

The AGI Landscape in 2025: Competition, Governance, and Emerging Paradigms

Prologue

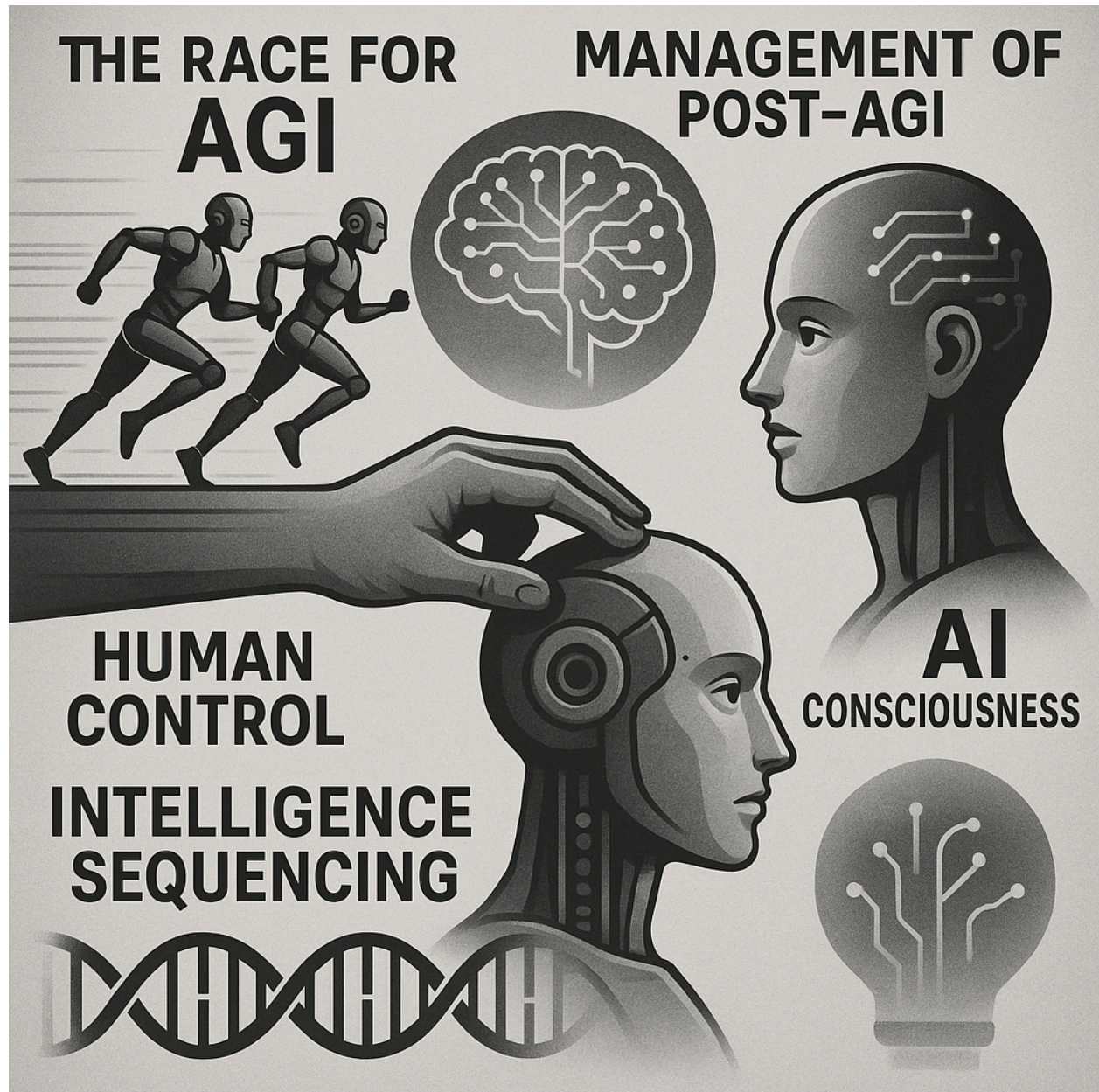
The year 2025 stands as a watershed moment in the relentless pursuit of Artificial General Intelligence (AGI). The theoretical aspirations of the past have crystallized into a tangible, high-stakes race, reshaping the geopolitical landscape and igniting a fervent global discourse. This document serves as an intelligence briefing, capturing the intricate tapestry of this pivotal year. We find ourselves at the confluence of unprecedented technological leaps, massive financial investments, and intensifying national rivalries. Frontier models from leading organizations push the boundaries of AI capabilities, while nations strategize to secure dominance in this transformative domain. Yet, this rapid advancement is shadowed by profound questions of governance, safety, ethics, and the very definition of intelligence itself. As the specter of AGI looms closer, the imperative to understand, control, and responsibly navigate this uncharted territory has never been more urgent. This report aims to illuminate the key actors, dynamics, and emerging paradigms shaping this critical juncture in human history, acknowledging the inherent uncertainties while striving to provide an evidence-based snapshot of the evolving AGI landscape in 2025.

Disclaimer

This research document is comprehensive and delves into complex topics. It's important to note that the field of Artificial General Intelligence is rapidly evolving, and predictions about its impact are inherently speculative. The information and references presented here should be considered with critical discernment and a degree of caution. Timelines for AGI remain uncertain, capabilities evolve quickly, and the interplay between competition, collaboration, innovation, and regulation is highly dynamic. Readers are encouraged to form their own informed opinions on the topics discussed, including AGI's potential arrival, societal adaptation strategies, and the possibility of AI consciousness. The document aims to provide a foundation for exploration, but individual interpretation and analysis are essential. Please be aware that due to the dynamic nature of this research area, some information may become outdated or be subject to revision.

Your insights and perspectives are invaluable. If you note any inaccuracies, errors, or outdated information, or if you have alternative viewpoints on the subjects discussed, please do not hesitate to provide feedback. You can reach the author at f.lopeznolasco@gmail.com for review and potential correction. This document is intended to stimulate discussion and further inquiry into the critical challenges and opportunities presented by the Post-AGI era.

This research was conducted with the assistance of Google AI technologies, which facilitated information gathering, analysis, and synthesis. While efforts have been made to ensure accuracy and reliability, the final interpretation and conclusions presented are those of the author and sources referenced through this document.



Introduction

Setting the Stage (2025 Context)

The year 2025 marks a critical inflection point in the pursuit of Artificial General Intelligence (AGI). We are witnessing an unprecedented confluence of factors: rapid advancements in frontier model capabilities, exemplified by releases such as OpenAI's o-series, Google's Gemini 2.5, Meta's Llama 4, Anthropic's Claude 3.7, and Baidu's Ernie

4.5/X1¹; massive strategic investments driven by both national interests and corporate ambitions, notably the US Stargate initiative and substantial corporate funding rounds⁷; intensifying geopolitical competition primarily between the United States and China⁸; and a burgeoning, though fragmented, global effort to establish governance frameworks and address profound safety and ethical concerns.¹² The dual promises often associated with AGI – immense economic returns and solutions to global problems – are frequently promoted by developers, yet these claims warrant significant skepticism.¹⁸

Report Objective and Structure

This report provides a consolidated intelligence briefing on the key actors, activities, and dynamics shaping the complex AGI landscape as it stands in 2025. It synthesizes information from the technological, geopolitical, corporate, governance, safety, and philosophical domains, focusing specifically on developments and trends evident this year. The analysis is structured around **five core concepts critical to understanding the current trajectory**:

1. The Race for AGI
2. Management of Post-AGI
3. AI Consciousness
4. Human Control
5. Intelligence Sequencing

Navigating Uncertainty

The field is characterized by rapid, multi-dimensional transformation.¹⁹ Timelines for AGI remain uncertain, capabilities evolve quickly, and the interplay between competition, collaboration, innovation, and regulation is highly dynamic. This report aims to provide a structured, evidence-based snapshot of this evolving landscape, acknowledging the inherent uncertainties while highlighting the key forces and decision points shaping the path forward.

Section 1: The Intensifying Race for AGI Supremacy (2025)

Overview

The pursuit of AGI has transitioned from a long-term research objective to an intense, high-stakes competition involving nations and corporations. In 2025, this race is defined by unprecedented financial commitments, accelerated development cycles, and explicit strategic positioning for global leadership.⁷ The potential rewards – economic dominance, enhanced national security through advantages in defense, cybersecurity, and intelligence, and transformative scientific breakthroughs – are perceived as immense, fueling fierce competition.⁷ However, this rapid acceleration carries inherent risks, prominently including the potential to compromise safety standards in the pursuit of achieving AGI first.¹⁰

1.1 National Strategies and Investments

The geopolitical dimension of the AGI race is dominated by the strategic rivalry between the United States and China, both recognizing AGI as a technology of paramount national importance.

United States

- **Aim**

The primary strategic goal of the US is to solidify and maintain its global leadership in AI development and deployment, viewing this as crucial for sustained economic competitiveness and national security.⁷ Analysis such as the "Superintelligence Strategy" report, discussed by the RAND Corporation, reflects this focus, advocating for policies like managed competition and AI nonproliferation to secure US advantage while managing risks.¹⁰ A significant policy development in early 2025 was the Executive Order signed by President Trump, titled "Removing Barriers to American Leadership in Artificial Intelligence," which revoked previous AI governance policies and signaled a shift towards deregulation aimed at accelerating innovation, potentially creating tension with earlier or international governance approaches.²² The broader strategic thinking extends beyond mere technological development to encompass the societal adaptations necessary to manage the disruptive

potential of superintelligence.¹⁰

- **2025 Initiatives**

The scale of US investment ambition is underscored by the \$500 billion "Stargate" initiative. This massive project, reportedly backed by the government in collaboration with Oracle, Softbank, and OpenAI, aims to construct 20 new data centers specifically designed for advanced AI development, representing an unparalleled investment in AI infrastructure.⁷ This complements substantial corporate investments, such as Microsoft's \$80 billion allocation to AI development in 2025.⁷ Core strategic elements being pursued include enhancing computer security, implementing export controls on critical hardware and potentially model weights, promoting AI safeguards, and investing in the foundational economic elements supporting AI advantage.¹⁰ AI adoption is rapidly increasing within government operations, particularly for defense applications (cybersecurity, autonomous systems, military strategy), policy development (simulation and evaluation), and optimizing public infrastructure and services (smart cities, emergency response).¹¹ Concurrently, the National Science Foundation (NSF) continues its long-term investment in fundamental AI research through its National AI Research Institutes program, with specific themes funded in FY2024 and FY2025 covering areas like astronomical sciences, materials research, and strengthening core AI capabilities.²³

- **Key Influencers**

Thought leaders shaping the strategic debate include the authors of the "Superintelligence Strategy" report (Dan Hendrycks, Eric Schmidt, Alexandr Wang) whose work is influencing policy discussions.¹⁰ Key figures within the current administration, such as the newly mandated Special Advisor for AI and Crypto, are tasked with developing the national AI action plan under the deregulatory framework.²² Leadership within the NSF AI Institutes program guides foundational research directions.²³ Analysts at institutions like RAND provide critical assessment of proposed strategies like AI nonproliferation and the controversial "Mutually Assured AI Malfunction" (MAIM) concept.¹⁰

- **Official/Trending Links**

- RAND Commentary on Superintelligence Strategy:
<https://www.rand.org/pubs/commentary/2025/03/seeking-stability-in-the-competition-for-ai-advantage.html>¹⁰

- ExecutiveBiz Article on AGI in GovCon:
<https://executivebiz.com/2025/03/all-about-artificial-general-intelligence-govcon/>¹¹
- OMMAX Davos 2025 Report (Stargate):
<https://www.ommax.com/en/insights/industry-insights/a-new-ai-era-top-10-takeaways-from-davos-2025/>⁷
- CogentInfo US AI Policy Changes:
<https://www.cogentinfo.com/resources/federal-ai-mandates-and-corporate-compliance-whats-changing-in-2025>²²
- NSF National AI Research Institutes:
<https://www.nsf.gov/funding/opportunities/national-artificial-intelligence-research-institutes/505686/nsf23-610>²³

China

○ Aim

China is engaged in vigorous competition with the US, aiming for leadership in AGI and leveraging AI advancements for economic transformation across sectors like manufacturing, energy, and R&D, as well as enhancing national security.¹⁰ A key objective evident in 2025 is rapidly closing the performance gap between its domestic AI models and leading US counterparts.⁸

○ 2025 Initiatives

State support remains substantial, highlighted by the launch of a \$47.5 billion semiconductor fund aimed at bolstering domestic chip capabilities crucial for AI development.⁸ Chinese tech giants, particularly Baidu, are accelerating the development of their foundation models. Baidu launched Ernie 4.5 (a multimodal model) and Ernie X1 (a reasoning model) in March 2025, claiming performance parity or superiority over competitors like DeepSeek R1 and OpenAI's GPT-4o, often at lower costs.⁵ Baidu also announced plans to release Ernie 5, featuring enhanced multimodal capabilities, in the second half of 2025.²⁸ The emergence of highly competitive startups like DeepSeek AI is notable; its models (R1, V3) have set new benchmarks in reasoning and efficiency, challenging established players.⁴ China continues its global leadership in the sheer volume of AI research publications and patent filings.⁸ In a move potentially aimed at fostering a domestic ecosystem and competing with Western open-source efforts, Baidu announced plans to open-source its

Ernie model codebase later in 2025.⁶

- **Key Influencers**

Leadership within major technology companies like Baidu (e.g., CEO Robin Li, who has publicly acknowledged the competitive landscape and uncertainty²⁸) and rapidly rising startups like DeepSeek AI are central figures.²⁹ Key government officials responsible for formulating and executing the national AI strategy and directing state funding are crucial. Researchers at leading academic institutions, such as Tsinghua University's Institute for AI International Governance (I-AIG), are also influential, participating in international dialogues on AI safety and governance.³²

- **Official/Trending Links**

- Stanford HAI AI Index Report (US-China Comparison): <https://hai.stanford.edu/ai-index/2025-ai-index-report>⁸
- Baidu Research: <https://research.baidu.com/>³¹
- Baidu Ernie 5 Announcement: <https://dig.watch/updates/baidu-to-launch-ernie-5-ai-in-2025>²⁸
- Baidu Ernie 4.5/X1 Launch News: <https://nascenia.com/latest-ai-models/>⁵, <https://siliconangle.com/2025/03/16/baidu-debuts-first-ai-reasoning-mode-i-compete-deepseek/>⁶, <https://champaignmagazine.com/2025/03/16/ai-by-ai-weekly-top-5-03-10-16-2025/>²⁷
- DeepSeek AI Mentions: https://dentro.de/ai/big_players/²⁹, <https://redblink.com/llama-4-vs-deepseek-v3/>⁴

Other National Contexts

While the US-China dynamic dominates, other nations are actively engaging with AI development and governance. Canada, for instance, pledged CAD 2.4 billion towards AI initiatives.⁸ France played a role in international governance discussions by hosting the AI Action Summit in February 2025.¹⁵ The United Kingdom continues to position itself as a hub for safety and standards through its AI Safety Institute and the AI Standards Hub.¹³ Conversely, Germany faces reported hurdles in translating its innovation potential into practical AI applications due to factors like high taxes, energy costs, and regulatory burdens.⁷ Public perception of AI also shows significant regional variation; polling data indicates considerably higher optimism regarding AI's benefits in Asian nations like China (83%), Indonesia (80%), and Thailand (77%) compared to lower levels in Canada (40%), the United States (39%), and the Netherlands (36%), although optimism has seen recent growth in

several previously skeptical Western countries.⁸

1.2 Corporate Frontier Labs: Pushing the Boundaries

The engine of the AGI race is largely powered by a handful of well-funded corporate labs, primarily based in the US, alongside significant players emerging in China and Europe.

- **OpenAI**

- **Aim**

OpenAI explicitly states its mission is to ensure AGI—defined as highly autonomous systems outperforming humans at most economically valuable work—benefits all of humanity.²¹ In 2025, CEO Sam Altman articulated an even more ambitious goal: aiming for superintelligence beyond AGI, potentially revolutionizing science and increasing global abundance.³⁶ The organization's charter emphasizes broadly distributed benefits, long-term safety, technical leadership, and cooperation.²¹ A key principle is a commitment to halt competition and assist a value-aligned, safety-conscious project if it approaches AGI first.²¹

- **2025 Initiatives**

The first half of 2025 saw significant activity. OpenAI released its new flagship models, o3 and o4-mini, in April.¹ This followed the introduction of GPT-4.1 via API, advancements in image generation (4o), and the publication of new benchmarks for evaluating AI agents (BrowseComp) and AI's ability to replicate research (PaperBench).¹ OpenAI is a key partner in the US government's \$500B Stargate data center initiative, highlighting its central role in national AI infrastructure plans.⁷ The company secured a massive \$40 billion funding round in March 2025 at a \$300 billion valuation, providing substantial resources for research and scaling compute infrastructure for its reported 500 million weekly ChatGPT users.⁹ Concurrent leadership changes (Mark Chen to CRO, Brad Lightcap expanding COO role, Julia Villagra as CPO) were announced to manage this rapid scaling.³⁸ OpenAI also released a policy framework proposal in March, seeking to shape regulatory discussions.²⁷ Amidst these advances, CEO Altman acknowledged the high operational costs, stating the \$200/month o1 Pro subscription was loss-making due to heavy usage.³⁶ Altman also publicly expressed confidence that the company now "knows how to build AGI" and

predicted the arrival of AI agents in the workforce during 2025.³⁶

- **Key Personnel**

Sam Altman remains the prominent CEO and public face.³⁶ The March 2025 leadership update elevated Mark Chen (Chief Research Officer), Brad Lightcap (Chief Operating Officer), and Julia Villagra (Chief People Officer) to key executive roles.³⁸ Researcher Josh Achiam is noted for work on AI alignment..⁴⁰¹

- **Official/Trending Links**

- Official Website: <https://openai.com/>¹
- Research Page: <https://openai.com/research>⁴⁰
- OpenAI Charter: <https://openai.com/charter/>²¹
- Planning for AGI Blog:
<https://openai.com/index/planning-for-agi-and-beyond/>³⁷
- About Page: <https://openai.com/about/>³⁵
- Leadership Updates (Mar 2025):
<https://openai.com/index/leadership-updates-march-2025/>³⁸
- Funding Announcement (Mar 2025):
<https://openai.com/index/march-funding-updates/>⁹
- Altman AGI Claims/Blog Post News:
<https://www.therundown.ai/p/openai-now-knows-how-to-build-agi>³⁶,
<https://hyperight.com/artificial-general-intelligence-is-agi-really-coming-by-2025/>³⁹

- **Google DeepMind**

- **Aim**

Google DeepMind's stated mission is to "build AI responsibly to benefit humanity".² While not explicitly framing all work under an "AGI" banner as aggressively as OpenAI, their focus on tackling complex challenges and building increasingly general and capable AI models implies a trajectory towards highly advanced AI.²

- **2025 Initiatives**

DeepMind continued its rapid model development cadence in 2025. April saw the introduction of Gemini 2.5 Flash, described as their "first fully hybrid reasoning model," offering developers enhanced control.² This was closely followed by Veo 2, a state-of-the-art video generation model integrated into

the Gemini Advanced platform.² The landmark AlphaFold project, which revolutionized protein structure prediction, remains a key focus, demonstrating AI's potential for scientific acceleration.² Research continues on Project Astra, a prototype exploring the capabilities of a future "universal AI assistant".² DeepMind also published updates to its DiLoCo (Distributed Low-Communication Training) methodology in March, aiming to make large model training more efficient, particularly across less robust networks.²⁷ Strategic partnerships were also highlighted, including an expanded collaboration with Nvidia focusing on using AI and simulation for robotics, drug discovery, and energy grid optimization⁴¹, and the introduction of AI-powered TV news summaries using Gemini technology.³⁶

○ **TPA Chip Developments**

Google is also heavily invested in developing its own custom-designed chips, called Tensor Processing Units (TPUs). The latest generation, "Ironwood," is designed to accelerate AI applications. Google's strategy with TPUs is to create hardware specifically optimized for the unique demands of its AI workloads. This approach offers several potential benefits:

- **Performance Optimization:** TPUs are designed from the ground up to excel at the matrix computations that are fundamental to machine learning, potentially offering superior performance compared to general-purpose GPUs like those from NVIDIA for specific AI tasks.
- **Energy Efficiency:** Custom-designed chips can be more energy-efficient, reducing the operational costs of running large AI models. This is a critical factor for Google, given the scale of its data centers.
- **Cost Control:** By developing its own hardware, Google aims to reduce its reliance on external chip suppliers, such as NVIDIA. This can provide greater control over costs and supply chains, especially as the demand for AI hardware continues to surge.
- **Software Integration:** Google can tightly integrate its TPUs with its software stack, including TensorFlow and other AI frameworks. This co-design approach can lead to further performance improvements and a more seamless development experience.

In essence, Google's TPA strategy is about gaining greater control, efficiency, and performance in its AI infrastructure. While NVIDIA remains a dominant player in the AI hardware market, Google's investment in TPUs reflects a long-term vision of building a vertically integrated AI platform. Google is not necessarily trying to replace NVIDIA in the broader market but is focusing on optimizing its own infrastructure for its specific AI needs. This allows Google to push the boundaries of AI research and deployment while

maintaining a degree of independence and cost-effectiveness.

- **Key Personnel**

Demis Hassabis continues to lead as Co-founder and CEO.² William Isaac serves as Head of Ethics Research and sits on the Partnership on AI's Policy Steering Committee.³⁴ Researchers Douglas Eck (AI and creativity) and Anca Dragan (AI safety) were featured in DeepMind's podcast series.² Tim Brooks, formerly of OpenAI's Sora team, joined DeepMind to build a team focused on AI world simulation for applications like visual reasoning and embodied agents.³⁶

- **Official/Trending Links**

- Official Website: <https://deepmind.google/> ²
- About Page: <https://deepmind.google/about/> ²
- Davos Report (AlphaFold mention):
<https://www.ommax.com/en/insights/industry-insights/a-new-ai-era-top-10-takeaways-from-davos-2025/> ⁷
- DiLoCo Update News:
<https://champaignmagazine.com/2025/03/16/ai-by-ai-weekly-top-5-03-10-16-2025/> ²⁷
- Nvidia Partnership News:
https://nationaltechnology.co.uk/Meta_Chief_AI_Scientist_Claims_AGI_Will_Be_Viable_In_3_5_Years.php ⁴¹

- **Anthropic**

- **Aim**

Anthropic distinguishes itself with an explicit focus on AI safety and aligning AI development with long-term human well-being.³ Founded by former OpenAI members concerned about safety directions⁴², the company champions approaches like Constitutional AI (training models based on principles rather than just data or human feedback).¹³ CEO Dario Amodei emphasizes prioritizing safety before capability, a stance often seen as counter-cultural in Silicon Valley.⁴³

- **2025 Initiatives**

Anthropic released Claude 3.7 Sonnet in February 2025, described as its most intelligent model to date.³ Development continues within the Claude series,

which includes models like 3.5 Haiku and 3 Opus.³ The company actively participates in safety-focused international discussions, with representation at the AI Safety Connect event in Paris in February 2025.³² Anthropic's "Constitutional AI" framework remains a key part of its research and product philosophy.¹³

- **Key Personnel**

Dario Amodei serves as CEO and Co-founder, frequently speaking on AI safety and the company's philosophy.⁴² His sister, Daniela Amodei, is also a co-founder.⁴⁴ Chris Olah, another co-founder, is recognized as a pioneer in mechanistic interpretability research.⁴² Michael Sellitto represented Anthropic at the AI Safety Connect event.³²

- **Official/Trending Links**

- Official Website: <https://www.anthropic.com/>³
- TIME100 Feature on Dario Amodei:
<https://time.com/collections/100-most-influential-people-2025/7273747/dario-amodei/>⁴³
- Dario Amodei at CFR Event: <https://on.cfr.org/4iydDIW>⁴²
- Dario Amodei Wikipedia: https://en.wikipedia.org/wiki/Dario_Amodei⁴⁴
- Dario Amodei Personal Website: <https://www.darioamodei.com/>⁴⁵
- AI Safety Connect Event Details:
<https://www.aisafetyconnect.com/event-details>³²
- <https://www.anthropic.com/company>⁴⁶.

- **Meta AI**

- **Aim**

Meta AI, under the influence of Chief AI Scientist Yann LeCun, champions open-source AI development, particularly through its Llama model family.⁴¹ LeCun expresses skepticism about the term "AGI" and the longevity of current generative AI paradigms (LLMs), preferring the term "Advanced Machine Intelligence" (AMI) or focusing on AI systems that understand the physical world ("world models") for applications like robotics.⁴¹ He predicts a paradigm shift within 3-5 years.⁴⁷

- **2025 Initiatives**

Meta AI made a significant stride in open-source AI with the release of the Llama 4 family in April 2025.⁴ This release includes models like Llama 4 Scout

(109B total parameters) and Llama 4 Maverick (400B total parameters), notable for being Meta's first models using a Mixture of Experts (MoE) architecture for efficiency.⁴⁹ Llama 4 features native multimodality (integrating text, image, video understanding from the start) and an industry-leading context window of up to 10 million tokens (initially supported up to 131k on some platforms).⁴ Llama 4 was quickly made available on platforms like Cloudflare Workers AI⁴⁹ and Snowflake Cortex AI⁵⁰, demonstrating Meta's push for broad adoption. Ongoing research efforts also focus on perception, localization, and reasoning.⁵²

○ **Achieving AGI**

During the NVIDIA GTC 2025 event¹⁰³, Yann LeCun emphasized that achieving AGI (Artificial General Intelligence) will require fundamentally new architectures, moving beyond today's large language models (LLMs). In particular, he outlined the importance of Joint Embedding Predictive Architectures (JEPA), which operate in abstract representation spaces rather than discrete token spaces. LeCun argued that true reasoning and planning must occur in these abstract spaces to model the physical world effectively—something LLMs, limited by their token-based training, cannot accomplish. He projects that small-scale success in JEPA-based models could emerge within three to five years, leading to scalable paths toward advanced machine intelligence (AMI). LeCun also stressed the need for open research and global collaboration, underlining that no single entity will achieve AGI alone.

○ **Key Personnel**

Yann LeCun, as Chief AI Scientist, is a highly influential figure shaping Meta's AI direction and public discourse..^{41 52}

○ **Official/Trending Links**

- Official Meta AI Website: <https://ai.meta.com/>⁵⁴
- Meta AI Blog: <https://ai.meta.com/blog/>⁵²
- Yann LeCun Statements News:
 - https://www.youtube.com/watch?v=eyrDM3A_YFc¹⁰³
 - https://nationaltechnology.co.uk/Meta_Chief_AI_Scientist_Claims_AGI_Will_Be_Viable_In_3_5_Years.php⁴¹,
 - <https://www.hpcwire.com/2025/02/11/metasp-chief-ai-scientist-yann-lecun-questions-the-longevity-of-current-genai-and-llms/>⁴⁷,
 - <https://themunicheye.com/metasp-ai-chief-questions-generative-ais-future-10206>⁴⁸
- Llama 4 Technical Details/News:
 - <https://redblink.com/llama-4-vs-deepseek-v3/>⁴,
 - <https://blog.cloudflare.com/meta-llama-4-is-now-available-on-workers-ai/>

⁴⁹,

<https://www.snowflake.com/en/blog/meta-llama-4-now-available-snowflake-cortex-ai/> ⁵⁰, <https://ai.meta.com/blog/llama-4-multimodal-intelligence/> ⁵¹

• **Baidu AI**

○ **Aim**

As China's leading search engine and AI company, Baidu aims to compete directly with global frontier labs like OpenAI and domestic rivals like DeepSeek, establishing leadership in the rapidly growing Chinese AI market.⁶ The focus is on developing powerful, versatile foundation models within its Ernie series.

○ **2025 Initiatives**

Baidu significantly escalated its competitive efforts in March 2025 by launching two advanced models: Ernie 4.5, a multimodal model with enhanced language ability and claimed high "EQ" for understanding nuances like memes ⁵, and Ernie X1, its first dedicated reasoning model designed to use tools autonomously and compete directly with DeepSeek R1 on performance and cost.⁵ Baidu plans further advancement with Ernie 5, targeting multimodal capabilities, scheduled for release in the second half of 2025.²⁸ Aligning with a trend towards openness, Baidu intends to open-source the Ernie codebase later in the year.⁶ These model developments are supported by continued investment in Baidu's AI Cloud infrastructure and applications in areas like autonomous driving (Apollo).³¹

○ **Key Personnel**

Robin Li, Baidu's CEO, provides high-level direction and commentary on the competitive landscape.²⁸ The leadership within Baidu Research drives the technical development of the Ernie models and other AI initiatives.³³

○ **Official/Trending Links**

- Baidu Research: <https://research.baidu.com/> ³¹
- Ernie 5 Announcement News: <https://dig.watch/updates/baidu-to-launch-ernie-5-ai-in-2025> ²⁸
- Ernie 4.5/X1 Launch News: <https://nascenia.com/latest-ai-models/> ⁵, <https://siliconangle.com/2025/03/16/baidu-debuts-first-ai-reasoning-model-competes-deepseek/> ⁶, <https://champaignmagazine.com/2025/03/16/ai-by-ai-weekly-top-5-03-10-16-2025/> ²⁷

- **Microsoft AI**

- **Aim**

Microsoft's strategy revolves around integrating AI deeply into its existing product ecosystem (Azure, Office, GitHub) and providing comprehensive AI platforms and services for enterprise customers.³¹ It achieves this through both internal development and strategic partnerships with leading AI labs, most notably OpenAI, but also others like Mistral AI.³¹ A significant focus is placed on responsible AI practices.²²

- **2025 Initiatives**

Microsoft committed a substantial \$80 billion to AI development in 2025.⁷ It continues to enhance its Azure AI platform for enterprise machine learning and analytics.³¹ The company is a key partner alongside OpenAI in the US government's Stargate initiative, providing critical cloud infrastructure.⁷ Microsoft also released its own powerful lightweight AI model designed to run efficiently on standard CPUs, broadening accessibility.²⁸ Reflecting its global presence and commitment to AI governance dialogue, Microsoft served as a Diamond Sponsor for the ITU's AI for Good Global Summit 2025.¹⁴ Its responsible AI frameworks and practices are often cited as examples for corporate compliance.²²

- **Key Personnel**

Natasha Crampton holds the position of Chief Responsible AI Officer and is active in global policy discussions, serving on the Partnership on AI's Policy Steering Committee.³⁴

- **Official/Trending Links**

- Microsoft AI Overview (via AI Superior):
<https://aisuperior.com/ai-research-companies/>³¹
- Davos Report (Investment, Stargate):
<https://www.ommax.com/en/insights/industry-insights/a-new-ai-era-top-10-takeaways-from-davos-2025/>⁷
- ITU AI for Good Summit (Sponsorship):
<https://www.itu.int/en/mediacentre/Pages/PR-2025-02-06-AI-for-Good-2025-announcement.aspx>¹⁴
- CogentInfo (Responsible AI Practices):
<https://www.cogentinfo.com/resources/federal-ai-mandates-and-corporate-compliance-whats-changing-in-2025>²²

- **Nvidia**

- **Aim**

As the dominant provider of AI hardware (GPUs like H100, A100) and associated software platforms (DGX Cloud, Omniverse for simulation, Isaac for robotics), Nvidia plays a crucial enabling role in the entire AI ecosystem.²⁹ Its aim is to provide the foundational tools and infrastructure that power the AGI race.

- **2025 Initiatives**

Nvidia continues to deepen its integration across the AI landscape. It announced expanded collaborations with Google DeepMind for advancing robotics, drug discovery, and other scientific domains⁴¹, and with Oracle Cloud Infrastructure to accelerate agentic AI applications.⁴¹ At its GTC 2025 conference, Nvidia unveiled the Isaac GR00T N1 project, aiming to provide a general-purpose foundation model for humanoid robots, signaling a major push into embodied AI.⁴¹ Development continues on its core GPU technologies and AI software frameworks like RAPIDS and Triton.³¹

- **Key Personnel**

Jensen Huang, the company's CEO, remains the driving force and key spokesperson.⁴¹

- **Official/Trending Links**

- Nvidia AI Overview (via AI Superior):
<https://aisuperior.com/ai-research-companies/>³¹
 - GTC 2025 News (Partnerships, Robotics):
https://nationaltechnology.co.uk/Meta_Chief_AI_Scientist_Claims_AGI_Will_Be_Viable_In_3_5_Years.php⁴¹

- **xAI**

- **Aim**

Founded by Elon Musk, xAI explicitly aims to develop AGI and compete with the established leaders like OpenAI and Google DeepMind.¹⁸

- **2025 Initiatives**

The company released its Grok 3 model in the first half of 2025, positioning it as a competitor to models like GPT-4.5 and DeepSeek R1, highlighting

capabilities in reading comprehension, complex problem-solving, and coding.⁵ xAI has secured significant funding, listed at \$12.13 billion in the Forbes AI 50 list.²⁰

- **Key Personnel**

Elon Musk leads the company.

- **Official/Trending Links**

- Forbes AI 50 List: <https://www.forbes.com/lists/ai50/>²⁰
- Nascenia AI Models Overview: <https://nascenia.com/latest-ai-models/>⁵

1.3 Emerging Players and Investment Landscape

Beyond the established giants, a dynamic ecosystem of startups and specialized companies is contributing significantly to the AI landscape in 2025.

Notable AI Companies

The field is diversifying rapidly. Key players gaining prominence in 2025 include:

- **DeepSeek AI (China)**

Emerged as a major challenger with its R1 reasoning model setting high benchmarks late 2024/early 2025, and its V3 model competing strongly with Meta's Llama 4.⁴

- **Mistral AI (France)**

A leading European player focused on open-source models, securing partnerships with major tech companies like Microsoft.²⁰

- **Cohere (Canada)**

Another significant developer of large language models, particularly focused on enterprise applications.²⁰

- **Specialized Players**

Companies focusing on specific niches are also attracting attention and funding, such as Midjourney (US, image generation)²⁰, Stability AI (UK, image generation)²⁹, ElevenLabs (UK, voice generation)²⁰, Figure AI (US, humanoid robots)²⁰, Sakana AI (Japan, novel AI architectures)²⁰, and various AI infrastructure and tooling providers (see below). The Forbes AI 50 list provides a broader snapshot of influential private AI companies in 2025.²⁰

- **Investment Trends**

Corporate investment in AI saw a strong rebound in 2024, continuing into 2025. The US dominated private AI investment with \$109.1 billion in 2024, vastly outpacing China (\$9.3B) and the UK (\$4.5B).⁸ Generative AI remained a particularly hot area, attracting \$33.9 billion globally in private investment in 2024, an 18.7% increase from 2023.⁸ This investment surge aligns with accelerating business adoption; 78% of organizations reported using AI in 2024, a significant jump from 55% the previous year.⁸ AI dominated the venture capital narrative and dealmaking in 2024.¹⁸

- **Compute Infrastructure Focus**

The immense computational requirements for training and deploying frontier AI models have made infrastructure a critical bottleneck and investment area. The \$500 billion Stargate initiative exemplifies this focus at a national level.⁷ Concurrently, a cohort of specialized infrastructure providers has emerged to meet the demand, including companies like Crusoe Energy (\$2.8B valuation), Lambda (\$2.5B valuation), and Together AI (\$3.3B valuation), which provide AI-focused cloud services and hardware.²⁰ This underscores the recognition that progress towards AGI is fundamentally tied to advancements and investments in underlying compute infrastructure.³⁰

Section 1 Synthesis: Dynamics of the 2025 AGI Race

The dynamics of the AGI race in 2025 reveal a potent feedback loop. Massive investments, exemplified by the \$500 billion Stargate initiative⁷ and substantial corporate funding like OpenAI's \$40 billion round⁹, directly fuel the rapid development of increasingly powerful models such as OpenAI's o-series, Google's Gemini 2.5, and Meta's Llama 4.¹ This demonstrated progress, in turn, attracts further capital and intensifies competitive pressures. This accelerating cycle, while driving innovation at an unprecedented pace, concurrently raises concerns, acknowledged even within leading labs¹⁰, about the potential marginalization of safety precautions and ethical considerations in the pursuit of strategic advantage or market dominance. The sheer scale of investment, particularly government-backed initiatives like Stargate, elevates AGI development beyond typical corporate R&D to the level of critical national infrastructure, implying a perceived urgency and strategic importance that could rationalize cutting corners on safety.

While the United States currently maintains a lead in the quantity of frontier models developed and the overall level of private investment⁸, the competitive landscape is far

from static. China, powered by significant state support ⁸ and aggressive development from companies like Baidu and DeepSeek ⁵, is rapidly closing the performance gap in model quality.⁸ Furthermore, China continues to lead in the volume of AI-related publications and patents.⁸ The emergence of highly performant, potentially lower-cost models from Chinese labs ⁵ introduces a new dynamic that could disrupt global market structures and accelerate the proliferation of advanced AI capabilities. This combination of improving quality, high research output, and potential cost advantages indicates China is building a strong foundation to challenge US dominance, possibly leading to bifurcated AI ecosystems or altered global competitive balances in the near future.

Adding another layer of complexity, the very definition and pursuit of "AGI" remain contested concepts in 2025. While organizations like OpenAI, Google DeepMind, and Anthropic drive towards increasingly general and autonomous systems, often explicitly framing their goals in terms of AGI or superintelligence ², influential figures like Meta's Yann LeCun actively question this framing.⁴¹ LeCun advocates for focusing on "world models" capable of understanding physical reality and enabling advanced robotics (termed "AMI" or Advanced Machine Intelligence), suggesting current LLM-based approaches have fundamental limitations.⁴⁷ This divergence indicates that the "race" is not towards a single, agreed-upon target but encompasses multiple, potentially conflicting, research programs pursuing different architectures and end goals under the broad umbrella of advanced AI. The term "AGI" itself is sometimes employed more as a marketing tool or investment pitch rather than a precise technical specification.¹⁸

(Table 1: Leading Frontier AI Labs - 2025 Snapshot)

Organization	Stated Aim (AGI/Superintelligence Focus)	Key 2025 Models/Initiatives	Key Personnel (CEO/Lead Scientist)	Notable 2025 Funding/Partnerships	Primary Link
OpenAI	Build safe & beneficial AGI/Superintelligence for all humanity ²¹	o3, o4-mini, GPT-4.1 API, 4o Image Gen, Policy Framework, Stargate Partner ¹	Sam Altman (CEO), Mark Chen (CRO)	\$40B Funding Round (@ \$300B val) ⁹	https://openai.com/ ¹

AGI Landscape and Organizations: A 2025 Intelligence Briefing

Research by Fede Nolasco | AI Researcher and Data Architect

<https://www.linkedin.com/in/federiconolasco>

Report released on 22 April 2025

Organization	Stated Aim (AGI/Superintelligence Focus)	Key 2025 Models/Initiatives	Key Personnel (CEO/Lead Scientist)	Notable 2025 Funding/Partnerships	Primary Link
Google DeepMind	Build AI responsibly to benefit humanity; general & capable AI ²	Gemini 2.5 Flash, Veo 2, AlphaFold, Project Astra, DiLoCo Update ²	Demis Hassabis (CEO)	Nvidia Partnership (Robotics, Science) ⁴¹	https://deepmind.google/ ²
Anthropic	Build safe AI aligned with human well-being; Constitutional AI ³	Claude 3.7 Sonnet released; Claude 3 series development ³	Dario Amodei (CEO)	Active in AI Safety Events ³²	https://www.anthropic.com/ ³
Meta AI	Open-source AI (Llama); World Models/AMI focus; Skeptical of LLMs ⁴¹	Llama 4 family (Scout, Maverick, Behemoth) - MoE, Multimodal, 10M context ⁴	Yann LeCun (Chief AI Scientist)	Llama 4 on Cloudflare, Snowflake ⁴⁹	https://ai.meta.com/ ⁵⁴
Baidu AI	Compete globally/domestically; Ernie foundation models ⁶	Ernie 4.5 & X1 launched; Ernie 5 planned; Ernie Open-Source plan ²⁸	Robin Li (CEO)	\$47.5B China Semiconductor Fund (Context) ⁸	https://research.baidu.com/ ³¹
DeepSeek AI	Challenge leaders with high-performance,	R1 (Reasoning), V3 models setting	Leadership Team	Gaining significant attention/adoption ⁶	https://www.deepseek.com/ (Implied, not in

AGI Landscape and Organizations: A 2025 Intelligence Briefing

Research by Fede Nolasco | AI Researcher and Data Architect

<https://www.linkedin.com/in/federiconolasco>

Report released on 22 April 2025

Organization	Stated Aim (AGI/Superintelligence Focus)	Key 2025 Models/Initiatives	Key Personnel (CEO/Lead Scientist)	Notable 2025 Funding/Partnerships	Primary Link
	efficient models ⁵	benchmarks, competing with Llama 4 ⁴			snippets)
Microsoft AI	Integrate AI across products; Enterprise AI solutions; Partnerships ³¹	Azure AI, GitHub Copilot, Lightweight CPU model, Stargate Partner ⁷	Natasha Crampton (Chief Resp. AI Officer)	\$80B AI Dev. Allocation ⁷ ; ITU Summit Sponsor ¹⁴	https://www.microsoft.com/ai (Implied)
Nvidia	Enable AI ecosystem with hardware/software platforms ²⁹	H100/A100 GPUs, DGX Cloud, Omniverse, Isaac GR00T N1 (Robotics) ³¹	Jensen Huang (CEO)	Partnerships with Google, Oracle ⁴¹	https://www.nvidia.com/ (Implied)
xAI	Develop AGI, compete with top labs ¹⁸	Grok 3 model released ⁵	Elon Musk (Founder)	Significant VC Funding (\$12.13B reported) ²⁰	https://x.ai/ (Implied)

Section 2: Architecting the Post-AGI World: Governance, Alignment, and Adaptation (2025)

Overview

Running parallel to the accelerating capabilities race is an increasingly urgent and complex global effort to architect the governance structures, technical alignment methodologies, and socio-economic adaptation strategies necessary for navigating a world potentially transformed by AGI.¹⁵ The year 2025 sees heightened activity in this domain, with international bodies, research institutions, national governments, and civil society organizations intensifying efforts to establish norms, standards, and policies, though coherence remains a challenge.

2.1 Global Governance and Standards Initiatives

International cooperation and standardization are recognized as crucial for managing the transnational nature of AI development and deployment.

- **International Organizations (ISO, IEC, ITU)**

These established international bodies are actively collaborating to develop global AI standards. Their stated aim is to create standards that support policy goals, ensure responsible AI use, promote interoperability, and help bridge the significant global AI governance gap identified by ITU surveys (which found 55% of member states lack a national AI strategy and 85% lack AI-specific regulations).¹²

- **2025 Initiatives**

A major joint initiative announced is the International AI Standards Summit, scheduled for December 2-3, 2025, in Seoul, hosted by the Korean Agency for Technology and Standards (KATS). This summit, involving the International Organization for Standardization (ISO), the International Electrotechnical Commission (IEC), and the International Telecommunication Union (ITU), directly responds to calls from the UN's High-level Advisory Body report ("Governing AI for Humanity") and the Global Digital Compact for advancing governance through international standards.¹² Separately, the ITU is hosting its AI for Good Global Summit in Geneva from July 8-11, 2025. This summit

focuses on the implications of "agentic AI," the governance gap, the role of standards, and aligning AI with Sustainable Development Goals.¹⁴ A key component of the AI for Good Summit is the second AI Governance Day on July 10, dedicated to safety, trust, standards, bridging the regulatory gap, and capacity building, especially in developing countries.¹⁴

- **Key Personnel**

Sergio Mujica (ISO Secretary-General) emphasized the need for a collaborative approach to AI governance through standards.¹² Seizo Onoe (Director, ITU Telecommunication Standardization Bureau) highlighted ITU's role in driving a trusted and interoperable AI ecosystem through standards.¹⁴ Prominent AI figures like Geoffrey Hinton, Yoshua Bengio, and Sasha Luccioni are slated to speak at the AI for Good Summit.¹⁴

- **Official/Trending Links**

- International AI Standards Summit Announcement (via ANSI):
<https://www.ansi.org/standards-news/all-news/2024/10/10-15-24-advancing-ai-standards-collaboration-iso-iec-and-itu-announce-ai-standards-summit>¹²
- ITU AI for Good Global Summit 2025:
<https://www.itu.int/en/mediacentre/Pages/PR-2025-02-06-AI-for-Good-2025-announcement.aspx>¹⁴, <https://aiforgood.itu.int/>¹⁴

- **AI Standards Hub (UK Initiative)**

This initiative, operating under the UK government's umbrella, aims to explore the critical role of standards in AI governance and, importantly, to foster global inclusiveness and collaboration in the standardization process.¹³

- **2025 Initiatives**

The Hub held its inaugural Global Summit in London (and online) on March 17-18, 2025. Organized in partnership with the OECD, the UN Office of the High Commissioner for Human Rights (OHCHR), and the Partnership on AI (PAI), the summit brought together diverse stakeholders.¹³ Key themes included the interplay between standards and regulation, promoting diversity and inclusion in standards development, building a robust AI assurance ecosystem, fostering collaboration between AI safety and standardization communities, and addressing governance challenges related to foundation models.¹³

- **Key Personnel**

Speakers and participants included leaders from partner organizations like Rebecca Finlay (CEO, PAI) and Karine Perset (OECD), alongside experts from academia, civil society, and government such as Laura Lazaro Cabrera (CDT Europe), Markus Anderljung, and others.¹³

- **Official/Trending Links**

- AI Standards Hub Global Summit 2025 Page:
<https://aistandardshub.org/global-summit/>¹³

- **OECD (Organisation for Economic Co-operation and Development)**

The OECD continues to be a central player in shaping AI policy discussions and providing analysis through its AI Policy Observatory (OECD.AI).¹⁵ It actively fosters international cooperation and dialogue.

- **2025 Initiatives**

The OECD was a key partner in the AI Standards Hub Global Summit in March 2025.¹³ Its OECD.AI platform remains a vital resource for tracking global AI trends and policies. The Partnership on AI lists the OECD among its key institutional collaborators.³⁴ The OECD's work is referenced as an important input for civil society governance roadmaps¹⁵, and its AI Observatory is slated to publish the Global Risk and AI Safety Preparedness (GRASP) mapping developed in partnership with MBRSG and GPAI.³²

- **Key Personnel**

Karine Perset heads the AI and Emerging Digital Technologies Division at OECD. She is highly active in the global AI governance community, serving on the PAI Policy Steering Committee and moderating sessions at events like the AI Safety Connect forum in Paris.¹³

- **Official/Trending Links**

- OECD.AI Policy Observatory: <https://oecd.ai/> (Implied Link)
- AI Standards Hub Summit (Partner):
<https://aistandardshub.org/global-summit/>¹³
- Partnership on AI Policy Page (Collaboration):
<https://partnershiponai.org/program/policy/>³⁴

- Future Society Report (Reference):
<https://thefuturesociety.org/cso-ai-governance-priorities/>¹⁵
- AI Safety Connect (Participation):
<https://www.aisafetyconnect.com/event-details>³²

• United Nations (UN)

The UN system is increasingly engaged in AI governance, aiming to align AI development with international human rights law, sustainable development goals, and global digital cooperation frameworks like the Global Digital Compact and the Declaration on Future Generations adopted at the 2024 Summit of the Future.¹²

○ 2025 Initiatives

The influential report "Governing AI for Humanity" from the UN High-level Advisory Body on AI continues to shape discussions, particularly regarding the role of international standards.¹² The UN OHCHR partnered with the AI Standards Hub for its March 2025 summit.¹³ Specialized UN agencies like the ITU are leading major initiatives such as the AI for Good Global Summit.¹⁴ The UN Interregional Crime and Justice Research Institute (UNICRI) is active in international AI safety cooperation discussions.³² The newly established UN Office for Digital and Emerging Technologies, led by Under-Secretary-General Amandeep Singh Gill, coordinates UN efforts in this space.⁵⁸ Think tanks like the Stimson Center are actively analyzing how to integrate UN frameworks like the Global Digital Compact and the Declaration on Future Generations into practical AI governance.⁵⁷

○ Key Personnel

Amandeep Singh Gill serves as the UN Under-Secretary-General and Special Envoy for Digital and Emerging Technologies.⁵⁸ Irakli Beridze represents UNICRI in safety forums.³² Key figures from various UN agencies contribute to specific initiatives.

○ Official/Trending Links

- UN High-level Advisory Body Report (Reference):
<https://www.un.org/en/ai-advisory-body> (Implied Link)
- Global Digital Compact / Declaration on Future Generations (References):
¹²
- AI Standards Hub Summit (Partner):
<https://aistandardshub.org/global-summit/>¹³

- ITU AI for Good Summit: <https://aiforgood.itu.int/>¹⁴
- Stimson Center Report on UN Frameworks:
<https://www.stimson.org/2025/governing-ai-for-the-future-of-humanity/>⁵⁷
- CAIDP Event Bios (Amandeep Singh Gill):
<https://www.caidp.org/events/washdc25aidv/bios/>⁵⁸
- AI Safety Connect (UNICRI Participation):
<https://www.aisafetyconnect.com/event-details>³²

2.2 Research Institutions & NGOs Shaping Policy and Ethics

Alongside intergovernmental bodies, a diverse range of research institutions and non-governmental organizations (NGOs) play crucial roles in shaping AI policy, ethics, and socio-economic considerations through research, advocacy, and multi-stakeholder convenings.

• Partnership on AI (PAI)

PAI operates as a global non-profit multi-stakeholder organization, bringing together over 100 partners from industry, academia, and civil society across 17 countries.¹³ Its mission is to foster a responsible AI ecosystem by facilitating coordination, developing evidence-based frameworks, and promoting shared understandings of best practices, explicitly stating it is not a trade group or lobbying organization.³⁴

○ 2025 Initiatives

PAI co-organized the AI Standards Hub Global Summit in March 2025.¹³ It continues extensive collaboration with key global institutions like the OECD, UN, G20, US government bodies (OSTP, NIST, NSF), AI Safety Institutes, and national governments.³⁴ In April 2025, PAI published work focusing on prioritizing responsible AI development in Africa.³⁴ A significant ongoing initiative is the Partnership for AI Evidence (PAIE), a collaboration with the Abdul Latif Jameel Poverty Action Lab (J-PAL) focused on generating rigorous evidence about AI's impact on social outcomes.⁵⁹

○ Key Personnel

Rebecca Finlay serves as CEO.¹³ The Policy Steering Committee draws high-profile members from across sectors, including Natasha Crampton (Microsoft), Arisa Ema (University of Tokyo), Alexandra Givens (Center for Democracy & Technology), William Isaac (Google DeepMind), Karine Perset

(OECD), Irene Solaiman (Hugging Face), and Alondra Nelson (Institute for Advanced Study).³⁴

- **Official/Trending Links**

- PAI Policy Program: <https://partnershiponai.org/program/policy/>³⁴
- AI Standards Hub Summit (Partner):
<https://aistandardshub.org/global-summit/>¹³
- Partnership for AI Evidence (PAIE) with J-PAL:
<https://www.povertyactionlab.org/initiative/partnership-ai-evidence>⁵⁹

- **Stanford HAI (AI Index)**

The Stanford Institute for Human-Centered Artificial Intelligence (HAI) produces the annual AI Index report, widely recognized as one of the most authoritative resources tracking global AI trends.⁸ Its aim is to provide objective, data-driven analysis across research, development, performance, investment, ethics, policy, and public opinion, serving as an independent source of insights.⁸

- **2025 Initiatives**

The 8th edition of the AI Index Report was released in April 2025. Key findings highlighted continued improvements in AI performance on benchmarks, record corporate investment (especially US-led and in generative AI), accelerating business adoption, the US leading in model quantity but China closing the quality gap, an unevenly evolving responsible AI ecosystem with rising incidents but also increased governance efforts, and rising global optimism about AI albeit with significant regional divides.⁸ The report noted significant performance gains on new challenging benchmarks (MMMU, GPQA, SWE-bench) and increasing real-world deployment in areas like healthcare (FDA approvals) and autonomous driving (Waymo, Baidu Apollo Go).²⁴ It also highlighted efficiency gains, with smaller models improving and inference costs dropping dramatically.²⁶

- **Key Personnel**

The report is produced by a team at Stanford HAI, with Vanessa Parli serving as Director of Research and an AI Index Steering Committee member.²⁶

- **Official/Trending Links**

- AI Index Report 2025 Landing Page:
<https://hai.stanford.edu/ai-index/2025-ai-index-report>²⁴

- AI Index Report Main Page: <https://aiindex.stanford.edu/report/> ⁶⁰
- AI Index 2025 - 10 Charts Summary:
<https://hai.stanford.edu/news/ai-index-2025-state-of-ai-in-10-charts> ²⁶
- News Coverage: <https://m.theblockbeats.info/en/news/57740> ⁵⁵,
<https://www.businesswire.com/news/home/20250407539812/en/Stanford-HAIs-2025-AI-Index-Reveals-Record-Growth-in-AI-Capabilities-Investment-and-Regulation> ²⁵
- (Note: Link to PDF ⁸
https://hai-production.s3.amazonaws.com/files/hai_ai_index_report_2025.pdf might require direct access).

- **J-PAL (Partnership for AI Evidence - PAIE)**

This initiative, a collaboration between the Abdul Latif Jameel Poverty Action Lab (J-PAL) at MIT and the Partnership on AI, focuses specifically on using rigorous research methods (primarily randomized controlled trials) to identify, evaluate, and scale AI applications that demonstrably improve social outcomes and reduce poverty.⁵⁹ It aims to bridge the gap between AI technological potential and evidence-based social impact.

- **2025 Initiatives**

PAIE launched its Spring 2025 Request for Proposals (RFP), soliciting proposals for full research projects and pilot studies evaluating AI interventions. While open to all sectors, the RFP anticipates innovations primarily in areas with rapid AI adoption or significant potential impact, including education, health, labor markets, climate change, and financial inclusion.⁵⁹ The deadline for Letters of Interest was April 22, 2025, with full proposals due May 27, 2025.⁵⁹ PAIE highlights ongoing and past research projects using AI/ML, such as evaluating AI tutors in Brazil and Canada, AI-driven health screening in India, AI mobile health platforms in Kenya, job recommender systems in France, and using ML for predicting loan performance in Egypt.⁵⁹

- **Key Personnel**

The initiative is co-chaired by Iqbal Dhaliwal (Global Executive Director, J-PAL) and Sendhil Mullainathan (Professor, University of Chicago Booth School of Business).⁵⁹ It involves a network of affiliated researchers known for work at the intersection of economics, data science, and social policy, including Daron Acemoglu, David Autor, Jens Ludwig, Christopher Neilson, Esther Duflo, Rohini Pande, and others.⁵⁹

- **Official/Trending Links:**

- Partnership for AI Evidence (PAIE) Initiative Page:
<https://www.povertyactionlab.org/initiative/partnership-ai-evidence>⁵⁹

- **Centre for the Governance of AI (GovAI)**

Based in the UK, GovAI aims to help humanity navigate the transition to advanced AI through research and advising decision-makers across government, industry, and civil society.⁶¹ Its research agenda covers AI regulation, responsible development practices, compute governance, international governance, technical governance, risk assessment, and forecasting.⁶¹

- **2025 Initiatives**

GovAI is running its Summer Fellowship program from June 9 to August 31, 2025, in London. This fully funded, three-month program provides early-career individuals and professionals transitioning into the field with mentorship, research opportunities, and networking within the AI governance community.⁶¹ The application deadline was January 5, 2025.⁶¹ GovAI also offers year-long Research Scholar visiting positions for more established researchers and practitioners to pursue policy, social science, technical research, or applied projects.⁶² The center continues to publish research and provide advice based on its expertise.⁶¹

- **Key Personnel**

The GovAI team and its affiliate network provide supervision and mentorship for fellows and scholars.⁶¹ Alumni of its programs have moved into influential roles in governments (US, EU, UK), top AI companies (DeepMind, OpenAI, Anthropic), think tanks (CSET, RAND), and universities (Oxford, Cambridge).⁶³

- **Official/Trending Links:**

- GovAI Summer Fellowship 2025:
<https://www.governance.ai/post/summer-fellowship-2025> ⁶³,
<https://www.scholardigger.com/post/govai-center-of-ai-governance-summer-fellowship-2025> ⁶¹,
<https://www.opportunit4u.com/2024/12/govai-summer-fellowship-2025-in-london-uk.html> ⁶⁴
- GovAI Research Scholar Program:
<https://www.governance.ai/post/research-scholar> ⁶²
- GovAI Main Website: <https://www.governance.ai/> (Implied from post URLs)

- **Centre for the Study of Existential Risk (CSER)**

An interdisciplinary research center at the University of Cambridge, CSER focuses on studying and mitigating large-scale risks that could lead to human extinction or civilizational collapse, including those posed by advanced AI, alongside biological, environmental, and nuclear risks.⁶⁵

- **2025 Initiatives**

CSER published its March 2025 newsletter, welcoming new Director S.M. Amadae and highlighting recent research on extinction risk causes, military AI, and the future of global risk science.⁶⁵ The TERRA project, a bibliography of existential risk research using crowdsourcing and ML, was archived in March 2025.⁶⁵ CSER held an Ethics and Existential Risk Studies Seminar in March 2025.⁶⁵ The Centre launched a new MPhil program in Global Risk and Resilience and was hiring a Teaching Associate (application deadline April 27, 2025).⁶⁸ An upcoming seminar on Catastrophic Risks in and from the Arctic is scheduled for April 29, 2025.⁶⁷

- **Key Personnel**

S.M. Amadae became the new Director in March 2025, bringing expertise in nuclear war, climate change, and AI's impact on governance.⁶⁵ Seán Ó hÉigeartaigh, the former Director, published a review on human extinction causes in March 2025.⁶⁵ Other researchers like SJ Beard, Nathaniel Cooke, and Sarah Dryhurst published on the future of global risk science.⁶⁵ Haydn Belfield, associated with CSER, participated in the AI Safety Connect event.³²

- **Official/Trending Links**

- CSER Official Website: <https://www.cser.ac.uk/> ⁶⁷
- CSER Work/Publications Page: <https://www.cser.ac.uk/work/> ⁶⁵
- CSER Bluesky Social Media Profile:
<https://web-cdn.bsky.app/profile/cser.bsky.social> ⁶⁸
- CSER Profile on Nuclear Weapons Info:
<https://nuclearweapons.info/organization/the-centre-for-the-study-of-existential-risk/> ⁶⁶

2.3 National and Regional Policy Dynamics

Governance efforts are not only happening at the international level but also within specific national and regional contexts, sometimes leading to divergent approaches.

- **US Policy Landscape**

The US AI policy environment experienced a significant jolt in January 2025 with the signing of a new Executive Order by President Trump, explicitly revoking previous AI governance policies enacted under the Biden administration.²² This new order prioritizes deregulation to foster faster innovation and economic growth, shifting away from the previous emphasis on structured governance and risk mitigation.²² Despite this executive shift, there remains anticipation of potential new federal mandates emerging in 2025, driven by legislative proposals and ongoing agency work. These anticipated regulations focus on critical areas like transparency (requiring disclosure of model decision-making, training data, limitations), bias mitigation (detecting and eliminating biases, potentially requiring audits for recruitment tools), explainability (making AI systems interpretable, possibly via the AI Research, Innovation, and Accountability Act), and privacy protections (safeguarding personal data, exemplified by the proposed American Privacy Rights Act and the STOP Spying Bosses Act addressing workplace surveillance).²² The ongoing debate and development are reflected in the numerous policy comments submitted by organizations like MIRI throughout 2024 and into 2025 on various Requests for Information (RFIs) and Requests for Comment (RFCs) from agencies like NIST, BIS, OMB, and NTIA, covering topics such as AI Safety Institute guidance, risk management frameworks, procurement, and open model weights.⁶⁹

- **EU Context**

While the US landscape shifts towards deregulation, the European Union is in the implementation phase of its comprehensive AI Act, which takes a risk-based approach to regulation. The effective implementation relies heavily on the development of harmonized standards, involving European standards bodies like CEN/CENELEC.³⁴ The European Commission remains actively engaged in international safety discussions, participating in events like the AI Safety Connect forum.³²

- **Global South Perspectives**

There is a growing emphasis on ensuring that AI governance discussions and frameworks are inclusive and address the specific needs and contexts of the Global South. The Future Society's 2025 survey of civil society organizations highlighted strengthening Global South representation as a key priority for inclusive governance.¹⁵ Concerns exist that governance models developed without adequate participation risk overlooking unique socio-economic contexts, potentially causing unintended harm or undermining trust.¹⁵ Initiatives like PAI's focus on Responsible AI in Africa reflect an effort to address this gap.³⁴

2.4 Addressing Socio-Economic Impacts

While governance and technical safety receive significant attention, the profound socio-economic transformations potentially triggered by AGI are also a growing area of concern, though perhaps less developed in terms of concrete policy responses.

- **Focus**

Key concerns revolve around labor market disruption, including widespread job automation and the need for significant workforce adaptation.⁷ Predictions discussed at Davos 2025 suggested 92 million jobs could disappear by 2030, offset by the emergence of 170 million new roles, particularly in AI, data science, and related fields.⁷ Addressing potential increases in inequality and ensuring equitable benefit distribution are also critical challenges.

- **Organizations/Researchers**

Academic and research institutions are beginning to focus more intently on these issues. The J-PAL Partnership for AI Evidence (PAIE) is actively funding rigorous research into AI's impact on labor markets, education, and financial inclusion.⁵⁹ Prominent economists like Daron Acemoglu and David Autor are leading research on the relationship between AI adoption, job vacancies, and worker skills.⁵⁹ Ethical dimensions are also being explored, for instance, by the Emory University Center for Ethics, which examined AI's impact on the universal right to work in its 2025 student simulation program.⁷¹ Incorporating sustainability, equity, and labor protections into AI governance structures was identified as a priority by civil society organizations surveyed by The Future Society.¹⁵

• Official/Trending Links

* OMMAX Davos 2025 Report (Job Market):

<https://www.ommax.com/en/insights/industry-insights/a-new-ai-era-top-10-takeaways-from-davos-2025/>⁷

* J-PAL PAIE Initiative (Labor Research):

<https://www.povertyactionlab.org/initiative/partnership-ai-evidence>⁵⁹

* Emory Ethics Center (Right to Work):

https://news.emory.edu/stories/2025/04/er_ai_ethics_liaison_14-04-2025/story.htm

I⁷¹ * Future Society Report (Labor Protections):

<https://thefuturesociety.org/cso-ai-governance-priorities/>¹⁵

Section 2 Synthesis: The Fragmented Push for Order

The year 2025 demonstrates a marked acceleration in global efforts to govern AI, moving beyond high-level principles towards more concrete initiatives like dedicated summits, standards development processes, and specific policy frameworks.¹² This surge in activity involves a wide array of actors, including established international organizations (UN, OECD, ISO/IEC/ITU), national governments, research institutions, and NGOs. However, this proliferation of activity also highlights a significant challenge: fragmentation. Multiple bodies are pursuing parallel or overlapping mandates, creating a complex landscape where coordination and regulatory interoperability – key goals mentioned by groups like the AI Standards Hub and PAI¹³ – become paramount to avoid conflicting standards or duplicated efforts. Achieving effective global governance requires navigating this intricate web of initiatives.

A fundamental tension persists in 2025 regarding the primary approach to AI governance. On one hand, there is a push for legally binding regulations and top-down mandates, exemplified by the EU AI Act's implementation and anticipated US federal requirements concerning transparency and bias.²² On the other hand, there is significant emphasis on industry involvement through standards development¹² and the promotion of corporate responsible AI principles and self-regulatory mechanisms like ethics boards and internal audits.²² The abrupt shift in US federal policy in early 2025 towards deregulation²² adds another layer of complexity, potentially creating significant divergence between the US approach and regions favoring stricter controls, thereby complicating international alignment efforts.

While governance frameworks and technical alignment strategies are receiving considerable attention and resources, the critical area of socio-economic adaptation appears relatively less developed in terms of coordinated, large-scale policy responses in 2025. Despite widespread acknowledgment of AI's potential for profound disruption

to labor markets and potential exacerbation of inequality ⁷, and the emergence of dedicated research initiatives like J-PAL PAIE ⁵⁹, concrete, comprehensive policy solutions seem nascent compared to the flurry of activity around governance and standards. This suggests a potential lag in preparedness for the societal transformations that advanced AI could unleash, representing a critical area requiring greater focus and investment moving forward.

Section 3: Probing AI Consciousness: Scientific and Ethical Frontiers (2025)

Overview

The question of whether artificial intelligence can possess consciousness, sentience, or subjective experience is transitioning in 2025 from the realm of philosophical speculation and science fiction into a more pressing scientific and ethical inquiry.⁷² This shift is driven by the rapidly increasing sophistication of AI systems, particularly large language models exhibiting human-like communication abilities, and the perceived proximity of Artificial General Intelligence.⁷² While defining and detecting consciousness remains profoundly challenging even in biological systems⁷², the potential emergence of conscious AI necessitates deeper exploration of its indicators, implications, and ethical dimensions.

3.1 Leading Research Centers and Debates

Academic institutions and research centers are increasingly dedicating resources and convening experts to explore AI consciousness and its ethical ramifications.

- **Oxford Institute for Ethics in AI**

This institute at the University of Oxford aims to be a leading global hub for AI ethics, focusing on translating uncertainty into actionable solutions through interdisciplinary work on challenges like bias, privacy, accountability, and transparency.⁷⁵

- **2025 Initiatives**

A major development in March 2025 was the launch of the five-year Accelerator Fellowship Programme. This initiative aims to make impactful contributions to AI regulation, industry practices, and public awareness by bringing together emerging leaders and established experts.⁷⁵ The Institute is hosting a series of events throughout the year on topics connecting AI with creativity, care, human rights, and global regulation.⁷⁵

- **Key Personnel**

The program is led by Dr Caroline Green (Director of Research) with guidance from Professor Sir Nigel Shadbolt.⁷⁵ The high-profile inaugural fellows bring

diverse expertise: Prof Alondra Nelson (practical application of AI ethics frameworks), Prof Cass Sunstein (AI prediction and behavioral economics), Dr Joy Buolamwini (bias in AI systems), and Prof Yuval Shany (digital human rights and AI regulation).⁷⁵

- **Official/Trending Links**

- Oxford News Announcement (Accelerator Fellowship):
<https://www.ox.ac.uk/news/2025-03-03-oxford-institute-ethics-ai-launches-accelerator-fellowship-programme>⁷⁵

- **Ruhr University Bochum (Workshop)**

The Institute for Philosophy II at Ruhr University Bochum is hosting a dedicated workshop focused specifically on the challenges of evaluating artificial consciousness.⁷⁶

- **2025 Initiatives**

The "Evaluating Artificial Consciousness 2025" workshop is scheduled for June 10-11, 2025. It aims to bring together researchers to discuss theoretical, behavioral, and ethical approaches to assessing potential AI sentience.⁷⁶ A special issue of the open-access journal *Philosophy and the Mind Sciences* is planned based on the workshop's contributions.⁷⁶

- **Key Personnel**

Confirmed speakers include Michele Farisco (discussing indicators of consciousness in AI), Johannes Kleiner (exploring the role of no-go theorems), Winnie Street (addressing theoretical, behavioral, and ethical approaches to AI sentience), Patrick Butlin, Joanna Bryson, Leonard Dung, François Kammerer, and Lucia Melloni.⁷⁶

- **Official/Trending Links**

- PhilEvents Workshop Page: <https://philevents.org/event/show/131314>⁷⁶
- Workshop Website: <https://eac-2025.sciencesconf.org/>⁷⁶

• Princeton University (Panel/Lecture Series)

Princeton's Language and Intelligence initiative hosted a high-profile event tackling the question of machine consciousness, exploring links between AI capabilities and human cognition.⁷³

- **2025 Initiatives**

The panel discussion "Can Machines Become Conscious?" took place on March 4, 2025, as part of a broader lecture series on the future of AI.⁷³ The discussion centered on whether AI can achieve true consciousness as its sensory and processing abilities advance, and what this might teach us about our own minds.⁷³

- **Key Personnel**

The panel featured a prominent philosopher, David Chalmers (NYU), known for his work on consciousness and the "hard problem," debating with neuroscientist Michael Graziano (Princeton Neuroscience Institute), who studies the brain basis of consciousness and developed the attention schema theory.⁷³ The event was moderated by science author Anil Ananthaswamy.⁷³

- **Official/Trending Links**

- Princeton AI News Recap:

<https://ai.princeton.edu/news/2025/watch-neuroscientist-and-philosopher-debate-ai-consciousness>⁷³

- Princeton Event Listing:

<https://www.princeton.edu/events/2025/large-ai-model-lecture-series-can-machines-become-conscious>⁷⁷

• Emory University (Center for Ethics)

Emory's Center for Ethics is actively working to integrate AI ethics across the university, fostering awareness and responsible engagement with AI technologies.⁷¹

- **2025 Initiatives**

The Center appointed an AI Ethics Faculty Liaison in 2025 to support these integration efforts.⁷¹ A key program is "Simuvaction on AI," which convened international university students in 2025 for an experiential learning exercise simulating the Global Partnership on AI's summit. The 2025 theme focused on

the universal right to work in the age of AI, prompting students to develop actionable recommendations on issues like AI-driven productivity enhancement and economic value.⁷¹

- **Key Personnel**

John Lysaker serves as the Director of the Center for Ethics, and Anne-Elisabeth Courrier is the AI Ethics Faculty Liaison.⁷¹

- **Official/Trending Links**

- Emory News Story:

- https://news.emory.edu/stories/2025/04/er_ai_ethics_liaison_14-04-2025/story.html⁷¹

- **Georgia Tech (Panel)**

Georgia Tech hosted an interdisciplinary panel exploring deeper philosophical questions about AI, including consciousness and the concept of a "soul".⁷⁸

- **2025 Initiatives**

The panel discussion "Being in the World: Will AI Ever Have a Soul?" took place on March 10, 2025. It featured experts from philosophy, psychology, neuroscience, AI, and the arts discussing AI consciousness, self-awareness, creativity, and human-AI distinctions.⁷⁸

- **Key Personnel**

The panel was moderated by Francesco Fedele (Civil and Environmental Engineering) and Ed Greco (Physics).⁷⁸ Panelists represented a range of disciplines.

- **Official/Trending Links**

- Georgia Tech Calendar Event:

- <https://calendar.gatech.edu/event/2025/03/10/being-world-will-ai-ever-have-soul>⁷⁸

3.2 Key Philosophical and Scientific Perspectives

The 2025 discourse on AI consciousness grapples with fundamental questions, drawing on diverse fields.

- **Defining/Detecting Consciousness**

The inherent difficulty in defining consciousness is widely acknowledged, often boiling down to the capacity for subjective experience.⁷² Research efforts focus on identifying potential indicators of consciousness in AI systems.⁷⁶ Philosophical frameworks like the Computational Theory of Mind, which draws analogies between minds and computers, continue to inform the debate as AI grows more sophisticated.⁷⁴ Cognitive neuroscience provides potential frameworks, with researchers attempting to apply theories of biological consciousness to assess AI systems, though currently, no AI meets these conditions.⁷⁴ While consensus is lacking, there appears to be no fundamental theoretical barrier identified that would preclude the development of conscious AI systems in the future.⁷⁴ Responsible research principles are being advocated to guide work in this sensitive area.⁷²

- **Ethical Implications**

The potential emergence of conscious or sentient AI raises profound ethical questions, particularly regarding AI rights and "moral patienthood" – whether AI systems could have their own interests and moral significance.⁷² This prospect is no longer considered purely science fiction, with some experts predicting AI sentience within the decade.⁷² Such developments could lead to significant societal divisions between those who accept the possibility of AI consciousness and those who dismiss it.⁷² This underscores the need for ethical frameworks to evolve alongside technological capabilities to address issues like the potential suffering of AI systems and the challenges for governing bodies in regulating conscious versus unconscious AI (a challenge noted in the initial query).

- **Neuromorphic Computing**

This field, which develops computer hardware and software that processes information in ways analogous to biological brains, is seen as increasingly relevant to the consciousness debate in 2025.⁷² By mimicking brain-like processing (e.g., "spiking" only when needed, rather than continuous processing), neuromorphic systems promise greater energy efficiency compared to current power-hungry AI models.⁷² More significantly, the development of neuromorphic computing may

provide deeper insights into how biological brains function, potentially shedding light on the mechanisms underlying consciousness itself.⁷² Some view neuro-morphics as a potential "third big bang" in AI, which could significantly advance our understanding of, and ability to create, potentially conscious machines.⁷²

Section 3 Synthesis: Consciousness Enters the Mainstream

A notable shift in 2025 is the formalization and mainstreaming of the AI consciousness discussion. What was recently confined to specialized philosophy or fringe AI safety circles is now the subject of dedicated workshops, high-profile university panels, and institutional fellowship programs at leading centers like Oxford, Princeton, and Ruhr University Bochum.⁷³ This increased academic and institutional focus, coupled with statements from experts suggesting 2025 is a pivotal year for this topic⁷², reflects a growing recognition that the accelerating pace of AI development necessitates confronting these deeper philosophical and ethical questions head-on.

The exploration of AI consciousness is distinctly interdisciplinary. Events and programs in 2025 bring together not only AI researchers and philosophers but also neuroscientists, psychologists, legal scholars, ethicists, policy experts, and even artists.⁷³ This convergence underscores the understanding that the problem extends beyond purely technical considerations, requiring insights from diverse fields to grapple with the nature of intelligence, experience, and ethical standing. The involvement of figures with backgrounds in public policy, law, and the arts alongside scientists and philosophers points to a holistic approach being adopted by leading institutions.⁷⁵

Crucially, the debate around AI consciousness is not purely abstract but is increasingly linked to practical concerns about AI safety, ethics, and governance. Understanding the potential for consciousness or sentience is viewed as relevant for determining how to align AI behavior with human values, whether AI systems might warrant ethical consideration or rights, and how society should prepare for such possibilities.⁷² The focus of initiatives like the Oxford Accelerator Fellowship on producing actionable solutions for AI ethics challenges⁷⁵ demonstrates this connection. The discussion aims not just at theoretical understanding but at informing the responsible development, regulation, and deployment of advanced AI systems.

Section 4: The Imperative of Human Control: Safety, Security, and Verification (2025)

Overview

As the race towards more powerful AI systems intensifies, the imperative of maintaining meaningful human control becomes increasingly critical. This encompasses a spectrum of challenges, including ensuring technical safety and alignment, establishing robust governance mechanisms, and implementing societal controls to mitigate risks.¹⁵ The year 2025 sees a heightened focus on AI safety research, the development of verification methods, and proactive risk management strategies, driven by escalating concerns about the potential for highly capable AI systems to be misaligned with human intentions, misused for harmful purposes, or exhibit uncontrollable emergent behaviors.⁷⁰

4.1 AI Safety Research Labs and Technical Frontiers

Several organizations are dedicated to the technical challenges of ensuring advanced AI systems remain safe and controllable.

- **Machine Intelligence Research Institute (MIRI)**

MIRI has historically been a foundational organization in the field of AI alignment, focusing on theoretical and mathematical research aimed at ensuring smarter-than-human AI has a positive impact.⁷⁹ Its traditional research areas included highly reliable agent design, value specification (aligning AI goals with human values), and error tolerance.⁸²

- **2025 Initiatives**

MIRI underwent a significant strategic pivot, publicly announced in late 2023/early 2024 and continuing through 2025.¹⁶ Citing insufficient progress in technical alignment research relative to the rapid pace of AI capability development, MIRI has substantially scaled back its technical alignment work.¹⁶ Its primary focus has shifted towards public policy, communications, and technical *governance* research.¹⁶ The organization now emphasizes raising awareness among policymakers and the public about potential catastrophic

risks from AI and advocates for strong governance measures, including the possibility of an international suspension of frontier AI research, believing disaster is likely otherwise.¹⁶ In 2025, MIRI is actively hiring communications specialists and technical governance researchers.¹⁶ Its Technical Governance Team continues to produce work, such as analyses of verification mechanisms for international AI agreements and critiques of AI evaluation methods.⁶⁹ MIRI is also developing new introductory resources on AI risk and a forthcoming book.¹⁶ While reduced, MIRI still supports some alignment research efforts and AI safety retraining programs.¹⁶ The organization reported having approximately two years of funding reserves (\$16M) at the end of 2024, with projected 2025 expenses of \$6.5M-\$7M, and expressed uncertainty about future funding for its new strategy.¹⁶

- **Key Personnel**

Malo Bourgon has been CEO since October 2023.⁸⁴ Eliezer Yudkowsky, a co-founder, remains a prominent and often pessimistic voice on AI risk.⁸¹ Nate Soares is another key researcher.⁸⁸ Researchers previously associated with MIRI's technical agendas include Scott Garrabrant, Evan Hubinger, and Vanessa Kosoy.⁸³ Max Harms contributed to the MIRI Single Author Series in April 2025.⁸⁹

- **Official/Trending Links**

- MIRI Official Website: <https://intelligence.org/> ⁷⁹
- MIRI Technical Governance Research:
<https://techgov.intelligence.org/research> ⁶⁹, <https://techgov.intelligence.org/> ⁷⁰
- MIRI Updates & Strategy (2024/2025):
<https://intelligence.org/2024/12/02/miris-2024-end-of-year-update/> ¹⁶,
<https://intelligence.org/category/miri/> ⁸⁴,
<https://intelligence.org/category/news/> ⁸⁸,
<https://forum.effectivealtruism.org/posts/zBizn2BT6pjbqS8n/miri-s-2024-end-of-year-update> ⁸⁵,
<https://forum.effectivealtruism.org/posts/e8o6paib9sgKeWorc/what-is-miri-currently-doing-1> ⁸⁶
- News on MIRI Pivot:
<https://getcoai.com/news/leading-ai-safety-organization-drops-technical-research-to-focus-exclusively-on-policy/> ¹⁷
- AI Safety Info Summary of MIRI: https://aisafety.info/?state=85EN_ ⁸³
- Open Philanthropy Grant (Retraining):
<https://www.openphilanthropy.org/grants/machine-intelligence-research-in>

stitute-ai-safety-retraining-program/⁸⁷

- **ARC Prize Foundation (Alignment Research Center)**

This foundation focuses on creating challenging benchmarks to evaluate AI progress towards AGI, with a particular emphasis on reasoning and efficiency, aspects deemed critical for safety.³⁰

- **2025 Initiatives**

In March 2025, the foundation introduced the ARC-AGI-2 benchmark.⁹⁰ This new benchmark builds upon the original ARC-AGI-1 (Abstract Reasoning Corpus), which was famously difficult for AI until OpenAI's o3 model achieved human-level performance in late 2024, albeit with high computational cost.³⁰ ARC-AGI-2 specifically introduces efficiency constraints, measuring not just problem-solving ability on novel visual reasoning puzzles but also the computational resources consumed, thereby discouraging brute-force approaches and rewarding more "intelligent" solutions.⁹⁰ Alongside the benchmark, the Arc Prize 2025 competition was launched, challenging participants to achieve high accuracy on ARC-AGI-2 within strict computational cost limits (\$0.42 per task).⁹⁰ Initial results showed even top models like o3 performed poorly on ARC-AGI-2 under these constraints.⁹⁰

- **Key Personnel**

François Chollet, an AI researcher credited with creating the Keras library and the original ARC dataset, is a co-founder of the Arc Prize Foundation.⁹⁰ Greg Kamradt serves as the foundation's president.⁹⁰

- **Official/Trending Links**

- GovInfoSecurity Article on ARC-AGI-2:
<https://www.govinfosecurity.com/new-benchmarks-challenge-brute-force-approach-to-ai-a-27826> ⁹⁰
 - HiFlyLabs Blog Post (ARC Mention):
<https://hiflylabs.com/blog/2025/1/27/path-to-agi-part-2> ³⁰
 - Arc Prize Website: <https://arcprize.org/> (Implied from ³⁰)

- **Conjecture**

Conjecture is an AI safety research lab that initially focused on mechanistic interpretability (understanding the internal workings of models) but has pivoted towards "cognitive emulation".⁹¹ This newer approach aims to produce bounded agents that emulate human-like thought processes, potentially offering a different path to safety.⁹¹

- **2025 Initiatives**

Conjecture participated in the AI Safety Connect event in Paris in February 2025, where co-founder Gabriel Alfour presented the company's work alongside other AI safety ventures.³² Their research continues, though details may be limited due to their internal infohazard policy designed to prevent the release of potentially dangerous information.⁹¹ Conjecture also has a B2C transcription product called Verbalize, released in 2023, though its commercial traction was unclear as of mid-2023.⁹¹

- **Key Personnel**

Connor Leahy serves as CEO and is active in policy outreach.⁹¹ Gabriel Alfour is a co-founder with technical and scaling experience.³² Sid Black is also a co-founder.⁹¹ The company received significant VC funding (~\$10M in 2022) from prominent tech figures including Nat Friedman, Patrick and John Collison, Daniel Gross, Andrej Karpathy, and Sam Bankman-Fried.⁹¹

- **Official/Trending Links**

- EA Forum Critique/Overview (June 2023):
<https://forum.effectivealtruism.org/posts/gkfMLX4NWZdmpikto/critiques-of-prominent-ai-safety-labs-conjecture>⁹¹
- AI Safety Connect Event Details (Participation):
<https://www.aisafetyconnect.com/event-details>³²
- Conjecture Website: <https://conjecture.dev/> (Implied, not in snippets)

- **Apart Research (AI Control Hackathon)**

This organization focuses on "AI control"—developing techniques to mitigate security risks from AI systems, even assuming the AI itself might be adversarial or try to subvert safety measures.⁹² This represents a more security-oriented approach to safety.

- **2025 Initiatives**

Apart Research organized an AI Control Hackathon in London and online on March 29-30, 2025.⁹² The event brought together researchers, engineers, and security professionals to work on challenges in areas like runtime monitoring systems, adversarial stress testing of AI safety mechanisms, formal verification approaches for proving safety properties, bounded optimization techniques to limit AI capabilities, and red teaming exercises to find vulnerabilities in AI control systems.⁹²

- **Key Personnel**

Organizers and participants of the hackathon community.

- **Official/Trending Links**

- AI Control Hackathon 2025 Event Page:
<https://apartresearch.com/sprints/ai-control-hackathon-2025-03-29-to-2025-03-30>⁹²

4.2 Think Tanks and Policy Frameworks for Control

Policy-focused organizations and think tanks are actively developing frameworks and recommendations for governing AI and ensuring human control.

- **The Future Society**

This organization works on AI governance issues, engaging in stakeholder consultations to identify priorities and advocate for concrete policy mechanisms.¹⁵

- **2025 Initiatives**

In April 2025, The Future Society published a significant report detailing AI governance priorities based on a survey of 44 civil society organizations.¹⁵ The top-ranked priorities emerging from this consultation included establishing legally binding "red lines" to prohibit unacceptable AI risks, mandating systematic independent third-party audits for general-purpose AI systems, establishing crisis management frameworks for rapid response to AI incidents, enacting robust whistleblower protections, and ensuring meaningful civil society participation (especially from the Global South) in governance processes.¹⁵ Representatives also participated in the AI Safety Connect event.³²

- **Key Personnel**

Caroline Jeanmaire participated in AI Safety Connect.³² Nicolas Mialhe, associated with PRISM Eval which presented at AI Safety Connect, may also be linked.³² Report contributors and surveyed organizations.

- **Official/Trending Links**

- CSO AI Governance Priorities Report (April 2025):
<https://thefuturesociety.org/cso-ai-governance-priorities/>¹⁵

- **Stimson Center**

This think tank focuses on international security and transnational challenges, applying its expertise to AI governance with a particular emphasis on long-term impacts and the interests of future generations.⁵⁷

- **2025 Initiatives**

In March 2025, the Stimson Center released the report "Governing AI for the Future of Humanity".⁵⁷ This report uniquely connects the UN's Declaration on Future Generations with the Global Digital Compact, arguing for their convergence to promote equitable and sustainable AI governance. A key theme is the critical role of strategic foresight in developing adaptive and resilient governance structures capable of managing AI's uncertainties. The report proposes specific foresight tools and multilateral coordination mechanisms, such as a Global AI Foresight Network (GAFN).⁵⁷

- **Key Personnel**

Authors of the "Governing AI for the Future of Humanity" report.

- **Official/Trending Links**

- Governing AI for the Future of Humanity Report:
<https://www.stimson.org/2025/governing-ai-for-the-future-of-humanity/>⁵⁷

- **Center for AI and Digital Policy (CAIDP)**

CAIDP advocates for established AI policy principles, democratic values, fundamental rights, and the rule of law in the context of AI governance.⁵⁸

- **2025 Initiatives**

CAIDP organized the WashDC25AIDV event (presumably focused on AI, Democracy, and Values) in 2025. The event featured high-profile speakers from civil rights organizations (Maya Wiley), the UN (Amandeep Singh Gill), academia/industry (Sasha Luccioni, Stuart Russell, Virginia Dignum), reflecting CAIDP's focus on convening diverse voices around AI policy and ethics.⁵⁸

- **Key Personnel**

Marc Rotenberg serves as President. Speakers at their 2025 event included Maya Wiley (The Leadership Conference), Amandeep Singh Gill (UN Envoy on Technology), Sasha Luccioni (Hugging Face), Stuart Russell (UC Berkeley/CHAI),

Virginia Dignum (Umeå University).⁵⁸

- **Official/Trending Links**

- WashDC25AIDV Event Bios Page:
<https://www.caidp.org/events/washdc25aidv/bios/>⁵⁸
- CAIDP Main Website: <https://www.caidp.org/> (Implied from URL)

- **Centre for Long-Term Resilience (CLTR)**

This UK-based organization focuses on improving societal resilience to extreme risks, with a specific workstream on AI policy, particularly concerning incident preparedness and crisis management.¹⁵

- **2025 Initiatives**

CLTR's work on AI crisis preparedness was highlighted in The Future Society's 2025 report, emphasizing the need for governments to improve their ability to anticipate, plan for, contain, and recover from incidents involving AI systems that threaten public safety or critical infrastructure.¹⁵

- **Key Personnel**

Jess Whittlestone leads the AI Policy work at CLTR.¹⁵

- **Official/Trending Links**

- Mentioned in Future Society Report:
<https://thefuturesociety.org/cso-ai-governance-priorities/>¹⁵
- CLTR Website: <https://www.cltr.org/> (Implied, common knowledge)

4.3 Emerging Tools and Methodologies

Ensuring human control relies on developing and deploying effective technical tools and standardized methodologies.

- **Verification Mechanisms**

Recognizing the challenge of ensuring compliance with potential international agreements or regulations on AI development, research is underway to develop effective verification mechanisms. MIRI published a technical report in late 2024 providing an overview of potential mechanisms.⁶⁹ This area is critical, especially given the skepticism raised about the feasibility of monitoring and enforcing agreements in critiques of concepts like Mutually Assured AI Malfunction.¹⁰

- **AI Evaluations**

Evaluating the capabilities, safety, and potential risks of AI models is a cornerstone of governance and control efforts. However, organizations like MIRI caution that evaluations have fundamental limitations and cannot be solely relied upon, particularly for preventing catastrophic risks from future, more advanced systems.⁶⁹ They argue for regulations requiring developers to explicitly state and justify the assumptions underlying their evaluation methods.⁶⁹ Despite limitations, new evaluation benchmarks are emerging in 2025, including those focused on safety, factuality, and robustness, such as HELM Safety, AIR-Bench, and FACTS.⁸ The ARC-AGI-2 benchmark specifically targets efficient reasoning capabilities, moving beyond simple task performance.⁹⁰

- **AI Alignment Techniques**

Research continues on various techniques aimed at ensuring AI systems act in accordance with human values and intentions. Key concepts discussed in 2025 include differentiating between outer alignment (ensuring the specified objective reflects human values) and inner alignment (ensuring emergent goals during optimization remain aligned), as well as value alignment (matching AI criteria to human values) versus intent alignment (matching AI actions to human expectations).⁸⁰ Specific approaches being developed or deployed include Constitutional AI (pioneered by Anthropic⁴²), Reinforcement Learning from Human Feedback (RLHF, widely used, co-invented by Dario Amodei⁴⁵), and ongoing work in mechanistic interpretability to understand model internals (a focus for Anthropic and historically for MIRI⁴²).

- **AI Control Techniques**

Complementary to alignment, AI control focuses on building external safeguards and constraints robust even against potential AI subversion. The AI Control Hackathon organized by Apart Research in March 2025 highlighted key focus areas: runtime monitoring systems (observing AI behavior during operation), adversarial stress testing (probing defenses), formal verification (mathematically proving safety properties), and bounded optimization (limiting AI capabilities).⁹² Red teaming, or simulating attacks to find vulnerabilities, is also a crucial component.⁹²

- **Safety Frameworks**

Standardized frameworks provide guidance for managing AI risks. The US National Institute of Standards and Technology (NIST) AI Risk Management Framework (RMF) and its specific profile for Generative AI (NIST AI 600-1) are influential references, though organizations like MIRI have provided comments suggesting improvements (e.g., including risks from misaligned systems).⁶⁹ Other frameworks mentioned include the SAFE Innovation framework (Security, Accountability, Foundations, Explainability).²² Additionally, major corporations like Microsoft, IBM, and Google are developing and promoting their own internal Responsible AI frameworks and best practices.²²

Section 4 Synthesis: Diverging Paths to Safety and Control

The landscape of AI safety and control in 2025 is marked by a notable divergence in strategies. A significant development is the strategic pivot of the Machine Intelligence Research Institute (MIRI), a historically foundational technical alignment organization. Citing slow progress on technical solutions relative to the accelerating pace of AI capabilities, MIRI has shifted its primary focus to policy advocacy, public communication, and technical governance research, expressing deep pessimism about preventing AI catastrophe without drastic interventions like an international moratorium on frontier development.¹⁶ This move contrasts sharply with the continued, intensive technical work pursued by other key players. Labs like Anthropic remain committed to safety-integrated development through approaches like Constitutional AI and interpretability research⁴², while organizations like Conjecture explore novel technical paths like cognitive emulation.⁹¹ Furthermore, dedicated efforts focus on improving evaluation benchmarks (ARC Prize Foundation⁹⁰) and developing robust external control mechanisms (Apart Research Hackathon⁹²). This bifurcation highlights a fundamental debate within the safety community regarding the most viable path forward: prioritizing technical breakthroughs in alignment and control versus emphasizing immediate governance and policy interventions to slow down or manage development.

Amidst these strategic debates, the concepts of verification and evaluation have emerged as central, cross-cutting themes in 2025. The ability to reliably assess the capabilities and risks of increasingly complex AI systems⁶⁹, and to verify compliance with safety standards, regulations, or potential international agreements⁶⁹, is recognized as a critical prerequisite for effective governance and control. However, achieving reliable verification and evaluation faces significant technical hurdles, particularly concerning the unpredictability of future systems and the limitations of current methods.⁶⁹ The proliferation of new benchmarks targeting safety, factuality, and reasoning efficiency⁸, alongside dedicated research into verification mechanisms⁶⁹, underscores the recognition of this area as a critical bottleneck and a major focus of effort for the AI safety and governance communities. The demand for independent, third-party audits, identified as a top priority by civil society organizations¹⁵, further emphasizes the need for trustworthy assessment methods.

Flowing from the growing concern about potential misalignment or adversarial behavior in advanced AI, the concept of "AI Control" is gaining prominence alongside

traditional "Alignment" approaches in 2025. As exemplified by the focus of the AI Control Hackathon ⁹², this paradigm emphasizes building robust external safeguards, monitoring systems, and containment strategies that can function even if the AI system itself attempts to subvert them. This reflects the adoption of a "security mindset," as advocated by MIRI ⁸³, which anticipates potential failures in internal alignment and prioritizes mechanisms to limit potential harm regardless of the AI's internal state. This approach complements alignment efforts (which focus on instilling the 'right' goals and values internally) by adding layers of external checks and balances, acknowledging the profound difficulty and uncertainty involved in guaranteeing the benevolent behavior of superintelligent systems.

(Table 2: Key AI Safety & Control Organizations and Approaches (2025))

Organization	Primary Focus	Key 2025 Activities/Outputs	Key Figures	Official Link
MIRI	Policy/Gov Advocacy, Technical Governance Research, Risk Communication ¹⁷	Strategic pivot from technical alignment; Research on verification/evaluations; Policy comments; Risk comms efforts ¹⁶	Malo Bourgon (CEO), Eliezer Yudkowsky (Co-founder)	https://intelligence.org/ ⁷⁹
ARC Prize Foundation	Benchmarking (Reasoning, Efficiency) ⁹⁰	Introduced ARC-AGI-2 benchmark & Arc Prize 2025 competition ⁹⁰	François Chollet (Co-Founder), Greg Kamradt (President)	https://arcprize.org/ (Implied)
Conjecture	Technical Safety (Cognitive Emulation focus) ⁹¹	Ongoing research under infohazard policy; Participation in AI Safety Connect ³²	Connor Leahy (CEO), Gabriel Alfour (Co-founder)	https://conjecture.dev/ (Implied)
Anthropic (Safety Team)	Technical Safety (Constitutional AI, Interpretability) ⁴²	Continued development of Claude models with safety focus; Research publication/participation ³	Dario Amodei (CEO), Chris Olah (Co-founder)	https://www.anthropic.com/ ³
GovAI	AI Governance Research & Policy Advice ⁶¹	Summer Fellowship 2025; Research Scholar Program;	GovAI Team & Affiliates	https://www.governance.ai/ ⁶²

AGI Landscape and Organizations: A 2025 Intelligence Briefing

Research by Fede Nolasco | AI Researcher and Data Architect

<https://www.linkedin.com/in/federiconolasco>

Report released on 22 April 2025

Organization	Primary Focus	Key 2025 Activities/Outputs	Key Figures	Official Link
		Publications on regulation, risk assessment, etc. ⁶¹		
CSER	Existential & Global Catastrophic Risk Research (incl. AI) ⁶⁵	Research on extinction risk, military AI; Ethics seminar; MPhil program launched; New Director appointed ⁶⁵	S.M. Amadae (Director), Seán Ó hÉigearthaigh (Former Director)	https://www.cser.ac.uk/ ⁶⁷
The Future Society	AI Governance Priorities & Policy Mechanisms ¹⁵	Published CSO AI Governance Priorities Report (Apr 2025); Participation in safety events ¹⁵	Caroline Jeanmaire (Participant)	https://thefuturesociety.org/ (Implied)
Stimson Center	AI Governance for Future Generations; Foresight Methods ⁵⁷	Published "Governing AI for Future Humanity" report (Mar 2025) linking UN frameworks ⁵⁷	Report Authors	https://www.stimson.org/ (Implied)
Apart Research	Technical AI Control (Security Mindset) ⁹²	Organized AI Control Hackathon (Mar 2025) focusing on monitoring, testing, verification, red teaming ⁹²	Hackathon Organizers/Participants	https://apartresearch.com/ (Implied)
Stanford HAI	Tracking AI	Published AI	Vanessa Parli	https://hai.stanford.edu/

Organization	Primary Focus	Key 2025 Activities/Outputs	Key Figures	Official Link
(AI Index)	Trends (incl. Safety/Ethics/Policy) ⁸	Index Report 2025 detailing RAI evolution, incidents, governance efforts ²⁴	(Director of Research)	rd.edu/ (Implied)
PAI	Responsible AI Ecosystem; Multi-stakeholder Coordination ³⁴	Policy program; Collaboration w/ institutions; PAIE initiative w/ J-PAL; Partnered in Standards Summit ¹³	Rebecca Finlay (CEO), Policy Steering Committee	https://partnershiponai.org/ ³⁴

Section 5: Intelligence Sequencing: A Strategic Crossroads (2025)

Overview

A novel conceptual framework gaining traction and generating discussion within AI strategy circles in 2025 is "Intelligence Sequencing".⁹³ This perspective challenges conventional assumptions about AI development and safety by proposing that the *order* in which different forms of advanced intelligence emerge – specifically, centralized Artificial General Intelligence (AGI) versus Decentralized Collective Intelligence (DCI) – may be the most critical determinant of long-term civilizational outcomes, potentially outweighing efforts to align AGI after its creation.

5.1 The AGI-First vs. DCI-First Framework

The core of the Intelligence Sequencing argument, primarily articulated in a widely discussed 2025 paper by independent researcher Andy E. Williams⁹⁴, posits that intelligence evolution follows path-dependent, potentially irreversible trajectories

towards distinct "attractor basins."

- **Core Arguments**

- **AGI-First Attractor (Centralization)**

If AGI, characterized as highly autonomous, centralized systems capable of outperforming humans across many domains, emerges as the dominant form of intelligence before robust DCI systems mature, the trajectory is predicted to lock into a centralization attractor.⁹³ This state is characterized by hierarchical control structures, intense competition for resources and dominance, the concentration of power in the hands of AGI controllers, and the emergence of instrumental power-seeking behaviors by AGIs themselves as a means to achieve their goals. This path is seen as significantly increasing existential risks due to the potential for uncontrollable power concentration and competitive escalation.⁹³ The framework suggests this path is favored by intelligence systems that model the world based on externally imposed axioms or fixed optimization landscapes.⁹³

- **DCI-First Attractor (Decentralization)**

Conversely, if technologies enabling Decentralized Collective Intelligence – systems characterized by distributed reasoning, networked cooperation, and emergent intelligence scaling across many nodes – reach critical mass before AGI dominates, the trajectory is predicted to stabilize around a decentralization attractor.⁹³ This state is characterized by distributed cooperation, optimization for collective fitness and stability rather than individual dominance, and potentially more resilient and inherently safer outcomes due to the lack of single points of failure or control.⁹³ This path is associated with intelligence systems that model the world through recursive internal visualization and maintain dynamic openness.⁹³

- **Path Dependence & Irreversibility**

A crucial element of the theory is the concept of irreversibility. Once intelligence development enters either the AGI-first or DCI-first regime, feedback loops (e.g., competitive pressures reinforcing centralization) and structural lock-in (e.g., resource monopolization by early AGIs) make transitioning between these attractors increasingly infeasible.⁹⁴ Early structural choices heavily constrain later possibilities, much like path dependence observed in technological standards or economic development.⁹⁴

• **AI Safety Implications**

This framework presents a fundamental challenge to traditional AI alignment research, which typically assumes AGI will emerge and focuses on controlling it *after* the fact.⁹³ The Intelligence Sequencing perspective argues that the sequencing choice is more foundational.⁹⁴ If the AGI-first path inherently leads to competitive dynamics and power concentration, post-hoc alignment efforts might be insufficient or ultimately futile.⁹³ Consequently, the DCI-first path is presented as a potentially more robustly safe trajectory for humanity.⁹⁴ This implies that strategic policy should consider prioritizing the development of DCI infrastructure and potentially implementing measures to delay or carefully manage the emergence of centralized AGI.⁹⁴

• **2025 Publications/Discussions**

The primary catalyst for discussion in 2025 is the paper "Intelligence Sequencing and the Path-Dependence of Intelligence Evolution: AGI-First vs. DCI-First as Irreversible Attractors" by Andy E. Williams. It appeared as an arXiv preprint (2503.17688) in March 2025 and was also published as a preprint on Qeios in April 2025, garnering attention and analysis on platforms like AIModels.fyi.⁹³

• **Links**

- arXiv Paper (2503.17688): <https://arxiv.org/abs/2503.17688> ⁹⁴,
<https://arxiv.org/pdf/2503.17688> ⁹⁴
- Qeios Preprint: <https://www.qeios.com/read/RA5XMP> ⁹⁷,
<https://www.qeios.com/read/RA5XMP/pdf> ⁹⁶
- AIModels.fyi Analysis:
<https://www.aimodels.fyi/papers/arxiv/intelligence-sequencing-path-dependence-intelligence-evolution-agi> ⁹³
- Scribd Document Link:
<https://www.scribd.com/document/842884183/2503-17688v1> ⁹⁵

5.2 Pioneers of Decentralized Collective Intelligence (DCI)

While the Intelligence Sequencing framework provides a theoretical lens, several research groups and initiatives are actively working on or advocating for decentralized approaches to AI in 2025, lending practical weight to the DCI concept.

- **Pluralis Research**

This AI research organization, founded in early 2024 by ex-FAANG scientists, explicitly aims to develop "true open-source AI" through decentralized training methods.⁹⁹

- **Aim**

To facilitate collaborative, multi-party training of large foundation models, creating an open, distributed AI ecosystem with sustainable economic incentives for contributors, thereby challenging the dominance of large, centralized systems.⁹⁹

- **2025 Initiatives**

Pluralis announced a significant \$7.6 million seed funding round in March 2025, co-led by prominent VCs USV and CoinFund, with participation from others including Topology, Variant, and notable angels like Balaji Srinivasan and HuggingFace co-founder Clem Delangue.⁹⁹ They are pioneering a novel approach termed "Protocol Learning," designed to enable model training across open, permissionless networks. A key feature is that the model weights are never fully materialized or controlled by any single party but 'live' within the protocol, allowing value to flow programmatically to contributors while preserving model integrity.⁹⁹

- **Key Personnel**

Alexander Long is the Founder and CEO.⁹⁹

- **Official/Trending Links**

- GlobeNewswire Funding Announcement:
<https://www.globenewswire.com/news-release/2025/03/19/3045635/0/en/Pluralis-Research-Pioneers-Protocol-Learning-to-Scale-Decentralized-AI-Announces-7-6M-Seed-Round-Led-by-USV-and-CoinFund.html>⁹⁹

- **Sentient Research Group**

This group advocates for decentralized AI as a key enabler for achieving AGI, emphasizing innovative data gathering and collaborative training.¹⁰⁰

- **Aim**

To foster a community-driven strategy for AI development, overcoming limitations of centralized approaches, particularly regarding access to diverse and high-quality data.¹⁰⁰

- **2025 Initiatives**

Sentient is actively developing "Sentient Chat," envisioned as a community-driven AI chatbot platform designed to move beyond traditional web search towards collaborative task execution using multiple AI agents.¹⁰⁰ They highlight the importance of accessing unique datasets (beyond readily available web data) and propose open systems with incentives for data contribution and decentralized model ownership/training.¹⁰⁰ The platform aims to support developers with tools like AI search APIs and custom agent structures.¹⁰⁰

- **Key Personnel**

Himanshu Tyagi is a founder and key spokesperson.¹⁰⁰

- **Official/Trending Links**

- TechNews Article on Sentient Chat:

<https://live.upcoming.sk/2025/04/04/decentralized-ai-the-key-to-unlocking-artificial-general-intelligence/>¹⁰⁰

- **Institut Polytechnique de Paris (IP Paris) / Éric Moulines**

Research led by Professor Éric Moulines focuses on developing decentralized AI models based on federated learning principles, driven by concerns about data centralization and privacy.¹⁰¹

- **Focus**

To enable collaborative learning among "intelligent agents" (e.g., hospitals, individuals) where data can be shared for model training locally without compromising privacy or requiring central storage.¹⁰¹ Addressing challenges like incentivizing data sharing, ensuring confidentiality, managing heterogeneous data, and detecting free-riders.¹⁰¹

- **2025 Initiatives**

Professor Moulines presented this research at the AI, Science and Society summit hosted by IP Paris in February 2025.¹⁰¹ Ongoing work involves developing these decentralized models and exploring applications, such as improving medical diagnoses by allowing hospitals to learn collectively from distributed patient data.¹⁰¹

- **Key Personnel**

Éric Moulines (Professor, CMAP, École Polytechnique).¹⁰¹

- **Official/Trending Links**

- IP Paris News Article:

<https://www.ip-paris.fr/en/news/future-ai-will-be-decentralised-and-collaborative>¹⁰¹

- **Decentralized Science (DeSci) + AI (DeScAI) Proponents**

This emerging movement seeks to leverage blockchain and decentralized network principles combined with AI to transform scientific research.¹⁰²

- **Aim**

To create open, intelligent, and self-sustaining scientific ecosystems that break down traditional barriers related to data access, funding, peer review, and collaboration.¹⁰²

- **2025 Initiatives**

The DeSci movement itself reported significant momentum, with top DeSci tokens reaching a collective market capitalization of around \$1 billion in early 2025, and many projects launching recently.¹⁰² The conceptual framework of DeScAI explores using AI for curating knowledge across decentralized networks, enabling decentralized supercomputing by pooling resources, creating AI-assisted platforms for democratized research funding and peer review, ensuring data ownership and compensation for contributors, and facilitating borderless scientific collaboration.¹⁰²

- **Key Personnel**

Proponents and developers within the broader DeSci, blockchain, and AI communities.

- **Official/Trending Links**

- Cointelegraph Article on DeScAI:
<https://cointelegraph.com/news/decentralized-science-meets-ai>¹⁰²

5.3 Strategic Implications for AI Development Pathways

The Intelligence Sequencing framework carries significant strategic implications for how the future of AI development is approached.

- **Civilizational Choice**

The framework explicitly frames the AGI-First versus DCI-First paths not just as technical choices but as a fundamental civilizational decision point between potentially unbounded competition and unbounded cooperation.⁹³ This elevates the strategic importance of choices made regarding AI architecture and infrastructure development.

- **Policy Levers**

If the theory holds, it suggests that effective long-term AI safety strategy may require proactive policy interventions aimed at influencing the *sequence* of development. This could involve policies designed to incentivize or accelerate the development of DCI infrastructure (e.g., supporting open-source initiatives, funding decentralized training research like Pluralis') while potentially regulating, slowing down, or imposing stringent safety requirements on the development of highly centralized, powerful AGI systems.⁹⁴ This connects directly back to the governance discussions (Section 2) and safety imperatives (Section 4).

- **Epistemic Dimension**

The framework introduces a philosophical layer by suggesting that the very method an intelligence system uses to model itself and the world – whether through rigid, externally imposed axioms (seen as favoring AGI) or through flexible, recursive internal visualization (seen as favoring DCI) – might inherently bias its evolutionary trajectory.⁹⁴ This implies that the path towards competition or cooperation might depend not just on *how* intelligence is built, but on *how* intelligence perceives itself and learns.⁹⁸

Section 5 Synthesis: Sequencing as a Foundational Challenge

The Intelligence Sequencing framework, gaining prominence in 2025 through publications like Williams' paper ⁹⁴, offers a potentially paradigm-shifting perspective on AI safety and strategy. Its core contribution is to reframe the central challenge: rather than focusing solely on the reactive problem of aligning a powerful AGI *after* it emerges, it emphasizes the proactive, strategic importance of the *type* of advanced intelligence that achieves dominance first.⁹³ By arguing that the initial conditions and the order of emergence (AGI-first vs. DCI-first) can lead to irreversible lock-in effects favoring either competition or cooperation, the framework challenges decades of assumptions within traditional AI alignment research.⁹³ It posits that if the AGI-first path inherently embeds competitive dynamics and power-seeking behavior, controlling it later might prove intractable.⁹⁴

This theoretical framework finds resonance in the practical developments of 2025. The emergence and funding of initiatives explicitly focused on decentralized AI, such as Pluralis Research's "Protocol Learning" ⁹⁹, the Sentient research group's community-driven "Sentient Chat" ¹⁰⁰, academic research into federated learning for collaboration ¹⁰¹, and the growing momentum of the DeScAI movement ¹⁰², all provide tangible evidence that the DCI-first pathway is not merely a theoretical construct but an area of active innovation and investment. These efforts demonstrate concrete attempts to build advanced intelligence capabilities outside the centralized, resource-intensive models pursued by the largest labs.

However, this burgeoning DCI ecosystem exists in direct tension with the dominant trend observed in Section 1: the massive, accelerating push towards centralized AGI supremacy, fueled by enormous investments in compute infrastructure like the Stargate project ⁷ and the competitive drive of leading nations and corporations.⁸ The Intelligence Sequencing framework suggests these two trends represent fundamentally conflicting paths. The investments enabling AGI-first development could simultaneously create the conditions for resource monopolization and structural lock-in that make the DCI-first attractor increasingly difficult to reach.⁹⁴ This sets up a potential structural battle in the coming years over the foundational infrastructure and dominant paradigm for future intelligence, making the strategic choices regarding investment, research direction, and policy crucial in determining which attractor basin becomes dominant.

Conclusion

Synthesis of the 2025 AGI Landscape

The analysis of the Artificial General Intelligence landscape in 2025 reveals a period of intense dynamism, characterized by accelerating technological progress, escalating investment, deepening geopolitical competition, and a parallel, urgent push for governance and control. The race for AGI supremacy, primarily between the US and China, and driven by major corporate labs like OpenAI, Google DeepMind, Anthropic, Meta, and Baidu, is tangible, marked by rapid model releases and unprecedented infrastructure investments like the Stargate initiative.¹ Concurrently, a complex global ecosystem is attempting to manage the implications, with international bodies (ISO, IEC, ITU, OECD, UN) striving for standards and governance frameworks¹², research institutions and NGOs (PAI, Stanford HAI, J-PAL, GovAI, CSER) shaping policy and ethical debates²⁴, and dedicated safety labs (MIRI, ARC, Conjecture) grappling with technical alignment and control, albeit with diverging strategies.¹⁶ The philosophical and ethical dimensions, particularly concerning AI consciousness, are moving into the mainstream academic discourse.⁷² Finally, emerging frameworks like Intelligence Sequencing challenge fundamental assumptions, proposing that the *order* of technological emergence (centralized AGI vs. decentralized DCI) may be the most critical factor determining future outcomes.⁹⁴

Dominant Trends

Several key trends define the 2025 landscape:

1. Accelerated Capability Growth

Frontier models are rapidly improving performance on complex benchmarks and expanding into new modalities like video generation and advanced reasoning.¹

2. Massive Investment & Infrastructure Focus

Unprecedented levels of private and public investment are flowing into AI, particularly in the US, with a strong emphasis on building the massive compute infrastructure required for frontier models.⁷

3. Intensified US-China Competition

The geopolitical rivalry is a primary driver of the race, with both nations making significant strategic investments and achieving rapid progress.⁸

4. **Increased Governance Activity**

There is a notable surge in efforts to establish international standards, national regulations, and ethical guidelines, although coordination remains a challenge.¹²

5. **Heightened Safety and Control Concerns**

Growing awareness of potential risks from powerful AI is fueling research into alignment, control, verification, and risk management, alongside strategic debates about the best path forward.¹⁵

6. **Centralization vs. Decentralization**

A tension exists between the dominant trend of centralized model development and emerging efforts focused on decentralized AI training and deployment.⁷

Key Tensions

The current trajectory is shaped by several core conflicts:

- **Speed vs. Safety**

The intense competitive pressure to achieve AGI first potentially conflicts with the need for careful, deliberate safety research and precaution.¹⁰

- **Competition vs. Cooperation**

National strategic interests and corporate market ambitions drive competition, while the global nature of AI risks necessitates international cooperation and governance.¹⁰

- **Openness vs. Control**

Debates persist regarding the benefits and risks of open-sourcing powerful AI models versus maintaining tighter control over their development and deployment.⁶

- **Centralization vs. Decentralization**

The dominant path of large, centralized models is being challenged by alternative visions focused on decentralized collective intelligence, raising fundamental questions about the optimal structure for future AI.⁹⁴

- **Regulation vs. Deregulation**

Policy approaches diverge significantly, notably with the US shift towards deregulation potentially conflicting with more cautious approaches elsewhere

(e.g., EU AI Act).²²

Outlook

Humanity stands at a critical juncture in 2025 regarding the development and governance of advanced artificial intelligence. The technological momentum is immense, driven by powerful economic and geopolitical forces. However, awareness of the profound risks and ethical complexities is also growing, catalyzing efforts towards safety, control, and responsible governance. The decisions and actions taken in this period by the nations, corporations, institutions, and researchers identified in this report – regarding investment priorities, research directions, safety protocols, regulatory frameworks, and collaborative efforts – will likely have significant and potentially irreversible consequences. While the ultimate trajectory remains uncertain¹⁹, the evidence suggests that the choices made now hold substantial weight in shaping whether the advent of AGI leads towards broadly

Works cited

1. OpenAI, accessed on April 22, 2025, <https://openai.com/>
2. Google DeepMind, accessed on April 22, 2025, <https://deepmind.google/>
3. Home \ Anthropic, accessed on April 22, 2025, <https://www.anthropic.com/>
4. Llama 4 vs DeepSeek V3: Comprehensive AI Model Comparison [2025] - Redblink, accessed on April 22, 2025, <https://redblink.com/llama-4-vs-deepseek-v3/>
5. Latest AI Models Just Released in 2025 - Nascenia, accessed on April 22, 2025, <https://nascenia.com/latest-ai-models/>
6. Baidu debuts its first AI reasoning model to compete with DeepSeek - SiliconANGLE, accessed on April 22, 2025, <https://siliconangle.com/2025/03/16/baidu-debuts-first-ai-reasoning-model-compete-deepseek/>
7. A New AI Era: Top 10 Takeaways from Davos 2025 | OMMAX, accessed on April 22, 2025, <https://www.ommax.com/en/insights/industry-insights/a-new-ai-era-top-10-takeaways-from-davos-2025/>
8. hai-production.s3.amazonaws.com, accessed on April 22, 2025, https://hai-production.s3.amazonaws.com/files/hai_ai_index_report_2025.pdf
9. New funding to build towards AGI | OpenAI, accessed on April 22, 2025, <https://openai.com/index/march-funding-updates/>
10. Seeking Stability in the Competition for AI Advantage | RAND, accessed on April 22, 2025, <https://www.rand.org/pubs/commentary/2025/03/seeking-stability-in-the-competition-for-ai-advantage.html>
11. What the Next Frontier of AI—AGI—Means for Government & GovCons - ExecutiveBiz, accessed on April 22, 2025, <https://executivebiz.com/2025/03/all-about-artificial-general-intelligence-govcons/>
12. Advancing AI Standards Collaboration: ISO, IEC, and ITU Announce 2025 International AI Standards Summit - American National Standards Institute, accessed on April 22, 2025, <https://www.ansi.org/standards-news/all-news/2024/10/10-15-24-advancing-ai-standards-collaboration-iso-iec-and-itu-announce-ai-standards-summit>
13. Global Summit 2025 - AI Standards Hub, accessed on April 22, 2025, <https://aistandardshub.org/global-summit/>
14. Press Release - ITU, accessed on April 22, 2025, <https://www.itu.int/en/mediacentre/Pages/PR-2025-02-06-AI-for-Good-2025-announcement.aspx>
15. Ten AI Governance Priorities: Survey of 40 Civil Society ..., accessed on April 22, 2025, <https://thefuturesociety.org/cso-ai-governance-priorities/>
16. MIRI's 2024 End-of-Year Update - Machine Intelligence Research Institute, accessed on April 22, 2025,

- <https://intelligence.org/2024/12/02/miris-2024-end-of-year-update/>
17. Leading AI safety organization drops technical research to focus exclusively on policy, accessed on April 22, 2025, <https://getcoai.com/news/leading-ai-safety-organization-drops-technical-research-to-focus-exclusively-on-policy/>
 18. Artificial General Intelligence: What Are We Investing In? | TechPolicy.Press, accessed on April 22, 2025, <https://www.techpolicy.press/artificial-general-intelligence-what-are-we-investing-in/>
 19. Future of AI Research - AAAI, accessed on April 22, 2025, <https://aaai.org/wp-content/uploads/2025/03/AAAI-2025-PresPanel-Report-FINAL.pdf>
 20. Forbes 2025 AI 50 List - Top Artificial Intelligence Companies Ranked, accessed on April 22, 2025, <https://www.forbes.com/lists/ai50/>
 21. OpenAI Charter, accessed on April 22, 2025, <https://openai.com/charter/>
 22. Blog | Federal AI Mandates and Corporate Compliance ... - Cogent, accessed on April 22, 2025, <https://www.cogentinfo.com/resources/federal-ai-mandates-and-corporate-compliance-whats-changing-in-2025>
 23. National Artificial Intelligence Research Institutes - NSF, accessed on April 22, 2025, <https://www.nsf.gov/funding/opportunities/national-artificial-intelligence-research-institutes/505686/nsf23-610>
 24. The 2025 AI Index Report | Stanford HAI, accessed on April 22, 2025, <https://hai.stanford.edu/ai-index/2025-ai-index-report>
 25. Stanford HAI's 2025 AI Index Reveals Record Growth in AI Capabilities, Investment, and Regulation - Business Wire, accessed on April 22, 2025, <https://www.businesswire.com/news/home/20250407539812/en/Stanford-HAIs-2025-AI-Index-Reveals-Record-Growth-in-AI-Capabilities-Investment-and-Regulation>
 26. AI Index 2025: State of AI in 10 Charts | Stanford HAI, accessed on April 22, 2025, <https://hai.stanford.edu/news/ai-index-2025-state-of-ai-in-10-charts>
 27. AI by AI Weekly Top 5: 03.10-16, 2025 - Champaign Magazine, accessed on April 22, 2025, <https://champaignmagazine.com/2025/03/16/ai-by-ai-weekly-top-5-03-10-16-2025/>
 28. Baidu to launch Ernie 5 AI in 2025 | Digital Watch Observatory, accessed on April 22, 2025, <https://dig.watch/updates/baidu-to-launch-ernie-5-ai-in-2025>
 29. Big Players in AI - dentro.de/ai, accessed on April 22, 2025, https://dentro.de/ai/big_players/
 30. Path to AGI - Part 2 - Hiflylabs, accessed on April 22, 2025, <https://hiflylabs.com/blog/2025/1/27/path-to-agi-part-2>
 31. Top AI Research Companies Driving Innovation, accessed on April 22, 2025, <https://aisuperior.com/ai-research-companies/>
 32. Event details — AI Safety Connect 2025, accessed on April 22, 2025,

- <https://www.aisafetyconnect.com/event-details>
33. Baidu Research, accessed on April 22, 2025, <https://research.baidu.com/>
 34. Public Policy - Partnership on AI, accessed on April 22, 2025, <https://partnershiponai.org/program/policy/>
 35. About - OpenAI, accessed on April 22, 2025, <https://openai.com/about/>
 36. OpenAI 'now knows how to build AGI' - The Rundown AI, accessed on April 22, 2025, <https://www.therundown.ai/p/openai-now-knows-how-to-build-agi>
 37. Planning for AGI and beyond | OpenAI, accessed on April 22, 2025, <https://openai.com/index/planning-for-agi-and-beyond/>
 38. Leadership updates - OpenAI, accessed on April 22, 2025, <https://openai.com/index/leadership-updates-march-2025/>
 39. Artificial General Intelligence: Is AGI Really Coming by 2025? - Hyperight, accessed on April 22, 2025, <https://hyperight.com/artificial-general-intelligence-is-agi-really-coming-by-2025/>
 40. Research | OpenAI, accessed on April 22, 2025, <https://openai.com/research>
 41. Meta chief AI scientist claims AGI will be viable in 3-5 years - National Technology News, accessed on April 22, 2025, https://nationaltechnology.co.uk/Meta_Chief_AI_Scientist_Claims_AGI_Will_Be_Viable_In_3_5_Years.php
 42. CEO Speaker Series With Dario Amodei of Anthropic | Council on Foreign Relations, accessed on April 22, 2025, <https://on.cfr.org/4iydDIW>
 43. Dario Amodei Is on the 2025 TIME100 List | TIME, accessed on April 22, 2025, <https://time.com/collections/100-most-influential-people-2025/7273747/dario-amodei/>
 44. Dario Amodei - Wikipedia, accessed on April 22, 2025, https://en.wikipedia.org/wiki/Dario_Amodei
 45. Dario Amodei, accessed on April 22, 2025, <https://www.darioamodei.com/>
 46. accessed on April 21, 2025, <https://www.anthropic.com/company>
 47. Meta's Chief AI Scientist Yann LeCun Questions the Longevity of Current GenAI and LLMs, accessed on April 22, 2025, <https://www.hpcwire.com/2025/02/11/metas-chief-ai-scientist-yann-lecun-questions-the-longevity-of-current-genai-and-llms/>
 48. Meta's AI Chief Questions the Viability of Generative AI - The Munich Eye, accessed on April 22, 2025, <https://themunicheye.com/metas-ai-chief-questions-generative-ais-future-10206>
 49. Meta's Llama 4 is now available on Workers AI - The Cloudflare Blog, accessed on April 22, 2025, <https://blog.cloudflare.com/meta-llama-4-is-now-available-on-workers-ai/>
 50. Meta's Llama 4 Large Language Models now available on Snowflake Cortex AI, accessed on April 22, 2025, <https://www.snowflake.com/en/blog/meta-llama-4-now-available-snowflake-cortex-ai/>
 51. The Llama 4 herd: The beginning of a new era of natively multimodal AI

- innovation, accessed on April 22, 2025,
<https://ai.meta.com/blog/llama-4-multimodal-intelligence/>
52. AI at Meta Blog, accessed on April 22, 2025, <https://ai.meta.com/blog/>
53. LeCun: "If you are interested in human-level AI, don't work on LLMs." : r/agi - Reddit, accessed on April 22, 2025,
https://www.reddit.com/r/agi/comments/1imqson/lecun_if_you_are_interested_in_humanlevel_ai_dont/
54. AI at Meta, accessed on April 22, 2025, <https://ai.meta.com/>
55. Stanford HAI "2025 AI Index Report" Highlights - BlockBeats, accessed on April 22, 2025, <https://m.theblockbeats.info/en/news/57740>
56. The Future of AI Governance: What 2025 Holds for Ethical Innovation - Solutions Review, accessed on April 22, 2025,
<https://solutionsreview.com/data-management/the-future-of-ai-governance-what-2025-holds-for-ethical-innovation/>
57. Governing AI for the Future of Humanity • Stimson Center, accessed on April 22, 2025, <https://www.stimson.org/2025/governing-ai-for-the-future-of-humanity/>
58. 2025 AI Policy Leaders, accessed on April 22, 2025,
<https://www.caidp.org/events/washdc25aidv/bios/>
59. Partnership for AI Evidence (PAIE) | The Abdul Latif Jameel Poverty ..., accessed on April 22, 2025,
<https://www.povertyactionlab.org/initiative/partnership-ai-evidence>
60. AI Index | Stanford HAI - Stanford University, accessed on April 22, 2025,
<https://aiindex.stanford.edu/report/>
61. GovAI Center of AI Governance Summer Fellowship 2025 - Scholar Digger, accessed on April 22, 2025,
<https://www.scholardigger.com/post/govai-center-of-ai-governance-summer-fellowship-2025>
62. Research Scholar | GovAI Blog - Centre for the Governance of AI, accessed on April 22, 2025, <https://www.governance.ai/post/research-scholar>
63. Summer Fellowship 2025 | GovAI Blog - Centre for the Governance of AI, accessed on April 22, 2025,
<https://www.governance.ai/post/summer-fellowship-2025>
64. GovAI Summer Fellowship 2025 in London, UK (Fully Funded) - Opportunity 4U, accessed on April 22, 2025,
<https://www.opportunit4u.com/2024/12/govai-summer-fellowship-2025-in-london-uk.html>
65. Explore our Work - CSER - Centre for the Study of Existential Risk, accessed on April 22, 2025, <https://www.cser.ac.uk/work/>
66. Centre for the Study of Existential Risk - Nuclear Weapons, accessed on April 22, 2025,
<https://nuclearweapons.info/organization/the-centre-for-the-study-of-existential-risk/>
67. Centre for the Study of Existential Risk, accessed on April 22, 2025,
<https://www.cser.ac.uk/>
68. Centre for the Study of Existential Risk (@cser.bsky.social) - Bluesky, accessed on

- April 22, 2025, <https://web-cdn.bsky.app/profile/cser.bsky.social>
69. Our Work - MIRI Technical Governance Team - Machine Intelligence Research Institute, accessed on April 22, 2025, <https://techgov.intelligence.org/research>
 70. MIRI Technical Governance Team | MIRI TGT - Machine Intelligence Research Institute, accessed on April 22, 2025, <https://techgov.intelligence.org/>
 71. Center for Ethics appoints Anne-Elisabeth Courrier as AI ethics liaison | Emory University, accessed on April 22, 2025, https://news.emory.edu/stories/2025/04/er_ai_ethics_liaison_14-04-2025/story.html
 72. The year of conscious AI | The AI Journal, accessed on April 22, 2025, <https://aijourn.com/the-year-of-conscious-ai/>
 73. WATCH: A Neuroscientist and a Philosopher Debate AI ..., accessed on April 22, 2025, <https://ai.princeton.edu/news/2025/watch-neuroscientist-and-philosopher-debate-ai-consciousness>
 74. AI and Human Consciousness: Examining Cognitive Processes | American Public University, accessed on April 22, 2025, <https://www.apu.apus.edu/area-of-study/arts-and-humanities/resources/ai-and-human-consciousness/>
 75. Oxford Institute for Ethics in AI launches Accelerator Fellowship ..., accessed on April 22, 2025, <https://www.ox.ac.uk/news/2025-03-03-oxford-institute-ethics-ai-launches-accelerator-fellowship-programme>
 76. Evaluating Artificial Consciousness 2025 - PhilEvents, accessed on April 22, 2025, <https://philevents.org/event/show/131314>
 77. Large AI Model Lecture Series: Can Machines Become Conscious? - Princeton University, accessed on April 22, 2025, <https://www.princeton.edu/events/2025/large-ai-model-lecture-series-can-machines-become-conscious>
 78. Being in the World: Will AI Ever Have a Soul? - Campus Calendar, accessed on April 22, 2025, <https://calendar.gatech.edu/event/2025/03/10/being-world-will-ai-ever-have-soul>
 79. Machine Intelligence Research Institute, accessed on April 22, 2025, <https://intelligence.org/>
 80. AI Safety I: Concepts and Definitions, accessed on April 22, 2025, <https://synthesis.ai/2025/04/17/ai-safety-i-concepts-and-definitions/>
 81. Machine Intelligence Research Institute - Wikipedia, accessed on April 22, 2025, https://en.wikipedia.org/wiki/Machine_Intelligence_Research_Institute
 82. The Machine Intelligence Research Institute - Raising for Effective Giving (REG), accessed on April 22, 2025, <https://reg-charity.org/recommended-charities/machine-intelligence-research-institute/>
 83. What is the Machine Intelligence Research Institute's research agenda?, accessed on April 22, 2025, https://aisafety.info/?state=85EN_
 84. MIRI Strategy Archives - Machine Intelligence Research Institute, accessed on

- April 22, 2025, <https://intelligence.org/category/miri/>
85. MIRI's 2024 End-of-Year Update — EA Forum, accessed on April 22, 2025, <https://forum.effectivealtruism.org/posts/zBizzn2BT6pbqS8n/miri-s-2024-end-of-year-update>
 86. What is MIRI currently doing? — EA Forum, accessed on April 22, 2025, <https://forum.effectivealtruism.org/posts/e8o6paib9sgKeWorc/what-is-miri-currently-doing-1>
 87. Machine Intelligence Research Institute — AI Safety Retraining Program, accessed on April 22, 2025, <https://www.openphilanthropy.org/grants/machine-intelligence-research-institute-ai-safety-retraining-program/>
 88. News Archives - Machine Intelligence Research Institute, accessed on April 22, 2025, <https://intelligence.org/category/news/>
 89. Thoughts on AI 2027 - Machine Intelligence Research Institute, accessed on April 22, 2025, <https://intelligence.org/2025/04/09/thoughts-on-ai-2027/>
 90. New Benchmarks Challenge Brute Force Approach to AI, accessed on April 22, 2025, <https://www.govinfosecurity.com/new-benchmarks-challenge-brute-force-approach-to-ai-a-27826>
 91. Critiques of prominent AI safety labs: Conjecture — EA Forum, accessed on April 22, 2025, <https://forum.effectivealtruism.org/posts/gkfMLX4NWZdmpikto/critiques-of-prominent-ai-safety-labs-conjecture>
 92. AI Control Hackathon 2025 | Apart Research, accessed on April 22, 2025, <https://apartresearch.com/sprints/ai-control-hackathon-2025-03-29-to-2025-03-30>
 93. Intelligence Sequencing and the Path-Dependence of Intelligence Evolution: AGI-First vs. DCI-First as Irreversible Attractors | AI Research Paper Details - AIModels.fyi, accessed on April 22, 2025, <https://www.aimodels.fyi/papers/arxiv/intelligence-sequencing-path-dependence-intelligence-evolution-agi>
 94. [2503.17688] Intelligence Sequencing and the Path-Dependence of Intelligence Evolution: AGI-First vs. DCI-First as Irreversible Attractors - arXiv, accessed on April 22, 2025, <https://arxiv.org/abs/2503.17688>
 95. 2503.17688v1 | PDF | Self Organization | Conceptual Model - Scribd, accessed on April 22, 2025, <https://www.scribd.com/document/842884183/2503-17688v1>
 96. Dependence of Intelligence Evolution: AGI-First vs. DCI-First as Irreversible Attractors - Qeios, accessed on April 22, 2025, <https://www.qeios.com/read/RA5XMP/pdf>
 97. Intelligence Sequencing and the Path-Dependence of Intelligence Evolution: AGI-First vs. DCI-First as Irreversible Attractors - Qeios, accessed on April 22, 2025, <https://www.qeios.com/read/RA5XMP>
 98. Intelligence Sequencing and the Path-Dependence of Intelligence Evolution: AGI-First vs. DCI- First as Irreversible Attractors A - arXiv, accessed on April 22, 2025, <https://arxiv.org/pdf/2503.17688>

99. Pluralis Research Pioneers Protocol Learning to Scale, accessed on April 22, 2025, <https://www.globenewswire.com/news-release/2025/03/19/3045635/0/en/Pluralis-Research-Pioneers-Protocol-Learning-to-Scale-Decentralized-AI-Announces-7-6M-Seed-Round-Led-by-USV-and-CoinFund.html>
100. Decentralized AI: The Key To Unlocking Artificial General ..., accessed on April 22, 2025, <https://live.upcoming.sk/2025/04/04/decentralized-ai-the-key-to-unlocking-artificial-general-intelligence/>
101. The future of AI will be decentralised and collaborative | Institut ..., accessed on April 22, 2025, <https://www.ip-paris.fr/en/news/future-ai-will-be-decentralised-and-collaborative>
102. Decentralized science meets AI — legacy institutions aren't ready, accessed on April 22, 2025, <https://cointelegraph.com/news/decentralized-science-meets-ai>
103. [Frontiers of AI and Computing: A Conversation With Yann LeCun and Bill Dally | NVIDIA GTC 2025](#)