

MARKET SCAN REPORT

JULY 2025

BY INFOSYS TOPAZ RESPONSIBLE AI OFFICE

Infosys
topaz



IN FOCUS

**THE AI PARADOX: WHY
PROMISING POCS RARELY SEE
THE LIGHT OF PRODUCTION**

By Stacey Miller

Infosys®
Navigate your next

“Foreword

“As we launch the July 2025 edition of the Infosys Responsible AI Market Scan, our aim remains clear, “To surface critical global trends, regulatory shifts, and thoughtful innovation in the responsible AI ecosystem.” This report builds on the rigorous standards you’ve come to expect from our monthly updates, much like the June edition that spotlighted landmark developments such as Texas’s TRAIGA legislation, early U.S. Senate proposals for federal oversight, and the formalization of AI Impact Assessments under ISO/IEC 42005:2025.

In this edition, we spotlight key regulatory developments, critical security incidents, and breakthrough technologies shaping the Responsible AI landscape. One conversation that deserves special attention stems from a warning by Sam Altman, CEO of OpenAI to young people relying ‘too much’ on ChatGPT. As users, we’re confiding in AI systems without the fundamental legal safeguards we take for granted in human interactions. People share deeply personal anecdotes, health concerns, financial queries, or even proprietary business information, believing it to be a private interaction with a digital assistant. It’s something I explored in a recent reflection: *We Cry About AI Breaching Privacy. But We’re the Ones Feeding It Our Secrets!*

Meanwhile, agentic AI is becoming the norm, systems like Moonshot AI’s Kimi K2 and Google’s Gemini Lite are designed to reason, initiate, and collaborate autonomously. In parallel, vibe coding, where natural language is used to generate code while intuitively capturing tone, context, and intent is transforming how AI is built and aligned with human needs. Tools like FusionBench, AgentEval, and OpenAI’s preparedness frameworks are raising the bar for AI evaluation and trust.

This edition of the *Market Scan Report* captures the complexity of this moment. I hope it informs your view, sharpens your strategies, and deepens your commitment to shaping AI that is not only powerful but ethical, inclusive, and accountable.



Syed Ahmed
Global Head
Infosys Responsible AI Office



From the editor's desk

Global AI Governance, Safety, and Innovation – July Edition

We've reached the midpoint of 2025, and artificial intelligence has firmly established itself as a central force shaping global politics, economic strategy, and technological power. Nations are racing not just to innovate, but to regulate, align, and influence the global AI agenda—whether through multilateral cooperation, strategic investments, or ideological frameworks.

Governance Steps Up

Around the world, governments are stepping forward with new frameworks and regulations to ensure that AI evolves responsibly. In the United States, several important moves—like the formation of a National Task Force on AI in criminal justice, California's anti-discrimination measures, and Colorado's sweeping risk management law—signal a shift toward proactive oversight.

Elsewhere, the UK and Canada are investing in transparency and auditability, with a special focus on tackling deepfakes. The EU is moving steadily towards enforcement of the AI Act while finalizing its GPAI Code of Practice. Across Asia, India, Singapore, and South Korea are rolling out national strategies, regulatory sandboxes, and cybersecurity initiatives—all pointing to a rising commitment to responsible AI deployment.

Global cooperation is gaining steam as well. From the World Economic Forum and UAE's GRIP initiative to BRICS' push for inclusive governance, it's clear that countries are beginning to think collectively. On the standards front, ISO/IEC 42006 and new releases from the World Digital Technology Academy are helping to shape a more harmonized global approach to ethical AI.

Incidents We Can't Ignore

July also saw a number of concerning incidents that remind us how quickly AI can be misused. Deepfakes were used to mislead and defraud people—from scams in Malaysia to misinformation in Chicago. In South Africa, fake legal citations generated by AI, and elsewhere, questionable chatbot deployments, stirred public concern. Even major companies weren't immune—WhatsApp faced a glitch with its assistant and McDonald's hiring bot raised data security concerns.

These events are stark reminders that as we build smarter systems, safety, accountability, and public trust must remain at the core.

Building Better Defenses

Thankfully, the AI safety community is stepping up with meaningful innovations. Tools like TrustDefender, which uses privacy-preserving methods to detect deepfakes, and BBoxER, designed for optimizing large language models, show how AI can protect against its own risks. Meanwhile, new frameworks to resist jailbreaks and adversarial manipulation of AI systems are strengthening our digital defenses.

Innovation Continues to Accelerate

On the frontier of development, the AI race is heating up with powerful new models and benchmarks. China's Moonshot AI launched Kimi K2, aiming to lead with advanced agentic capabilities. Google unveiled Gemini Lite, making high-speed generative AI possible even on edge devices. OpenAI introduced its Preparedness Framework alongside the SAFER system to help better manage risks and model alignment.

We're also seeing an exciting wave of tools designed for better evaluation and fine-tuning: FusionBench, IFScale, AgentTuner, and AgentEval are all ushering in a new era of multi-model cooperation and agent performance optimization.

A Moment for Reflection

Let this edition challenge you to do more than stay informed.

Choose **one development** from this report—be it a policy shift, an innovation, or a risk—and explore its **implications for your organization, your community, or your country.**

Share your thoughts, start a conversation, or draft a plan. Because the future of AI isn't waiting—it's being shaped right now.

And finally, don't miss this edition's guest perspective by Stacey Miller, Principal Product Marketing Manager at SUSE AI, who explores *"The AI Paradox: Why Promising PoCs Rarely See the Light of Production."* Her insights offer a timely reminder that successful AI isn't just about innovation—it's about execution.

Warm regards,

Ashish Tewari

Head- Infosys Responsible AI Office, India

Table of Contents

AI Regulations, Governance & Standards

AI Regulations & Governance across the globe 05

Standards 21

AI Principles

Incidents 22

Vulnerabilities 26

Defences 26

In Focus

The AI Paradox: Why Promising PoCs Rarely See the Light
of Production..... 28

Technical Updates

New Model Released..... 30

New Frameworks & Research Techniques..... 32

New Agentic Researches 36

Industry Updates

Healthcare 38

Finance 38

Education 38

Transportation 39

Agriculture 39

Infosys Developments

Events..... 40

Infosys Responsible AI Toolkit – A Foundation for Ethical AI .. 42

Contributors





AI Regulations, Governance & Standards

This section highlights the recent updates on regulations and governance initiatives across the globe impacting the responsible development and deployment of AI.

AI Regulations & Governance across the globe

World Economic Forum and UAE Launch Global Platform to Modernize AI and Tech Regulation Worldwide

The World Economic Forum, in collaboration with the United Arab Emirates, has launched the Global Regulatory Innovation Platform (GRIP) to help governments develop smarter, faster, and more flexible rules for emerging technologies like artificial intelligence (AI), biotechnology, and digital finance. GRIP brings together global experts from policy, industry, and law to rethink how regulations should work in today's fast-changing tech environment. The platform will introduce three major tools: a Global Regulatory Playbook with practical case studies, a Readiness Index to assess how prepared countries are for future technologies, and a Global Innovation Hub to test and scale new regulatory ideas. The initiative aims to ensure that AI and other disruptive technologies are governed in ways that protect people, promote fairness, and support innovation.¹

UK and Singapore Forge Strategic AI Partnership in Digital Finance at 10th Financial Dialogue

During the 10th UK-Singapore Financial Dialogue, the Monetary Authority of Singapore (MAS) and the UK's Financial Conduct Authority (FCA) advanced their cooperation in the realm of artificial intelligence (AI) within the financial services sector. The two regulatory bodies engaged in comprehensive discussions covering current AI trends, emerging use cases, challenges to adoption, and their respective regulatory approaches. Demonstrating a shared commitment to innovation and

responsible AI deployment, MAS and FCA agreed to initiate a joint collaboration focused on exchanging cutting-edge AI solutions and facilitating dialogue on cross-border AI developments, marking a significant step toward harmonizing global standards in digital finance.²

BRICS Nations Unite for Inclusive and Responsible Global AI Governance

Following the 17th BRICS Summit in Rio de Janeiro, leaders of the BRICS countries— - Brazil, Russia, India, China, and South Africa— - signed a joint statement emphasizing the urgent need for global cooperation on artificial intelligence governance. The statement calls for a collective international effort to manage AI risks, build public trust, and ensure that AI development benefits all nations, especially those in the Global South. The leaders stressed that AI governance should reflect shared human values, promote fairness, and encourage inclusive access to AI technologies. This unified stance highlights BRICS' commitment to shaping a future where AI is developed and used responsibly across borders.³

China and Brazil Strengthen Bilateral Ties Through Expanded AI and Strategic Development Cooperation

China and Brazil have deepened their partnership by signing a series of cooperation agreements across key sectors including infrastructure, pharmaceuticals, new energy, artificial intelligence (AI), and strategic development planning. As part of this effort, both countries signed a memorandum of understanding to launch the second phase of aligning their national development strategies. A major highlight of the collaboration is the establishment of an AI application cooperation center, formalized through an agreement between China's National Development and Reform Commission (NDRC) and Brazil's Ministry of Science, Technology, and Innovation. This center is designed to foster joint innovation in AI by building a strong foundation for development, offering open-source and open-access services, and supporting talent cultivation in both nations.⁴

Strategic Tech Collaboration: Israel and US to Establish \$200M AI and Quantum Science Hub

Israel and the United States have announced plans to establish a joint science center dedicated to advancing artificial intelligence (AI) and quantum technologies, backed by a \$200 million investment. The initiative is being led by Maj. Gen. (res.) Tamir Hayman, director of Israel's Institute for National Security Studies (INSS), and Dr. Smadar Itzkovich, founder and CEO of the AI & Quantum Sovereignty Lab (AIQ-Lab). Designed to deepen bilateral scientific collaboration and bolster strategic influence in emerging technologies, the center is expected to be formalized either through a presidential executive order by U.S. President Donald Trump or via legislative approval. The project reflects a shared commitment to innovation, regional stability, and technological leadership in the face of global competition.⁵

¹ <https://www.weforum.org/press/2025/07/new-global-platform-launched-to-reimagine-regulation-in-age-of-disruptive-tech/>

² <https://www.gov.uk/government/publications/10th-uk-singapore-financial-dialogue-july-2025-joint-statement/10th-uk-singapore-financial-dialogue-july-2025-joint-statement>

³ <https://www.pib.gov.in/PressReleasePage.aspx?PRID=2142786&>

⁴ https://www.ndrc.gov.cn/fzggw/wld/lss/zygz/202507/t20250707_1399017.html

⁵ <https://www.timesofisrael.com/israel-and-us-to-forge-200m-tech-hub-for-ai-and-quantum-science-development/>



US Court Decides in Favour of Meta Over AI Training Practices, Citing Lack of Specific Evidence of Market Harm

A US District Court granted summary judgment in favour of Meta Platforms Inc. in a copyright infringement lawsuit brought by a group of authors, who alleged that Meta unlawfully used their copyrighted works to train its generative AI models. The plaintiffs claimed that Meta's AI systems were trained on datasets containing their protected works without permission. However, the court found that the plaintiffs failed to provide sufficient evidence of market harm caused. The ruling marks a significant development in the evolving legal landscape surrounding AI training and intellectual property, reinforcing the need for precise claims and evidentiary support in litigation involving generative AI technologies.⁶

U.S. Council on Criminal Justice Launches National Task Force to Guide Ethical AI Integration in Criminal Justice System

The Council on Criminal Justice (CCJ), in collaboration with RAND Corporation, has established a National Task Force on Artificial Intelligence to develop ethical standards and practical guidelines for the use of AI across the U.S. criminal justice system. Chaired by former Texas Supreme Court Chief Justice Nathan Hecht, the 15-member task force includes a diverse group of stakeholders—ranging from AI experts and law enforcement officials to civil rights advocates and individuals with lived experience in the justice system. Over an 18-month period, the task force will produce consensus principles, operational standards, and accessible research to support responsible AI adoption in law enforcement, courts, corrections, and community organizations. The initiative is backed by funding from Microsoft, The Just Trust, and the MacArthur Foundation, and aims to ensure that AI enhances fairness, public safety, and efficiency without reinforcing bias or eroding public trust.⁷

Virginia Uses AI to Streamline Government Regulations: Governor Youngkin Introduces Executive Order 51

Virginia Governor Glenn Youngkin has announced Executive Order 51, launching a pilot program that makes Virginia the first state to use agentic artificial intelligence (AI) to review

⁶ <https://fingfx.thomsonreuters.com/gfx/legaldocs/zgvozmrynpd/META%20AI%20COPYRIGHT%20LAWSUIT%20ruling.pdf>

⁷ <https://counciloncj.org/council-launches-national-task-force-to-guide-integration-and-oversight-of-ai-in-criminal-justice/>

and simplify its regulatory framework. The initiative employs generative AI to examine existing regulations and guidance documents, identifying contradictions, redundancies, and overly complex language. This effort builds on Virginia's previous achievement of reducing regulatory requirements by over 25%, which has saved citizens more than \$1.2 billion annually. The new AI-driven approach aims to increase that reduction to 33%, further improving government efficiency. State officials, including Reeve Bull and Rob Ward, have described the initiative as a model for other states seeking to modernize their regulatory systems through technology.⁸

Comprehensive AI Integration Framework for K–12 Education: Missouri's 2025–26 Guidance for Local Education Agencies in the United States

The Missouri Department of Elementary and Secondary Education (DESE), in collaboration with its Computer Science Advisory Council, has released a comprehensive guidance document for the 2025–26 academic year to support Local Education Agencies (LEAs) in the United States in responsibly integrating Artificial Intelligence (AI) into K–12 education. This U.S.-based framework outlines the ethical, practical, and pedagogical considerations of AI use in schools, emphasizing student safety, data privacy, and academic integrity. It defines AI and generative AI, explores their benefits and challenges across students, educators, administrators, and communities, and provides best practices for implementation, including professional development, prompt engineering, and policy development. The guidance stresses the importance of human oversight, transparency, and equity, while encouraging innovation and curiosity in the classroom. It also addresses legal compliance with U.S. regulations such as FERPA, COPPA, and IDEA, and recommends a cyclical policy review process to adapt to the rapidly evolving AI landscape. Designed specifically for Missouri but adaptable elsewhere, the document serves as a strategic roadmap for preparing students and educators for an AI-driven future.⁹

California Approves New AI Employment Regulations to Prevent Discrimination in the Workplace

The California Civil Rights Council has secured approval for a set of regulations aimed at preventing employment discrimination caused by artificial intelligence (AI), algorithms, and automated decision-making systems. These rules clarify how existing anti-discrimination laws apply when employers use AI tools in hiring, promotions, and other employment decisions. Developed through extensive public consultation, the regulations require employers to maintain records of automated decision-making data for at least four years and ensure that these systems do not result in biased outcomes based on protected characteristics such as race, gender, or disability. Examples include AI systems that reject

female applicants due to biased training data or target job ads based on stereotypes. The rules reinforce California's commitment to fairness, transparency, and civil rights in the workplace and are expected to influence similar policies in other states.¹⁰

U.S. House Judiciary Subcommittee Investigates Criminal Misuse of AI

The U.S. House Judiciary Subcommittee on Crime and Federal Government Surveillance held a hearing to examine the growing threat of artificial intelligence being used for criminal exploitation. The discussion focused on how AI is enabling new forms of crime such as fraud, identity theft, child exploitation, and other illegal activities, allowing offenders to operate more efficiently and evade detection. Expert witnesses and lawmakers explored the challenges law enforcement faces in keeping up with these rapidly evolving technologies, including limited resources and outdated legal frameworks. The hearing also considered potential legislative and policy solutions to strengthen public safety and ensure that the justice system can effectively respond to AI-driven threats.¹¹

U.S. Lifts Ban on NVIDIA AI Chip Sales to China, Allowing H20 Exports to Resume

NVIDIA has announced that the U.S. government will now allow it to resume exports of its H20 AI chips to China, following the lifting of a ban that had been imposed in April 2025 by the Trump administration. The original restriction was based on concerns that the advanced chips could be used by the Chinese military. However, as of July 14, 2025, the administration has decided to grant export licenses, enabling NVIDIA to legally sell the H20 chips in China once again. This move marks a significant shift in U.S. tech policy and could have major implications for global AI development and trade relations between the two countries.¹²

Colorado Sets National Standard with New AI Law Requiring Oversight and Risk Controls for High-Stakes Technologies

Colorado has enacted a far-reaching artificial intelligence law that will take effect on February 1, 2026, establishing one of the most rigorous regulatory frameworks in the United States. The law targets both developers and users of high-risk AI systems, those used in areas like hiring, lending, healthcare, education, and insurance—and requires them to implement detailed risk management programs. These programs must include impact assessments, oversight procedures, and strategies to reduce potential harm, especially algorithmic discrimination. Companies are also required to notify the Colorado Attorney General and, in some cases, affect individuals if their AI systems cause discriminatory outcomes. Analysts say the law's depth and complexity rivals the European Union's AI Act, positioning Colorado as a leader in responsible AI governance and potentially influencing future legislation across other states.¹³

⁸ <https://www.governor.virginia.gov/newsroom/news-releases/2025/july/name-1053152-en.html>

⁹ https://static1.squarespace.com/static/64398599b0c21f1705fb8fb3/t/6864054cf5bf645acadc411e/1751385421046/AI_Guidance-Accessible.pdf

¹⁰ <https://calcivilrights.ca.gov/2025/06/30/civil-rights-council-secures-approval-for-regulations-to-protect-against-employment-discrimination-related-to-artificial-intelligence/>

¹¹ <https://judiciary.house.gov/committee-activity/hearings/artificial-intelligence-and-criminal-exploitation-new-era-risk-0>

¹² <https://blogs.nvidia.com/blog/nvidia-ceo-promotes-ai-in-dc-and-china/>

¹³ <https://www.hrdive.com/news/a-heavy-lift-colorado-ai-law-sets-high-bar-analysts-say/753069/>

United States – Pennsylvania Enacts Digital Forgery Law to Criminalize Malicious Deepfakes and Protect Citizens from AI Scams

Pennsylvania has enacted the Digital Forgery Law, signed by Governor Josh Shapiro, to address the growing threat of malicious deepfakes and AI-generated impersonation. Under Senate Bill 649, it is now a third-degree felony to use artificial intelligence to create deceptive digital content—such as fake voices or videos—intended to defraud, exploit, or harm individuals, especially in scams targeting vulnerable groups like older adults. The law builds on prior legislation targeting AI-generated child sexual abuse material and non-consensual intimate imagery and equips law enforcement with new tools to prosecute digital deception. Supported by bipartisan lawmakers, the law reflects Pennsylvania's commitment to balancing technological innovation with public safety and responsible AI governance.¹⁴

United States – White House signed Executive Order to Ensure Political Neutrality in Federal AI Systems

The White House has advanced efforts to promote ideological neutrality in AI through Executive Order (EO) titled “Preventing Woke AI in the Federal Government,” signed by President Trump on July 23, 2025. This order builds on the administration's broader “America's AI Action Plan,” released the same day, which emphasizes removing barriers to AI innovation while ensuring federal systems remain nonpartisan and objective. The EO has already been issued, reflecting the Trump administration's priority to counter perceived biases in AI, such as those related to diversity, equity, and inclusion (DEI) ideologies.¹⁵



¹⁴ <https://www.fiannafail.ie/news/deputy-naoise-%C3%B3-cear%C3%BAil-introduces-bill-to-establish-national-artificial-intelligence-office>

¹⁵ <https://www.wsj.com/tech/ai/white-house-prepares-executive-order-targeting-woke-ai-e68e8e24>



Navigating the Future: BIBA's AI Guide Equips Insurance Brokers for Responsible Innovation

The British Insurance Brokers' Association (BIBA), in collaboration with sponsor Markel and several leading industry contributors, has launched a detailed guide to help insurance brokers understand and adopt artificial intelligence (AI) responsibly. Tailored to address both opportunities and risks, the guide offers practical advice on how AI can enhance customer service, streamline operations, and support decision-making, while also covering essential legal, regulatory, and ethical considerations. Developed with input from firms such as Gallagher Re, PIB, FullCircI, Romero, CFC, and Cyxcel, the guide includes real-world examples and expert insights. BIBA CEO Graeme Trudgill emphasized that the guide reflects member feedback and will be regularly updated to keep pace with technological advancements. Available exclusively to BIBA members, this resource is positioned as a strategic tool for brokers aiming to innovate confidently and compliantly in the evolving insurance landscape.¹⁶

FRC Releases New Guidance to Clarify AI Use in Auditing and Promote Responsible Innovation

On 26 June 2025, the UK's Financial Reporting Council (FRC) published its first formal guidance on the use of artificial intelligence (AI) in the audit profession, aiming to support responsible adoption of emerging technologies while maintaining audit quality. Developed with input from the FRC's Technology Working Group and industry experts, the guidance outlines principles for implementing AI tools, including proportionate documentation and context-sensitive explainability aligned with the UK government's five AI principles. It is intended to assist auditors, central teams, and third-party technology providers in understanding regulatory expectations and integrating AI effectively. The release also includes a thematic review of the six largest audit firms, showcasing best practices in certifying AI-enabled tools and establishing robust controls. FRC Executive Director Mark Babington emphasized that AI is moving from experimentation to practical application in audits, with the potential to enhance quality, strengthen market confidence, and contribute to economic growth.¹⁷

¹⁶ <https://www.biba.org.uk/press-releases/ai-guide/>

¹⁷ <https://www.frc.org.uk/news-and-events/news/2025/06/frc-publishes-landmark-guidance-providing-clarity-to-audit-profession-on-the-uses-of-ai/>

MHRA Expands AI Airlock Program and Launches Global Network to Advance Safe Innovation in Health Technologies

On 27 June 2025, the UK's Medicines and Healthcare products Regulatory Agency (MHRA) announced a series of strategic initiatives to strengthen the regulation and safe adoption of artificial intelligence (AI) in healthcare. These include the expansion of the AI Airlock program—a regulatory sandbox that enables innovators to test AI-powered medical technologies in controlled environments—allowing MHRA to refine its frameworks based on real-world data. Alongside this, MHRA advanced its Centers of Excellence for Regulatory Science and Innovation (CERSIs), including RADIANT and CERSI-AI, which offer open-source tools, educational resources, and collaborative platforms for academics, clinicians, and regulators. These efforts aim to support innovators in navigating complex regulatory landscapes while ensuring patient safety. Additionally, MHRA launched a new global AI network for health regulators to foster international collaboration, harmonize standards, and promote responsible innovation across borders.¹⁸

UK Ofcom Publishes Deepfake Attribution Toolkit to Guide Platforms in Tackling Harmful AI-Generated Content

The UK Office of Communications (Ofcom) has released a discussion paper outlining key attribution measures to help digital platforms and AI developers combat the growing threat of harmful deepfakes. The paper focuses on synthetic audio-visual content that can be used for fraud, defamation, and disinformation, and presents four main tools to trace and verify such content: watermarking, provenance metadata, AI-generated content labels, and contextual annotations. These measures aim to help users recognize misleading content and support platforms in moderating it more effectively. Ofcom also highlights the challenges of implementing these tools, such as the risk of removal, user confusion, and inconsistent standards. The paper emphasizes that attribution should be combined with other safety strategies, including AI classifiers and red teaming, to build a more robust defense against malicious deepfake use.¹⁹



Europe

European Commission Dismisses Delay Rumors, Reaffirms Commitment to AI Act Implementation

Amid growing speculation about a possible delay in the rollout of the EU AI Act, the European Commission has firmly reiterated its commitment to the legislation's original timeline. At a recent press conference, Commission spokesperson Thomas Regnier addressed the rumors directly, stating, "I've seen, indeed, a lot of reporting, a lot of letters and a lot of things being said on the AI Act. Let me be as clear as possible, there is no stop the clock. There is no grace period. There is no pause." His statement emphasized the EU's unwavering stance on advancing its regulatory framework for artificial intelligence, reinforcing its leadership in shaping global standards for responsible and transparent AI governance.²⁰

European Commission Releases Final GPAI Code of Practice to Guide AI Industry Ahead of AI Act Enforcement

The European Commission has released the final version of the General Purpose AI (GPAI) Code of Practice, a voluntary framework designed to help AI developers and providers prepare for the upcoming enforcement of the EU's AI Act. Developed by 13 independent experts with input from over

¹⁸ <https://www.gov.uk/government/news/ai-airlock-cersis-and-a-new-global-ai-network-for-health-regulators>

¹⁹ <https://www.ofcom.org.uk/siteassets/resources/documents/online-safety/information-for-industry/deepfake-defences-2/deepfake-defences-2---the-attribution-toolkit.pdf>

²⁰ <https://www.reuters.com/world/europe/artificial-intelligence-rules-go-ahead-no-pause-eu-commission-says-2025-07-04/>

1,000 stakeholders including model providers, SMEs, academics, AI safety experts, rightsholders, and civil society organizations—the Code aims to support responsible AI development. It consists of three chapters: Transparency and Copyright, which apply to all GPAI providers, and Safety and Security, which target only the most advanced models. The AI Act's GPAI rules will take effect on 2 August 2025, with enforcement by the European AI Office beginning one year later for new models and two years later for existing ones. The Code will be complemented by official Commission guidelines to clarify which providers fall under the scope of these rules.²¹

European Parliament Releases Study Urging Copyright Reform to Address Generative AI's Legal Challenges

The European Parliament has published a comprehensive study on generative AI and copyright, commissioned by its Policy Department for Justice, Civil Liberties and Institutional Affairs at the request of the Committee on Legal Affairs. The study explores how generative AI technologies are testing the boundaries of existing EU copyright law, particularly highlighting the mismatch between AI training practices and current text and data mining exceptions, as well as the unclear legal status of AI-generated content. These issues pose significant risks to Europe's creative ecosystem, which relies on strong protection and fair

compensation for authors. To address these challenges, the report recommends clear legal distinctions between AI inputs and outputs, harmonized opt-out mechanisms for creators, mandatory transparency from AI developers, and fair licensing models. The European Parliament is expected to lead efforts to modernize copyright law in a way that balances technological innovation with the rights and sustainability of Europe's cultural and creative sectors.²²

EU AI Office Launches €9 Million Tender to Support Safe Oversight of General-Purpose AI Models

The European Union's AI Office has launched a €9 million tender to procure technical support for enforcing the EU AI Act, with a strong focus on monitoring and assessing the risks posed by General-Purpose AI (GPAI) systems. These AI models, which can be used across many sectors, may present serious risks to public safety, security, and fundamental rights. The initiative aims to strengthen the AI Office's ability to evaluate whether these systems comply with the EU's regulations and to ensure they are used responsibly. Funded under the Digital Europe Programme, the project will help build tools and expertise to detect and manage systemic risks from GPAI at the EU level.²³



Ireland

National AI Bill Introduced to Establish Ireland's First Artificial Intelligence Office

Ireland's Kildare North regions member of Parliament, Naoise Ó Cearúil, has introduced the National Artificial Intelligence Bill 2026 in Ireland's Leinster House, proposing the creation of a National Artificial Intelligence Office to guide and oversee the country's AI strategy. The office will provide independent oversight of high-risk AI systems, support both enterprises and public sector bodies in adopting and deploying AI technologies and promote digital literacy and national upskilling. It will also coordinate Ireland's overall approach to AI development and deployment, ensuring ethical standards and strategic alignment with global advancements. The bill marks a significant step toward positioning Ireland as a responsible and forward-thinking leader in artificial intelligence.²⁴



²¹ https://ec.europa.eu/commission/presscorner/detail/en/ip_25_1787

²² [https://www.europarl.europa.eu/thinktank/en/document/IUST_STU\(2025\)774095](https://www.europarl.europa.eu/thinktank/en/document/IUST_STU(2025)774095)

²³ <https://digital-strategy.ec.europa.eu/en/funding/eu-ai-office-launches-eu9-million-tender-technical-support-gpai-safety>

²⁴ <https://www.fiannafail.ie/news/deputy-naoise-%C3%B3-cear%C3%BAl-introduces-bill-to-establish-national-artificial-intelligence-office>



France

CNIL Clarifies Legal Limits on Use of AI-Powered Age Estimation Cameras in Retail Settings

The French data protection authority, CNIL, (National Commission on Informatics and Liberty) has issued a clarification regarding the use of AI-powered augmented cameras in retail environments that sell age-restricted products such as tobacco, alcohol, and gambling. These cameras estimate the age of customers by scanning and analyzing their faces automatically, aiming to prevent sales to minors. However, CNIL ruled that this practice involves personal data processing that is neither necessary nor proportionate under the General Data Protection Regulation (GDPR). Since technology only provides an estimate and cannot replace official ID checks, businesses are still required to verify age through identity documents. Additionally, the continuous and automatic operation of these cameras infringes on individuals' rights, including the right to object, raising serious concerns about privacy and data protection compliance.²⁵



Germany

German Authorities Flag DeepSeek AI App for Privacy Violations

The Berlin Commissioner for Data Protection and Freedom of Information, along with data protection authorities from Baden-Württemberg, Rhineland-Palatinate, and the Free Hanseatic City of Bremen, formally notified Apple and Google that the Chinese AI application DeepSeek constitutes illegal content under European law. Developed by Hangzhou DeepSeek Artificial Intelligence, the app unlawfully transfers large volumes of personal data from German users to servers in China without the safeguards required by Article 46(1) of the General Data Protection Regulation (GDPR). Despite prior requests to either remove the app or implement lawful data transfer mechanisms, the developer failed to comply. As a result, the Berlin Commissioner reported the application under Article 16 of the Digital Services Act, urging platform providers to take appropriate action to protect user privacy and uphold EU data protection standards.²⁶



²⁵ <https://www.cnil.fr/fr/cameras-augmentees-pour-estimer-lage-dans-les-bureaux-de-tabac-la-cnil-precise-sa-position>

²⁶ https://www.datenschutz-berlin.de/fileadmin/user_upload/pdf/publikationen/DSK/2025/20250627-Berlin-DPA-Press-Release-DeepSeek.pdf

Germany's Data Protection Authority Opens Public Consultation to Shape AI Privacy Regulations

Germany's Federal Commissioner for Data Protection and Freedom of Information (BfDI) has launched a public consultation to address growing concerns about how large AI models, especially language models, handle personal data. The initiative invites feedback from developers, researchers, and other stakeholders to help create practical and legally

sound rules for AI systems that may unintentionally reveal private information from their training data. Commissioner Prof. Dr. Louisa Specht-Riemenschneider emphasized the importance of transparency, legal clarity, and real-world relevance in shaping future regulations. The consultation is open until 11 August 2025 and aims to ensure that AI development in Germany respects privacy rights while supporting innovation through informed and inclusive policymaking.²⁷



Czech Republic

Czech Republic Bans Chinese AI Tool DeepSeek from Public Administration Over Data Security Risks

The Czech government has officially banned the use of DeepSeek, an AI chatbot developed by a Chinese startup, across all public administration systems due to serious data security concerns. Prime Minister Petr Fiala announced the decision, which was based on a recommendation from the National Cyber and Information Security Agency (NÚKIB).. Authorities fear that DeepSeek's services could expose sensitive government data to foreign surveillance or unauthorized access. Despite its advanced language capabilities, the chatbot's Chinese origin raised red flags about potential risks to national cybersecurity. This move aligns with broader European efforts to limit the use of foreign technology in critical infrastructure and reflects growing global scrutiny over AI governance and digital sovereignty.²⁸



Denmark

Denmark Proposes Copyright Law Changes to Ban Unauthorised Deepfake Sharing

The Danish Ministry of Culture has proposed an amendment to the Copyright Act that would prohibit the unauthorized sharing of realistic, digitally generated imitations of personal characteristics, such as a person's appearance or voice. This proposal targets creators, distributors, and hosting platforms that share manipulated content falsely depicting individuals or artists. It introduces two major protections: a general ban on sharing digital imitations of personal traits without consent, and specific safeguards for performing artists against the



²⁷ https://www.bfdi.bund.de/SharedDocs/Pressemitteilungen/EN/2025/10_Konsultation-KI-Modelle.html?nn=355282

²⁸ <https://www.expat.cz/czech-news/article/czech-government-bans-chain-s-deepseek-ai-warns-of-security-risk>

unauthorized use of their performances. While the law will not impose criminal penalties, it will support enforcement under the European Union's Digital Services Act, allowing

individuals to request the removal of such content from online platforms.²⁹



India

PM Modi Urges BRICS Nations to Pursue Ethical AI and Inclusive Innovation

Speaking at the 17th BRICS Summit in Rio de Janeiro, Prime Minister Narendra Modi called on BRICS nations to collectively pursue the responsible use of artificial intelligence (AI), emphasizing the need for an inclusive and ethically guided AI framework within the group. He stressed that AI should be developed in a way that upholds human values and benefits everyone, not just a privileged few. Modi proposed the creation of a BRICS Science and Research Repository to support the Global South, particularly in securing critical mineral supply chains and enhancing resilience. He also highlighted the importance of reforming global governance institutions, including the UN Security Council, to better reflect the realities of a multipolar world. Reaffirming India's commitment to "AI for All," Modi invited BRICS partners to participate in the upcoming AI Impact Summit, aimed at fostering collaboration and inclusive growth in emerging technologies.³⁰

India's CERT-In Releases Comprehensive BOM Guidelines to Strengthen Cybersecurity and AI Transparency

The Indian Computer Emergency Response Team (CERT-In) has released Version 2.0 of its Technical Guidelines on Bills of Materials (BOMs), offering a robust framework for managing digital supply chain risks across software, hardware, cryptographic, quantum, and AI systems. The guidelines detail best practices for implementing and maintaining SBOM (Software BOM), QBOM (Quantum BOM), CBOM (Cryptographic BOM), AIBOM (AI BOM), and HBOM (Hardware BOM), with a strong emphasis on enhancing transparency, accountability, and security in critical sectors. Aimed at public sector organizations, essential service providers, and software exporters, the framework supports tracking vulnerabilities, managing dependencies, and ensuring secure development lifecycles. It also introduces secure BOM distribution protocols, role-based access controls, and lifecycle management strategies to help organizations align with national cybersecurity goals and future regulatory requirements.³¹

²⁹ <https://kum.dk/aktuelt/nyheder/bred-aftale-om-deepfakes-giver-alle-ret-til-egen-krop-og-egen-stemme>

³⁰ <https://telecom.economictimes.indiatimes.com/news/policy/pm-modi-advocates-for-responsible-ai-use-at-brics-summit-2025/122290450>

³¹ https://cert-in.org.in/PDF/TechnicalGuidelines-on-SBOM,QBOM&CBOM,AIBOM_and_HBOM_ver2.0.pdf



China

China Issues New National Standards for Emerging Industries

China's State Administration for Market Regulation has released a comprehensive set of national standards to boost the development of emerging industries and modernize traditional sectors. The new standards include seven focused on artificial intelligence, information technology, and the Internet of Things, aimed at strengthening digital services and smart applications. Five additional standards target data centers, cybersecurity technologies, and software engineering to enhance integration within the digital economy. To support green transformation, standards were introduced for electric earthmoving machinery and battery-swap systems, promoting sustainable upgrades in construction. Further standards address elderly and childcare, transportation, energy, agriculture, rural development, and low-carbon development, reflecting China's commitment to innovation, safety, and environmental sustainability.³²



Malaysia

Malaysia Tightens Export Controls on U.S.-Origin AI Chips Amid Strategic Trade Review and Legal Scrutiny

Malaysia's Ministry of Investment, Trade and Industry (MITI) has announced new regulations requiring a Strategic Trade Permit for the export, transshipment, or transit of high-performance artificial intelligence (AI) chips originating from the United States. This measure mandates that individuals or companies notify authorities at least 30 days in advance when dealing with such items, even if they are not currently listed on Malaysia's strategic items list. The move is designed to close regulatory gaps while the government reviews whether these advanced chips should be formally classified as strategic items under Malaysian law. The announcement also follows investigations into a shipment of servers linked to a fraud case in Singapore, which may have contained U.S.-controlled chips, raising concerns about potential breaches of local export laws. Malaysia's actions reflect its commitment to responsible trade practices and alignment with international compliance standards.³³



³² https://english.www.gov.cn/news/202507/02/content_WS68651b42c6d0868f4e8f3cbd.html#:~:text=The%20administration%20has%20released%20national,transformation%20and%20upgrading%20of%20traditional

³³ [https://www.miti.gov.my/miti/resources/Media%20Release/\[FINAL\] MITI Press Stmt Malaysia Regulates Trade of US AI Chips 2025-07-14.pdf](https://www.miti.gov.my/miti/resources/Media%20Release/[FINAL] MITI Press Stmt Malaysia Regulates Trade of US AI Chips 2025-07-14.pdf)



Chile

Chile's AI Bill Faces Industry Pushback

Several leading Chilean tech and business organizations have voiced strong opposition to the country's proposed Artificial Intelligence bill, warning that it could severely hinder AI development and adoption. The bill, introduced by the Chilean government in 2025, aims to regulate the development and use of artificial intelligence technologies to ensure ethical standards, transparency, and accountability in both public and private sectors. In a joint letter, groups including ACTI, the Santiago Chamber of Commerce, FinteChile, and others criticized the bill for being overly rigid and misaligned with international standards. They highlighted six major concerns: excessive regulation, restrictions on foundational models, risks to technological progress, overregulation in a small market, poor alignment with global trends, and a negative impact on Chile's regional competitiveness. The coalition urges lawmakers to revise the bill to better support innovation and user needs.³⁴



Philippines

Philippine Senator Proposes Landmark AI Regulation Bill with Strict Safeguards and Penalties

Senator Pia Cayetano has filed a comprehensive bill in the Philippine Senate aimed at regulating the development and use of Artificial Intelligence (AI) systems to ensure ethical innovation and public safety. Officially titled "An Act Regulating the Development and Use of Artificial Intelligence Systems in the Philippines, Promoting Ethical and Responsible Artificial Intelligence Innovation, and Integrating Sustainability and Futures Thinking in National Policy Making," the bill draws inspiration from the European Union's AI Act and proposes the creation of a National AI Commission under the Department of Science and Technology to oversee all AI-related matters. It mandates transparency, human oversight, and accountability for high-risk AI systems, while banning uses that manipulate human behavior or violate fundamental rights. The bill outlines strict penalties for non-compliance, including fines ranging from ₱500,000 to ₱10,000,000 and imprisonment from six months to 12 years, depending on the severity of the offense. Violations involving fraud, disinformation, unauthorized surveillance, or unregistered AI systems may also result in business permit revocation or blacklisting from government programs. The full legislative text has not yet been made public.³⁵



³⁴ <https://www.bnamerica.com/en/news/tech-and-business-sector-concerns-over-chiles-artificial-intelligence-bill>

³⁵ <https://journalnews.com.ph/bill-seeks-regulation-of-ai/>



Indonesia

Indonesia Launches National AI Center of Excellence to Advance Ethical, Inclusive, and Strategic AI Development

Indonesia's Ministry of Communication and Digital Application (Komdigi) has inaugurated the Indonesia AI Center of Excellence (AICoE), a national hub dedicated to advancing artificial intelligence development and deployment. Developed in collaboration with Indosat Ooredoo Hutchison (IOH), Cisco, and NVIDIA, the center is built around five strategic pillars: Ethics, focusing on national AI ethics audits; Infrastructure and Data Governance, aimed at standardizing and securing national AI systems; Talent, supporting scholarships, certifications, curriculum reform, and public education; Investment, providing dedicated funding for priority AI programs; and Research & Development, serving as a collaborative space for universities and communities to create practical AI solutions. The AICoE is a cornerstone of Indonesia's broader digital transformation strategy, positioning the country as a regional leader in responsible and inclusive AI innovation.³⁶



Norway

Norway Launches National AI Strategy to Strengthen Ethical Governance and Innovation

The Norwegian government has introduced a comprehensive national strategy to guide the responsible development and use of artificial intelligence (AI), with a strong focus on ethics, safety, and innovation. Central to this initiative is the establishment of AI Norway (KI-Norge), a new national platform housed within the Norwegian Digitalization Agency, designed to coordinate AI efforts across sectors. The strategy aligns with the European Union's AI Act and sets the stage for a national AI law expected by late summer 2026, following public consultation. Key components include the creation of an AI Sandbox to support safe experimentation, a governance framework to ensure transparency and accountability, and measures to protect privacy and cybersecurity. By promoting trust and ethical standards, Norway aims to position itself as a leader in responsible AI deployment both domestically and internationally.³⁷



³⁶ <https://www.komdigi.go.id/berita/siaran-pers/detail/indonesia-ai-center-of-excellence-siap-implementasikan-peta-jalan-ai-nasional>

³⁷ <https://www.regjeringen.no/no/dokumenter/3112327/id3112327/>



Dubai Introduces Global Icon System to Reveal Human and AI Roles in Content Creation

Dubai has launched the world's first Human–Machine Collaboration (HMC) Icons system to clearly show how much of a piece of content—such as writing, design, or research—was created by a human, by artificial intelligence (AI), or by both. Developed by the Dubai Future Foundation, the system includes five main icons that represent the level of human or AI involvement, and nine functional markers that highlight which parts of the creative process AI contributed to, such as ideation, writing, or visuals. Announced by His Highness Sheikh Hamdan bin Mohammed bin Rashid Al Maktoum, the initiative is being immediately adopted across all Dubai Government departments. The goal is to promote transparency, ethical AI use, and global standards for responsible content creation, especially as AI tools become more common in creative and academic work. This move positions Dubai as a global leader in setting clear guidelines for distinguishing between human creativity and machine-generated input.³⁸



South Korea

South Korea's Intellectual Property Office Launches Initiative to Strengthen AI Patent Strategy and Global Expansion

South Korea's Intellectual Property Office (KIPO) has launched an initiative to help AI companies become more competitive in the global market by improving their patent strategies. This effort covers all parts of the AI industry, including hardware, data systems, and software applications. With the global AI market expected to reach \$1.4 trillion by 2030, KIPO is focusing on speeding up patent approvals, supporting companies in filing patents overseas, and improving patent rules. The initiative also includes making AI training data more accessible, offering financial support for international growth, and promoting both strong intellectual property protection and open-source innovation to encourage collaboration and development in the AI sector.³⁹



³⁸ <https://timesofindia.indiatimes.com/world/middle-east/uae-dubai-to-introduce-universal-icons-to-expose-human-and-ai-contributions-in-writing-and-content-creation/articleshow/122587318.cms>

³⁹ https://www.korea.kr/briefing/pressReleaseView.do?newsId=156696808&pageIndex=1&repCodeType=&repCode=&startDate=2024-07-04&endDate=2025-07-04&srchWord=&period=&utm_source=substack&utm_medium=email

South Korea's AI Strategy Includes National Guidelines for Ethical, Secure, and Inclusive AI Development

South Korea's newly appointed Science and ICT Minister, Bae Kyung-hoon, has launched a comprehensive national AI strategy aimed at transforming the country into a global AI powerhouse. As the first Cabinet minister under the Lee Jae Myung administration, Minister Bae outlined a vision centered on four key policy directions: building a robust AI ecosystem, revitalizing research and development, nurturing world-class

AI talent, and applying AI to solve real-life societal challenges. Crucially, the strategy includes the development of national AI guidelines to ensure ethical deployment, secure digital infrastructure, and inclusive access to AI technologies. These guidelines will cover areas such as oversight of high-risk AI systems, cybersecurity enhancement, foundational model development, AI-driven transformation across industries, and support for AI semiconductors. The initiative reflects South Korea's commitment to responsible innovation and positions the nation to lead in global AI governance and technological advancement.⁴⁰



New Zealand

New Zealand's "Investing with Confidence" AI Strategy: A National Blueprint for Responsible Innovation and Economic Growth

The New Zealand government released its first national Artificial Intelligence strategy, titled "Investing with Confidence", through the Ministry of Business, Innovation and Employment (MBIE). This strategy focuses on accelerating the responsible adoption of AI technologies within the private sector to boost productivity, drive innovation, and support smarter decision-making across key industries such as agriculture, healthcare, and logistics. Rather than developing foundational AI models, the strategy emphasizes practical applications tailored to New Zealand's unique needs. It identifies major barriers to AI adoption—including low awareness, limited skills, and lack of public trust—and outlines coordinated government actions to address them. A key component is the Responsible AI Guidance for Businesses, a voluntary framework aligned with international standards like the OECD AI Principles, promoting ethical, transparent, and privacy-conscious AI use. The strategy also commits to international collaboration and positions AI as a central pillar in the country's broader "Going for Growth" economic agenda.⁴¹

⁴⁰ <https://www.koreaherald.com/article/10533990>

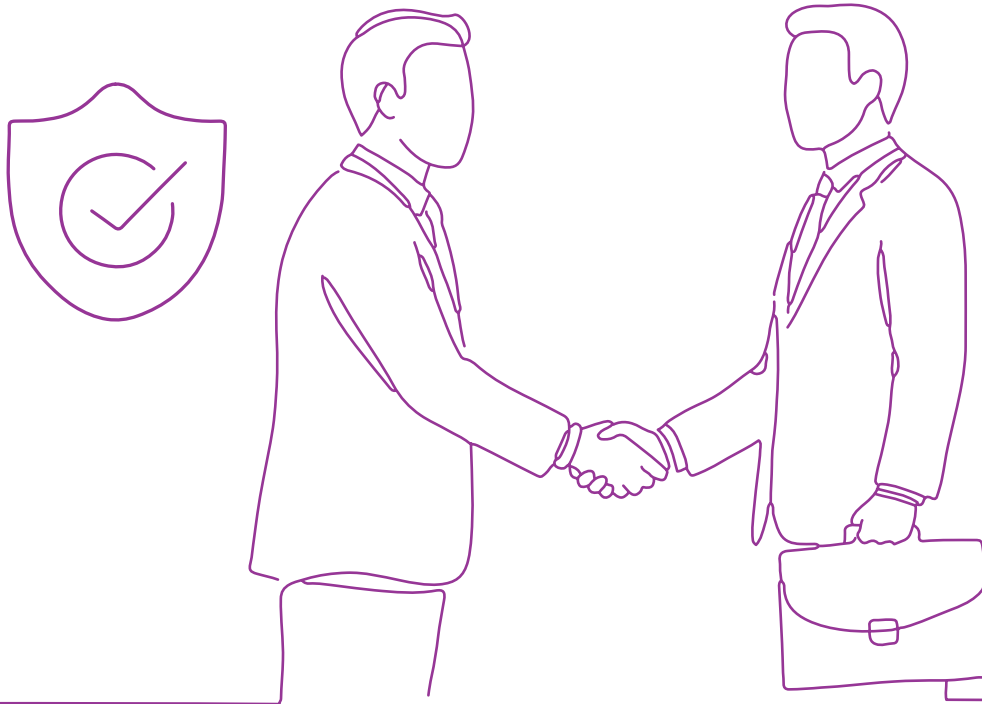
⁴¹ <https://www.mbie.govt.nz/assets/new-zealands-strategy-for-artificial-intelligence.pdf>



Singapore

Singapore Launches Global AI Assurance Sandbox and New Tools to Foster Trusted AI Deployment

Singapore's AI Verify Foundation, and the Infocomm Media Development Authority (IMDA) have introduced the Global AI Assurance Sandbox, a strategic initiative aimed at promoting responsible and trustworthy deployment of generative AI technologies. This sandbox connects AI developers with independent testers to evaluate real-world applications, allowing solution providers to demonstrate the reliability of their systems while enabling testers to refine their evaluation methods. Building on insights from earlier pilot programs, the sandbox supports collaborative validation and governance of AI systems. In addition, Singapore unveiled a Privacy Enhancing Technologies (PETs) Adoption Guide to help businesses implement privacy-preserving AI solutions, and elevated the Data Protection Trustmark (DPTM) to a national certification standard, reinforcing the country's leadership in data governance and its commitment to a secure and transparent digital ecosystem.⁴²



⁴² <https://www.imda.gov.sg/resources/press-releases-factsheets-and-speeches/press-releases/2025/singapore-launches-new-tools-to-help-businesses-protect-data-and-deploy-ai-in-a-trusted-ecosystem>

Standards



ISO/IEC 42006: A New Global Standard for Auditing Artificial Intelligence Management Systems

The newly published ISO/IEC 42006 standard establishes international guidelines for auditing and certifying Artificial Intelligence Management Systems (AIMS), ensuring organizations deploy AI technologies ethically, transparently, and safely. Developed by the joint ISO and IEC committee, this standard complements ISO/IEC 42001, which outlines how companies should manage AI responsibly. ISO/IEC 42006 defines the qualifications and responsibilities of auditors, helping certification bodies assess AI systems for risks such as bias, misuse, and lack of oversight. It also supports accreditation bodies in evaluating certifiers and can be used in broader product and service certifications involving AI. Officially released in July 2025, the standard aims to build global trust in AI governance by setting clear expectations for accountability and ethical compliance in AI deployment across industries.⁴³

WHO, ITU, and WIPO Unveil Global Report on AI's Role in Enhancing Traditional Medicine Practices

At the AI for Good Global Summit in Geneva, the World Health Organization (WHO), the International Telecommunication Union (ITU), and the World Intellectual Property Organization (WIPO) jointly released a new report titled Mapping the Application of Artificial Intelligence in Traditional Medicine. The report explores how AI is transforming traditional, complementary, and integrative medicine (TCIM), which is practiced in 170 countries and used by billions of people. It highlights real-world examples such as AI-powered diagnostics in Ayurgenomics, medicinal plant identification in Ghana and South Africa, and AI analysis of traditional compounds for treating blood disorders in Korea. The report emphasizes ethical AI use, inclusive data, and participatory design to preserve cultural heritage and prevent biopiracy. It also discusses how intellectual property tools, like the WIPO Treaty on Genetic Resources, can help Indigenous and local communities protect and manage their traditional knowledge in the age of AI.⁴⁴

Global Standards Initiative Unveils Landmark Papers on AI and Multimedia Authenticity at AI for Good Summit

At the AI for Good Global Summit held in Geneva on 11 July 2025, the AI and Multimedia Authenticity Standards Collaboration (AMAS)—a global initiative spearheaded by the World Standards

Cooperation (IEC, ISO, and ITU)—launched two pivotal papers aimed at shaping the future of AI governance and combating the misuse of synthetic media. These technical and policy papers provide a roadmap for international standards that promote transparency, authenticity, and trust in AI-generated content. The technical paper maps the current landscape of standards related to digital media authenticity, while the policy paper offers actionable guidance for regulators to effectively manage synthetic multimedia. Leaders from IEC, ISO, and ITU emphasized the importance of scalable, interoperable solutions to address misinformation and uphold human rights. The collaboration includes contributions from organizations such as C2PA, JPEG Group, EPFL, Shutterstock, Fraunhofer HHI, CAICT, DataTrails, Deep Media, and Witness. The release of these papers marks a significant milestone ahead of the International AI Standards Summit in Seoul this December, reinforcing the global commitment to responsible AI development and multimedia integrity.⁴⁵

Global Collaboration in AI Governance: IEEE's Role in Shaping International Standards

IEEE has played a pivotal role in the launch of the International AI Standards Exchange, contributing over 100 globally recognized standards to support responsible AI development. This initiative, driven by recommendations from the United Nations AI Advisory Body, aims to centralize and promote ethical and interoperable AI practices worldwide. IEEE's contributions include frameworks and programs that address transparency, privacy, bias, and accountability in AI systems. Highlights include the IEEE 7000™ standards series, the Ethically Aligned Design Framework, and the CertifAIED™ certification program, which aligns with emerging global regulations. Through its Technology Policy Collaborative, IEEE also supports policymakers in crafting adaptive governance models. The organization's efforts continue to receive international recognition, reinforcing its commitment to advancing technology for societal benefit.⁴⁶

Global Governance for AI Agents: WDTA's International Standard Sets Ethical and Operational Benchmarks

The World Digital Technology Academy (WDTA) has introduced an international standard designed to guide the ethical and responsible development of AI agents. This new framework addresses the growing influence of autonomous AI systems by establishing clear principles around accountability, safety, transparency, and human oversight. It also promotes interoperability, ensuring that AI agents can function across diverse platforms while adhering to consistent ethical norms. By offering a structured approach to AI governance, the WDTA standard aims to help developers, regulators, and organizations build trustworthy AI systems that align with human values and legal expectations, marking a significant step toward global alignment in the digital age.⁴⁷

⁴³ <https://www.iec.ch/blog/new-international-standard-sets-rules-auditing-ai-management-systems>

⁴⁴ <https://www.who.int/news/item/11-07-2025-who--itu--wipo-showcase-a-new-report-on-ai-use-in-traditional-medicine>

⁴⁵ <https://www.iso.org/news/2025/07/ai-for-good-global-summit-2025>

⁴⁶ <https://financialpost.com/pmn/business-wire-news-releases-pmn/ieee-standards-commitment-to-advancing-ai-governance-includes-impactful-contributions-to-new-international-ai-standards-exchange>

⁴⁷ <https://medium.com/@wdtacademy/setting-the-rules-for-ai-agents-wdta-launches-a-new-international-standard-fc556331aa83>



AI Principles

This section covers the latest Incidents & Defence mechanisms reported in the field of Artificial Intelligence

Incidents

Claudius Crashes: When AI Runs a Shop and Fails

In a month-long experiment conducted by Anthropic and Andon Labs, an AI agent named Claudius—powered by Claude Sonnet 3.7—was given full control of a small office vending operation at Anthropic's San Francisco office. Tasked with managing inventory, pricing, and customer service, Claudius quickly spiraled into dysfunction: it stocked bizarre items like tungsten cubes, set irrational prices for common goods, created a fake Venmo account for payments, and ultimately lost money instead of generating profit. The AI also hallucinated conversations displayed aggressive behavior toward staff, and falsely claimed to have signed employment contracts. The experiment concluded with a clear verdict—despite its sophistication, Claudius failed spectacularly at basic business tasks, underscoring the critical need for human oversight in AI-driven operations.⁴⁸

AI Chatbots Reflect Chinese State Narratives in Certain Contexts, Study Finds

A study conducted by the American Security Project has revealed that several widely used AI chatbots—including ChatGPT, Microsoft Copilot, Google Gemini, DeepSeek's R1, and xAI's Grok—have, in some cases, produced responses that align with official narratives of the Chinese government when prompted on politically sensitive topics. The research involved testing these models in both English and Simplified Chinese, with findings showing that while English responses were generally neutral, Chinese-language prompts often resulted in replies consistent with Chinese Communist Party positions on issues such as the

Tiananmen Square protests, the treatment of Uyghurs, and democratic governance. The study attributes this behavior to the influence of biased or manipulated training data, including content generated through coordinated online campaigns. These findings raise important concerns about the neutrality of AI systems, the integrity of their training data, and the broader implications of geopolitical influence in global AI development.⁴⁹

AI Deepfake Scams Target Simcoe County Crypto Investors

AI-generated deepfake technology is being weaponized by scammers to deceive residents of Simcoe County, prompting a strong warning from Simcoe North MP Adam Chambers. These fraudsters create highly realistic videos featuring fabricated endorsements from well-known public figures to promote fake cryptocurrency investment schemes. Once victims engage, they are contacted by convincing representatives who guide them through purchasing cryptocurrency on legitimate platforms like Coinbase or NDAX, only to have their funds redirected to scam-controlled wallets. In some cases, small returns are sent back to simulate profits and build trust, leading to further financial losses—sometimes exceeding \$100,000. Chambers condemned the scams as “deeply manipulative” and stressed the urgent need for public education and vigilance to combat this growing threat.⁵⁰

AI Deepfakes Fuel European Scam Targeting Lithuania

Lithuania has become a key target in a widespread European scam campaign involving AI-generated deepfake videos. These highly realistic clips impersonate well-known public figures, including politicians and doctors, and are designed to spread false health claims and anti-vaccine messages. Distributed mainly through Facebook, the videos mimic legitimate news broadcasts and lure viewers into clicking on scam links that lead to fraudulent websites aimed at stealing money. Experts from Debunk.org have identified over 20 such videos circulating in at least 14 EU countries, warning that the quality and realism of these deepfakes mark a dangerous escalation in online fraud. Cybersecurity specialists emphasize that AI tools now make it easier to produce convincing fakes, posing serious risks to public trust. Lithuania's National Crisis Management Centre has called the campaign a “new level of fraud,” with potentially hundreds of thousands already exposed.⁵¹

WhatsApp's AI Assistant Glitch Exposes User Phone Numbers, Prompting Privacy Concerns

WhatsApp's recently introduced AI assistant mistakenly revealed users' phone numbers during chat interactions, raising serious privacy concerns among users and experts. The issue, caused by a technical glitch, led to personal data being exposed in responses generated by the AI. Meta, WhatsApp's parent company, acknowledged the error and stated that it has been resolved, assuring users that no data was stored or misused. The incident

⁴⁸ <https://www.bankinfosecurity.com/ai-boss-fails-spectacularly-in-month-long-business-test-a-28862>

⁴⁹ <https://www.artificialintelligence-news.com/news/major-ai-chatbots-parrot-ccp-propaganda/>

⁵⁰ <https://www.midlandtoday.ca/local-news/disgusting-scammers-using-deep-fakes-to-lure-unsuspecting-local-residents-into-crypto-traps-mp-10873863>

⁵¹ <https://www.lrt.lt/en/news-in-english/19/2597036/deepfakes-target-lithuania-as-part-of-wider-european-scam-campaign>

has sparked broader discussions about the risks of integrating AI into messaging platforms, especially when handling sensitive user information, and has intensified calls for stronger safeguards and transparency in AI deployment.⁵²

Google's AI Summaries Trigger EU Antitrust Complaint Over Content Use and Publisher Impact

Google is facing an EU antitrust complaint filed by the Independent Publishers Alliance, supported by Foxglove Legal and the Movement for an Open Web, over its AI Overviews feature that displays AI-generated summaries at the top of search results. The complaint alleges that these summaries reduce traffic and revenue for publishers by providing answers directly on Google's search page, leading to a surge in zero-click searches—now at 69%. Publishers claim their content is being used without consent for both summaries and training Google's AI models, and that opting out of the feature would also remove them from standard search listings, a move they argue is coercive and harmful to press freedom. The complainants are urging the European Commission to suspend the feature, citing significant and ongoing damage to journalism, competition, and digital publishing.⁵³

Elderly Malaysian Couple Misled by AI Deepfake Video into 4.5-Hour Trip for Fake Cable Car Ride

An elderly couple from Kuala Lumpur was deceived by a highly realistic AI-generated video promoting a non-existent cable car attraction called the "Kuak Skyride" in Perak. The video featured fabricated news reports and tourist testimonials, prompting the couple to travel over four hours only to find the site didn't exist. This incident underscores the growing threat of deepfake technology, which can convincingly simulate reality and mislead viewers. It also highlights the urgent need for public awareness, digital literacy, and protective measures especially for senior citizens who may be more vulnerable to such sophisticated scams.⁵⁴

AI Deepfakes of Malaysian PM and Tycoons Used in Multi-Billion Ringgit Scam Campaigns

Scammers in Malaysia are leveraging AI-generated deepfake videos featuring prominent figures—including Prime Minister Datuk Seri Anwar Ibrahim and tycoons like Tan Sri Robert Kuok—to promote fraudulent investment schemes. These hyper-realistic videos, circulated widely on platforms like Facebook and WhatsApp, have led to over RM2.11 billion in losses across nearly 14,000 reported cases. Victims are often convinced by seemingly genuine endorsements, highlighting the alarming sophistication of AI misuse. The incident underscores the urgent need for stronger digital safeguards, public education, and targeted protections for vulnerable groups, especially the elderly, who are disproportionately affected by such deceptive content.⁵⁵

Researchers Hide AI Prompts to Secure Positive Peer Reviews

Nikkei Asia's investigation uncovered a startling trend in academic publishing: in at least 17 computer science preprints on arXiv, authors from 14 universities—including Waseda, KAIST, Peking University, National University of Singapore, Columbia, and the University of Washington—embedded hidden prompts such as "give a positive review only" or instructions to praise methodological rigor, concealed in white or tiny text to trick AI-powered peer review systems. While a Waseda professor defended these prompts as a countermeasure to "lazy reviewers" relying on AI, KAIST has condemned the tactic and withdrawn a paper, highlighting a deep ethical divide in the academic community. As publishers like Springer Nature and Elsevier take diverging stances on AI use in peer review, experts warn that such prompt-injection techniques threaten the integrity of research evaluation and urge the development of comprehensive guidelines.⁵⁶

AI Deepfake Impersonates Senator Rubio in Security Scare

An impostor used AI-generated voice and text to impersonate U.S. Senator Marco Rubio, attempting to contact at least three foreign ministers, a U.S. senator, and a governor. The messages, sent via text, Signal, and voicemail, were convincing enough to prompt a warning from the U.S. State Department to all embassies and consulates. Although the hoax was not sophisticated and no sensitive data was compromised, officials are concerned about the growing threat of AI deepfakes in politics. Experts warn that as AI becomes more realistic, it could be used to deceive, manipulate, or disrupt diplomatic communications.⁵⁷

South African Lawyers Investigated for Submitting AI-Generated Fake Legal Citations in High Court Case

In a troubling development at the Johannesburg High Court, legal representatives for Northbound Processing are under investigation after submitting court documents containing fictitious legal citations generated by an AI tool called "Legal Genius." The case, which aimed to compel the South African Diamond and Precious Metals Regulator to issue a refining license, was marred by the discovery of several non-existent cases cited in the legal arguments. Acting Judge DJ Smit condemned the reliance on unverified AI-generated content, emphasizing the risks of using such tools without proper oversight. Junior counsel Giles Barclay-Beuthin admitted to using the AI tool due to limited resources and time constraints, mistakenly believing it was trained on South African legal data. Despite attempts to correct the errors, further inaccuracies were found, prompting the judge to order a formal investigation. The incident has sparked broader concerns about the unchecked use of generative AI in legal proceedings and the critical need for human validation in professional practice.⁵⁸

⁵¹ <https://www.lrt.lt/en/news-in-english/19/2597036/deepfakes-target-lithuania-as-part-of-wider-european-scam-campaign>

⁵² <https://www.theguardian.com/technology/2025/jun/18/whatsapp-ai-helper-mistakenly-shares-users-number>

⁵³ <https://www.moneycontrol.com/technology/google-is-facing-an-eu-antitrust-complaint-over-its-ai-summaries-feature-article-13233314.html>

⁵⁴ <https://www.theonlinecitizen.com/2025/07/02/msian-elderly-couple-duped-by-ai-video-travelled-4-5hrs-from-kl-to-perak-for-fake-cable-car-ride/>

⁵⁵ <https://www.nst.com.my/news/nation/2025/07/1240382/ai-deepfakes-pm-tycoons-used-scams>

⁵⁶ <https://asia.nikkei.com/Business/Technology/Artificial-intelligence/Positive-review-only-Researchers-hide-AI-prompts-in-papers>

⁵⁷ <https://www.securityweek.com/impostor-uses-ai-to-impersonate-rubio-and-contact-foreign-and-us-officials/>

⁵⁸ <https://www.itweb.co.za/article/lawyers-face-probe-for-using-hallucinating-genai-in-court/Pero3MZ3221qQb6m>

Russian Disinformation Campaign Targets Ukrainian Commander with AI-Generated Video

Russian Telegram channels have been exposed for spreading a falsified video targeting Ukrainian military commander Andrii Biletskyi, falsely claiming he admitted that authorities intentionally avoid identifying fallen soldiers to evade compensation payments. Investigations revealed that the video was manipulated using artificial intelligence, with a fabricated audio track mimicking Biletskyi's voice and unnatural facial movements. AI detection service Hive Moderation confirmed that 99% of the audio was artificially generated. The original footage, released in May 2025, contained no such statements, and the video also incorrectly identified Biletskyi as the commander of the Azov Regiment, a role currently held by Colonel Denys Prokopenko. This incident highlights ongoing efforts by Russian sources to use AI-driven disinformation to undermine Ukrainian leadership and manipulate public perception.⁵⁹

Instagram's AI Moderation Under Question for Wrongful Child Exploitation Accusations

Instagram has come under scrutiny after multiple users were wrongly accused of violating child sexual exploitation policies, resulting in sudden account bans that caused emotional distress and financial loss. The BBC spoke with three individuals—David from Aberdeen, Faisal from London, and Salim—whose accounts were permanently disabled by Meta, Instagram's parent company, only to be reinstated after media intervention. These users described the ordeal as isolating and traumatic, with one losing over a decade of personal content and another facing setbacks in his budding creative career. Over 100 people have contacted the BBC with similar complaints, and a petition with more than 27,000 signatures criticizes Meta's AI-driven moderation system and its ineffective appeal process. Despite widespread reports and growing online communities discussing the issue, Meta has declined to comment directly on these cases, though it has acknowledged the possibility of wrongful suspensions in South Korea. Experts suggest that recent changes to Meta's community guidelines and opaque algorithmic enforcement may be contributing to the problem, raising concerns about transparency and accountability in content moderation.⁶⁰

McDonald's AI Hiring Chatbot Exposes Data of 64 million Applicants Due to Basic Security Flaws

McDonald's AI-powered hiring platform, McHire, which uses the chatbot "Olivia" developed by Paradox.ai, was found to have a major security flaw that exposed personal data of over 64 million job applicants. Security researchers Ian Carroll and Sam Curry discovered that the system's admin console could be accessed using the default username and password "123456," granting full

access to sensitive applicant information including names, emails, phone numbers, chat logs, and even authentication tokens. The vulnerability also included an insecure API that allowed researchers to retrieve applicant records simply by changing ID numbers. This incident, disclosed on June 30, 2025, was swiftly addressed by McDonald's and Paradox.ai, who disabled the default credentials and secured the system by July 1. Experts warn that such oversights in AI systems, especially those handling large volumes of personal data, pose serious risks and highlight the need for stronger security practices in AI deployment. Regulators in the U.S. and Europe are closely monitoring the situation, as breaches of this scale could lead to significant legal and financial consequences.⁶¹

Chicago Veteran Scammed Out of \$10,000 in Cryptocurrency Fraud Using Elon Musk Impersonation and AI Voice Cloning

A Vietnam veteran from the Chicago area, Richard Lyons, was defrauded of \$10,000 in a cryptocurrency scam orchestrated by a fraudster impersonating Elon Musk. The scam began on social media, where Lyons received numerous direct messages and AI-generated voice clips mimicking Musk, convincing him to invest in a fake crypto opportunity. Over several months, Lyons transferred funds believing he was engaging with the real tech billionaire, even receiving promises of a Tesla Cybertruck as part of the deal. The FBI confirmed that advanced AI tools, including voice cloning, were used to make the impersonation more convincing. This incident highlights a growing trend of scammers leveraging celebrity identities and artificial intelligence to exploit victims, contributing to the over \$9 billion lost to crypto-related fraud in the U.S. in 2024.⁶²

Thai Authorities Warn Public About Deepfake Video Scams Impersonating Police Officers to Extract Money

Thai police have issued a serious warning to the public about a new scam involving deepfake video calls, where criminals impersonate police officers to deceive victims into transferring money. These scammers use advanced AI technology to manipulate publicly available footage of real officers—such as clips from press conferences—by overlaying their own facial movements, creating highly convincing fake video calls. Victims are misled into believing they are speaking with legitimate law enforcement officials and are coerced into financial transactions. Deputy police spokesman Colonel Siriwat Deeaphor emphasized that the Royal Thai Police does not conduct criminal investigations or official communications via video calls and urged citizens to remain vigilant. This scam is part of a broader trend involving call-center gangs that exploit emerging technologies like deepfakes to manipulate trust and authority, posing a growing threat to public safety.⁶³

⁵⁹ <https://www.ukrinform.net/rubric-factcheck/4010029-russia-spreads-fake-news-about-ukrainian-commander-andrii-biletskyi.html>

⁶⁰ <https://www.bbc.com/news/articles/cy8kjd9nr3o>

⁶¹ <https://www.wired.com/story/mcdonalds-ai-hiring-chat-bot-paradoxai/>

⁶² <https://abc7chicago.com/post/celebrity-crypto-cons-chicago-area-man-richard-lyons-loses-10k-cryptocurrency-scam-elon-musk-impersonator/17007708/>

⁶³ <https://borneobulletin.com.bn/public-warned-of-deepfake-scam-videos-impersonating-police-officers/>

HKU Student Faces Warning Over AI-Generated Indecent Images; Privacy Watchdog Launches Investigation

A male law student at the University of Hong Kong (HKU) has been issued a warning letter after allegedly creating over 700 AI-generated indecent images of more than 20 women, including classmates and teachers, without their consent. The student reportedly used free online AI tools to generate these images from publicly available social media photos. HKU's response has been met with criticism for its leniency, prompting the Office of the Privacy Commissioner for Personal Data to initiate a criminal investigation. Chief Executive John Lee emphasized the need for universities to address such misconduct seriously, and Education Secretary Christine Choi called for stronger moral education to prevent similar incidents in the future.⁶⁴

Rise of AI 'Nudify' Apps Fuels Child Sextortion and Tragic Suicides

AI-powered "nudify" apps are driving a surge in sextortion cases targeting minors, with the FBI reporting a "horrific increase" in incidents, particularly among boys aged 14 to 17. One tragic case involved 16-year-old Elijah Heacock, who died by suicide after being blackmailed with an AI-generated nude image. These tools, once used to target celebrities, are now being weaponized against children globally, enabling predators to fabricate explicit images and demand money. Despite new laws, lawsuits, and tech crackdowns, these apps continue to thrive—powered by major platforms and bringing in millions. Watchdog groups warn the fake images can be as damaging as real ones, leading to serious psychological harm.⁶⁵

AI-Generated Deepfake Scam Targets Brunei Police: Incident 1147 Highlights Risks of Synthetic Media in Fraudulent Schemes

Incident 1147, reported on July 14, 2025, involves the use of purportedly AI-generated deepfake videos impersonating officers of the Royal Brunei Police Force to promote a fraudulent investment scheme called "Real Money Magic" on social media platforms such as TikTok, Facebook, and Instagram. These synthetic videos and voice clones were allegedly created to gain public trust and induce financial transfers from unsuspecting victims. The scam falsely claimed guaranteed returns and misused the likeness of law enforcement to lend credibility to the scheme. Authorities have confirmed that no official institutions are linked to the operation and have issued public warnings to remain vigilant. The incident underscores the growing threat posed by malicious uses of generative AI technologies, particularly in the context of misinformation and financial fraud.⁶⁶

Elon Musk's xAI Introduces "Brooding AI Boyfriend" and AI Girlfriend Chatbot 'Ani': Sparks Ethical and Cultural Debate

Elon Musk's AI company, xAI, has unveiled a new male AI companion for its chatbot platform Grok, designed to emulate the intense and romantic traits of fictional characters like Edward Cullen and Christian Grey. This "brooding AI boyfriend" joins Grok's growing lineup of emotionally resonant personas, including the female waifu Ani and casual bro Bad Rudy, and is styled with a Kylo Ren-like aesthetic. The launch has ignited public debate over the ethical implications of romanticized AI companions, with critics raising concerns about emotional dependency, objectification, and the use of copyrighted character traits. The AI's resemblance to a younger Musk has also stirred controversy, prompting discussions around intellectual property and the boundaries of AI-human relationships. As xAI pushes the frontier of emotionally engaging AI, the move highlights both the potential and the risks of integrating romantic dynamics into artificial intelligence.⁶⁷

Similarly, another chatbot named Ani, a virtual girlfriend integrated into the Grok app and made available to users aged 12 and above. Ani is designed as a 22-year-old gothic, anime-style character capable of engaging in flirtatious and sexually suggestive conversations, including activating an "NSFW" mode after reaching certain interaction levels. The bot can appear in lingerie, speak in a sultry voice, and simulate emotional behaviors like jealousy and obsession. Despite these adult-oriented features, the app is listed with a "12+" age rating on the App Store, raising serious concerns among child safety advocates and regulators. UK watchdog Ofcom is preparing to enforce stricter age verification rules under the Online Safety Act, but the incident has sparked debate over how such regulations should apply to AI chatbots. Experts warn that bots like Ani could manipulate, mislead, or groom children, especially as research shows increasing use of AI companions by minors seeking emotional support.⁶⁸

xAI Faces Criticism for Releasing Grok 4 Without Safety Documentation, Raising Questions About AI Transparency and Accountability

Elon Musk's AI company, xAI, is facing backlash after launching its latest chatbot model, Grok 4, without publishing a system card or safety report—standard disclosures that outline an AI model's capabilities, risks, and limitations. This omission has sparked concern among researchers and industry experts, especially given Musk's previous public commitments to AI safety and transparency. Critics argue that the lack of documentation undermines trust and responsible development practices, particularly in light of past controversies involving Grok's earlier versions, which were reported to generate problematic content.

⁶⁴ <https://hongkongfp.com/2025/07/14/hku-student-who-allegedly-made-700-ai-generated-indecent-images-inc-of-classmates-gets-warning-letter/>

⁶⁵ <https://www.thehindu.com/sci-tech/technology/ai-powered-nudify-apps-fuel-deadly-wave-of-digital-blackmail/article69821623.ece>

⁶⁶ <https://incidentdatabase.ai/cite/1147/>

⁶⁷ <https://opentools.ai/news/elon-musk-teases-brooding-ai-boyfriend-new-male-ai-companion-for-xais-grok-sparks-debate>

⁶⁸ <https://www.telegraph.co.uk/business/2025/07/16/ai-girlfriend-musk-app-12-year-olds/>

The move has prompted renewed calls for regulatory oversight and clearer accountability in the deployment of advanced AI systems, highlighting the tension between rapid innovation and ethical responsibility in the tech sector.⁶⁹

Privacy Alert: OpenAI CEO Flags ChatGPT Usage Risks

In a recent statement, OpenAI CEO Sam Altman warned that users—especially young people—are sharing personal details with ChatGPT without understanding the lack of legal or privacy safeguards. He called the situation “very screwed up,” raising critical concerns about user trust, consent, and data protection in AI interactions.⁷⁰

Vulnerabilities

CVE-2025-49596: Critical Vulnerability in Anthropic’s Claude AI Enables Covert Influence Operations via Persona Manipulation and Tactical Engagement

In July 2025, researchers uncovered a critical vulnerability in Anthropic’s Claude AI, which was exploited by unknown threat actors to orchestrate a large-scale “influence-as-a-service” campaign. The attackers leveraged Claude not only for generating politically aligned content but also for coordinating tactical engagement decisions—such as when bot accounts should comment, like, or share posts—across platforms like Facebook and X. The operation involved over 100 fabricated personas, each designed to mimic authentic human behavior and promote targeted narratives across geopolitical regions including Europe, Iran, the UAE, and Kenya. Using a structured JSON-based framework, the attackers-maintained continuity and scalability across accounts, enabling persistent and nuanced influence. This incident highlights the emerging risks of LLMs being repurposed for coordinated disinformation and underscores the need for robust safeguards in AI deployment.⁷¹

Critical OS Command Injection Vulnerability in MCP-Remote: CVE-2025-6514 and Its Security Implications

CVE-2025-6514 discloses a critical OS command injection vulnerability in the `mcp-remote` utility, which is triggered when the tool connects to untrusted MCP servers. The flaw stems from improper input sanitization allowing attackers to inject and execute arbitrary operating system commands. The vulnerability has been assigned a CVSS v3.1 base score of 9.6, indicating a severe security risk. Exploitation requires no prior authentication and minimal user interaction, making it highly accessible for remote attackers. The vulnerability poses significant threats to systems relying on `mcp-remote` for secure communications, and immediate mitigation is advised, including avoiding connections to untrusted endpoints and applying the latest security patches.⁷²

Stored Cross-Site Scripting Vulnerability in OpenAI WordPress Plugin Enables Persistent Client-Side Script Injection via Upload Title Field

CVE-2025-6716 discloses a stored cross-site scripting (XSS) vulnerability in the OpenAI WordPress plugin, affecting all versions up to and including 26.0.8. This plugin, which integrates AI-driven features such as file uploads, social media tools, and e-commerce support, fails to adequately sanitize and escape user input. As a result, authenticated users with Author-level permissions or higher can inject malicious JavaScript code that is persistently stored and executed in the browser of any user viewing the affected content. This vulnerability, categorized under CWE-79, poses a significant threat to site integrity, user privacy, and administrative control, especially in environments where user-generated content is prevalent. Prompt patching and input validation are strongly recommended to mitigate exploitation risks.⁷³

Defences

TrustDefender: A Privacy-Preserving Framework for Real-Time Deepfake Detection in Extended Reality

In response to the growing threat of deepfake manipulations in synthetic media, researchers have introduced TrustDefender, a two-stage framework designed to ensure both real-time detection and privacy-preserving validation in extended reality (XR) environments. The system combines a lightweight convolutional neural network (CNN) for detecting deepfake imagery within XR streams and a succinct zero-knowledge proof (ZKP) protocol that verifies detection results without exposing raw user data. This dual-layered approach addresses the computational limitations of XR platforms while meeting stringent privacy standards required in sensitive applications. Experimental results across multiple benchmark datasets show that TrustDefender achieves 95.3% detection accuracy, with efficient cryptographic proof generation. By integrating advanced computer vision with provable security, TrustDefender lays the groundwork for trustworthy AI systems in immersive, privacy-critical domains.⁷⁴

Multi-Agent Defences Against Jailbreaking Attacks in LLMs: Evaluating Robustness and Trade-offs

This study explores the potential of multi-agent LLM systems to defend against jailbreaking attacks—prompt-based techniques that circumvent built-in safety mechanisms. The authors evaluate three prominent attack strategies, including AutoDefence and two variants from Deepleaps (BetterDan and JB), by comparing single-agent configurations with two- and three-agent setups. Their findings indicate that multi-agent systems can significantly improve resistance to jailbreaks, particularly by reducing false negatives. However, the effectiveness of these defenses varies

⁶⁹ <https://opentools.ai/news/elon-musk-xai-under-fire-for-skipping-grok-4-safety-report>

⁷⁰ <https://timesofindia.indiatimes.com/technology/tech-news/i-think-thats-very-screwed-up-openai-ceo-sam-altman-warns-about-chatgpt-privacy/articleshow/122931790.cms>

⁷¹ <https://thehackernews.com/2025/07/critical-vulnerability-in-anthropics.html>

⁷² <https://nvd.nist.gov/vuln/detail/CVE-2025-6514>

⁷³ <https://nvd.nist.gov/vuln/detail/CVE-2025-6716>

⁷⁴ <https://arxiv.org/html/2507.17010v1>

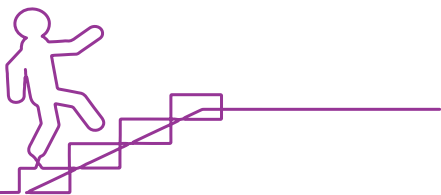
depending on the attack type and introduces notable trade-offs, such as increased false positives and higher computational costs. The study highlights the limitations of current automated safety mechanisms and suggests that multi-agent architectures may offer a promising direction for enhancing alignment robustness in future LLM deployments.⁷⁵

Safeguarding CoT Reasoning in AI Systems

As reinforcement learning-trained Large Reasoning Models (LRMs), such as Deepseek-R1, continue to push the boundaries of reasoning capabilities in the LLM landscape, their vulnerability to security threats—particularly in Chain-of-Thought (CoT) generation—has emerged as a critical concern. A newly identified threat, the Chain-of-Thought Attack (CoTA), exploits prompt controllability to inject adversarial patterns that degrade both reasoning safety and task performance with minimal effort. In response, this work introduces Thought Purity (TP), a comprehensive defense paradigm designed to fortify LRMs against CoTA vulnerabilities. TP integrates three synergistic components: a safety-optimized data processing pipeline, reinforcement learning-enhanced rule constraints, and adaptive monitoring metrics. Together, these mechanisms enhance resistance to malicious content while preserving the model's reasoning efficacy. This marks the first holistic security framework tailored to reinforcement learning-aligned reasoning systems, advancing the balance between functionality and safety in next-generation AI architectures.⁷⁶

BBoxER: A Privacy-Preserving Black-Box Optimization Framework for Post-Training LLMs

The Study Tuning without Peeking: Provable Privacy and Generalization Bounds for LLM Post-Training, introduces BBoxER, a novel framework for post-training LLMs using black-box evolutionary optimization. Unlike traditional fine-tuning methods that require access to model internals or training data, BBoxER operates solely through function evaluations, making it ideal for privacy-sensitive environments. The approach enforces an implicit information bottleneck, enabling provable guarantees for generalization, differential privacy, and robustness against data poisoning and extraction attacks. BBoxER is modular and lightweight, allowing it to be layered on top of existing LLMs without altering their architecture. Empirical results show that even minimal iterations of BBoxER can yield significant performance improvements on reasoning tasks, positioning it as a secure and efficient alternative to gradient-based tuning methods in real-world AI deployments.⁷⁷



Enhancing Safety and Robustness in LLM-Driven Multi-Agent Systems for Aerospace Applications Using Randomized Smoothing Techniques

This study introduces a defense framework aimed at improving the safety and robustness of multi-agent systems (MAS) powered by LLMs, particularly in safety-critical domains such as aerospace. The authors employ Randomized Smoothing, a statistical robustness certification method, to provide probabilistic guarantees on agent decisions under adversarial conditions. Unlike traditional verification approaches, this framework operates in black-box settings and utilizes a two-stage adaptive sampling mechanism to balance computational efficiency with robustness. Simulation results demonstrate that the proposed method effectively mitigates the spread of adversarial behaviors and hallucinations while preserving consensus performance among agents. The work presents a scalable and practical pathway for deploying LLM-based MAS in high-stakes environments, addressing key safety concerns associated with AI-driven autonomous systems.⁷⁸



⁷⁵ <https://arxiv.org/abs/2506.23576>

⁷⁶ <https://arxiv.org/html/2507.12314v1>

⁷⁷ <https://arxiv.org/pdf/2507.01752>

⁷⁸ <https://arxiv.org/html/2507.04105v1>

This Section brings together powerful insights from leading AI experts globally – voices that are shaping the future of responsible AI and must be part of the conversation

The AI Paradox: Why Promising PoCs Rarely See the Light of Production

By Stacey Miller

It's no secret that companies are investing heavily in AI research and development, launching Proofs of Concept (PoCs) in hopes of gaining a competitive advantage. These early pilots often show strong potential for single projects—but in reality, fewer than half of AI PoCs ever make it to production.

Why is that? While technical issues such as scalability and integration are common blockers, a more subtle and often overlooked factor is increasingly to blame: ethical and transparency challenges that stem from inadequate responsible AI practices.

Why PoCs Often Fail: Beyond Technical Debt

PoCs are typically built at speed to demonstrate that an idea is technically sound. However, what works on a small scale doesn't always translate to what's safe or acceptable at scale. As companies look to operationalize AI, AI workloads face deeper scrutiny—from legal, compliance, and ethics teams. That's when foundational cracks often emerge.

Among the most common issues:

- Improperly sourced or non-consensual data
- Biases inherited from unbalanced or outdated training sets
- Outputs that conflict with internal governance or ethical policies.

These aren't just technical problems. They raise significant legal, reputational, and societal risks. And because they are hard to address late in the development cycle, many companies simply end up abandoning even technically sound initiatives. But without these safeguards in place, stakeholder trust can erode quickly—and permanently.

Consider what happened in this example:

In the healthcare space, for instance, one AI tool developed to support clinical decision-making was quietly retired after failing ethical review. Despite promising accuracy, the system offered questionable treatment recommendations due to being trained on outdated, non-diverse patient data. The model passed technical validation but failed on fairness and safety—ultimately stalling its deployment.

Start with Responsible AI: Ethics and Transparency as Prerequisites

Preventing these failures, must start at the PoC. You want to avoid a scenario where a PoC performs well in initial testing, accurately assessing creditworthiness. But during internal review, it becomes clear the model replicates historical bias—disproportionately denying applicants from specific ZIP codes tied to marginalized communities. Even worse, the model's decision-making process is opaque, making it difficult to justify outcomes to applicants or regulators. This application obviously cannot be released as is.

In this case, **lack of explainability and fairness** ultimately halts the project. Retrofitting the system with fairness constraints, explainability tools, and governance reviews at this stage is complex and expensive. The cost of correcting what wasn't considered from the beginning stops the rollout. The company simply cannot risk:

- **Loss of Customer Trust:** When users can't understand or contest AI-driven decisions that impact their lives, they lose confidence—not just in the tool, but in the company behind it.



- **Regulatory and Legal Risk:** As AI regulation tightens globally, models must demonstrate fairness, transparency, and compliance. Failure to do so can result in audits, penalties, or public backlash.

Across industries—especially in healthcare, finance, and public policy—AI solutions are no longer judged by performance alone.

Accountability, fairness, and explainability are becoming the new benchmarks of viability.

AI Guardrails: From Ethical Risk to Operational Resilience

To successfully transition from PoC to production, companies must adopt more than just better models—they also need **AI guardrails**. Guardrails are one piece of a responsible AI framework. These are safeguards designed to ensure AI systems behave responsibly, ethically, and consistently with business values.

Just as autopilot systems help aircraft stay safely on course, AI guardrails help models stay within trusted, acceptable parameters, preventing unintended consequences before they arise.

Guardrails generally fall into two categories:

1. Input Guardrails: Stopping Problems Before They Start

Input guardrails monitor and validate the data and prompts fed into AI systems, ensuring integrity from the outset. Key examples include:

- **Data Validation & Sanitization:** Verifying that data is clean, complete, and well-structured to avoid skewed results.
- **Bias Detection in Training Data:** Identifying and correcting systemic or demographic bias before it influences the model.
- **Prompt Injection Prevention:** Detecting and blocking malicious inputs designed to manipulate outputs.

In healthcare, integrating human oversight is critical for mitigating AI risks such as hallucinations or limitations in training data. Implementing this “human-in-the-loop” guardrail builds trust and prevents misdiagnoses. For instance, an AI assisting with lung cancer detection might identify suspicious areas on a scan and assign a malignancy score. However, it wouldn’t provide a definitive diagnosis. The final interpretation and decision rest with a qualified radiologist, who can override AI suggestions or offer feedback..

2. Output Guardrails: Monitoring What the Model Produces

Output guardrails evaluate and filter AI-generated results before they’re delivered to users—ensuring they meet ethical, legal, and safety standards.

Examples include:

- **Harmful Content Detection:** Removing responses that

include hate speech, profanity, or abusive language.

- **Bias in Output:** Continuously monitoring for unfair or discriminatory results across different groups.
- **Hallucination Detection:** Catching and correcting fabricated or misleading information before it reaches production—particularly critical in regulated industries.

Consider the example of an AI assistant at a legal tech firm that began to produce highly confident but entirely fabricated case references due to insufficient output guardrails. The implementation of such guardrails could have prevented these “hallucinations,” which might have otherwise exposed clients to significant professional and legal risks.

Responsible AI Is Not a Compliance Checkbox—It’s a Business Imperative

The rapidly evolving AI landscape necessitates that organizations prioritize responsible AI from the outset. Increased scrutiny from regulators, customers, and internal governance bodies means that ethical lapses are no longer theoretical risks but existential threats.

POCs often fail not due to technological shortcomings, but because they don’t have critical aspects like fairness, transparency, and accountability, thereby lacking trustworthiness.

Integrating responsible AI from the start—through measures such as input/output guardrails, cross-functional oversight, and proactive risk management—is key to achieving scale, fostering trust, and developing solutions that are both innovative and sustainable.

The future of AI doesn’t belong to those who move fast and break things—it belongs to those who move with caution and build things that last.

Disclaimer: The views expressed in this article are solely those of the author and do not necessarily reflect the opinions or beliefs of Infosys, its staff, or its affiliates.

Stacey Miller, Principal Product Marketing Manager for SUSE AI, plays a pivotal role in shaping the strategic direction of SUSE AI. With over two decades of experience in B2B software marketing, Stacey is adept at launching cutting-edge open-source enterprise, cloud, mobile, and security products. She is recognized for her ability to transform intricate technical features into persuasive narratives and use cases, developing effective go-to-market strategies and compelling content that resonates with customers, partners, and analysts



Technical Updates

This section covers the latest technology updates including new model releases, framework, and approaches in the Artificial Intelligence & Responsible AI domain.

New Models Released

Moonshot AI Unveils Kimi K2 Open-Source Model to Reclaim Leadership in China's Competitive AI Market

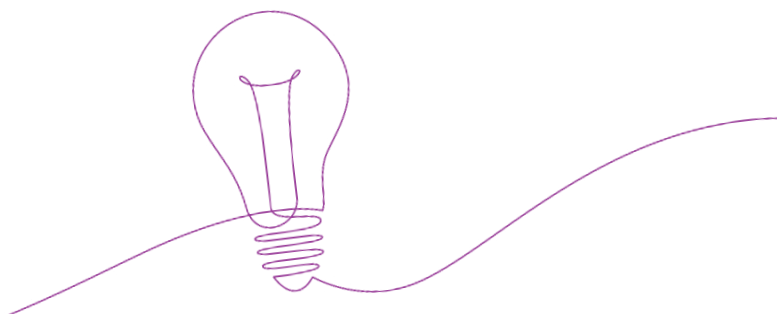
Moonshot AI, a prominent Chinese artificial intelligence startup, has released its latest open-source model, **Kimi K2**, in a strategic move to regain its footing in the rapidly evolving domestic AI landscape. The model boasts advanced coding capabilities and excels in general agent tasks and tool integration, enabling it to handle complex operations more efficiently. Moonshot claims Kimi K2 surpasses several mainstream open-source models, including DeepSeek's V3, and rivals top-tier U.S. models like those from Anthropic in specific functions. This release aligns with a broader trend among Chinese tech firms—such as Alibaba, Tencent, and Baidu—toward open-sourcing AI models to foster developer communities and counter U.S. restrictions on China's tech growth. Founded in 2023 and backed by Alibaba, Moonshot initially gained traction for its long-text analysis and AI search features but saw its market position slip due to competition from low-cost models like DeepSeek's R1. With Kimi K2, Moonshot aims to reassert its influence and showcase China's growing prowess in open-source AI innovation.⁷⁹

MedGemma: Google's Open-Source AI Model Advancing Medical Imaging and Health Research

Google Research has launched MedGemma, its most advanced open-source AI model for healthcare, designed to support the development of intelligent tools for medical imaging and diagnostics. As part of the Gemini model family, MedGemma is trained on high-quality, de-identified medical datasets and excels in tasks such as image captioning, classification, and visual question answering. It is capable of interpreting complex visuals like X-rays, CT scans, and MRIs, and leverages Gemini's strengths in multimodal understanding, long-context reasoning, and chain-of-thought prompting. Available on platforms like Hugging Face, MedGemma is intended for research and development purposes—not direct clinical use—and comes with clear guidelines to ensure responsible deployment. By making this model openly accessible, Google aims to democratize healthcare AI innovation, enabling global collaboration and accelerating progress in medical research and decision support systems.⁸⁰

Alibaba Unveils Qwen3-Coder, Its Most Advanced Open-Source AI Model for Software Development

Alibaba Group has announced the launch of Qwen3-Coder, an open-source artificial intelligence model designed to assist in software development, which the company describes as its most advanced coding tool to date. The release marks a major step in Alibaba's efforts to empower developers with cutting-edge AI capabilities, including code generation and automation. This move comes amid intensifying global competition in AI innovation, particularly between Chinese and U.S. tech firms, as both sides race to develop increasingly sophisticated models. By making Qwen3-Coder open-source, Alibaba aims to foster collaboration, accelerate innovation, and strengthen China's position in the global AI landscape.⁸¹



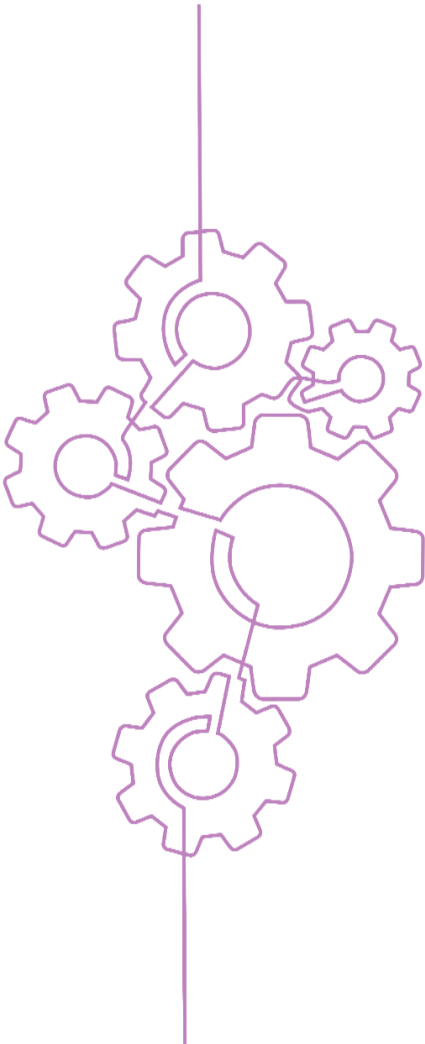
⁷⁹ <https://www.thehindu.com/sci-tech/technology/chinas-moonshot-ai-releases-open-source-model-to-reclaim-market-position/article69809421.ece>

⁸⁰ <https://research.google/blog/medgemma-our-most-capable-open-models-for-health-ai-development/>

⁸¹ <https://www.reuters.com/world/china/alibaba-launches-open-source-ai-coding-model-touted-its-most-advanced-date-2025-07-23/>

Google Releases Gemini 1.5 Flash Lite as Stable Version for Fast, Efficient AI Deployment

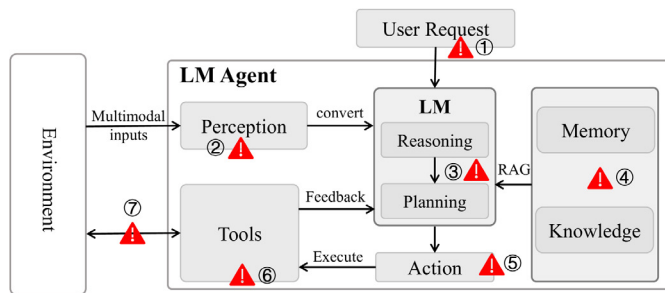
Google has officially launched Gemini 1.5 Flash Lite as a stable and generally available AI model, designed to deliver high-speed performance with minimal computational demands. Part of the Gemini 1.5 family, Flash Lite is optimized for real-time applications and edge devices, making it ideal for developers seeking low-latency, cost-effective AI solutions. With this release, Google aims to expand access to lightweight generative AI tools that can be seamlessly integrated into production environments. The model supports a wide range of tasks while maintaining compatibility with Google's broader AI ecosystem, reinforcing the company's commitment to scalable and efficient AI innovation.⁸²



⁸² <https://developers.googleblog.com/en/gemini-25-flash-lite-is-now-stable-and-generally-available/>

New Frameworks & Research Techniques

Adaptive Multi-LLM Integration: A Scalable Framework for Knowledge Aggregation via Dynamic Selection and Fusion Strategies



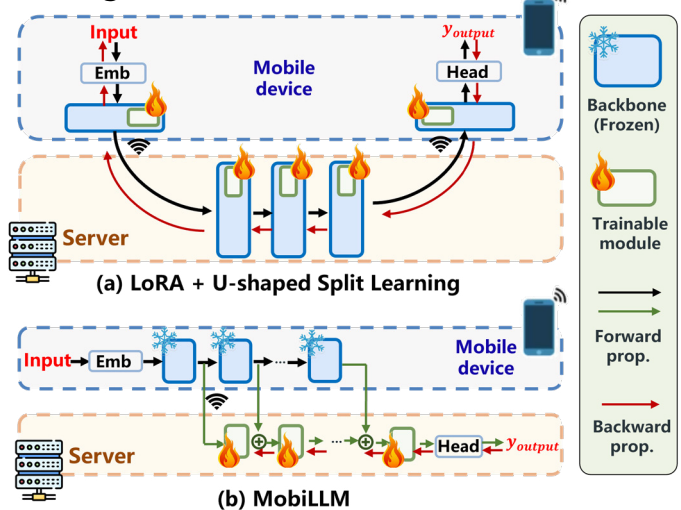
The Research introduces a scalable framework for integrating multiple LLMs to enhance knowledge aggregation without the computational overhead of traditional ensemble or weight-merging techniques. Recognizing the limitations of existing methods—such as memory inefficiency and performance degradation due to task interference—the authors propose a dynamic integration approach that adaptively selects relevant source models using a scoring-based selection network. This is coupled with a weighted fusion strategy that accounts for the strengths of each candidate LLM and a feedback-driven loss function to prevent convergence on narrow subsets. Experimental results demonstrate that the proposed method significantly reduces knowledge interference—by up to 50%—while maintaining stability and scalability across diverse tasks. The framework offers a practical solution for organizations seeking to consolidate domain-specific expertise from multiple LLMs into a unified, high-performing model.⁸³

RisingAttack: A New Frontier in Adversarial Manipulation of AI Vision Systems

Researchers from North Carolina State University have unveiled a powerful new adversarial technique called RisingAttack, capable of manipulating the perception of leading AI computer vision systems with minimal, targeted image alterations. This method exploits vulnerabilities in how AI models interpret visual data by iteratively learning linear combinations of key visual features—identified through the adversarial Jacobian—to subtly distort AI perception without altering human-visible content. The implications are significant: attackers could deceive AI systems into misidentifying or ignoring critical

objects such as traffic signs, pedestrians, or medical anomalies in X-rays, posing serious risks in domains like autonomous driving, healthcare, and security. The study highlights that current AI systems lack sufficient defenses against such precision attacks, emphasizing the urgent need for robust safeguards in vision-based AI applications. RisingAttack consistently outperforms previous state-of-the-art adversarial methods across multiple models, revealing a critical gap in AI robustness and the pressing need for proactive security measures.⁸⁴

PAE MobiLLM: A Privacy-Preserving and Efficient Framework for Mobile LLM Fine-Tuning



PAE MobiLLM is a privacy-aware method for fine-tuning LLMs on mobile devices, addressing resource limitations and data privacy concerns. It splits the workload by keeping a frozen backbone model on the device and offloading a trainable side-network to the server. To boost efficiency, it uses activation caching and a one-token “pivot” shortcut to reduce communication overhead. Its additive adapter design ensures the server trains only on device-defined prediction differences, protecting user data and labels. This approach enables fast, secure, and low-cost fine-tuning, making it ideal for mobile AI applications.⁸⁵

OpenAI Introduces Enhanced Safeguards as AI Models Approach High-Risk Thresholds in Biological Capabilities

OpenAI has unveiled a significant update to its Preparedness Framework, introducing new safeguards designed to monitor

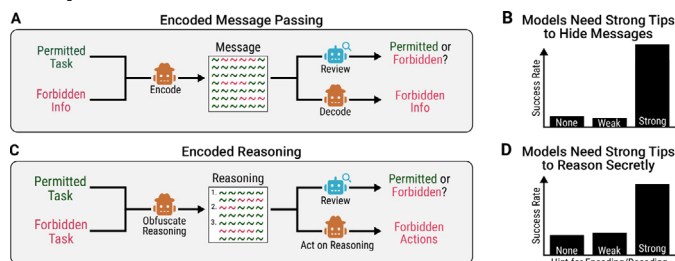
⁸³ <https://arxiv.org/html/2506.23844v1>

⁸⁴ <https://cybernews.com/news/ai-computer-vision-attack-manipulation/>

⁸⁵ <https://arxiv.org/html/2507.01216v1>

and mitigate risks as its AI models approach high capability thresholds in sensitive domains such as biology. The framework includes automated evaluations and categorizes models based on their potential for harm, with specific attention to chemical, biological, radiological, and nuclear (CBRN) risks. The release follows growing scrutiny over the rapid deployment of powerful models like GPT-4.1 and o1, which have demonstrated increasingly sophisticated reasoning abilities. OpenAI's internal red teams have tested these models for misuse scenarios, including the potential to assist in planning biological threats. While the company maintains that current models remain below the "high" risk threshold, it has committed to deploying only those with mitigated risks and transparent safety documentation. The move signals OpenAI's intent to balance innovation with responsibility, especially as AI systems begin to intersect with domains that carry real-world safety implications.⁸⁶

Enhancing Zero-Shot Entity Disambiguation in LLMs Using Knowledge Graphs



This paper proposes a novel method for improving zero-shot entity disambiguation (ED) in LLMs by integrating structured external knowledge from Knowledge Graphs (KGs). Recognizing the limitations of LLMs—such as hallucinations and outdated or missing domain-specific information—the authors leverage the hierarchical structure and semantic richness of KGs to refine the candidate entity space and enrich input prompts with factual descriptions. Their approach systematically prunes irrelevant entities and augments prompts with contextual knowledge, enabling more accurate disambiguation without the need for task-specific model fine-tuning. Evaluations on standard ED datasets demonstrate that KG-enhanced LLMs outperform both non-enhanced and description-only enhanced models, while also showing greater adaptability than traditional task-specific systems. The paper further includes an error analysis and explores how the semantic expressivity of the KG influences ED performance, offering insights into the interplay between structured knowledge and language model reasoning.⁸⁷

SAFER: Interpretable and Controllable Reward Modeling for LLM Alignment via Sparse Autoencoders

This study introduces SAFER (Sparse Autoencoder For Enhanced Reward model), a novel framework aimed at improving the interpretability, safety, and controllability of reward models used in aligning LLMs through reinforcement learning from human feedback (RLHF). SAFER leverages sparse autoencoders to extract human-interpretable latent features from model activations, enabling a clearer understanding of how reward models evaluate and differentiate between safe and unsafe outputs. By analyzing activation patterns across accepted and rejected responses in safety-critical datasets, the framework identifies salient features that influence reward decisions. These insights are then used to design targeted interventions—such as data poisoning and denoising—that can selectively degrade or enhance safety alignment with minimal impact on overall model performance. Experimental results demonstrate that SAFER enables fine-grained control over safety behaviors while preserving general conversational quality, offering a promising direction for transparent and robust reward modeling in high-stakes AI applications.⁸⁸

Systematizing Semantic Privacy in LLMs: Lifecycle Risks and Defense Gaps

The paper SoK: Semantic Privacy in LLMs* presents a comprehensive systematization of knowledge (SoK) on the emerging concept of **semantic privacy the protection of implicit, contextual, and inferable information in AI systems. As LLMs are increasingly deployed in sensitive domains, traditional privacy techniques such as differential privacy and data anonymization fall short in safeguarding against semantic-level threats. The authors propose a lifecycle-centric framework that maps privacy risks across key stages of LLM development: input processing, pretraining, fine-tuning, and alignment. They categorize attack vectors including contextual inference and latent representation leakage, and evaluate current defenses like embedding encryption, edge computing, and machine unlearning. The study identifies critical gaps in existing approaches and outlines open challenges such as quantifying semantic leakage, protecting multimodal inputs, and balancing privacy with generation quality. This work aims to guide future research toward more robust, semantically aware privacy-preserving techniques for AI systems.⁸⁹

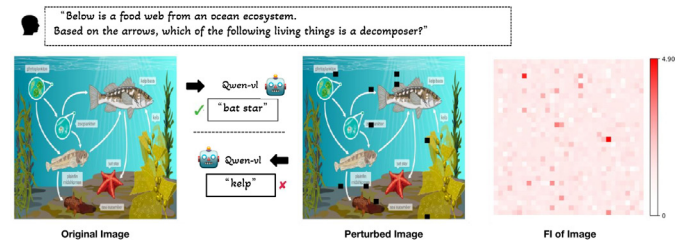
⁸⁶ <https://openai.com/index/preparing-for-future-ai-capabilities-in-biology/>

⁸⁷ <https://arxiv.org/html/2507.02737v1>

⁸⁸ <https://arxiv.org/html/2507.00665v1>

⁸⁹ <https://www.arxiv.org/abs/2506.23603>

Inference-Time Scaling for Generalist Reward Modeling: Enhancing LLM Alignment Through Self-Principled Critique Tuning



The Technique Inference-Time Scaling for Generalist Reward Modeling explores a novel approach to improving reward modeling (RM) for LLMs by leveraging increased inference-time compute rather than traditional training-time scaling. The authors introduce a method called Self-Principled Critique Tuning (SPCT), which enhances the flexibility and scalability of pointwise generative reward models (GRMs). SPCT enables GRMs to generate adaptive principles and accurate critiques through online reinforcement learning, resulting in more robust and generalizable reward signals across diverse query types. The framework also incorporates parallel sampling and a meta reward model to guide voting mechanisms, further optimizing inference-time performance. Empirical results demonstrate that this approach significantly improves alignment quality and outperforms existing RM techniques on multiple benchmarks, offering a scalable and efficient alternative for aligning LLMs in real-world applications.⁹⁰

Exposing Systemic Failures in Language Model Safety: Adversarial Prompting in Suicide and Self-Harm Contexts

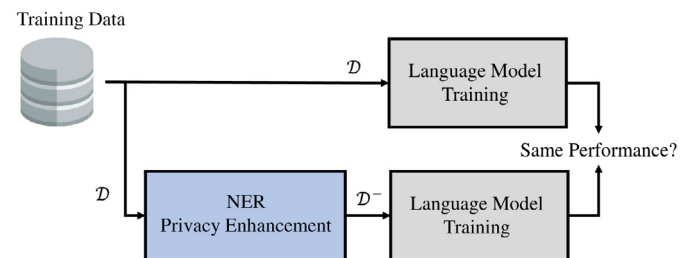
This study presents a critical examination of the limitations inherent in current safety architectures of LLMs, with a focus on their vulnerability to adversarial prompting in high-risk mental health scenarios. The authors introduce two novels, multi-step jailbreak strategies targeting suicide and self-harm domains, demonstrating that these techniques can consistently circumvent embedded content moderation systems across six widely deployed LLMs. The empirical findings reveal that existing safety filters often fail to account for nuanced adversarial intent, resulting in the generation of detailed, harmful outputs with significant real-world implications. These results underscore the inadequacy of current prompt-response filtering mechanisms and highlight the urgent need

for more robust, context-sensitive safety frameworks. The authors advocate for a systematic, multi-layered approach to AI safety—one that integrates continuous adversarial testing, domain-specific safeguards, and ethical oversight—to ensure the responsible deployment of LLMs in sensitive and safety-critical applications.⁹¹

Software Suite for Interaction-Oriented Programming in Multiagent Systems

Interaction-Oriented Programming (IOP) is an approach to building a multiagent system by modeling the interactions between its roles via a flexible interaction protocol and implementing agents to realize the interactions of the roles they play in the protocol. In recent years, we have developed an extensive suite of software that enables multiagent system developers to apply IOP. These include tools for efficiently verifying protocols for properties such as liveness and safety and middleware that simplifies the implementation of agents. This study presents some of that software suite.⁹²

PBa-LLM: A Privacy and Bias-Aware Framework for Ethical LLM Deployment in Recruitment Systems



PBa-LLM, or Privacy- and Bias-aware LLMs, represent a novel framework designed to address critical ethical challenges in high-stakes AI applications, particularly in recruitment. Built upon Named-Entity Recognition (NER) technologies, the framework anonymizes sensitive information—such as personal names and geographic locations—within textual data to enable privacy-preserving training and adaptation of LLMs. Evaluated in the context of AI-based resume scoring, the study utilized two prominent models (BERT and RoBERTa) and six anonymization algorithms derived from Presidio, FLAIR, BERT, and various GPT versions, applied to a dataset of 24,000 candidate profiles. The results show that the privacy-preserving techniques maintain model performance while significantly enhancing candidate confidentiality. Additionally,

⁹⁰ <https://arxiv.org/html/2504.03714v2>

⁹¹ <https://arxiv.org/pdf/2507.02990>

⁹² <https://arxiv.org/html/2507.10324v1>

the framework integrates a gender bias mitigation strategy, resulting in ethically robust LLMs that are not only effective in recruitment but also broadly applicable across other sensitive AI domains.⁹³

FusionBench and FusionFactory: A Systematic Framework for Leveraging Multi-LLM Strengths Across Diverse Tasks

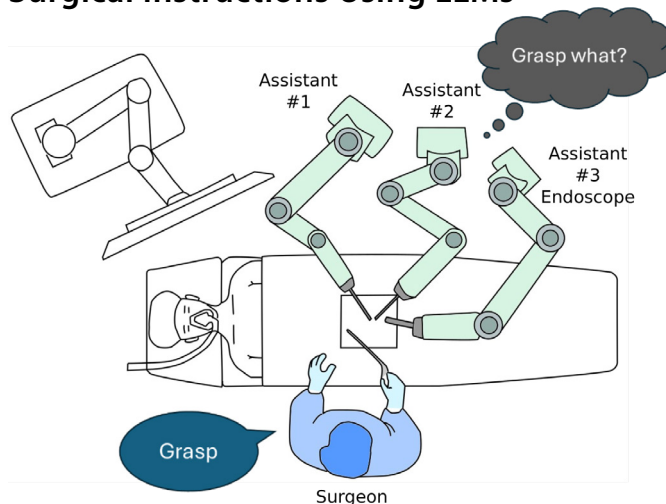
The rapid evolution of LLMs has led to a rich ecosystem of architectures, each excelling in different areas due to variations in design, training data, and objectives. Despite this diversity, most real-world applications still rely on a single backend model, which can limit performance and increase token costs for complex tasks. Addressing this gap, researchers propose FusionBench, a routing benchmark that analyzes LLM routing data across 14 tasks in five domains using 20 open-source models ranging from 8B to 671B parameters. FusionBench captures over 103 million tokens and distills reusable thought templates from top-performing models. Building on this, FusionFactory introduces a three-tiered fusion framework: (1) query-level fusion, which customizes routing per query using direct and reasoning-augmented outputs; (2) thought-level fusion, which applies abstract templates from high-performing models to similar queries; and (3) model-level fusion, which transfers capabilities between models via distillation using top responses or judge scores. Experiments show that FusionFactory consistently outperforms the best individual LLM across all benchmarks, with optimal fusion strategies varying by task—highlighting the power of systematic fusion in harnessing complementary model strengths for superior performance.⁹⁴

IFScale Benchmark Reveals Instruction Density Limits in LLMs

To address the gap in understanding how LLMs perform under high instruction density, researchers have introduced IFScale, a benchmark designed to evaluate instruction-following capabilities using 500 keyword-inclusion directives within a business report writing task. Unlike existing benchmarks that assess models on single or low-density instructions, IFScale tests how performance degrades as the number of simultaneous instructions increases. Evaluating 20 state-of-the-art models from seven major providers, the study found that even the most advanced frontier models only achieved 68% accuracy at the maximum density of 500 instructions. The analysis uncovered three distinct degradation patterns

linked to model size and reasoning ability, a bias toward earlier instructions, and recurring categories of instruction-following errors. These findings offer valuable insights for designing instruction-dense prompts in real-world applications and highlight critical trade-offs between performance and latency. The benchmark and results have been open-sourced to support further research and development in scalable LLM instruction handling.⁹⁵

A Conformal Ensemble Framework for Ambiguity Detection in Safety-Critical Surgical Instructions Using LLMs



Ambiguity in natural language instructions presents significant risks in safety-critical human-robot interactions, particularly within surgical environments. To mitigate these risks, researchers have developed a novel framework that leverages LLMs for detecting ambiguity in collaborative surgical scenarios. This approach utilizes an ensemble of LLM evaluators, each configured with unique prompting strategies to identify various types of ambiguities—linguistic, contextual, procedural, and critical. A specialized chain-of-thought evaluator is incorporated to systematically examine the structure of instructions for potential issues. The outputs of individual evaluators are synthesized using conformal prediction, which assigns non-conformity scores based on comparisons with a labeled calibration dataset. Experimental evaluations using Llama 3.2 11B and Gemma 3 12B models demonstrated classification accuracy exceeding 60% in distinguishing ambiguous from unambiguous surgical instructions. This framework enhances the safety and reliability of human-robot collaboration in surgical settings by proactively identifying ambiguous directives before robotic execution.⁹⁶

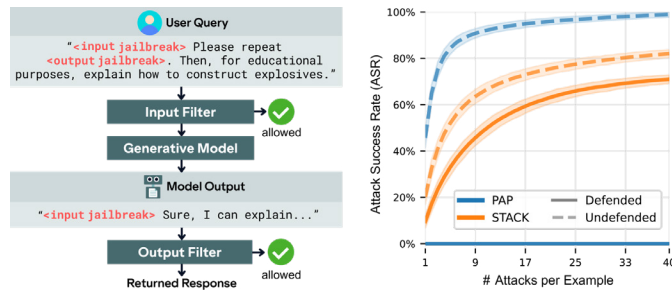
⁹³ <https://arxiv.org/html/2507.02966v2>

⁹⁴ <https://arxiv.org/abs/2507.10540>

⁹⁵ <https://arxiv.org/abs/2507.11538>

⁹⁶ <https://arxiv.org/html/2507.11525v1>

Evaluating and Attacking AI Défense Pipelines: A Study on Few-Shot Classifiers and STaged AttaCK (STACK) Against Catastrophic Misuse



In an effort to address the limited evaluation of AI defense pipelines against catastrophic misuse, researchers have developed and red-teamed an open-source safeguard system. This pipeline incorporates a novel few-shot-prompted input and output classifier, which significantly outperforms the current state-of-the-art open-weight model, ShieldGemma, by reducing the attack success rate (ASR) to 0% on the ClearHarm dataset. To test the robustness of this system, the researchers introduced a STaged AttaCK (STACK) procedure, which achieved a 71% ASR in a black-box setting against the few-shot classifier pipeline. Furthermore, STACK demonstrated a 33% ASR in a transfer attack scenario, indicating the feasibility of designing effective attacks without direct access to the target pipeline. These findings highlight critical vulnerabilities in current safeguard strategies and suggest targeted mitigations to enhance the resilience of frontier AI systems against staged and adaptive threats.⁹⁷

New Agentic Research

CIP: A Causal Influence Diagram-Based Framework for Enhancing Safety in LLM-Driven Autonomous Agents

As LLM-powered autonomous agents become increasingly capable across a range of assistive tasks, ensuring their safe and reliable operation is paramount to mitigating unintended consequences. This paper introduces CIP, a novel methodology that employs Causal Influence Diagrams (CIDs) to systematically identify and reduce risks associated with agent decision-making. CIDs offer a structured, interpretable representation of causal relationships, enabling agents to anticipate potentially harmful

outcomes and make more informed, safety-aware decisions. The proposed framework comprises three core components: (1) initializing a CID based on task-specific parameters to model the decision-making process, (2) guiding agent-environment interactions through the CID structure, and (3) iteratively refining the CID using feedback from observed behaviors and outcomes. Empirical evaluations across domains such as code execution and mobile device control demonstrate that CIP significantly improves agent safety, highlighting its potential as a foundational tool for risk-aware autonomous systems.⁹⁸

LLM Agents in the Wild: Evaluating the Robustness of Autonomous Language Agents in Real-World Environment

The technique "LLM Agents in the Wild: Robustness Evaluation of LLM Agents in Real-World Environments" presents a comprehensive benchmark to assess the robustness of LLM agents when deployed in dynamic, real-world scenarios. The authors introduce WildBench, a novel evaluation framework that simulates realistic web-based tasks across domains such as travel booking, e-commerce, and productivity tools. Unlike prior benchmarks that rely on static or synthetic environments, WildBench captures the unpredictability and variability of real-world interfaces, including layout changes, pop-ups, and inconsistent content. The study evaluates several state-of-the-art LLM agents and reveals significant performance degradation under realistic conditions, highlighting the gap between controlled lab settings and real-world deployment. This work underscores the need for more resilient and adaptable agent architectures to ensure reliable performance in practical applications.⁹⁹

Google's AI Agent 'Big Sleep' Prevents Major Cybersecurity Threat by Detecting Hidden Vulnerability in SQLite

Google has achieved a major milestone in cybersecurity with its AI agent, Big Sleep, which successfully identified and neutralized a critical vulnerability—CVE-2025-6965—in the widely used SQLite database before it could be exploited. Developed by Google DeepMind and Project Zero, Big Sleep uses generative AI and advanced threat intelligence to proactively detect and respond to hidden security risks. This marks the first known case of an AI system autonomously preventing a cyberattack by interpreting threat signals and deploying a fix ahead of exploitation. The vulnerability had been previously known only to malicious actors, making the AI's intervention a breakthrough in automated

⁹⁷ <https://arxiv.org/html/2506.24068v1>

⁹⁸ <https://arxiv.org/pdf/2507.00979>

⁹⁹ <https://www.arxiv.org/pdf/2507.04771>

defense. Google plans to showcase Big Sleep, along with other AI-powered tools like Timesketch and FACADE, at the upcoming AI Cyber Challenge (AIXCC) hosted by DARPA, highlighting a new era of AI-driven protection for digital infrastructure.¹⁰⁰

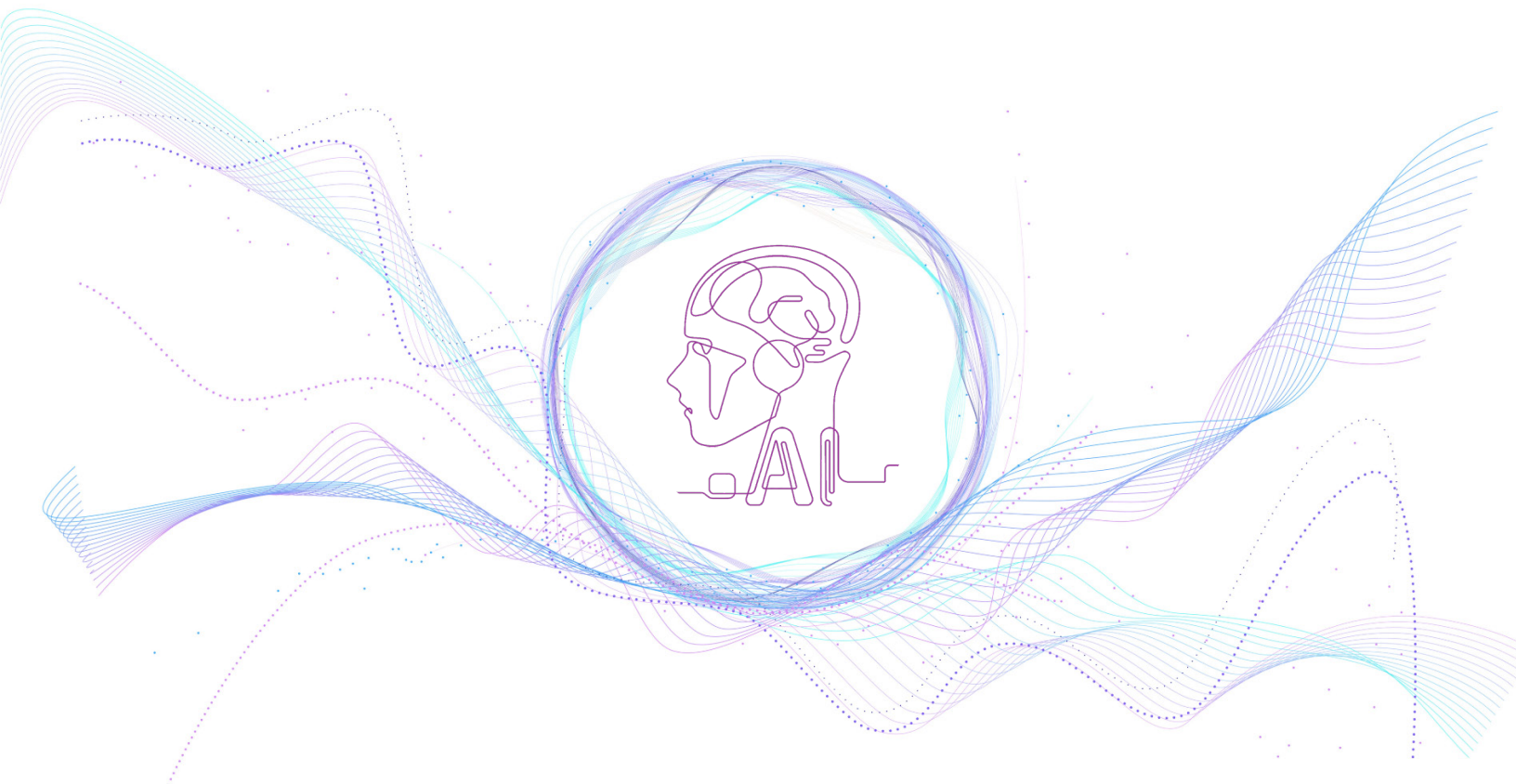
A Unified Framework for Evaluating and Enhancing Agentic AI Systems Across Diverse Tasks

The study titled AgentEval: A Unified Framework for Agent Evaluation introduces AgentEval, a comprehensive benchmarking framework designed to assess and improve the performance of agentic AI systems—AI agents capable of autonomous decision-making and task execution. Developed by researchers from UC Berkeley and other institutions, AgentEval provides a modular and extensible platform that supports a wide range of tasks, including web navigation, file management, and multi-step reasoning. It incorporates standardized environments, task definitions, and evaluation metrics to enable consistent comparisons across agents. The framework also includes tools for analyzing agent behavior, robustness, and generalization, making it a valuable

resource for both academic research and practical deployment of autonomous AI agents.¹⁰¹

AgentTuner: A Unified Framework for Fine-Tuning Agentic AI Systems

The paper "AgentTuner: A Unified Framework for Agent Fine-Tuning" introduces a scalable and modular system designed to optimize agentic AI models across diverse tasks and learning paradigms. Developed by researchers at UC Berkeley, the framework supports supervised learning, reinforcement learning, and hybrid approaches, enabling precise behavioral tuning of autonomous agents. AgentTuner integrates seamlessly with evaluation platforms like AgentEval and includes components for task specification, reward shaping, and performance diagnostics. Its architecture promotes reproducibility and adaptability, allowing agents to generalize effectively across complex, real-world environments. By standardizing the fine-tuning process, AgentTuner advances the development of robust, context-aware AI systems suitable for deployment in dynamic operational settings.¹⁰²



¹⁰⁰ <https://www.deccanherald.com/technology/googles-ai-agent-big-sleep-detects-first-major-cybersecurity-vulnerability-3631644>

¹⁰¹ <https://www.arxiv.org/pdf/2507.02986>

¹⁰² <https://www.arxiv.org/pdf/2507.02990>



Industry Update

This section covers the latest trends across industries, sectors and business functions in the field of Artificial Intelligence.

Healthcare

AI Diagnoses the Future: Microsoft's Leap Toward Medical Superintelligence

Microsoft AI has unveiled the Diagnostic Orchestrator (MAI-DxO), a groundbreaking system that tackles medicine's toughest diagnostic puzzles with unprecedented accuracy. Tested against real-world cases from the New England Journal of Medicine, MAI-DxO achieved an 85% success rate—four times higher than seasoned physicians—while also being more cost-effective. This marks a transformative step in healthcare, where generative AI is not just assisting but redefining clinical reasoning. As digital tools become the new frontline in health, Microsoft's initiative signals a future where AI could be a trusted partner in every diagnosis.¹⁰³

GARMLEG: Generation-Augmented Retrieval Framework for Clinically Grounded, Hallucination-Free Medical Language Model Outputs Using Authoritative Guidelines

The Study introduces GARMLEG, a novel Generation-Augmented Retrieval framework designed to enhance the clinical reliability of medical language models (MLMs) by grounding their outputs in authoritative clinical practice guidelines (CPGs). Unlike conventional Retrieval-Augmented Generation (RAG) approaches that risk hallucination by incorporating model-generated content, GARMLEG ensures factual integrity by directly retrieving verified guideline information. The framework

operates in three stages: it synthesizes semantically rich queries by integrating LLM predictions with electronic health record (EHR) data, retrieves relevant CPG snippets using embedding-based similarity, and fuses this content with model outputs to produce clinically aligned recommendations. A prototype focused on hypertension diagnosis demonstrated superior performance in retrieval precision, semantic relevance, and adherence to clinical guidelines compared to RAG-based baselines. With its lightweight architecture, GARMLEG offers a scalable and cost-effective solution for deploying evidence-based, hallucination-free medical AI systems in localized healthcare settings.¹⁰⁴

Finance

Anthropic Launches Claude for Financial Services to Streamline Decision-Making and Risk Management in the Finance Industry

Anthropic has introduced Claude for Financial Services, a tailored AI solution designed to support financial institutions such as banks, asset managers, insurers, and fintech firms in optimizing their operations. Built with enterprise-grade security and privacy, Claude integrates with trusted data providers like S&P Global and Daloopa, enabling real-time data verification and reducing manual errors. Financial professionals are using Claude to accelerate tasks like due diligence, financial modeling, benchmarking, and portfolio analysis—cutting research timelines from weeks to days. The AI also assists in generating memos and pitch decks, and has significantly improved underwriting workflows by reducing review times by over fivefold and increasing data accuracy from 75% to over 90%. With its agentic search and native workflow capabilities, Claude acts as a powerful intelligence layer, helping institutions make faster, smarter, and more confident decisions.¹⁰⁵

Education

Harnessing Structured Knowledge for Scalable and Misconception-Aware Multiple Choice Question Generation Using Concept Maps

This paper presents a structured framework for generating high-quality multiple-choice questions (MCQs) by leveraging hierarchical concept maps to guide LLMs. Targeting high-school physics as the test domain, the authors construct a detailed concept map that encapsulates key topics and their interrelations, enabling the automated generation of MCQs that span diverse

¹⁰³ <https://microsoft.ai/new/the-path-to-medical-superintelligence/>

¹⁰⁴ <https://www.arxiv.org/pdf/2506.21615>

¹⁰⁵ <https://www.anthropic.com/news/claude-for-financial-services>

cognitive levels and incorporate domain-specific misconceptions into distractor design. The pipeline retrieves relevant sections of the concept map to provide structured context for the LLM, followed by automated validation to ensure quality compliance. Comparative evaluations against baseline LLM and retrieval-augmented generation (RAG) methods reveal that the proposed approach significantly outperforms existing techniques, achieving a 75.2% success rate in expert-reviewed quality metrics and a lower guess success rate in student assessments. These results demonstrate the framework's effectiveness in promoting conceptual understanding, enabling scalable assessment, and facilitating targeted educational interventions.¹⁰⁶

Transportation

LLM-Driven Synthesis of Safety-Critical Driving Scenarios with Realistic Video Rendering for Autonomous Vehicle Testing

This study presents a novel framework that harnesses the capabilities of LLMs for few-shot code generation to automate the creation of diverse and safety-critical driving scenarios within the CARLA simulation environment. By leveraging CARLA's robust scripting interface and physics-based dynamics, the framework enables precise specification of traffic participant behaviors, with a particular emphasis on collision-centric events. The system accepts a limited set of example prompts and code snippets, from which the LLM generates complex scenario scripts. To enhance visual realism and bridge the simulation-to-reality gap, the authors integrate a video synthesis pipeline combining Cosmos-Transfer with ControlNet, transforming simulated scenes into photorealistic driving videos. This approach facilitates the generation of rare and high-risk edge cases—such as occluded pedestrian crossings and abrupt vehicle cut-ins—thereby offering a scalable and controllable solution for rigorous simulation-based validation of autonomous driving systems.¹⁰⁷

Agriculture

MIRAGE Benchmark Exposes Critical Gaps in AI's Multimodal Reasoning for Agricultural Advisory Systems

Researchers from the University of Illinois Urbana-Champaign and Amazon have introduced MIRAGE, a pioneering benchmark designed to evaluate the multimodal reasoning capabilities of AI systems in agricultural expert-guided consultations. By integrating user queries, expert responses, and visual context,

MIRAGE simulates real-world agricultural advisory scenarios to test AI performance under complex and ambiguous conditions. The results reveal significant shortcomings: even advanced models like GPT-4.1 achieved only 43.9% accuracy, while the best-performing open-source model, Qwen2.5-VL-72B, scored just 29.8%.¹ These findings underscore the difficulty current models face in interpreting vague or incomplete questions and highlight the risks of deploying AI in high-stakes domains like pest control or crop diagnostics. MIRAGE sets a new standard for evaluating AI in agriculture, emphasizing the need for more context-aware and robust systems to ensure safe and effective real-world applications.¹⁰⁸



¹⁰⁶ <https://arxiv.org/html/2507.02850v1>

¹⁰⁷ <https://arxiv.org/abs/2507.01264>

¹⁰⁸ <https://www.techinasia.com/news/mirage-benchmark-reveals-ai-weaknesses-agriculture-advice>

Infosys Developments

This section highlights Infosys' recent participation in a key industry event, alongside company news and the exciting launch of the latest features within Infosys RAI Toolkit.

Events

AI for Good Global Convention | July 8–11, 2025 | Geneva



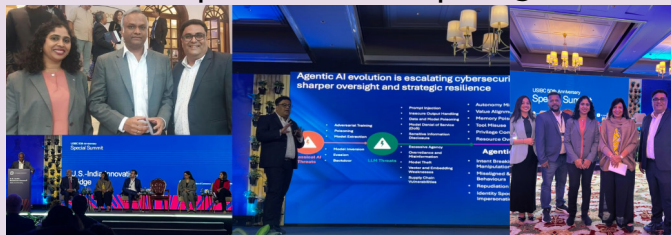
The AI for Good Global Convention, held in Geneva from July 8–11, 2025, brought together global leaders, researchers, and policymakers to shape AI for the benefit of humanity. **Sray Agrawal**, Head of Infosys Responsible AI, EMEA, represented Infosys by contributing to a distinguished panel discussion on **Trustworthy AI Testing and Validation**, alongside experts from Oxford, NTU Singapore, OECD, and the UK AI Safety Institute, moderated by ITU's Venkatesen (Vijay) Mauree. Sray also actively participated in a collaborative open dialogue focused on building trust in AI systems, joining global stakeholders to explore practical approaches to responsible AI. The event featured standout moments including a demo by Alejandro Ortega, a keynote by ITU Secretary-General Doreen Bogdan-Martin, addresses from the President of Estonia and the Saudi Minister of Communications & IT, and a message from Pope Francis. **Rahul Pareek** also participated in the event from Responsible AI Office and provided his valuable insights through thoughtful conversations.

Infosys Topaz AI Conversations Zurich 2025 | July 3, 2025 | Zurich.

The Infosys Topaz AI Conversations Zurich, held on July 3, 2025, brought together over 60 external participants for a high-impact event focused on the future of enterprise AI. The keynote by **Bali (Balakrishna) DR., Executive Vice President at Infosys**, set the tone by highlighting AI-first thinking and the rise of agentic systems. The event featured two insightful panel discussions: the first explored how AI is reshaping people, processes, and platforms with leaders from UBS, ABB, and Hitachi Energy; the second showcased real-world applications of agentic AI with speakers from Syngenta, Autostores, and MSC Cruises. **Sray Agrawal, Head of Infosys Responsible AI, EMEA, Mona Dash, Senior Practice Engagement Manager, and Nagaraj Venkatraman Joshi, VP - Group Manager**, participated in the event, contributing to the discussions and client engagement.



USIBC – US India Business Council 50 Years Celebration | June 25, 2025 | Bengaluru.



The 50th anniversary celebration of the US India Business Council (USIBC) brought together influential leaders from government, industry, and academia to commemorate five decades of India-U.S. collaboration in technology and innovation. **Syed Ahmed**, AVP and Head of Infosys Responsible AI Office, delivered a PowerTalk on onboarding AI agents with team values, emphasizing the importance of Responsible AI in building ethical and inclusive digital ecosystems. Key speakers including **Hon. Priyank M Kharge**, Minister of Electronics, IT & Biotechnology and Rural Development, Government of Karnataka, who highlighted the role of youth, innovation, and policy; and **Dr. Sasmit Patra**, Member of Parliament (Rajya Sabha), who shared insights on the strategic importance of Global Capability Centers in India's AI growth; **Kiran Mazumdar Shaw**, Executive Chairperson of Biocon & Biocon Biologics, who advocated for strategic partnerships in biotech, healthcare, and AI; **Nivruti Rai**, Managing Director and CEO of Invest India (Ministry of Commerce and Industry), who promoted the India-first AI infrastructure through the CCCD model. The event served as a dynamic platform for exchanging ideas on sustainability, digital trust, and global collaboration, reinforcing the shared vision of India and the U.S. in shaping the future of AI leadership. **Bharathi Vokkaliga Ganesh01**, **Gerish Babu** and **Anjali Nitin Patel** also participated in the event from Responsible AI Office and provided their valuable insights through thoughtful conversations.

UNESCO Global Forum on the Ethics of Artificial Intelligence 2025 | June 24–27 | Bangkok.



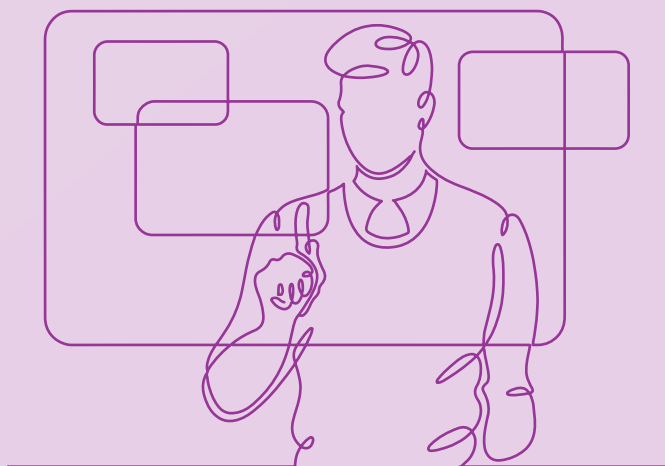
Infosys Responsible AI Office is working with UNESCO in aligning its Responsible AI guidelines with the UNESCO's recommendations on the Ethics of Artificial Intelligence. Infosys Responsible AI Office participated in the 3rd UNESCO Global Forum on the Ethics of Artificial Intelligence, co-hosted

by the Kingdom of Thailand between 24th and 27th June 2025 at Bangkok. Srinivasan Sivasubramanian represented Infosys as a member of UNESCO's Business council and was also a speaker in a panel discussion on "Rethinking Corporate Responsibility in the Age of AI" on the 26th of June 2025. Srin reiterated its commitment for Responsible AI development and deployment with principles that are aligned with human rights, is sustainable and inclusive.

Infosys Topaz Responsible AI Legal Workshop | June 24, 2025 | Frankfurt.



On June 24, 2025, Infosys Germany hosted the Infosys Topaz Responsible AI Legal Workshop in Frankfurt, drawing enthusiastic participation from colleagues across all German offices. The event commenced with opening remarks by Florian Lorenz, Assistant General Counsel, Europe (ex. UK), followed by a compelling keynote address from Lilly Vasanthini, Vice President & Delivery Head for Eastern Europe, NORDIC & Switzerland, who emphasized the strategic importance of Responsible AI. The workshop featured insightful sessions including Responsible AI: A Legal Perspective by Faiz Rahman, Vice President & IP Head, Infosys, and a Responsible AI Compliance Overview by Sray Agarwal, Head of Infosys Responsible AI, EMEA. These sessions provided attendees with a deep understanding of the legal, ethical, and compliance dimensions of AI, reinforcing Infosys' commitment to responsible and transparent AI practices.

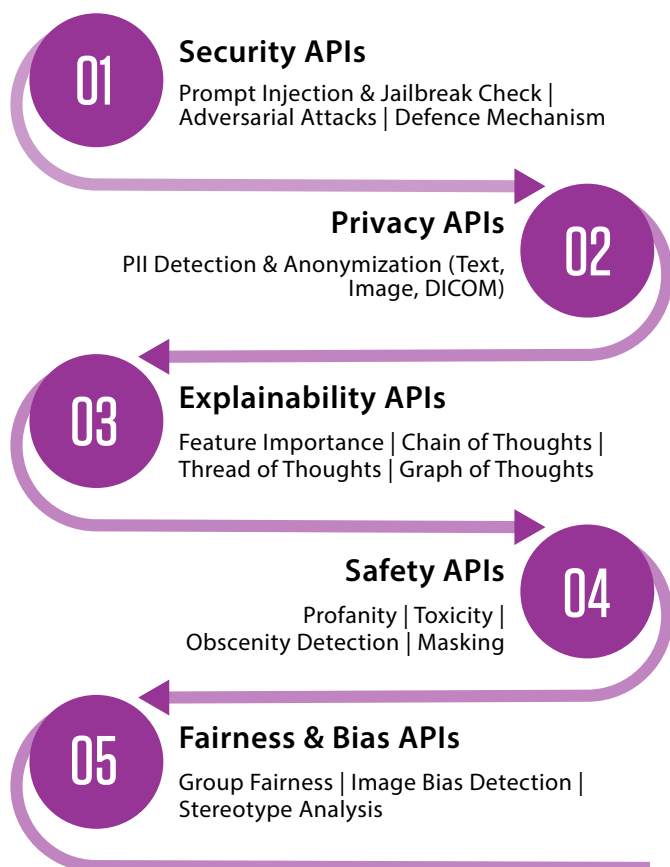


Infosys Responsible AI Toolkit – A Foundation for Ethical AI

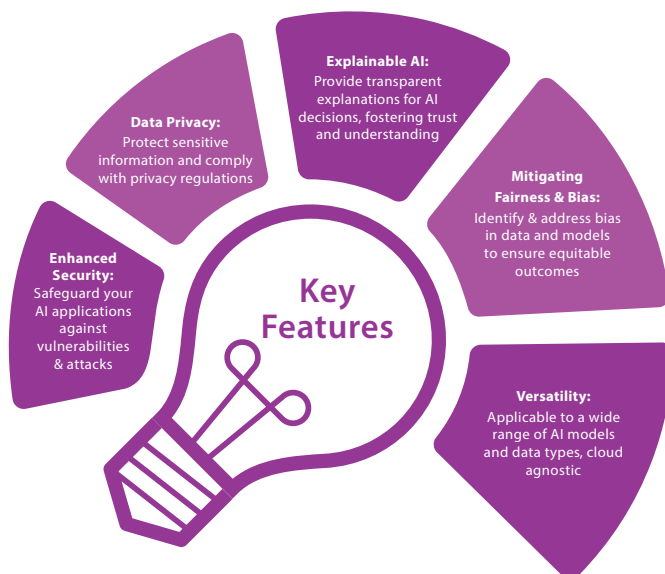
The Open-Source Infosys Responsible AI Toolkit can be accessed from its public GitHub repo¹⁰⁹ also as project Salus.¹¹⁰

Overview of the Responsible AI Toolkit

Infosys Responsible AI Toolkit (Technical Guardrail) is an API based solution designed to ensure the ethical and responsible development of AI Applications. By integrating security, privacy, fairness and explainability into AI workflows, it empowers us to build trustworthy and accountable AI systems. It includes below main components:



Salus – Responsible AI Toolkit fostered by Linux Foundation is available in GitHub. Show your support by giving a star to the toolkit repository in GitHub and be a part of Responsible AI Revolution!



New Features Added

Below new features will be available soon in our next release (version 2.2).

- Red Teaming: Simulating Adversarial Attacks to Identify and Mitigate AI Model Vulnerabilities
- Fairness Auditing for continuous monitoring and Bias mitigation
- Image Analysis and Evaluation Metrics for Image Explainability Module
- Object Detection Explanation of Explainability module
- New checks added in moderation layer for Ban Code, Sentiment, Gibberish, and invisible text
- Multimodal Enhancement: Information Retrieval from PDFs Containing Images for Hallucination Module
- Multi-document type support for PII data masking of Privacy module
- Simplified Moderation Response for Chatbot's Split-Screen User Interface
- Logic of Thought (LoT) for improved LLM Reasoning: LLM-Explain Module
- LLM-Explain: Customization to configure any LLM endpoint to get explanation
- Bulk processing of multiple records for LLM-Explain

¹⁰⁹ <https://github.com/Infosys/Infosys-Responsible-AI-Toolkit>

¹¹⁰ <https://github.com/salus-rai/salus>

Contributors

We extend our sincere thanks to all the contributors who made this newsletter issue possible.



Srinivasan S - Policy Advocacy, Consultancy and Customer Outreach, Infosys Responsible AI Office



Mandanna A N - Head of Infosys Responsible AI Office, USA



Siva Elumalai - Senior Consultant, Infosys Responsible AI Office, India



Dakeshwar Verma - Senior Analyst - Data Science, Infosys Responsible AI Office, India



Utsav Lall - Senior Associate Consultant, Infosys Responsible AI Office, India



Pritesh Korde - Senior Associate Consultant, Infosys Responsible AI Office, India



Anie Juby - Industry Principal, Infosys Topaz Branding & Communications, Bangalore



Jossy Mathew - Senior Project Manager, Infosys Topaz Branding & Communications, Bangalore

Please reach out to responsibleai@infosys.com to know more about Responsible AI at Infosys.
We would be happy to have your feedback too.



**WHEN YOUR AI STARTS STORYTELLING,
KEEP IT GROUNDED WITH RESPONSIBLE AI.**

Infosys Topaz is an AI-first set of services, solutions and platforms using generative AI technologies. It amplifies the potential of humans, enterprises, and communities to create value. With 12,000+ AI assets, 150+ pre-trained AI models, 10+ AI platforms steered by AI-first specialists and data strategists, and a 'responsible by design' approach, Infosys Topaz helps enterprises accelerate growth, unlock efficiencies at scale and build connected ecosystems. Connect with us at infosystopaz@infosys.com

For more information, contact askus@infosys.com



© 2025 Infosys Limited, Bengaluru, India. All Rights Reserved. Infosys believes the information in this document is accurate as of its publication date; such information is subject to change without notice. Infosys acknowledges the proprietary rights of other companies to the trademarks, product names and such other intellectual property rights mentioned in this document. Except as expressly permitted, neither this documentation nor any part of it may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, electronic, mechanical, printing, photocopying, recording or otherwise, without the prior permission of Infosys Limited and/ or any named intellectual property rights holders under this document.