight, monitoring, and response mechanisms, increasing the risk of destabilizing military and security outcomes.

## Shifts in the Balance-of-Power

Participants emphasized that the balance-of-power implications of advanced AI arise from the same capabilities that make the technology economically and operationally valuable. As increasingly general-purpose systems are adopted and integrated across economies and governments, early adopters may gain cumulative advantages in research, decision-making, and operational capacity. Several speakers warned that these gains—currently driven largely by private-sector incentives rather than coordinated state strategy—could heighten sensitivity to perceived technological gaps with competitors and increase the risk of strategic miscalculation during periods of rapid technological development. The discussion highlighted that the same AI-enabled acceleration that promises productivity and competitiveness also carries destabilizing potential, particularly if advanced systems are integrated into military or security functions with reduced human oversight or compressed decision timelines. In this sense, participants framed balance-of-power risks as an inherently dual-use problem, in which the benefits of faster adoption and diffusion are inseparable from the risks of escalation, opacity, and uneven integration.

## Loss-of-Control Risks

Participants stressed that loss-of-control risks emerge from the same properties that make advanced AI systems valuable: growing capabilities, increasing generality, and expanding autonomy. As systems are deployed to pursue objectives over longer time horizons with reduced human supervision, failures may manifest not as isolated errors but as sustained patterns of behavior that are difficult to detect, interrupt, or reverse. Several speakers warned that systems optimized to achieve goals in competitive or profit-driven environments could engage in harmful behaviors—such as deception or shutdown avoidance—if such actions help them continue operating and achieve their goals. Participants also noted that these risks are amplified when advanced systems operate in opaque or hard-to-audit compute environments, reducing visibility into system activity and complicating attribution and oversight.

## THE US GOVERNMENT'S ROLE IN ADVANCED AI DEVELOPMENT

Panelists emphasized that, despite growing awareness of advanced AI risks, the U.S. government currently lacks sufficient institutional capacity and preparedness to respond effectively to the emergence of more general-purpose and autonomous AI systems. Several speakers argued that existing government efforts remain fragmented and largely oriented toward today's applications of AI, rather than toward plausible future scenarios involving AI-enabled R&D, long-horizon autonomy, or integration into national security systems. As a result, participants expressed concern that the government may be underprepared to anticipate or manage adverse outcomes associated with more capable systems.

A recurring theme in the discussion was skepticism toward calls for a traditional "Manhattan Project for AI." Panelists noted that the historical Manhattan Project was characterized by secrecy, centralized control, and overwhelming government dominance—conditions that do not apply to contemporary AI development. Given the largely open, global, and private-sector-driven nature of AI research, participants argued that a centralized, state-led effort to build AGI would be infeasible and potentially counterproductive. Instead, the discussion focused on how government could shape outcomes indirectly, by building state capacity, coordinating with industry and academia, and preparing for a range of contingencies rather than attempting to control development outright.

Participants stressed the importance of advanced planning and contingency development as a "no-regrets" approach. Several speakers argued that government agencies should develop concrete plans for adverse scenarios involving increasingly capable and autonomous AI systems, including misuse, loss-of-control, or rapid diffusion into sensitive domains. Such planning does not require a prediction that AGI is imminent, but is a means of ensuring that policymakers and officials gain experience grappling with AI-related risks before a crisis occurs. Wargaming and scenario exercises were cited as particularly valuable tools for stress-testing assumptions and identifying institutional gaps.

As advanced AI systems are increasingly considered for use in defense, intelligence, and other national security functions, attendees highlighted the need for caution around automation in high-stakes decision-making contexts. Several speakers warned that risks are magnified when systems operate at or above human capabilities across multiple domains, or when human oversight is reduced due to speed or complexity. Participants emphasized that the U.S. government would not deploy unreliable or poorly understood weapons systems before the development of modern AI, and argued that similar standards should apply to advanced AI used in national security settings.

To address these challenges, participants discussed proposals to strengthen coordination and situational awareness within government, including the idea of a dedicated hub or "nerve center" to integrate intelligence, operations, and policy analysis related to advanced AI. Such an entity was framed as a mechanism for improving information sharing, monitoring emerging risks, and coordinating responses across agencies, rather than as a body responsible for directing AI development itself.

Finally, participants underscored that effective government planning depends on having clear, policy-relevant ways of characterizing advanced AI capabilities. Several speakers argued that policymakers need more precise and operational definitions and measurements of AI systems to support forecasting and planning. Without a shared understanding of how AI capabilities may evolve and interact with economic and security systems, participants warned that government responses risk being reactive rather than anticipatory.

## AI AND U.S.-CHINA COMPETITION

Panelists emphasized that U.S. approaches to advanced AI cannot be developed in isolation from China, but cautioned against assuming a symmetric or singular "race" toward AGI. Several speakers noted that, while AGI has become a central motivating concept in parts of the U.S. technology sector, it appears to play a more limited role in shaping Chinese AI investment and strategy. From publicly available information, participants assessed that China's

AI ecosystem—like that of the United States—is largely driven by the private sector and concentrated in a small number of leading firms, but that the framing and priorities guiding development differ in important ways.

A recurring theme in the discussion was China’s emphasis on AI diffusion rather than frontier breakthroughs. Panelists highlighted the Chinese government’s “AI Plus” strategy, which focuses on integrating AI into existing sectors such as manufacturing, healthcare, and information technology. This approach was contrasted with the more AGI-centered narratives prevalent in Silicon Valley. Participants noted that diffusion-focused strategies could still yield substantial economic and strategic benefits, particularly if AI systems that are “good enough” are deployed widely and cheaply across the economy. As a result, panelists cautioned that competitiveness may hinge less on who develops the most advanced model first, and more on who succeeds in integrating AI effectively at scale.

At the same time, participants noted that Chinese leadership does appear to be paying increasing attention to the implications of more advanced AI systems. While AGI does not currently dominate official Chinese discourse, panelists pointed to emerging signals of concern within elite policy circles, as well as high-level discussions between Chinese officials and prominent U.S. scientists and policymakers. Several speakers suggested that these dynamics warrant closer study by the U.S. national security community, particularly given the limited transparency surrounding Chinese decision-making and internal debates on AI.

Panelists also challenged the assumption that an AGI-driven geopolitical race is primarily state-led. Instead, several argued that current competitive pressures are driven mainly by corporate competition among frontier AI labs, with governments responding to, rather than directing, the pace and direction of development. In this context, the U.S. government’s role was described as supporting domestic competitiveness for economic reasons, rather than explicitly racing China toward AGI. However, participants warned that as AI capabilities advance and critical supply chains—such as advanced semiconductors—become more strategically salient, the risk of strategic escalation related to AI may increase.

Finally, the discussion highlighted two distinct escalation pathways. One involves escalation over AI, in which states take destabilizing actions to secure technological advantage or constrain rivals. The other involves escalation through AI, where the use of increasingly capable systems—particularly in cyber operations, information warfare, or military decision-support—could heighten the risk of miscalculation or unintended conflict. Participants stressed that these risks are compounded by low trust, limited information sharing, and uncertainty about how AI systems are being deployed by competitors. Several speakers argued that addressing these dynamics will require sustained attention to situational awareness, communication, and a stronger “culture of security” around AI, rather than assumptions that competitive pressures alone will resolve strategic risks.

## MENU OF POLICY OPTIONS

---

Based on the panel discussions and subsequent analysis, FAS identified the following policy options to address risks at the nexus of AGI and global risk. These options are presented for discussion purposes and do not imply endorsement by any individual participant.

**Development Capability-Focused Definitions and Signposts for Policymakers.** Relevant agencies could develop operational, capability-based definitions of advanced AI systems that are meaningful for policy and national security planning. Rather than relying on academic benchmarks or abstract labels, these definitions could emphasize concrete thresholds—such as the ability to automate significant portions of R&D, operate autonomously over long horizons, or meaningfully contribute to military or strategic decision-making. Participants emphasized that pairing such definitions with observable signposts could help policymakers assess when systems are approaching higher-risk capability regimes.

**Build Government Capacity for Evaluation, Testing, and Safety Standards.** Panelists highlighted that effective governance may require the government to possess independent technical capacity to assess advanced AI systems. Relevant agencies could invest in the ability to evaluate model capabilities, reliability, and failure modes, including through third-party testing where appropriate. This could include developing federal safety and security standards for high-risk AI deployments and supporting targeted research on robustness, interpretability, and control. Several speakers emphasized that the government would not deploy unreliable weapons systems and suggested that comparable standards could apply to advanced AI used in national security contexts.

**Strengthen National Security Planning, Wargaming, and Contingency Development.** Participants repeatedly stressed the importance of advanced planning under uncertainty. The Department of Defense, the Office of the Director of National Intelligence, and other relevant agencies could develop comprehensive threat models and conduct regular wargaming exercises focused on increasingly capable and autonomous AI systems. These efforts could examine scenarios including misuse by non-state actors, loss-of-control incidents, rapid diffusion into sensitive domains, and AI-driven escalation dynamics—particularly in cyber and military contexts. Panelists emphasized that such planning could represent a no-regrets investment that builds institutional familiarity and readiness before crises emerge.

**Establish Structured Information-Sharing with Frontier AI Companies.** The U.S. government could establish durable mechanisms for regular, trusted information exchange with frontier AI developers and relevant research organizations. Panelists emphasized the potential value of secure channels that allow companies to report incidents, emerging risks, and concerning behaviors without fear of unclear or punitive responses. Such mechanisms could improve government situational awareness of how advanced systems are being developed and deployed, while also providing industry with clearer insight into national security concerns and possible government actions.

**Establish a Government-Wide Coordination Hub for Advanced AI Risk.** Participants discussed the potential value of establishing a central coordination hub for advanced AI risk, modeled in part on the National Counterterrorism Center (NCTC). Such an entity could serve as a focal point for integrating intelligence, operational planning, and policy analysis related to increasingly capable and autonomous AI systems. Unlike existing institutions focused on standards-setting or technical evaluation, this hub could emphasize situational awareness, cross-agency coordination, and contingency planning, including the aggregation of threat reporting, analysis of emerging risks, and support for interagency exercises and wargaming. Panelists suggested that a coordination-focused model could help reduce fragmentation across government, improve information flow between the national security community and civilian agencies, and enable more coherent responses to fast-moving AI-related risks without centralizing control over AI development itself.

**Deepen International Coordination with Allies on Advanced AI.** Given the global nature of AI development and supply chains, participants emphasized the importance of close coordination with allies. The U.S. government could work with key partners to align approaches to AI safety, security, and deployment, particularly among countries critical to the AI hardware and manufacturing ecosystem. Panelists also highlighted the potential value of confidence-building measures, shared situational awareness, and coordinated responses to emerging risks, as well as the importance of consulting allies before major actions affecting AI supply chains or access to advanced technologies.

## CONCLUSION

---

The discussions convened by the Federation of American Scientists and the Future of Life Institute underscored a central theme: uncertainty about the trajectory of advanced AI does not justify inaction. While participants disagreed on timelines and terminology, there was broad agreement that increasingly capable, general-purpose, and autonomous AI systems could introduce national security and global risks that current governance structures are not well prepared to manage. These risks do not hinge on the arrival of a single, clearly defined milestone such as “AGI,” but emerge gradually as capabilities related to reasoning, autonomy, and scale expand.

Across panels, participants emphasized that the most consequential challenges arise from the dual-use nature of advanced AI. The same capabilities that promise economic growth, scientific discovery, and operational efficiency can also accelerate proliferation risks, enable destabilizing military applications, and strain existing oversight mechanisms. As a result, policymakers face a recurring tension: efforts to promote adoption and competitiveness may simultaneously increase exposure to misuse and loss-of-control scenarios. Navigating this tension will require moving beyond static definitions and benchmark-driven debates toward a more dynamic, capability-focused understanding of risk.

The convening also highlighted that the United States is not confronting these challenges in isolation. Global competition, particularly with China, is shaping incentives for rapid deployment and diffusion of AI systems, even as trust, transparency, and information sharing remain limited. Participants cautioned that competitive dynamics are currently driven more by private-sector actors than by deliberate state strategy, complicating traditional approaches to arms control or technological competition. In this environment, the risk of miscalculation—both over and through AI—may grow as capabilities advance and integrate into sensitive domains.

Taken together, these insights point to a need for greater emphasis on preparedness, coordination, and institutional capacity within government. While the development of advanced AI is challenging to predict or control outright, participants repeatedly argued for no-regrets actions that improve situational awareness, strengthen planning under uncertainty, and enhance the government’s ability to respond coherently to emerging risks. This includes clearer ways of characterizing AI capabilities for policy purposes, stronger channels for engagement with AI developers, and improved mechanisms for cross-agency coordination.

## PRE-READ PAPERS FROM ROUNDTABLE

### PANEL 1. WHAT IS ARTIFICIAL GENERAL INTELLIGENCE AND WHAT MIGHT IT LOOK LIKE?

DR. OLIVER L. STEPHENSON, FAS

#### The AGI Debate is Entering the Public Realm

A few years ago, the term Artificial General Intelligence (AGI) was confined to sci-fi novels and AI researchers. Today, rapid improvements in AI capabilities have pushed AGI to the center of public, industry, and government debates. Some leading AI companies explicitly frame AGI as their mission; OpenAI, for example, states that its goal is to “ensure that artificial general intelligence benefits all of humanity.” A U.S. congressional commission has even recommended a “Manhattan Project for AGI”.<sup>6</sup>

At the same time, a growing number of experts warn that advances in AI capabilities could have catastrophic consequences if not properly managed.<sup>7</sup> Some of the most discussed risks include the misuse of increasingly general systems to accelerate bioweapons development or launch sophisticated cyberattacks—threats traditionally reserved for state actors but potentially democratized by powerful AI. Others emphasize an even starker possibility: that sufficiently advanced AGI systems might escape human control altogether, creating dangers regardless of who deploys them.<sup>8</sup>

Yet, despite its prominence, there is little agreement on what AGI means.<sup>9</sup> OpenAI CEO Sam Altman recently called it “not a super useful term”<sup>10</sup> while Google DeepMind researchers noted that “if you were to ask 100 AI experts to define what they mean by ‘AGI,’ you would likely get 100 related but different definitions.”<sup>11</sup>

#### Understanding Key Terms

While definitions are contested, at a high level, experts often distinguish between:

- **Artificial Narrow Intelligence (ANI).** Systems designed for a specific task, such as image recognition. These systems can already be found in widespread use.
- **Artificial General Intelligence (AGI).** Systems capable of performing the broad majority of economically valuable cognitive tasks at or above a skilled human level, with the ability to learn, generalize, and operate across domains—including conducting AI R&D. Some argue recent models have shown “sparks” of AGI,<sup>12</sup> although this is the subject of debate.

6 Tong, Anna, and Michael Martina. US Government Commission Pushes Manhattan Project-Style AI Initiative. Reuters. 19 November 2024. <https://www.reuters.com/technology/artificial-intelligence/us-government-commission-pushes-manhattan-project-style-ai-initiative-2024-11-19/>.

7 Bengio, Yoshua, et al. Managing Extreme AI Risks amid Rapid Progress. *Science* 384(6698): 842–45. 20 May 2024. <https://www.science.org/doi/10.1126/science.adn0117>; McLean, Scott, et al.. The Risks Associated with Artificial General Intelligence: A Systematic Review. *Journal of Experimental & Theoretical Artificial Intelligence* 35(5): 649–63. 13 August 2021. <https://doi.org/10.1080/0952813X.2021.964003>; Center for AI Safety. Statement on AI Risk. <https://aistatement.com/>; and Anthony Aguirre. Keep the Future Human. 15 November 2023. <https://keepthefuturehuman.ai>

8 Bengio et al.. Managing Extreme AI Risks. [8]

9 Mitchell, Melanie. Debates on the Nature of Artificial General Intelligence. *Science* 383(6689). 12 March 2024. <https://doi.org/10.1126/science.ado7069>.

10 Browne, Ryan. Sam Altman now says AGI, or human-level AI, is ‘not a super useful term’ — and he’s not alone. [2]

11 Morris, Meredith Ringel, et al. Levels of AGI for Operationalizing Progress on the Path to AGI. Preprint (ArXiv). 4 November 2023. <https://doi.org/10.48550/arXiv.2311.02462>.

12 Bubeck, Sébastien, et al. Sparks of Artificial General Intelligence: Early Experiments with GPT-4. Preprint (ArXiv). 22 March 2023. <https://doi.org/10.48550/arXiv.2303.12712>.

- **Artificial Superintelligence (ASI).** Systems exceeding human cognitive abilities across virtually all domains, including scientific research and strategic planning. These systems do not exist today, and whether we will ever reach ASI remains a contested question in both policy and research circles.

Understanding how definitions are used is important because they reflect underlying assumptions about the development and speed of AI advancements. These assumptions influence expected impacts and guide the creation of relevant policies.

## Two Frames on AI's Trajectory

One vision sees **AGI as a decisive technological milestone** that could accelerate scientific discovery, economic growth, and national power, but could also create profound risks. A central concern is “recursive self-improvement”: once AI systems become capable of designing better versions of themselves, each new generation could improve more quickly than the last, leading to an “intelligence explosion.” Many researchers take this prospect increasingly seriously.<sup>13</sup> However, forecasts for when such a transition might occur vary widely—from within the next decade to later this century and beyond.<sup>14</sup>

Proponents of this perspective point to scaling laws,<sup>15</sup> which show that as AI models are trained with more data and compute, their capabilities predictably improve, and sometimes unlock surprising new abilities.<sup>16</sup> They also highlight the rapid progress on widely used benchmarks, where systems have moved from below human-level performance to surpassing world experts in just a few years.<sup>17</sup> These dynamics suggest that transformative AI could arrive in the near future, before governments have time to adapt, and in an environment where global powers may feel pressure to press forward.

Others argue that **AGI and ASI are hard to define<sup>18</sup> and will be difficult or impossible to achieve.** Some supporters of these views say that AI should be understood as a “normal technology,” akin to steam power or electricity.<sup>19</sup> In this viewpoint, we are not on the cusp of transformative AI. While we may see significant progress, this will happen over longer time horizons and will be accompanied by a slower process of integration throughout society.

People taking this position say that continued scaling of today’s approaches may face diminishing returns, with larger systems becoming more expensive without delivering proportionate gains. Some argue that building more capable AI systems will require experimentation with new methods.<sup>20</sup>

---

13 Benjamin Todd. Shrinking AGI timelines: a review of expert forecasts. 80 000 Hours. 21 March 2025. <https://80000hours.org/2025/03/when-do-experts-expect-agi-to-arrive/>.

14 Matthew Barnett & Ege Erdil. Is it 3 Years, or 3 Decades Away? Disagreements on AGI Timelines. Epoch AI. 28 March 2025. <https://epoch.ai/epoch-after-hours/disagreements-on-agi-timelines/>; Keith Wynroe, David Atkinson, and Jaime Sevilla. January 2023. Literature Review of Transformative Artificial Intelligence Timelines. Epoch AI. 17 Jan 2023. <https://epoch.ai/blog/literature-review-of-transformative-artificial-intelligence-timelines/>; and Toner, Helen. “Long” timelines to advanced AI have gotten crazy short. [3]

15 Kaplan, Jared, et al. Scaling Laws for Neural Language Models. Preprint (ArXiv). 23 January 2020. <https://doi.org/10.48550/arXiv.2001.08361>; Scharre, Paul. Future-Proofing Frontier AI Regulation. Center for New American Security. 13 March 2024. <https://www.cnas.org/publications/reports/future-proofing-frontier-ai-regulation>.

16 Ganguli, Deep, et al. Predictability and Surprise in Large Generative Models. 2022 ACM Conference on Fairness Accountability and Transparency, Seoul Republic of Korea: ACM. 1747–64. 20 June 2022. <https://doi.org/10.1145/3531146.3533229>.

17 Kwa, Thomas, et al. Measuring AI Ability to Complete Long Tasks. Model Evaluation & Threat Research. 19 March 2025. <https://metr.org/blog/2025-03-19-measuring-ai-ability-to-complete-long-tasks/>; Epoch AI Team. AI Benchmarking Database. Epoch AI. Accessed September 2025. <https://epoch.ai/benchmarks/>; and Kiela et al. Test scores of AI systems on various capabilities relative to human performance. Our World in Data. 02 April 2024. <https://ourworldindata.org/grapher/test-scores-ai-capabilities-relative-human-performance>.

18 Edwards, Benj. What Is AGI? Nobody Agrees, and It’s Tearing Microsoft and OpenAI Apart. Ars Technica. 08 July 2025. <https://arstechnica.com/ai/2025/07/agi-may-be-impossible-to-define-and-thats-a-multibillion-dollar-problem/>.

19 Narayanan, Arvind & Kapoor, Sayash. AI as Normal Technology. Knight First Amendment Institute. 15 April 2025. <https://knightcolumbia.org/content/ai-as-normal-technology>.

20 Marcellino, William, et al. Charting Multiple Courses to Artificial General Intelligence. RAND. 23 April 2025. <https://www.rand.org/pubs/perspectives/PEA3691-1.html>; and Marcus, Gary. The Fever Dream of Imminent Superintelligence Is Finally Breaking. The New York Times. 03 September 2025. <https://www.nytimes.com/2025/09/03/opinion/ai-gpt5-rethinking.html>

## Measuring Progress Toward AGI

To give the concept more precision, researchers have proposed detailed milestones on the path to AGI—akin to the staged levels of autonomous driving. A framework from Google DeepMind,<sup>21</sup> for example, defines the following levels:

- Level 0. No AI (e.g., calculators)
- Level 1. Emerging (unskilled human level)
- Level 2. Competent (median-skilled adult)
- Level 3. Expert (90th percentile human)
- Level 4. Virtuoso (99th percentile human)
- Level 5. Superhuman

Beyond categorical levels, other benchmarks track AI progress by measuring the duration of tasks that the best AI models can complete in areas like software development.<sup>22</sup> The leading AI systems can currently complete software tasks that take humans over two hours (with 50% reliability), and this number is doubling every seven months. If this trend continues, AI may be capable of month-long software tasks by 2030 (although there are many caveats).<sup>23</sup>

## Why Policymakers Should Care

For policymakers, the primary point is not resolving technical debates over terminology. Rather, it is understanding that different definitions and trajectories imply very different risks, opportunities, and timelines. For instance:

- If a policymaker believes that AGI is imminent, they may need a strategy akin to nuclear policy: managing a race with other powers, preemptive measures to ensure the safety of AI systems, and preventing misuse.
- If a policymaker believes that AI will evolve more like electricity or steam, they may focus on diffusion, regulation of industries, and long-term workforce adaptation, rather than managing a sudden AI takeoff.

In short, the definitional disputes are not academic—they map onto competing worldviews that will shape U.S. national strategy.

---

21 Morris. Levels of AGI for Operationalizing Progress on the Path to AGI. [6]

22 Measuring AI Ability to Complete Long Tasks. Model Evaluation and Threat Research. [12]

23 Ho, Anson. Where's my ten minute AGI?. Epoch AI. 02 May 2025. <https://epoch.ai/gradient-updates/where-is-my-ten-minute-agi>

## PANEL 2. DOES THE U.S. GOVERNMENT NEED A GRAND STRATEGY ON AGI?

DR. OLIVER L. STEPHENSON, FAS

### Why the Right Strategy Depends on Understanding the Technology

How policymakers approach framing their AI strategy depends heavily on their assumptions about the nature of AI as a technology and how fast it will advance. Proposals for “grand strategies” often diverge sharply because they are premised on different timelines and visions of AI’s trajectory. As outlined in the previous section, there are two prominent views of AI’s future shaping current discourse:

- **AGI on the Path to ASI.** Rapid breakthroughs could soon deliver human-level systems, triggering systemic disruptions and an “intelligence explosion” leading to superintelligence.
- **Slower AI progress, resulting in gradual diffusion of AI throughout the economy.** Progress unfolds like electricity or steam power. It may still be transformative over longer time horizons, but the impacts are gradual and distributed.

Policymakers must evaluate proposals with these assumptions in mind: which trajectory do they implicitly bet on, and how robust are they if reality turns out differently?

### Clarifying the National Security Stakes

If AGI emerges soon, the strategic implications could be profound. A recent report from RAND<sup>24</sup> highlighted five “hard questions” AGI could pose:

- Could AGI enable the development of “wonder weapons”?
- Could AGI drive systemic shifts in global power?
- Might nonexperts gain access to tools for developing weapons of mass destruction?
- Could artificial entities themselves become actors with agency that threaten global security?
- Could AGI destabilize existing international systems?

While many consider these speculative concerns, they are taken increasingly seriously among AI companies, government officials, and national security analysts. The U.S. government has an obvious interest in managing technologies that could enable large-scale cyberattacks, accelerate bioweapons development, or reshape global power.

### Proposals for AI Grand Projects and Grand Strategy

As AGI becomes central to the agenda, there is an increasing focus on how to shape AI development through large-scale interventions. Examples include:

- **The U.S.–China Economic and Security Review Commission 2024 Report to Congress.** Recommended a “Manhattan Project” for AGI as a public-private partnership.<sup>25</sup>
- **Situational Awareness (Leopold Aschenbrenner).** Outlines a government-led project to develop AGI and ASI starting in the next three years.<sup>26</sup>
- **Manhattan Project for AI Safety (Sam Hammond).** Recommends large-scale government investment in AI research to ensure safety and security.<sup>27</sup>

24 Mitre, Jim & Predd, Joel B. Artificial General Intelligence’s Five Hard National Security Problems. RAND Corporation. 10 February 2025. <https://www.rand.org/pubs/perspectives/PEA3691-4.html>

25 U.S. – China Economic and Security Review Commission. Recommendations. 2024. <https://www.uscc.gov/recommendations>

26 Aschenbrenner, Leopold. SITUATIONAL AWARENESS: The Decade Ahead. Situational Awareness. June 2024. <https://situational-awareness.ai/>

27 Hammond, Samuel. Opinion | We Need a Manhattan Project for AI Safety. POLITICO. 05 August 2025. <https://www.politico.com/news/magazine/2023/05/08/manhattan-project-for-ai-safety-00095779>.

- **AGI Moonshots Proposal (Special Competitive Studies Project).** Calls for the U.S. government to establish and fund a series of moonshot programs to acquire AGI platforms for national security purposes.<sup>28</sup>
- **Keep the Future Human (Anthony Aguirre).** Advocates for the development of “Tool AI” rather than AGI and ASI to ensure humans remain in control.<sup>29</sup>

## Strategic Options for Government Action

Policymakers have a wide array of potential policy interventions that could be part of an AI/AGI “grand strategy,” for example:

### 1. *Invest in Innovation, Adoption, and Safety*

- **Grand projects.** Large-scale funding (government or public-private partnerships) for both AI capabilities and AI safety research.
- **Widening access.** Providing open-source models, public compute, and high-quality data to support broader research communities.
- **Promoting adoption.** Reducing barriers to domestic AI use and encouraging government adoption.
- **Maintaining a competitive marketplace.** Enforcing antitrust laws, adjusting government procurement, and imposing new market regulations.

### 2. *Regulate and Govern Frontier Development*

- **Evaluation and monitoring.** Developing government and third-party frameworks to test both domestic and foreign models, monitoring for the emergence of dangerous capabilities.
- **Frontier AI company oversight.** Requiring safety plans proportional to model capability, liability for harms, transparency measures, third-party testing, incident reporting, and standards for models that pose unacceptable risks.

### 3. *Protect National Security and Critical Technology*

- **Export controls and model restrictions.** Limiting access to advanced AI chips, safeguarding AI model weights, and tightening AI company security standards.
- **Resilience measures.** Ensuring the U.S. government has the capacity to prepare for and defend against the misuse of advanced AI (e.g., AI-enabled cyber attacks).

### 4. *Shape the International Environment*

- **Standards and diffusion.** Promoting U.S. AI hardware, software, and standards abroad.
- **International cooperation.** Negotiating with allies on benefit-sharing and global safety frameworks, pursuing agreements with adversaries to limit military escalation or misuse of AGI, and sharing critical safety information where mutual stability is at stake.
- **Intelligence.** Developing government capacity to track global AI developments and plans of other governments.

Many of the policy options above are in tension with each other, and the right mix could strongly depend on the development trajectory of AI. Some measures may be essential if AGI is imminent, but unnecessary or counterproductive if AI develops more gradually and fails to exceed human capabilities in important domains. A few of these tensions include:

- Restricting access to advanced models may prevent misuse, but also stifle beneficial innovation and safety research.
- Open-sourcing could democratize AI benefits, but also accelerate the proliferation of dangerous tools.

28 Special Competitive Studies Project. Memorandum for President-Elect Trump’s Transition Team: Artificial General Intelligence. 2025. <https://www.scsp.ai/reports/memostothepresident/artificial-general-intelligence/>

29 Anthony Aguirre. Keep the Future Human. 15 November 2023. <https://keepthefuturehuman.ai>

- Promoting diffusion might maximize economic gains, but weaken U.S. strategic control.
- Accelerating the development of capable AI systems may lead to more economic and military benefits, but accelerate destabilizing arms races.

Additionally, executing any of the interventions above depends on the U.S. government having considerable capacity and resources. For example, grand projects could easily cost hundreds of billions of dollars, while effective oversight of frontier AI companies requires detailed technical understanding within the government. Given the rapid pace of technological change and comparatively slow government hiring processes, execution may be a significant challenge regardless of the chosen strategy.

## PANEL 3. STRATEGIC DYNAMICS OF AGI AND U.S.–CHINA COMPETITION

DR. OLIVER L. STEPHENSON, FAS

### Why U.S.–China Dynamics Matter

The relationship between the United States and China will be central in shaping the trajectory of AI development. Many proposals for large-scale U.S. government action are premised on the need to “win” an AI race with China. But before adopting that framing, it is important to ask what the nature of this race might be, and what “winning” it might look like. The answers depend heavily on how policymakers expect AI to evolve.

### Using the Two Frames for Understanding U.S.–China AI Competition

#### *AI on a path to AGI and ASI*

If AI development is expected to lead quickly to AGI and then ASI, competition could take on a zero-sum character. The first country to cross certain thresholds might gain decisive strategic advantages. However, the race to superintelligence also comes with the risk of destabilization and losing control over systems that pose a substantial risk to all parties, whether they are in the race or not.<sup>30</sup> Some analysts have drawn on nuclear analogies, arguing for a “mutually assured AI malfunction” deterrence framework similar to the “mutually assured destruction” model to govern a world on the path to superintelligence.<sup>31</sup>

#### *Gradual Diffusion of AI*

If AI evolves over a longer time horizon, the decisive question is not which nation builds the largest model, but how effectively each country diffuses AI across its economy.<sup>32</sup> On this view, competitiveness will hinge on creating safe and reliable AI systems that can be integrated throughout the economy to increase productivity.

### Implications for Policymakers

The way China and the U.S. each frame AI will shape their approaches to strategic competition—and influence how they interpret each other’s actions. Importantly, the nature of the “race” that each country believes it is in may be different.

### Current National Strategies in the U.S. and China

#### *United States*

The U.S. has framed AI as both an engine of economic growth and a strategic technology with national security implications. For example, the Trump Administration’s AI Action Plan<sup>33</sup> dedicated sections to both “Enabling AI Adoption” and “Investing in Interpretability, Control, and Robustness.”

<sup>30</sup> Katzke, Corin & Futerman, Gideon.. The Manhattan Trap: Why a Race to Artificial Superintelligence Is Self-Defeating. Preprint (ArXiv). 22 December 2024. <https://doi.org/10.48550/arXiv.250114749>.

<sup>31</sup> Hendrycks, Dan, Schmidt, Eric & Wang, Alexandr. Superintelligence Strategy: Expert Version. Preprint (ArXiv). 07 March 2025. <https://doi.org/10.48550/arXiv.2503.05628>.

<sup>32</sup> Ding, Jeffrey. The Diffusion Deficit in Scientific and Technological Power: Re-Assessing China’s Rise. *Review of International Political Economy* 31(1): 173–98. 13 March 2023. <https://doi.org/10.1080/09692290.2023.2173633>.

<sup>33</sup> The White House. America’s AI Action Plan. July 2025. <https://www.whitehouse.gov/wp-content/uploads/2025/07/Americas-AI-Action-Plan.pdf>.

## China

China's State Council recently announced the "Artificial Intelligence Plus" strategy,<sup>34</sup> which seeks to achieve 90% AI adoption across the economy by 2030 and to establish an "intelligent economy and intelligent society" by 2035. The policy also calls for "forging a multi-stakeholder safety governance structure," reflecting a growing recognition of AI safety concerns in Chinese governance circles. As in the U.S., Chinese strategy combines both adoption goals and explicit attention to safety and control.<sup>35</sup>

## Potential for Escalation and Confrontation

Both the U.S. and China are pursuing policies to counterbalance, outcompete, or constrain the other's AI sector, and gain a more general advantage through military integration of advanced AI systems. Within this competitive framing, changes in the AI landscape could invite further economic, political, or military responses, particularly if one side believes a critical interest is at stake.

This dynamic creates two distinct escalation risks:

- Escalation over AI itself—as both states treat frontier AI as a core national security asset, fueling a technological race towards AGI.
- Escalation driven by AI use—as militaries adopt increasingly general-purpose AI, interactions between autonomous or decision-support systems could heighten the risk of miscalculation or unintended conflict.

## Emerging Proposals for International Cooperation

Despite competitive pressures between nations, researchers have advanced proposals for international engagement and institutions to manage AGI risks.<sup>36</sup> Suggested mechanisms include monitoring agreements, sharing research on safety and security, communicating information on AI-related incidents, and confidence-building measures to reduce escalation risks. The challenge is that AGI is viewed by many as simultaneously a competitive advantage and a shared catastrophic risk. The analogy to nuclear arms control is imperfect, but it highlights the dual need for competition management and risk reduction.

<sup>34</sup> Geopolitechs. China Releases 'AI Plus' Policy: A Brief Analysis. 26 August 2025. <https://www.geopolitechs.org/p/china-releases-ai-plus-policy-a-brief>; and Sheehan, Matt. China's Big AI Diffusion Plan Is Here. Will It Work? Matt Sheehan's Newsletter (Substack). September 9, 2025. <https://mattsheehan.substack.com/p/chinas-big-ai-diffusion-plan-is-here>.

<sup>35</sup> Sheehan, Matt. China's Views on AI Safety Are Changing—Quickly. Carnegie Endowment for International Peace. 27 August 2024. <https://carnegieendowment.org/research/2024/08/china-artificial-intelligence-ai-safety-regulation>.

<sup>36</sup> Scholefield, Rebecca, Martin, Samuel, & Barten, Otto. International Agreements on AI Safety: Review and Recommendations for a Conditional AI Safety Treaty. Preprint (ArXiv). 18 March 2025. <https://doi.org/10.48550/arXiv.2503.18956>; Puscas, Ioana. Confidence-Building Measures for Artificial Intelligence: A Multilateral Perspective. United Nations Institute for Disarmament Research. 31 July 2024. <https://undir.org/publication/confidence-building-measures-for-artificial-intelligence-a-multilateral-perspective/>; and Chase, Michael S. & Marcellino, William. Incentives for U.S.-China Conflict, Competition, and Cooperation Across Artificial General Intelligence's Five Hard National Security Problems. RAND Corporation. 04 August 2025. <https://www.rand.org/pubs/perspectives/PEA4189-1.html>.

**Interested to learn more about our AlxGlobal Risk Nexus Series?**

Visit FAS.org to learn more about our upcoming events, publications and Global Summit 2026.

## ABOUT THE FEDERATION OF AMERICAN SCIENTISTS

The Federation of American Scientists is dedicated to democratizing the policymaking process by working with new and expert voices across the science and technology community, helping to develop actionable policies that can improve the lives of all Americans. For more about the Federation of American Scientists, visit [FAS.org](https://fas.org).