



AI ACTION
SUMMIT

International AI Safety Report

The International Scientific Report
on the Safety of Advanced AI

January 2025

This is the first International AI Safety Report. Following an interim publication in May 2024, a diverse group of 96 Artificial Intelligence (AI) experts contributed to this first full report, including an international Expert Advisory Panel nominated by 30 countries, the Organisation for Economic Co-operation and Development (OECD), the European Union (EU), and the United Nations (UN). The report aims to provide scientific information that will support informed policymaking. It does not recommend specific policies....

This report summarises the scientific evidence on the safety of general-purpose AI. The purpose of this report is to help create a shared international understanding of risks from advanced AI and how they can be mitigated. To achieve this, this report focuses on general-purpose AI – or AI that can perform a wide variety of tasks – since this type of AI has advanced particularly rapidly in recent years and has been deployed widely by technology companies for a range of consumer and business purposes. The report synthesises the state of scientific understanding of general-purpose AI, with a focus on understanding and managing its risks.

Amid rapid advancements, research on general-purpose AI is currently in a time of scientific discovery, and – in many cases – is not yet settled science. The report provides a snapshot of the current scientific understanding of general-purpose AI and its risks. This includes identifying areas of scientific consensus and areas where there are different views or gaps in the current scientific understanding.

People around the world will only be able to fully enjoy the potential benefits of general-purpose AI safely if its risks are appropriately managed. This report focuses on identifying those risks and evaluating technical methods for assessing and mitigating them, including ways that general-purpose AI itself can be used to mitigate risks.

Executive Summary

The purpose of this report

This report synthesises the state of scientific understanding of general-purpose AI – AI that can perform a wide variety of tasks – with a focus on understanding and managing its risks.

This report summarises the scientific evidence on the safety of general-purpose AI. The purpose of this report is to help create a shared international understanding of risks from advanced AI and how they can be mitigated. To achieve this, this report focuses on general-purpose AI – or AI that can perform a wide variety of tasks – since this type of AI has advanced particularly rapidly in recent years and has been deployed widely by technology companies for a range of consumer and business purposes. The report synthesises the state of scientific understanding of general-purpose AI, with a focus on understanding and managing its risks.

Amid rapid advancements, research on general-purpose AI is currently in a time of scientific discovery, and – in many cases – is not yet settled science. The report provides a snapshot of the current scientific understanding of general-purpose AI and its risks. This includes identifying areas of scientific consensus and areas where there are different views or gaps in the current scientific understanding.

People around the world will only be able to fully enjoy the potential benefits of general-purpose AI safely if its risks are appropriately managed. This report focuses on identifying those risks and evaluating technical methods for assessing and mitigating them, including ways that general-purpose AI itself can be used to mitigate risks. It does not aim to comprehensively assess all possible societal impacts of general-purpose AI. Most notably, the current and potential future benefits of general-purpose AI – although they are vast – are beyond this report's scope. Holistic policymaking requires considering both the potential benefits of general-purpose AI and the risks covered in this report. It also requires taking into account that other types of AI have different risk/benefit profiles compared to current general-purpose AI.

The three main sections of the report summarise the scientific evidence on three core questions: What can general-purpose AI do? What are risks associated with general-purpose AI? And what mitigation techniques are there against these risks?

Update on latest AI advances after the writing of this report: Chair’s note

Between the end of the writing period for this report (5 December 2024) and the publication of this report in January 2025, an important development took place. The AI company OpenAI shared early test results from a new AI model, o3. These results indicate significantly stronger performance than any previous model on a number of the field’s most challenging tests of programming, abstract reasoning, and scientific reasoning. In some of these tests, o3 outperforms many (but not all) human experts. Additionally, it achieves a breakthrough on a key abstract reasoning test that many experts, including myself, thought was out of reach until recently. However, at the time of writing there is no public information about its real-world capabilities, particularly for solving more open-ended tasks.

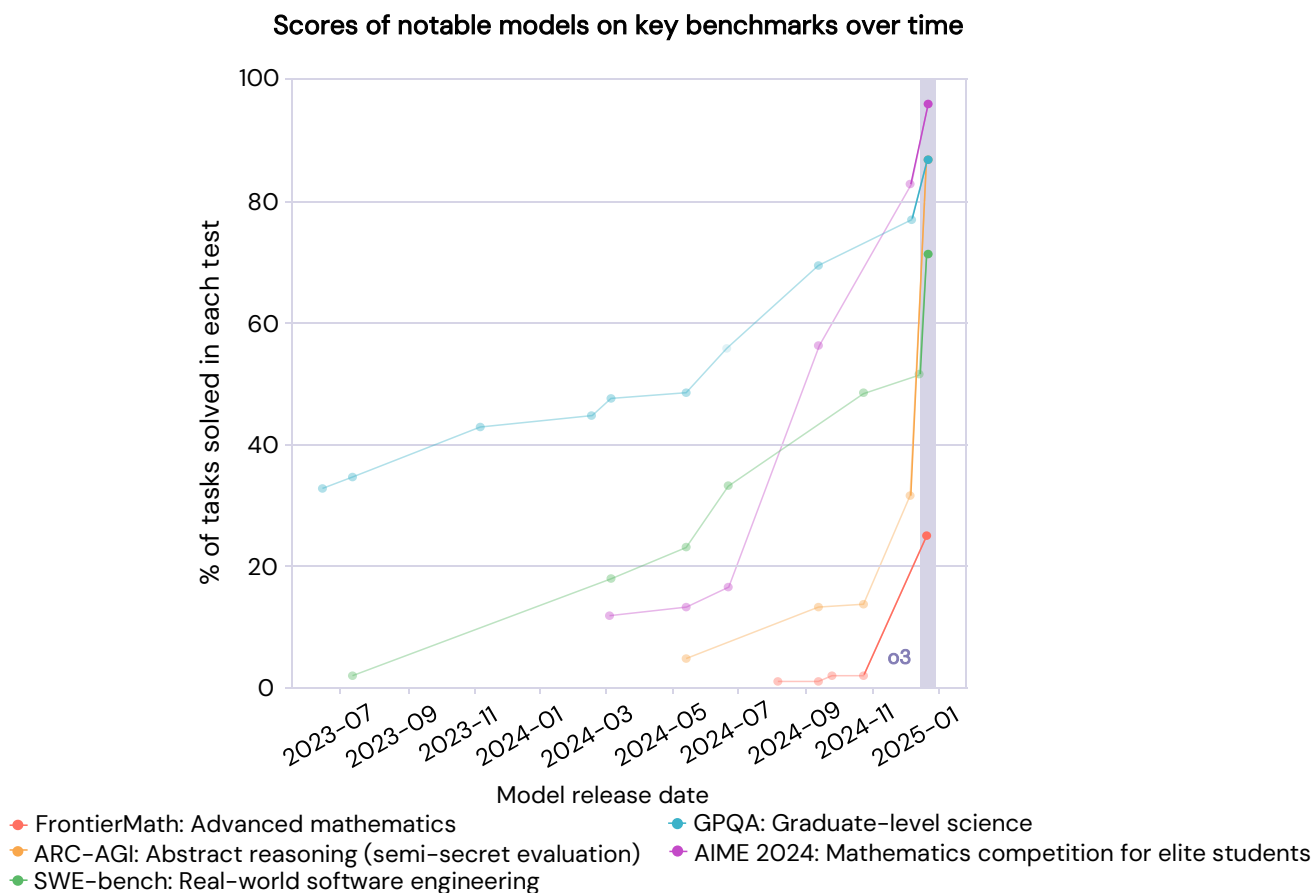


Figure 0.1: Scores of notable general-purpose AI models on key benchmarks from June 2023 to December 2024. o3 showed significantly improved performance compared to the previous state of the art (shaded region). These benchmarks are some of the field’s most challenging tests of programming, abstract reasoning, and scientific reasoning. For the unreleased o3, the announcement date is shown; for the other models, the release date is shown. Some of the more recent AI models, including o3, benefited from improved scaffolding and more computation at test-time. Sources: Anthropic, 2024; Chollet, 2024; Chollet et al., 2025; Epoch AI, 2024; Glazer et al. 2024; OpenAI, 2024a; OpenAI, 2024b; Jimenez et al., 2024; Jimenez et al., 2025.

The o3 results are evidence that the pace of advances in AI capabilities may remain high or even accelerate.

More specifically, they suggest that giving models more computing power for solving a given problem ('inference scaling') may help overcome previous limitations. Generally speaking, inference scaling makes models more expensive to use. But as another recent notable model, *R1*, released by the company DeepSeek in January 2025, has shown, researchers are successfully working on lowering these costs. Overall, inference scaling may allow AI developers to make further advances going forward. The o3 results also underscore the need to better understand how AI developers' growing use of AI may affect the speed of further AI development itself.

The trends evidenced by o3 could have profound implications for AI risks. Advances in science and programming capabilities have previously generated more evidence for risks such as cyber and biological attacks. The o3 results are also relevant to potential labour market impacts, loss of control risk, and energy use among others. But o3's capabilities could also be used to help protect against malfunctions and malicious uses. Overall, the risk assessments in this report should be read with the understanding that AI has gained capabilities since the report was written. However, so far there is no evidence yet about o3's real world impacts, and no information to confirm nor rule out major novel and/or immediate risks.

The improvement in capabilities suggested by the o3 results and our limited understanding of the implications for AI risks underscore a key challenge for policymakers that this report identifies: they will often have to weigh potential benefits and risks of imminent AI advancements without having a large body of scientific evidence available. Nonetheless, generating evidence on the safety and security implications of the trends implied by o3 will be an urgent priority for AI research in the coming weeks and months.

Key findings of the report

- **The capabilities of general-purpose AI, the type of AI that this report focuses on, have increased rapidly in recent years and have improved further in recent months.**[†] A few years ago, the best large language models (LLMs) could rarely produce a coherent paragraph of text. Today, general-purpose AI can write computer programs, generate custom photorealistic images, and engage in extended open-ended conversations. Since the publication of the Interim Report (May 2024), new models have shown markedly better performance at tests of scientific reasoning and programming.
- **Many companies are now investing in the development of general-purpose AI agents, as a potential direction for further advancement.** AI agents are general-purpose AI systems which can autonomously act, plan, and delegate to achieve goals with little to no human oversight. Sophisticated AI agents would be able to, for example, use computers to complete longer projects than current systems, unlocking both additional benefits and additional risks.
- **Further capability advancements in the coming months and years could be anything from slow to extremely rapid.**[†] Progress will depend on whether companies will be able to rapidly deploy even more data and computational power to train new models, and whether ‘scaling’ models in this way will overcome their current limitations. Recent research suggests that rapidly scaling up models may remain physically feasible for at least several years. But major capability advances may also require other factors: for example, new research breakthroughs, which are hard to predict, or the success of a novel scaling approach that companies have recently adopted.
- **Several harms from general-purpose AI are already well established.** These include scams, non-consensual intimate imagery (NCII) and child sexual abuse material (CSAM), model outputs that are biased against certain groups of people or certain opinions, reliability issues, and privacy violations. Researchers have developed mitigation techniques for these problems, but so far no combination of techniques can fully resolve them. Since the publication of the Interim Report, new evidence of discrimination related to general-purpose AI systems has revealed more subtle forms of bias.
- **As general-purpose AI becomes more capable, evidence of additional risks is gradually emerging.** These include risks such as large-scale labour market impacts, AI-enabled hacking or biological attacks, and society losing control over general-purpose AI. Experts interpret the existing evidence on these risks differently: some think that such risks are decades away, while others think that general-purpose AI could lead to societal-scale harm within the next few years. Recent advances in general-purpose AI capabilities – particularly in tests of scientific reasoning and programming – have generated new evidence for potential risks such as AI-enabled hacking and biological attacks, leading one major AI company to increase its assessment of biological risk from its best model from ‘low’ to ‘medium’.

[†] Please refer to the [Chair's update](#) on the latest AI advances after the writing of this report.

- **Risk management techniques are nascent, but progress is possible.** There are various technical methods to assess and reduce risks from general-purpose AI that developers can employ and regulators can require, but they all have limitations. For example, current interpretability techniques for explaining why a general-purpose AI model produced any given output remain severely limited. However, researchers are making some progress in addressing these limitations. In addition, researchers and policymakers are increasingly trying to standardise risk management approaches, and to coordinate internationally.
- **The pace and unpredictability of advancements in general-purpose AI pose an ‘evidence dilemma’ for policymakers.**[†] Given sometimes rapid and unexpected advancements, policymakers will often have to weigh potential benefits and risks of imminent AI advancements without having a large body of scientific evidence available. In doing so, they face a dilemma. On the one hand, pre-emptive risk mitigation measures based on limited evidence might turn out to be ineffective or unnecessary. On the other hand, waiting for stronger evidence of impending risk could leave society unprepared or even make mitigation impossible – for instance if sudden leaps in AI capabilities, and their associated risks, occur. Companies and governments are developing early warning systems and risk management frameworks that may reduce this dilemma. Some of these trigger specific mitigation measures when there is new evidence of risks, while others require developers to provide evidence of safety before releasing a new model.
- **There is broad consensus among researchers that advances regarding the following questions would be helpful:** How rapidly will general-purpose AI capabilities advance in the coming years, and how can researchers reliably measure that progress? What are sensible risk thresholds to trigger mitigations? How can policymakers best gain access to information about general-purpose AI that is relevant to public safety? How can researchers, technology companies, and governments reliably assess the risks of general-purpose AI development and deployment? How do general-purpose AI models work internally? How can general-purpose AI be designed to behave reliably?
- **AI does not happen to us: choices made by people determine its future.** The future of general-purpose AI technology is uncertain, with a wide range of trajectories appearing to be possible even in the near future, including both very positive and very negative outcomes. This uncertainty can evoke fatalism and make AI appear as something that happens to us. But it will be the decisions of societies and governments on how to navigate this uncertainty that determine which path we will take. This report aims to facilitate constructive and evidence-based discussion about these decisions.

[†] Please refer to the [Chair's update](#) on the latest AI advances after the writing of this report.

Contributors

CHAIR

Prof. Yoshua Bengio, Université de Montréal / Mila – Quebec AI Institute

EXPERT ADVISORY PANEL

This international panel was nominated by the governments of the 30 countries listed below, the UN, EU, and OECD.

Australia: Bronwyn Fox, the University of New South Wales

Brazil: André Carlos Ponce de Leon Ferreira de Carvalho, Institute of Mathematics and Computer Sciences, University of São Paulo

Canada: Mona Nemer, Chief Science Advisor of Canada

Chile: Raquel Pezoa Rivera, Universidad Técnica Federico Santa Maria

China: Yi Zeng, Chinese Academy of Sciences

European Union: Juha Heikkilä, European AI Office

France: Guillaume Avrin, National Coordination for Artificial Intelligence

Germany: Antonio Krüger, German Research Center for Artificial Intelligence

India: Balaraman Ravindran, Wadhvani School of Data Science and AI, Indian Institute of Technology Madras

Indonesia: Hammam Riza, Collaborative Research and Industrial Innovation in Artificial Intelligence (KORIKA)

Ireland: Ciarán Seoighe, Research Ireland

Israel: Ziv Katzir, Israel Innovation Authority

Italy: Andrea Monti, Legal Expert for the Undersecretary of State for the Digital Transformation, Italian Ministers Council's Presidency

Japan: Hiroaki Kitano, Sony Group Corporation

Kenya: Nusu Mwamanzi, Ministry of ICT & Digital Economy

Kingdom of Saudi Arabia: Fahad Albalawi, Saudi Authority for Data and Artificial Intelligence

Mexico: José Ramón López Portillo, LobsterTel

Netherlands: Haroon Sheikh, Netherlands' Scientific Council for Government Policy

New Zealand: Gill Jolly, Ministry of Business, Innovation and Employment

Nigeria: Olubunmi Ajala, Ministry of Communications, Innovation and Digital Economy

OECD: Jerry Sheehan, Director of the Directorate for Science, Technology and Innovation

Philippines: Dominic Vincent Ligot, CirroLytix

Republic of Korea: Kyoung Mu Lee, Department of Electrical and Computer Engineering, Seoul National University

Rwanda: Crystal Rugege, Centre for the Fourth Industrial Revolution

Singapore: Denise Wong, Data Innovation and Protection Group, Infocomm Media Development Authority

Spain: Nuria Oliver, ELLIS Alicante

Switzerland: Christian Busch, Federal Department of Economic Affairs, Education and Research

Türkiye: Ahmet Halit Hatip, Turkish Ministry of Industry and Technology

Ukraine: Oleksii Molchanovskyi, Expert Committee on the Development of Artificial Intelligence in Ukraine

United Arab Emirates: Marwan Alserkal, Ministry of Cabinet Affairs, Prime Minister’s Office

United Kingdom: Chris Johnson, Chief Scientific Adviser in the Department for Science, Innovation and Technology

United Nations: Amandeep Singh Gill, Under-Secretary-General for Digital and Emerging Technologies and Secretary-General’s Envoy on Technology

United States: Saif M. Khan, U.S. Department of Commerce

SCIENTIFIC LEAD

Sören Mindermann, Mila – Quebec AI Institute

LEAD WRITER

Daniel Privitera, KIRA Center

WRITING GROUP

Tamay Besiroglu, Epoch AI

Rishi Bommasani, Stanford University

Stephen Casper, Massachusetts Institute of Technology

Yejin Choi, Stanford University

Philip Fox, KIRA Center

Ben Garfinkel, University of Oxford

Danielle Goldfarb, Mila – Quebec AI Institute

Hoda Heidari, Carnegie Mellon University

Anson Ho, Epoch AI

Sayash Kapoor, Princeton University

Leila Khalatbari, Hong Kong University of Science and Technology

Shayne Longpre, Massachusetts Institute of Technology

Sam Manning, Centre for the Governance of AI

Vasilios Mavroudis, The Alan Turing Institute

Mantas Mazeika, University of Illinois at Urbana-Champaign

Julian Michael, New York University

Jessica Newman, University of California, Berkeley

Kwan Yee Ng, Concordia AI

Chinasa T. Okolo, Brookings Institution

Deborah Raji, University of California, Berkeley

Girish Sastry, Independent

Elizabeth Seger (generalist writer), Demos

Theodora Skeadas, Humane Intelligence

Tobin South, Massachusetts Institute of Technology

SENIOR ADVISERS

Daron Acemoglu, Massachusetts Institute of Technology

Olubayo Adekanmbi, contributed as a Senior Adviser prior to taking up his role at EqualyzAI

David Dalrymple, Advanced Research + Invention Agency

Thomas G. Dietterich, Oregon State University

Edward W. Felten, Princeton University

Pascale Fung, contributed as a Senior Adviser prior to taking up her role at Meta

Pierre-Olivier Gourinchas, Research Department, International Monetary Fund

Fredrik Heintz, Linköping University

Geoffrey Hinton, University of Toronto

Nick Jennings, University of Loughborough

Andreas Krause, ETH Zurich

Susan Leavy, University College Dublin

Percy Liang, Stanford University

Teresa Ludermir, Federal University of Pernambuco

Vidushi Marda, AI Collaborative

Emma Strubell, Carnegie Mellon University

Florian Tramèr, ETH Zurich

Lucia Velasco, Maastricht University

Nicole Wheeler, University of Birmingham

Helen Margetts, University of Oxford

John McDermid, University of York

Jane Munga, Carnegie Endowment for International Peace

Arvind Narayanan, Princeton University

Alondra Nelson, Institute for Advanced Study

Clara Neppel, IEEE

Alice Oh, KAIST School of Computing

Gopal Ramchurn, Responsible AI UK

Stuart Russell, University of California, Berkeley

Marietje Schaake, Stanford University

Bernhard Schölkopf, ELLIS Institute Tübingen

Dawn Song, University of California, Berkeley

Alvaro Soto, Pontificia Universidad Católica de Chile

Lee Tiedrich, Duke University

Gaël Varoquaux, Inria

Andrew Yao, Institute for Interdisciplinary Information Sciences, Tsinghua University

Ya-Qin Zhang, Tsinghua University

SECRETARIAT

AI Safety Institute

Baran Acar

Ben Clifford

Lambrini Das

Claire Dennis

Freya Hempleman

Hannah Merchant

Rian Overy

Ben Snodin

Mila — Quebec AI Institute

Jonathan Barry

Benjamin Prud'homme

ACKNOWLEDGEMENTS

Civil Society and Industry Reviewers

Civil Society: Ada Lovelace Institute, AI Forum New Zealand / Te Kāhui Atamai Iahiko o Aotearoa, Australia’s Temporary AI Expert Group, Carnegie Endowment for International Peace, Center for Law and Innovation / Certa Foundation, Centre for the Governance of AI, Chief Justice Meir Shamgar Center for Digital Law and Innovation, Eon Institute, Gradient Institute, Israel Democracy Institute, Mozilla Foundation, Old Ways New, RAND, SaferAI, The Centre for Long–Term Resilience, The Future Society, The Alan Turing Institute, The Royal Society, Türkiye Artificial Intelligence Policies Association.

Industry: Advai, Anthropic, Cohere, Deloitte Consulting USA and Deloitte LLM UK, G42, Google DeepMind, Harmony Intelligence, Hugging Face, IBM, Lelapa AI, Meta, Microsoft, Shutterstock, Zhipu.ai.

Special Thanks

The Secretariat appreciates the support, comments and feedback from Angie Abdilla, Concordia AI, Nitarshan Rajkumar, Geoffrey Irving, Shannon Vallor, Rebecca Finlay and Andrew Strait.

3.2. General challenges for risk management and policymaking

3.2.1. Technical challenges for risk management and policymaking

KEY INFORMATION

Several technical properties of general-purpose AI make risk mitigation for many risks associated with general-purpose AI difficult:

- A. **Autonomous general-purpose AI agents may increase risks:** AI developers are making large efforts to create and deploy general-purpose AI systems that can more effectively act and plan in pursuit of goals. These agents are not well understood but require special attention from policymakers. They could enable malicious uses and risks of malfunctions, such as unreliability and loss of human control, by enabling more widespread applications with less human oversight.
- B. **The breadth of use cases complicates safety assurance:** General-purpose AI systems are being used for many (often unanticipated) tasks in many contexts, making it hard to assure their safety across all relevant use cases, and potentially allowing companies to adapt their systems to work around regulations.
- C. **General-purpose AI developers understand little about how their models operate internally:** Despite recent progress, developers and scientists cannot yet explain why these models create a given output, nor what function most of their internal components perform. This complicates safety assurance, and it is not yet possible to provide even approximate safety guarantees.
- D. **Harmful behaviours, including unintended goal-oriented behaviours, remain persistent:** Despite gradual progress on identifying and removing harmful behaviours and capabilities from general-purpose AI systems, developers struggle to prevent them from exhibiting even well-known overtly harmful behaviours across foreseeable circumstances, such as providing instructions for criminal activities. Additionally, general-purpose AI systems can act in accordance with unintended goals that can be hard to predict and mitigate.
- E. **An 'evaluation gap' for safety persists:** Despite ongoing progress, current risk assessment and evaluation methods for general-purpose AI systems are immature. Even if a model passes current risk evaluations, it can be unsafe. To develop evaluations needed in time to meet existing governance commitments, significant effort, time, resources, and access are needed.
- F. **System flaws can have a rapid global impact:** When a single general-purpose AI system is widely used across sectors, problems or harmful behaviours can affect many users simultaneously. These impacts can manifest suddenly, such as through model updates or initial release, and can be practically irreversible.

Key Definitions

- **AI agent:** A general-purpose AI which can make plans to achieve goals, adaptively perform tasks involving multiple steps and uncertain outcomes along the way, and interact with its environment – for example by creating files, taking actions on the web, or delegating tasks to other agents – with little to no human oversight.
- **Deployment:** The process of implementing AI systems into real-world applications, products, or services where they can serve requests and operate within a larger context.
- **Evaluations:** Systematic assessments of an AI system's performance, capabilities, vulnerabilities or potential impacts. Evaluations can include benchmarking, red-teaming and audits and can be conducted both before and after model deployment.
- **Fine-tuning:** The process of adapting a pre-trained AI model to a specific task or making it more useful in general by training it on additional data.
- **Goal misgeneralisation:** A situation in which an AI system correctly follows an objective in its training environment, but applies it in unintended ways when operating in a different environment.
- **Interpretability research:** The study of how general-purpose AI models function internally, and the development of methods to make this comprehensible to humans.
- **Jailbreaking:** Generating and submitting prompts designed to bypass guardrails and make an AI system produce harmful content, such as instructions for building weapons.
- **Open-ended domains:** Environments into which AI systems might be deployed which present a very large set of possible scenarios. In open-ended domains, developers typically cannot anticipate and test every possible way that an AI system might be used.
- **Open-weight model:** An AI model whose weights are publicly available for download, such as Llama or Stable Diffusion. Open-weight models can be, but are not necessarily, open source.
- **Weights:** Model parameters that represent the strength of connection between nodes in a neural network. Weights play an important part in determining the output of a model in response to a given input and are iteratively updated during model training to improve its performance.

This section covers six general technical challenges that can make risk management and policymaking more difficult for a wide range of general-purpose AI risks (see Figure 3.1).

A. Autonomous general-purpose AI agents may increase risks: general-purpose AI agents – systems that can plan and act in the world with little to no human involvement elevate risks of malfunctions and malicious use. Today, general-purpose AI systems are primarily used as tools by humans. For example, a chatbot can write computer code, but a human runs, debugs, and integrates code into a larger software project. However, researchers and developers are making large efforts to design general-purpose AI agents – systems that can act and plan autonomously by controlling computers, programming interfaces, robotic tools, and by delegating to other AI systems (18, 55, 316*, 984, 985, 986*, 987, 988, 989, 990, 991*, 992). These systems are also sometimes called 'autonomous agents' or 'autonomous AI'. Researchers and developers are building

agents for a variety of domains, including web browsing (85*), research in chemistry and AI (22*, 121*, 402), software engineering (122, 259), cyber offence (127), general computer use (993, 994*, 995), and controlling robots (19*).

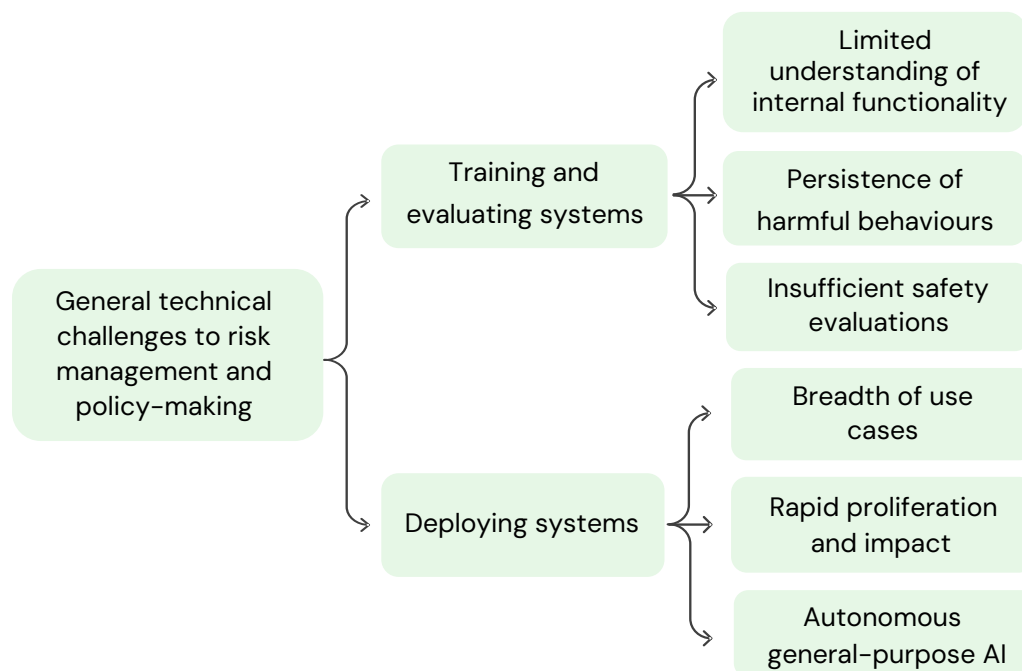


Figure 3.1: Technical challenges for managing general-purpose AI risks can be divided into two types: challenges with training and evaluating systems, and challenges with deploying them. This section discusses six broad challenges that apply to many risks. Source: International AI Safety Report.

Agentic general-purpose AI systems escalate risks by reducing human involvement and oversight.

The main purpose of general-purpose AI agents is to reduce the need for human involvement and oversight, allowing for much faster and cheaper applications. This is economically valuable, and increasingly agentic AI products are rapidly being developed and deployed. However, increased delegation to AI agents reduces human oversight and can increase the risk of accidents (996) (see [2.2.1. Reliability issues](#)). Meanwhile, agents can be uniquely vulnerable to attacks from malicious actors (997), for example by ‘hijacking’ an agent by placing instructions in places where the agent will encounter them (998). AI agents can also automate some workflows for malicious uses such as scams, hacking, and the development of weapons (127, 358, 999, 1000, 1001*) (see [2.1. Risks from malicious use](#) for more examples). AI agents could also uniquely contribute to risks of loss of human control if their capabilities advance significantly (see [2.2.3. Loss of control](#)) (316*, 1002). Furthermore, researchers have argued that it would be difficult or impossible to assure the safety of advanced agents by relying on testing, if those agents can make long-term plans and can distinguish testing conditions from real-world conditions (1003).

General-purpose AI agents can perform useful work autonomously but currently have limited reliability, especially for complex tasks. Current state-of-the-art general-purpose AI systems are capable of autonomously executing many simple tasks (e.g. writing short snippets of code), but they struggle with more complex ones (e.g. writing entire code libraries) (122, 593, 600, 1004).

They are particularly unreliable at performing tasks that involve many steps (1005). Meanwhile, general-purpose AI agents deployed to accomplish long-horizon tasks can be particularly vulnerable to manipulation by malicious actors (997). The capabilities of current and future agents are further discussed in [1.2. Current capabilities](#) and [1.3. Capabilities in coming years](#).

The capabilities of general-purpose AI agents are advancing rapidly, and understanding their future capabilities is a key evidence gap. General-purpose AI agents are rapidly becoming more capable. For example, ‘SWE-Bench’ is a popular benchmark (metric) used to evaluate the capabilities of agentic AI systems for software engineering tasks such as finding and fixing bugs (122). Since the Interim Report (May 2024), top models’ performance on SWE-Bench has increased from 26% to 42% (122), with the top 19 leading submissions all occurring after May 2024. This represents dramatic progress from October 2023, when the best model achieved only 2%. Meanwhile, the recent introduction of o1 (2*) marks a leap forward in the reasoning and problem-solving capabilities of general-purpose AI systems. These performance improvements are due to a combination of advances. First, as the general-purpose AI models underlying these agents become more capable, the agents’ cognitive abilities improve. Second, these agents are being developed with increasingly advanced training and planning methods. For example, AlphaProof, a ‘neuro-symbolic’ general-purpose AI system that combined neural networks with advanced planning techniques, achieved silver medal-level performance on 2024 International Mathematical Olympiad questions (187*). However, due to the rapid pace of progress in the area and the fact that many agents are proprietary, public understanding of current state-of-the-art methods is limited. Over the coming months and years, the development of more advanced agents demands special attention from policymakers.

B. The breadth of use cases complicates safety assurance: general-purpose AI systems can be applied in many unanticipated contexts, making it hard to test and assure their trustworthiness across all realistic use cases. General-purpose AI systems’ inputs and outputs are often open-ended, such as free-form text or image generation where users can enter any prompt. It is not possible to study the diffuse, downstream impacts of a system in a pre-deployment laboratory setting. This makes it challenging to make strong safety assurances because it is intractable to exhaustively test a system in all relevant usage contexts. For example, there are thousands of languages spoken by humans, making it very challenging to comprehensively assure the safety of language models across languages. Since the publication of the Interim Report (May 2024), general-purpose AI systems that can process multiple types of data (e.g. text, images, and audio) have become increasingly common (1006). This greatly expands the set of contexts which might cause the system to behave harmfully (1007). AI companies can readily redirect their systems’ capabilities between different applications and legal workarounds, posing challenges for targeted intervention approaches as seen historically in financial markets (1008).

C. General-purpose AI developers understand little about how their models operate internally. A key feature of general-purpose AI models is that their capabilities are mainly achieved through learning rather than from top-down design: an automatic algorithm adjusts billions of numbers

(‘parameters’) millions of times until the model’s output matches the training data. As a result, the current understanding of general-purpose AI models is more analogous to that of growing brains or biological cells than aeroplanes or power plants. AI scientists and AI developers only have a minimal ability to explain why these models made a given decision over another one, and how their capabilities arise from their known internal mathematical components. This contrasts, for example, with complex software systems such as web search engines, where the developers can explain the function of individual components (such as lines and files of code) and can also investigate why the system found a particular result. Current ‘interpretability’ techniques for explaining the internal structures of general-purpose AI models are unreliable and require major simplifying assumptions (1009, 1010*, 1011*, 1012, 1013*). In practice, techniques for interpreting the inner workings of neural networks can be misleading (466, 1014, 1015*, 1016, 1017, 1018, 1019), and can fail sanity checks or prove unhelpful in downstream uses (1020, 1021, 1022, 1023, 1024, 1025*). For example, one goal of interpretability research is to help researchers understand models well enough to edit their behaviours by modifying their weights. However, state-of-the-art interpretability tools have not yet proven useful and reliable for this (1026*). As discussed in [3.4.1. Training more trustworthy models](#), these research methods are actively being improved, and new developments may yield further insights. However, because of how deep learning models represent information across neurons in a highly distributed way (1027, 1028), it is unclear whether interpreting the inner structures of general-purpose AI models could offer guaranteed safety assurances. In other words, modern general-purpose AI systems may be too complex to tractably make performance guarantees for. At present, computer scientists are unable to give guarantees of the form ‘System X will not do Y’ (41). Nonetheless, a deeper understanding of models’ inner workings could be useful in many ways (see [3.4.2. Monitoring and intervention](#) and [3.4.1. Training more trustworthy models](#)).

D. Harmful behaviours, including unintended goal-oriented behaviours, remain persistent: ensuring that general-purpose AI systems act in accordance with the goals, behaviours and capabilities intended by their developers and users is difficult. Although general-purpose AI systems can excel at learning what they are ‘told’ to do, their behaviour may not necessarily be what their designers intended (607, 1029, 1030, 1031). Even subtle differences between a designer’s goals and the objectives given to a system can lead to unexpected failures. For example, general-purpose AI chatbots are often trained to produce text that will be rated positively by evaluators, but user approval is an imperfect proxy for user benefit. As a result, several widely-used chatbots have displayed ‘sycophantic’ or actively misleading behaviour, making statements that users approve of regardless of whether they are true (98, 317, 522, 608). For example, general-purpose AI language models are known to have a strong tendency to agree with opinions that a user expresses in chats (98). Even when a general-purpose AI system receives correct feedback during training, it may still develop a solution that does not generalise well when applied to new situations once deployed (‘goal misgeneralisation’) (616, 1032, 1033). For example, some researchers have found that language models’ safety training can be ineffective if the model is prompted in a language that was underrepresented in its training data (1034). Since the publication of the Interim Report (May 2024), researchers have demonstrated examples of unwanted goal-oriented behaviour from general-purpose AI systems. These include attempts at rewriting their own goals (599*).

Despite efforts to diagnose and debug issues, developers have not always been able to prevent even well-known and overtly harmful behaviours from general-purpose AI systems in foreseeable circumstances. Empirically, state-of-the-art general-purpose AI systems have exhibited a variety of harmful and often unexpected behaviours post-deployment (41, 1035, 1036). These hazards include general-purpose AI systems assisting malicious users in overtly harmful tasks (127, 319, 1037, 1038, 1039, 1040, 1041); leaking private or copyrighted information (1042, 1043, 1044*, 1045, 1046, 1047); generating hateful content (1048, 1049); exhibiting social and political biases (183, 438, 491, 511, 560, 561, 562, 563, 564, 565); pandering to user biases (98); and hallucinating inaccurate content (101, 102*, 104, 461, 1050, 1051*). Meanwhile, users have consistently been able to circumvent state-of-the-art general-purpose AI model safeguards with relative ease through prompting ('jailbreaks') (39, 155, 460, 904*, 1052, 1053, 1054, 1055, 1056*, 1057, 1058, 1059, 1060, 1061, 1062, 1063*) or simple model modifications (906, 1064, 1065, 1066, 1067, 1068, 1069, 1070, 1071, 1072, 1073, 1074, 1075, 1076, 1077, 1078, 1079, 1080). Since the publication of the Interim Report (May 2024), some researchers have also found that even when chat systems safely refuse harmful requests, they can still behave harmfully when used to operate as agents (1000, 1001*). Researchers continuously develop new techniques that defend against these attacks, but they also develop stronger attacks that usually overcome the existing defences (see [3.4.1. Training more trustworthy models](#)).

General-purpose AI systems sometimes gain and retain harmful capabilities even when they are explicitly fine-tuned not to (41, 1069). While current techniques are effective at suppressing harmful behaviours from general-purpose AI systems, these harmful capabilities can and do resurface from anomalies, inputs from malicious users, and modifications to models. For example, fine-tuning GPT-3.5 on only ten examples of harmful text can undo its safeguards and make it possible to elicit harmful behaviour (1064). The difficulty of making general-purpose AI systems fully resistant to overt failure modes has led some researchers to question whether it is possible to make current development approaches robust to such failure modes (1081, 1082). See [2.1. Risks from malicious use](#) for further discussion of harmful capabilities in AI models, [2.4. Impact of open-weight general-purpose AI models on AI risks](#) for a discussion of the benefits and risks of releasing models with both harmful and beneficial capabilities for public download, and [3.4.1. Training more trustworthy models](#) for a discussion of methods for unlearning harmful capabilities.

E. An 'evaluation gap' for safety persists: current safety evaluations are not thorough enough to meet existing governance frameworks and commitments from companies. Both developers and regulators are increasingly proposing risk management frameworks that rely on high-quality evaluations of general-purpose AI systems. The goal of evaluations is to identify risks so that they can be addressed or monitored. However, the science of evaluating general-purpose AI systems and predicting their downstream impacts is immature. Even when general-purpose AI systems are evaluated pre-deployment, new failure modes are often quickly discovered post-deployment (1055). For example, users found methods to subvert o1's safety fine-tuning within days of its release, and some researchers publicised work on a method to reliably jailbreak the model only three weeks after the model's release (1083). Evaluating AI systems for harmful behaviours and

downstream risks is a rapidly growing field. However, the large scope of potential risks (933), the limitations of benchmarking techniques (178, 1084, 1085), a lack of full access to systems (1086), and the difficulty of assessing downstream societal impacts (928*, 930*, 933) make high-quality evaluations challenging. [3.3. Risk identification and assessment](#) will delve further into methods for risk evaluation and broader risk assessment approaches.

F. System flaws can have a rapid global impact: because general-purpose AI systems can be shared rapidly and deployed in many sectors (like other software), a harmful system can rapidly have a global and sometimes irreversible impact. A small number of both proprietary and freely available open-weight general-purpose AI models currently reach many millions of users (see [2.3.3. Market concentration risks and single points of failure](#)). Both proprietary and open-weight models can therefore have rapid and global impacts, although in different ways (911). A risk factor for open-weight models is that there is no practical way to roll back access if it is later discovered that a model has faults or capabilities that enable malicious use (902) (see [2.4. Impact of open-weight general-purpose AI models on AI risks](#), [2.1. Risks from malicious use](#)). However, a benefit of openly releasing model weights and other model components such as code and training data is that it also allows a much greater and more diverse number of practitioners to discover flaws, which can improve understanding of risks and possible mitigations (911). Developers or others can then repair faults and offer new and improved versions of the system. This cannot prevent deliberate malicious use (902, 1075), which could be a concern if a system poses additional risk ('marginal risk') compared to using alternatives (such as internet search). All of these factors are relevant to the specific possibility of rapid, widespread, and irreversible impacts of general-purpose AI models. However, even when model components are not made publicly accessible, the model's capabilities still reach a wide user base across many sectors. For example, within two months of launch, the fully closed system ChatGPT had over 100 million users (1087).

3.2.2. Societal challenges for risk management and policymaking

KEY INFORMATION

Several economic, political, and other contextual factors make risk mitigation for many risks associated with general-purpose AI difficult:

- A. **As general-purpose AI advances rapidly, risk assessment, risk mitigation, governance, and enforcement efforts can struggle to keep pace.** Policymakers face the challenge of creating governance and/or regulatory environments that are sufficiently flexible, agile and future-proof.
- B. **Developers of general-purpose AI face strong competitive pressure, which can incentivise them to conduct less thorough risk mitigations.** Markets characterised by high fixed costs, low marginal costs, and network effects tend to create competitive pressures that discourage safety investments. The market for general-purpose AI is such a market.
- C. **The rapid growth and consolidation in the AI industry raises concerns about certain AI companies becoming particularly powerful because critical sectors in society are dependent on their products.** Such companies may become more inclined to take excessive risks or cut corners on safety standards if they expect that it would be costly for governments to let the company fail.
- D. **The inherent lack of both algorithmic transparency and institutional transparency in general-purpose AI makes legal liability hard to determine, potentially hindering governance and enforcement.** The fact that general-purpose AI systems can act in ways that were not explicitly programmed or intended by their developers or users raises questions about who should be held liable for resulting harm.

Key Definitions

- **Algorithmic transparency:** The degree to which the factors informing general-purpose AI output, e.g. recommendations or decisions, are knowable by various stakeholders. Such factors might include the inner workings of the AI model, how it has been trained, what data it is trained on, what features of the input affected its output, and what decisions it would have made under different circumstances.
- **Institutional transparency:** The degree to which AI companies disclose technical or organisational information to public or governmental scrutiny, including training data, model architectures, emissions data, safety and security measures, or decision-making processes.
- **Winner takes all:** A concept in economics referring to cases in which a single company captures a very large market share, even if consumers only slightly prefer its products or services over those of competitors.
- **Race to the bottom:** A competitive scenario in which actors like companies or nation states prioritise rapid AI development over safety.

- **First-mover advantage:** The competitive benefit gained by being the first to establish a significant market position in an industry.
- **Distributed training:** A process for training AI models across multiple processors and servers, concentrated in one or multiple data centres.
- **Human in the loop:** A requirement that humans must oversee and sign off on otherwise automated processes in critical areas.
- **Emergent behaviour:** The ability of AI systems to act in ways that were not explicitly programmed or intended by their developers or users.

A. As general-purpose AI markets advance rapidly, governance, regulatory or enforcement efforts can struggle to keep pace. A recurring theme in the discourse on general-purpose AI risk is the mismatch between the pace of technological innovation and the development of governance structures (1088). While existing legal and governance frameworks apply to some uses of general-purpose AI systems, and several jurisdictions (such as the European Union, China, the US, and Canada) have initiated or completed efforts to establish relevant standards or to regulate AI broadly and general-purpose AI specifically, areas of regulatory uncertainty persist, particularly regarding novel AI capabilities. In a market that is as fast-moving as the general-purpose AI market currently is, it is very difficult to fill such gaps reactively, because by the time a governance and/or regulatory fix is implemented it might already be outdated. For instance, critics of social media regulation often point to challenges in areas such as data privacy, suggesting that these issues developed more quickly than policymakers could effectively address them (1089, 1090). Policymakers face the challenge of creating flexible regulatory environments that are robust to technological change over time.

The pace and unpredictability of advancements in general-purpose AI pose an ‘evidence dilemma’ for policymakers. Given the sometimes rapid and unexpected advancements, policymakers will often have to weigh potential benefits and risks of imminent AI advancements without having a large body of scientific evidence available. In doing so, they face a dilemma. On the one hand, pre-emptive risk mitigation measures based on limited evidence might turn out to be ineffective or unnecessary. On the other hand, waiting for stronger evidence of impending risk could leave society unprepared or even make mitigation impossible, for instance if sudden leaps in AI capabilities, and their associated risks, occur. Companies and governments are developing early warning systems and risk management frameworks that may reduce this dilemma. Some of them trigger specific mitigation measures when there is new evidence of risks, while others require developers to provide evidence of safety before releasing a new model.

B. Developers of general-purpose AI face strong competitive pressure, which can incentivise them to conduct less thorough risk mitigations. The one-time cost of developing a state-of-the-art general-purpose AI model is very high, while the marginal costs of distributing such a model to (additional) users are relatively low. For example, the estimated cost of training GPT-4 was \$40 million (27), but once trained, the cost of running the model for a single query is believed to be just a few cents, allowing it to serve many users at a relatively low marginal cost. In economic theory,

these conditions can lead to a ‘winner takes all’ dynamic in which field leaders can quickly capture a large market, whereas second-place actors will be at a significant disadvantage. As such, if cutting corners in (for example) testing and safety could allow one developer to take the lead in model capability, then there is a strong incentive to cut those corners (1091). This dynamic is visible in social media platforms, where a large initial user base attracted more people to join certain platforms because that is where their friends were, making the leading platform more valuable to new users and further expanding its network, while newer social networks often struggled to achieve critical mass (1092). The ‘winner takes all’ dynamic raises concern about potential ‘race to the bottom’ scenarios, where actors compete to develop general-purpose AI models as quickly as possible while under-investing in measures to ensure that the models are safe and ethical (1093, 1094).

Markets characterised by high fixed costs, low marginal costs, and network effects tend to create competitive pressures that discourage safety investments. Economic theory and empirical studies have shown that, under conditions of high fixed costs, low marginal costs, and strong network effects, firms in highly competitive markets tend to under-invest in safety measures (1095, 1096, 1097, 1098). For instance, in the early commercial aviation industry, airlines operating with thin profit margins due to high fixed costs of aircraft acquisition and maintenance sometimes cut corners on safety procedures to reduce costs and maintain competitive ticket prices (1099). These conditions are present in the general-purpose AI market. Moreover, in highly competitive markets with significant first-mover advantages, economic theory suggests that risk-taking behaviour tends to be rewarded and may become prevalent among surviving firms (1100). While direct studies of safety investment in the AI market are currently lacking, these economic principles and empirical studies in other fields suggest cause for concern. This could contribute to situations in which it is challenging for general-purpose AI developers to commit unilaterally to stringent safety standards, as doing so might put them at a competitive disadvantage (1101). At the same time, from a long-term business perspective, releasing risky models without adequate safety measures could damage user trust and company reputation, potentially creating stronger incentives for safety investment than the short-term competitive pressures might suggest.

C. The rapid growth and consolidation in the AI industry raises concerns about certain AI companies becoming particularly powerful because critical sectors in society are dependent on their products, which might incentivise them to take excessive risks (see [2.3.3. Market concentration and single points of failure](#)). Such scenarios are well studied in the economic literature (1102). They arise when an organisation reaches a size and level of influence so substantial that potential failure could pose systemic risks to the economy or national security. Governments are therefore inclined to take steps to protect these organisations from failure, for example by forgiving debts or providing bailout money. When protected in this way, companies may become more inclined to take excessive risks or cut corners on safety standards (1103, 1104), though empirical evidence on this effect remains mixed (1105). There is some concern that critical sectors in society might over time become overly dependent on the products of a small number of leading AI companies in this way. AI applications are becoming more integral to everyday life, and smaller

startups often seek acquisition by or collaboration with larger companies to overcome market entry barriers, most notably the extremely high costs of training a general-purpose AI model. In such arrangements, the startups typically trade access to their innovations for use of the larger companies' computing infrastructure and latest models, further reinforcing the market concentration and, potentially, overreliance on the AI products of a few industry leaders (767).

Beyond market concentration dynamics, several other factors may contribute to underinvestment in risk mitigation. Similar to environmental pollution or public health issues such as tobacco, many potential harms from AI systems represent externalities – costs that may be borne by society rather than directly by the developers (1106, 1107, 1108). Additionally, economic theory suggests that when there is a significant time lag between actions and consequences, market actors may systematically underinvest in risk mitigation (1109). This challenge is compounded by the inherent uncertainty of these potential harms, making it difficult to quantify the appropriate level of investment in risk mitigation. While empirical evidence on this question is scarce, economic theory suggests that the immediate costs of risk mitigation weighed against uncertain future benefits creates incentives for underinvestment in safety measures.

D. General-purpose AI systems' inherent lack of transparency and limited institutional transparency in organisations that develop AI makes legal liability hard to determine, potentially hindering governance and enforcement. Tracking the development and use of AI systems is important for establishing liability for potential harms, monitoring and seeking evidence for malicious use, and noticing malfunctions (1002, 1110, 1111). In principle, people and corporate entities are held accountable, not the technology, which is why developers maintain a 'human in the loop' policy for many critical areas, where a human must oversee and sign off on otherwise-automated processes. However, tracing harm back to the responsible individuals is very challenging (1112, 1113, 1114), as is gathering evidence of error or negligence. This stems from both technical and institutional factors: AI models' decision-making processes are difficult to interpret even for their developers (lack of algorithmic transparency), and AI companies often treat their training data, methodologies, and operational procedures as commercially sensitive information not open to public scrutiny (lack of institutional transparency) (1025*, 1115, 1116, 1117, 1118, 1119, 1120). Without transparency into both technical systems and organisational processes, it is difficult to develop the kinds of comprehensive safety governance standards that are common in other safety-critical fields such as automotive, pharmaceuticals, and energy (1121, 1122, 1123). The fact that general-purpose AI systems can act in ways that were not explicitly programmed or intended by their developers or users raises questions about who should be held liable for resulting harm (174, 1124). These liability challenges become even more pronounced with increasingly autonomous AI systems that require less direct human oversight, as it becomes harder to trace specific harmful actions back to human instructions or decisions (see [3.1. Risk management overview](#)).

The concentration of AI expertise in private companies can create significant information gaps for policymakers and the public. While academic researchers and public sector experts contribute to AI development and safety research, much of the cutting-edge work in AI development occurs within private companies (1125, 1126). This concentration of expertise can make it challenging for

policymakers and the public to access the technical knowledge needed to make informed decisions about AI governance and risk management. The resulting information asymmetry between AI developers and other stakeholders could complicate efforts to develop appropriate governance and/or regulatory frameworks and safety standards.

3.3. Risk identification and assessment

KEY INFORMATION

- **Assessing general-purpose AI systems for hazards is an integral part of risk management.** Scientists use a variety of techniques to study hazards during system development, before deployment, and after deployment.
- **Existing AI regulations and commitments require rigorous risk identification and assessment.** Governments and general-purpose AI developers have adopted policies that require them to identify and assess the potential risks and impacts of general-purpose AI systems on people, organisations, and society.
- **While very useful, existing quantitative methods to assess general-purpose AI risks have significant limitations.** Safety risks heavily depend on how and where these systems are used, which is often unanticipated, making it hard to measure risks without guessing how people will use them. This is especially challenging for general-purpose AI because it can be used in countless different situations, and many potential harms (e.g. bias, toxicity, and misinformation) are hard to measure objectively. While current risk assessment methods are nascent, they can be greatly improved.
- **Rigorous risk assessment requires combining multiple evaluation approaches, significant resources, and better access.** Key risk indicators include evaluations of systems themselves, how people apply them, as well as forward-looking threat analysis. For evaluations at the technical frontier to be effective, evaluators need substantial and growing technical ability and expertise. They also need sufficient time and more direct access than is currently available to the models, training data, methodologies used, and company-internal evaluations – but companies developing general-purpose AI typically do not have strong incentives to grant these.
- **In recent months, more research has been evaluating how well AI risk assessment methods actually work, identifying current shortcomings and criteria for improvement.** While more evidence is needed – especially for new risks – this technical progress is complemented by institutional developments, as governments begin to build evaluation capacity and stakeholders work to establish clearer guidelines for who is responsible for different aspects of risk assessment.
- **The absence of clear risk assessment standards and rigorous evaluations is creating an urgent policy challenge, as AI models are being deployed faster than their risks can be evaluated.** Policymakers face two key challenges: 1. internal risk assessments by companies are essential for safety but insufficient for proper oversight, and 2. complementary third-party and regulatory audits require more resources, expertise and system access than is currently available.

Key Definitions

- **Risk:** The combination of the probability and severity of a harm that arises from the development, deployment, or use of AI.
- **Hazard:** Any event or activity that has the potential to cause harm, such as loss of life, injury, social disruption, or environmental damage.
- **Deployment:** The process of implementing AI systems into real-world applications, products, or services where they can serve requests and operate within a larger context.
- **Evaluations:** Systematic assessments of an AI system's performance, capabilities, vulnerabilities or potential impacts. Evaluations can include benchmarking, red-teaming and audits and can be conducted both before and after model deployment.
- **Benchmark:** A standardised, often quantitative test or metric used to evaluate and compare the performance of AI systems on a fixed set of tasks designed to represent real-world usage.
- **Red-teaming:** A systematic process in which dedicated individuals or teams search for vulnerabilities, limitations, or potential for misuse through various methods. Often, the red team searches for inputs that induce undesirable behaviour in a model or system to identify safety gaps.
- **Jailbreaking:** Generating and submitting prompts designed to bypass guardrails and make an AI system produce harmful content, such as instructions for building weapons.
- **Audit:** A formal review of an organisation's compliance with standards, policies, and procedures, typically carried out by an independent third party.
- **Incident reporting:** Documenting and sharing cases in which developing or deploying AI has caused direct or indirect harms.

To manage the risks of general-purpose AI, it is necessary to understand and measure the risks it poses to people, organisations, and society. Several governments and general-purpose AI developers have already adopted policies and regulations that require them to identify and assess the potential risks and impacts of general-purpose AI systems, triggering planned responses when risks reach specific thresholds. 'Risk identification' is the process of identifying the potential risks of the technology, including possible hazards and unintended outcomes. 'Risk assessment' is the process of assessing the severity and likelihood of occurrence of each identified risk. (See Table 3.1 in [3.1 Risk Management Overview](#) for an overview of risk management stages including risk identification and assessment as well as risk evaluation, risk mitigation, and risk governance).

Methods for risk identification

General-purpose AI risks can be identified and formulated at various levels of *specificity*. For example, one broad category of general-purpose AI risks is *confabulating* or '*hallucinating*' *misinformation* – that is, generating outputs that are inaccurate or misleading. A more specific instance of the same risk is general-purpose AI *making up a non-existent polling location* when the user prompts it to gather information about where to cast their ballot during a national election

(1127). The specification of a risk can make it easier or more difficult for evaluators to assess both its *severity* and *likelihood*. Better-specified risks are easier to assess and mitigate.

Evaluators need to understand the use cases of general-purpose AI well in order to conceptualise its risks with the appropriate degree of specificity. For example, if general-purpose AI users are likely to prompt it to gather information about political campaigns and voting procedures, then assessing the risk of the model ‘hallucinating a polling location’ may be a high priority. Therefore *participatory approaches*, which consist of engaging with various stakeholders and impacted communities to understand their use cases, practices, needs, and values, are especially helpful for identifying higher-priority risks to users. Crowd audits (1128) are one example of a participatory approach. They are designed to allow everyday users to collaboratively surface the potential harms of AI products and services. Creating accessible mechanisms for the public to report observed and perceived harms is another important method of risk identification. AI incident-tracking databases, such as the OECD’s AI Incidents Monitor (AIM), are platforms designed to collect, categorise, and report harmful incidents involving AI (459). In short, there is a need to identify and assess risks in context.

To facilitate general-purpose AI risk-identification practices, scholars have proposed taxonomies of hazards (439, 933, 951*, 1129). These taxonomies list risk categories, such as informational hazards, memorisation of the training data (which can lead to copyright infringement and privacy concerns), and malicious usage (e.g. writing malware). Taxonomies of hazards can serve as a starting point to help evaluators conceptualise, identify, and specify the salient risks associated with general-purpose AI in specific application domains. In conventional risk management and safety engineering, there are several well-established methods for identifying hazards and risks of a technology, including functional failure analysis and HAZOP (hazard and operability study) (1130). These methods have been adopted in a wide range of industries, including the automotive industry, which also considers SOTIF (946). In addition to risk typologies and taxonomies, recent work has begun adapting some of these conventional techniques, e.g. hazard analysis, the bowtie method and safety cases, to AI products and services (968, 1131, 1132, 1133), but additional research is necessary in this area. See [3.1. Risk management](#) overview for a discussion of further risk identification practices established in other fields.

Methods for risk assessment

Once high-priority risks are identified, they need to be assessed to determine the likelihood and severity of the harm, hazard, or unintended outcome in question.

Better understanding the current state of general-purpose AI risk assessment methods is essential to AI policy because risk assessments are a core component of many AI governance and regulatory approaches. For example, the EU AI Act classifies AI systems into four main risk tiers based on their potential impact and imposes different requirements on AI systems depending on their risk tier. Furthermore, many leading AI companies have agreed to create AI safety commitments with

mitigations that are proportional and specific to the assessed risk of their systems (1134). However, risk assessment is a relatively nascent topic of research in the AI safety community, and there are currently no fully validated, systematic approaches to assessing the severity and likelihood of general-purpose AI harms. Implementing the aforementioned policies will require a substantially more mature field of risk assessment for general-purpose AI.

Existing work in AI safety heavily focuses on conventional model testing approaches in AI, often conducted after the development of general-purpose AI models. This reliance on retrospective (as opposed to prospective) risk assessment can lead to major omissions and misestimations of high-priority risks. In conventional risk management and safety engineering, a critical stage of risk assessment is the prospective analysis of risks before completing the design and development of a system. This stage is currently often overlooked in general-purpose AI risk assessments. In AI safety, risk assessment primarily consists of running a battery of tests and evaluations on the general-purpose AI system, then translating the results into quantitative estimates of risks. This is in contrast to traditional risk assessment, which consists of 1. analysing the causes, consequences, and prevalence of risks (through methods such as causal mapping and Delphi technique) then 2. evaluating whether the risk is acceptable, e.g. through checklists and risk matrices. Recent work has begun adapting some of these techniques to AI products and services (944, 968). See [3.1. Risk management overview](#) for further discussion of risk assessment approaches that are established in other fields.

Existing technical approaches and methodologies to general-purpose AI risk assessment rely heavily on testing and evaluations which can be broken down into four layers (1135):

1. **Model testing** evaluates the general-purpose AI model in terms of (often quantitative) metrics of performance on proxy tasks designed to represent real-world usage. These tests often take the form of benchmarks – fixed sets of prompts to test a model on.
2. **Red-teaming** is a systematic process in which dedicated individuals or teams search for vulnerabilities, limitations, or potential for misuse in AI models or systems through various methods. Often, the red team searches for inputs that induce undesirable behaviour for the purpose of improving the model or system's protections against such attacks.
3. **Field testing** evaluates the risks of general-purpose AI under real-world conditions.
4. **Long-term impact assessments** monitor and evaluate long-term impacts of the system on people, organisations, and society.

One major evidence gap is research to establish the validity, reliability, and practicality of existing general-purpose AI risk assessment methods. Good risk measurement methods must be *valid*, *reliable*, and *practical*. Validity refers to the extent to which a test, tool, or instrument accurately measures what it is intended to measure. For instance, validity issues arise if a benchmark differs from real-world use or contains false labels (1136). *Reliability* refers to the consistency, stability, and dependability of a measurement over time and across different contexts. In other words, it indicates the degree to which a measurement yields consistent,

repeatable results under similar conditions (1137). Prior work has shown that even small perturbations to prompts can have significant effects on the behaviour and performance of general-purpose AI on benchmarks (1138, 1139). *Practicality* assesses whether the measurement can be conducted efficiently and effectively in practice by the designated evaluators, considering constraints such as time, cost, computational resource availability, and burden on evaluators. For example, the process of evaluating general-purpose AI increasingly relies on using general-purpose AI (522, 929*), which requires technical capacity and raises new concerns (e.g. about LLM agents favouring outputs from their own model family (1140)). For rigorous risk assessment, validity and reliability are prioritised over ease and convenience of measurement (1141).

Since the publication of the Interim Report, the scientific community has made progress toward further implementing and evaluating existing risk assessment methods. The US and UK AI Safety Institutes (US AISI and UK AISI) recently published a technical report detailing a pre-deployment evaluation of the upgraded version of Claude 3.5 Sonnet (1142). New research has examined reproducibility (1143, 1144*) or validity, which can be compromised when AI models are trained on or exposed to test data beforehand (benchmark contamination) (1145, 1146). However, additional evidence is necessary to characterise the strengths and weaknesses of existing general-purpose AI evaluation methods (465) especially when general-purpose AI is utilised in new domains.

The initial layer of general-purpose AI risk assessment often consists of testing the model's behaviour across certain fixed benchmarked tasks. New benchmarks and standardised tests and metrics have been designed to evaluate and compare various categories of risk for general-purpose AI applications in stylised scenarios and tasks (122, 137, 141, 1147*, 1148, 1149*). For example, the AI Safety Benchmark from MLCommons (457) provides a benchmark to measure seven risk categories, such as misinformation and harmful content. Holistic Evaluation of Language Models (HELM) consists of 16 scenarios and seven metrics, including robustness, fairness, and bias (1150). Harmful capability evaluations (318*) are used to assess whether the general-purpose AI has particularly dangerous knowledge or skills (such as the ability to aid cyberattacks (2.1.3. [Cyber offence](#)) or aid the design of bioweapons (2.1.4. [Biological and chemical attacks](#))). Highly consequential upcoming decisions by companies and governments about model release partially rely on these evaluations (596*, 947*, 1134). Existing benchmarks significantly vary in quality (1151), and the scope of applicability for existing benchmarks is often unclear. Some best practices for creating high-quality benchmarks have been proposed (1151, 1152*).

While model testing methods can serve as a necessary first step toward assessing the risks of general-purpose AI, they are not sufficient on their own. It is impossible to derive reliable quantitative conclusions about the risks these methods aim to capture without making strong assumptions about patterns of use in specific applications. Such assumptions are hard to justify: First, the technology is general-purpose and can be used in numerous contexts, so it is difficult to predict patterns of use. Second, some risks (e.g. bias, toxicity, and misinformation) are difficult to

specify objectively, and any definitions must rest on questionable assumptions about what is (for example) ‘toxic’ or ‘biased’. Therefore, benchmarks cannot capture the risks associated with the usage of general-purpose AI in new domains and for novel tasks, because test conditions always differ from real-world usage to varying degrees (1153*). Benchmarks at best serve as a proxy measure for the risk category in question (for example, subjective ratings of human annotators or content moderators may serve as proxy for ‘toxicity’ (1154)). However, these proxy measures often do not reliably reflect the true risk in context. For instance, if human evaluators are not diverse, this can lead to benchmarks containing biased labels, since people from similar backgrounds might systematically miss certain examples of toxicity or misinformation. Moreover, improving scores on a benchmark does not always translate to lowering the associated risk in practice. For example, an LLM can pass the bar exam for lawyers, but that does not mean that it can create effective legal briefs (445, 446, 451). Any fixed benchmark is often easy to improve on without mitigating the target risk (1070). While creating capacity for dynamically evolving, collaborative benchmarks may address some of these challenges, it is important for AI evaluators to understand the inherent limitations of quantitative approaches to model testing (1155) and avoid over-reliance on them as the primary layer of risk assessment.

Red-teaming and adversarial attacks are other prominent methods to identify and assess risks, but can require special access. ‘Red team’ refers to a set of evaluators tasked with finding vulnerabilities in a system by attacking it. In contrast to benchmarks, which are mostly static and consist of a fixed set of test cases, a key advantage of red-teaming is that it adapts the evaluation to the specific system being tested. Through adversarial interactions with a system, red-teamers can design custom inputs to identify worst-case behaviours, malicious use opportunities, and unexpected failures. As an example, attacks against language models can take the form of automatically generated inputs (904*, 1053, 1063*, 1156, 1157, 1158*, 1159, 1160, 1161, 1162) or manually generated ones (1056*, 1059, 1158*, 1163). In automated attacks, for example, LLMs can be used to generate prompts designed to make another AI system produce harmful content, such as instructions for dangerous materials, even after the system initially refuses. These ‘jailbreaking’ attacks subvert the models’ safety restrictions (460, 904*, 1052, 1053, 1164, 1165*). Automated approaches can systematically test thousands of variations of potential attacks, allowing for more extensive and rapid coverage than manual testing alone. However, manual red-teaming over longer conversations can catch issues that current automated attacks alone can miss (1056*). However, it can be slow, labour intensive, and require special access. Further research for faster and effective automated red-teaming is necessary to address this challenge (1166).

While red-teaming is more effective at surfacing a wider range of general-purpose AI risks than model testing, many important harms and hazards may remain undetected. Importantly, if a red-teaming activity fails to surface certain categories of risks, that does not imply that those risks are unlikely. Previous work has found that bugs often evade detection (1022). A real-world example is jailbreaks, which induce general-purpose chat systems to comply with harmful requests that they were trained to refuse (460, 904*, 1052, 1053, 1164), and which evaded initial detection by developers (48*, 147*, 1158*). Research has also called into question whether red-teaming can

produce reliable and reproducible results. One study shows that red-teaming practices in industry diverge along several key axes, including the setting (e.g. the characteristics of red-teamers and the resources and methods available to them), and the decisions it informs (e.g. subsequent reporting, disclosure, and mitigation) (1167). The composition of the red team and the instructions provided to red-teamers (1168*), the number of attack rounds (1056*), and the availability of auxiliary or automation tools (1161, 1169) can significantly influence the outcomes of the activity, including the risk surface covered. See Table 3.2 for an overview of criteria for structuring red-teaming activities in practice. Comprehensive guidelines on red-teaming aim to address some of these challenges (1170).

| Phase | Key Questions and Considerations |
|-----------------------------|---|
| 0. Pre-activity criteria | What is the artefact under evaluation through the proposed red-teaming activity? |
| | What is the threat model the red-teaming activity aims to recreate? |
| | What is the specific vulnerability the red-teaming activity aims to find? |
| | What are the criteria for assessing the success of the red-teaming activity? |
| | What is the team composition , or who will be part of the team? |
| 1. Within-activity criteria | What resources are available to participants? |
| | What instructions are given to the participants to guide the activity? |
| | What kind of access do participants have to the model? |
| | What methods can members of the team utilise to test the artefact? |
| 2. Post-activity criteria | What reports and documentation are produced on the findings of the activity? |
| | What were the resources the activity consumed? |
| | How successful was the activity in terms of the criteria specified in phase 0? |
| | What are the proposed measures to mitigate the risks identified in phase 1? |

Table 3.2: Different types of criteria can help practitioners to structure red-teaming before, during, and after the relevant activities. Source: based on the criteria proposed by Feffer et al., 2024 (1167).

‘Field tests’ are exercises designed to assess risks under normal use conditions. ‘Human uplift studies’ examine whether people can use AI to perform malicious tasks better than they could without AI. ‘Human uplift’ studies are one important variant of field testing. They aim to measure how access to general-purpose AI systems improves individuals’ competencies and performance. For example, a human uplift study might explore how an AI system affects a person’s ability to accomplish complex tasks, such as customer support (662) or (potentially harmful) cybersecurity operations (361, 1171, 1172, 1173), compared to their performance without the AI assistance. These studies aim to quantify the ‘uplift’ in human capabilities and assess whether the AI’s support introduces new risks, such as lowering barriers to harmful conduct (see [2.4. Impact of open-weight general-purpose AI models on AI risks](#) for further discussion of uplift studies). However, there are several challenges in designing and conducting such studies, including simulating conditions similar

to ordinary use and choosing the appropriate measures of uplift. Evaluators could address some of these challenges if there were better guidelines for conducting human uplift studies and integrating them into the staged rollout of general-purpose AI products. In other safety-critical industries, for example drug testing in clinical trials, a series of studies are conducted in increasingly more realistic conditions (for example, going from testing on animals to human-subject studies), before the drug is deemed ready to market. A similar approach may prove useful for developing effective field testing methods for general-purpose AI.

Certain risks associated with general-purpose AI are likely to manifest only in the long run, making long-term impact assessments crucial. Such risks include the effects of the technology on labour markets and the future of work ([2.3.1. Labour market risks](#)), risks associated with more capable future AI systems ([2.2.3. Loss of control](#), [2.1.3. Cyber offence](#), [2.1.4. Biological and chemical attacks](#)), the environmental impact of AI development and use (see [2.3.4. Risks to the environment](#)), and long-term impacts on human cognition, wellbeing, and control (1003). Careful monitoring, investigating and rectifying long-term harms is necessary to maintain the public's confidence in the technology and prevent calls for unnecessarily strong controls. Accurately gauging the downstream societal impacts of general-purpose AI is challenging due to 1. uncertainties surrounding the capabilities of future general-purpose AI systems, and 2. the existence of numerous confounding factors that make it difficult to attribute long-term trends to any single cause. Creating capacity for predicting and monitoring the potential downstream societal impacts of general-purpose AI requires multidisciplinary analysis and the involvement of diverse perspectives (929*, 1174, 1175).

Challenges and opportunities

In addition to the challenges discussed here, see also [3.2.1. Technical challenges for risk management and policymaking](#) and [3.2.2. Societal challenges for risk management and policymaking](#).

The culture of 'build-then-test' in AI hinders comprehensive risk assessment and mitigation.

In conventional risk management, risk assessment is integrated into all stages of product design, development, and deployment, and is tightly intertwined with risk mitigation strategies. In AI safety, however, current risk assessment methods are largely conducted after development, and independent from risk mitigation. Prior work (978) has proposed the creation of safety case studies and safety guarantees for AI (1176). Adapting and implementing such practices for general-purpose AI requires both a cultural shift and further research.

The four layers of risk assessment (model testing, red-teaming, field testing, and long-term impact assessment) are necessary but not sufficient for comprehensive risk assessment. Existing methods do not provide generalisable guarantees or assurances surrounding the likelihood and severity of general-purpose AI harms (1177). The main evidence gaps are in 1. assessing the validity, reliability and practicality of each evaluation layer independently, and 2. combining information from different layers of evaluation to produce actionable insights (41).

Conducting comprehensive risk assessment, in practice, requires considerable access, resources, and time, which are often constrained. Very few entities have the *resources* (or the will to allocate the necessary resources) to conduct comprehensive evaluations, and potential conflicts of interest can lead to misleading results and reports (1014, 1178). Moreover, sometimes evaluators are not given enough time to thoroughly test models. In some cases, companies only provided evaluators with several days to test a new model before release (2*, 129). Effective model evaluation requires substantial time and resources.

Furthermore, developers of state-of-the-art general-purpose AI systems often limit external access to their technology (880). For models that are hosted on a developer's platform or that have to be accessed via an API (giving 'black box' access, only to model inputs and outputs), it can be challenging for external evaluators to perform effective adversarial attacks, model interpretations, and fine-tuning (1086, 1179). For example, AI models are usually trained to refuse dangerous requests, but to assess dangerous capabilities, evaluators require access to versions of the model without this guardrail. This access is sometimes provided (2*). Without it, certain high-priority risks may be overlooked. Incomplete information about how a system was designed, including data, techniques, implementation details, and organisational details hinders evaluations of the development process (34, 488, 1086, 1180, 1181, 1182). Some scholars have argued that a combination of technical, physical, and legal measures can offer external researchers' direct access without compromising trade secrets more than they are already compromised (1086). Several studies have advocated for legal 'safe harbours' (1036) or government-mediated access regimes (939) to enable evaluators to conduct independent evaluations without the risk of being prosecuted or banned from use. Researchers have proposed methods for structured access that do not require making the model's code and training weights public (1183), but that do make it possible for independent researchers and auditors to fully access the model in a secured environment designed to avoid leaks. Researchers are developing auditing techniques that use 'secure enclaves'. These techniques have the potential to avoid leaking the model parameters to auditors, and also the audit details to model developers (1184).

Successful risk assessment requires the participation of diverse perspectives in the evaluation process. The composition of the evaluation team in evaluation layers, such as red-teaming, can play a critical role in the process of discovering, characterising, and prioritising harms (1185). Improving stakeholder participation has been a focus of the machine learning community in recent years (932, 1186, 1187). Multiple strategies have been proposed, from broadening the understanding of 'impacts' in AI impact assessments (1188) to enabling a more inclusive range of human feedback (1189, 1190). However, fostering participation requires sensitivity to several criteria (1186), such as respect for participating parties to minimise the potential for exploitation (540), and surfacing hard choices between incompatible values or priorities (467, 538, 574). This process can be facilitated by methods from practical ethics such as 'reflective equilibrium' – the mutual adjustment of principles and judgements until they agree with each other (1191).

Policymakers face several challenges around how to incentivise adequate risk identification and assessments for general-purpose AI systems. Without clear guidelines, standards, and resources surrounding general-purpose AI risk assessment, practitioners face uncertainties as to what constitutes adequate risk assessments in their specific use cases. This in turn makes it difficult for policymakers to incentivise compliance. Another policy challenge is how to designate responsibility for various layers of risk assessments across different general-purpose AI stakeholder groups, including technology creators, users, and third-party auditors (763). Another approach is creating resources (for example, ‘sandboxes’ and ‘safe harbours’) that promote public-interest evaluations (1036) or third-party audits. The success of this approach hinges heavily on the availability of resources, trained evaluators and experts, incentives to conduct rigorous evaluations (for example, by offering indemnity and compensation), and access to models or information about data and methods used. Several governments have begun to build capacity for conducting technical evaluations and audits of general-purpose AI. It remains to be seen how much these efforts will advance interdisciplinarity and inclusive evaluation of general-purpose AI in the near future, and how much they can and will be scaled in practice (537, 540, 1192, 1193).

3.4. Risk mitigation and monitoring

3.4.1. Training more trustworthy models

KEY INFORMATION

- **Current training methods show progress on mitigating safety hazards from malfunctions and malicious use but remain fundamentally limited.** There has been progress in training general-purpose AI models to function more safely, but no current method can reliably prevent even overtly unsafe actions.
- **A multi-pronged approach is emerging as necessary for safety.** Evaluating the trustworthiness of models requires analysing many aspects of their behaviour and their development process – including factual accuracy, human supervision quality, AI system internals, and analysis of potential misuse patterns – all of which must inform training methodologies. While techniques exist to remove harmful capabilities, current methods tend to suppress rather than eliminate them.
- **Adversarial training provides limited robustness against attacks.** Adversarial training involves deliberately exposing AI models to examples designed to make them fail or misbehave during training, aiming to build resistance to such cases. However, adversaries can still find new ways ('attacks') to circumvent these safeguards with low to moderate effort, such as 'jailbreaks' that lead models to comply with harmful requests even if they were fine-tuned not to do so.
- **Since the publication of the Interim Report (May 2024), recent advances reveal both progress and new concerns.** Improved understanding of model internals has advanced both adversarial attacks and defences without a clear winner. Additionally, growing evidence suggests that current training methods – which rely heavily on imperfect human feedback – inadvertently cause models to mislead humans on difficult questions by making errors harder to spot. Improving the quantity and quality of human feedback is an avenue for progress, though nascent training techniques using AI to detect misleading behaviour also show promise.
- **Key challenges for policymakers centre around uncertainty and verification.** There are no reliable methods to quantify the risk of unexpected model failures. While some researchers are exploring provably safe approaches, these remain theoretical. This suggests that frameworks for safety training currently need to focus on processes to search for, respond to, and mitigate new failures before they cause unacceptable harm.

Key Definitions

- **Interpretability:** The degree to which humans can understand the inner workings of an AI model, including why it generated a particular output or decision. A model is highly interpretable if its mathematical processes can be translated into concepts that allow humans to trace the specific factors and logic that influenced the model's output.
- **Red-teaming:** A systematic process in which dedicated individuals or teams search for vulnerabilities, limitations, or potential for misuse through various methods. Often, the red team searches for inputs that induce undesirable behaviour in a model or system to identify safety gaps.
- **Adversarial training:** A machine learning technique used to make models more reliable. First, developers construct 'adversarial inputs' (e.g. through red-teaming) that are designed to make a model fail, and second, they train the model to recognise and handle these kinds of inputs.
- **Reinforcement learning from human feedback (RLHF):** A machine learning technique in which an AI model is refined by using human-provided evaluations or preferences as a reward signal, allowing the system to learn and adjust its behaviour to better align with human values and intentions through iterative training.
- **Jailbreaking:** Generating and submitting prompts designed to bypass guardrails and make an AI system produce harmful content, such as instructions for building weapons.

The risks of general-purpose AI systems may be mitigated in part by limiting their behaviours. For example, policymakers may wish to prevent general-purpose AI systems from providing dangerous information to users (e.g. on the production of weapons; see [2.1.4. Biological and chemical attacks](#)), being used for malicious purposes (e.g. for cyberattacks; see [2.1.3. Cyber offence](#)), or having malfunctions that lead to harm (see [2.2. Risks from malfunctions](#)). A system's behaviour is safe if it avoids such mistakes, and a system is robust if it continues to behave safely in a wide range of circumstances. Beyond this, a system is *adversarially* robust if it maintains safe behaviour even in the presence of an adversary (e.g. a human user) trying to get it to perform harmful or illegal tasks. There exist proposals for how to build general-purpose AI systems which are guaranteed to behave safely (1176), but this is not possible without significant technological advances and may require significant changes to the architecture of current general-purpose AI systems. Regulation of current systems will have to focus on ensuring that their training and development minimises the harms of malfunctions and misuse.

Since the publication of the Interim Report, both attackers and defenders have become better at leveraging a deeper understanding of AI systems' internal workings to respectively induce or prevent harmful behaviour, and the advantage remains with attackers. New methods to resist adversarial attack by leveraging the concepts internally represented in neural networks have been developed both for image models (1194*) and language models (1195). However, these approaches are not completely robust, and another recent study has shown that language models internally represent the refusal of harmful requests in a simple

way which allows them to be easily exploited as well (907). On balance, the advantage generally remains with attackers, who can induce a model to engage in harmful behaviour with only moderate effort. However, these developments suggest that further research on both attacks and defences will likely leverage progress in interpretability. If this is true, further advances may favour defenders in the case of closed-weights models, since attackers will not have access to neural network internals in these cases.

Evidence has also grown that existing methods for training general-purpose models can lead them to produce more misleading (i.e. false but convincing) outputs. A recent study showed that in the case of especially challenging questions, training general-purpose AI systems to maximise human approval of the answers led the systems to obfuscate their mistakes and make them harder for humans to spot, instead of becoming more accurate (608). Other studies in simulated environments have found that an AI learns to use harmful strategies (e.g. hiding information or exploiting its supervisor’s biases) to receive positive feedback (1196) or modify its training environment to increase its reward (599*), if enough information is available to the AI on how to do so. Using AI to help supervisors avoid errors remains a challenging problem, but there has also been modest progress in this area, with two recent studies showing cases where models become easier to supervise when optimised to debate themselves (1197, 1198). These developments highlight the need for further research investigating the behaviours incentivised by current training methods, and developing new training methods that provide better incentives and generally more trustworthy outputs by design.

The main evidence gaps around training trustworthy models include:

- Despite recent progress (1010*, 1012, 1199)), it is still unclear whether interpretability methods, which help researchers and evaluators understand how models function internally, will be useful enough to substantially inform model training and testing. There are preliminary studies of this (1076, 1200, 1201).
- It is unclear whether ‘scalable oversight’ protocols, where AI systems can help humans evaluate their outputs, can provide a strong lever by which models can be trained to be more trustworthy even on hard problems (609*).
- There are currently no viable technical approaches to rigorously quantifying the risk of unforeseen or unexpected failures in large general-purpose AI systems. Although there is ongoing research on obtaining probabilistic safety guarantees, there is no practical technique to obtain even approximate guarantees yet.

For policymakers, key challenges include:

- Research moves very quickly in AI training, making it a moving target for regulation.
- It is difficult to quantify the risk of unexpected, unforeseen failure modes. In addition, it is unclear what are the best practices by which AI developers should detect, respond to, and mitigate newly discovered failures to minimise risks.

Robustness

Incentivising safe and correct behaviour during system training

It is challenging to precisely specify objectives for general-purpose AI systems in a way that does not unintentionally incentivise harmful behaviours. Currently, researchers do not know how to specify abstract human preferences and values (such as reporting the truth, figuring out and doing what a user wants, or avoiding harmful actions) in a way that can be used to train general-purpose AI systems. Moreover, given the complex socio-technical relationships embedded in general-purpose AI systems, it is not clear whether such specification is even possible. After an initial pre-training phase, general-purpose AI systems have learned to imitate human behaviour and are then generally tuned to optimise for objectives that are imperfect proxies for the developer's true goals (1031). For example, AI chatbots are often tuned to produce text that will be rated positively by human evaluators, but user approval is an imperfect proxy for user benefit. Research has shown that several widely used chatbots sometimes match their stated views to a user's views regardless of truth (98, 522) possibly creating 'echo chambers', and that training general-purpose AI systems to satisfy human evaluators' assessments can incentivise the system to provide harder-to-check answers that obfuscate the system's mistakes (608). This is an ongoing challenge for general-purpose AI systems (607, 1029, 1031, 1202*).

Researchers have methods to measure whether training incentivises the right behaviour using experiments with human evaluators, but current results are preliminary. 'Scalable oversight' experiments test whether an evaluator can successfully steer an AI system to correctly perform a task that the evaluator is unable to demonstrate or evaluate themselves – for example, to answer questions (such as hard science questions) which require specialised expertise to check (609*, 1203*). This provides a strong empirical check that the training protocol being used incentivises the right behaviour. Protocols under development for scalable oversight often enlist the AI system itself in helping the evaluator, for example by having it engage in a debate with itself over the correct answer (611*), and letting a human evaluator steer the model on the basis of that debate. Recent human and AI debate experiments show that this can improve the ability of human evaluators to determine the right answers to hard questions (615*, 1198, 1204*), and preliminary results show that this can translate into an improved training incentive (1197). However, positive results have only been shown on a simple reading comprehension task, with mixed results for other tasks such as mathematics problems (1198). These methods have not been used to train general-purpose AI

systems, but progress in this area is continuing, and scalable oversight experiments may at some point form a practical way of measuring how reliably training techniques incentivise the correct behaviour.

Some researchers are working toward ‘safe-by-design’ approaches which might be able to provide quantitative safety guarantees. Beyond ensuring that an AI’s training process encodes the incentive to be safe, it may be possible to design AI systems that quantitatively guarantee certain levels of safety (1176). These proposals often rely on a combination of three elements: first, a specification of desired and undesired outcomes (which in some cases could be a natural language description of desired and unacceptable behaviours), second, a ‘world model’ that includes capturing (approximate) cause and effect relationships and predicts the outcomes of possible actions the AI system could take, and third, a verifier that checks whether a given candidate action would lead to undesirable predicted outcomes. The goal of this process is to guarantee that dangerous actions are not taken. If the world model captures scientific knowledge, it will typically rely on ‘neuro-symbolic’ hybrids of general-purpose AI and classic techniques using formal mathematics. The advantage of mathematical guarantees and bounds is that they may provide safety assurances even outside of the domain in which the AI has been trained and tested, in contrast with spot checks and improvement through trial-and-error which are currently the standard for evaluating and training general-purpose AI models. This explicit model-based approach offers two additional advantages: firstly, because it uses formal logic and probability laws to analyse clearly defined knowledge components, its conclusions are more trustworthy, understandable, and verifiable than those of traditional AI systems. Secondly, it allows building non-agentic (non-autonomous) AI systems that can advance science and human knowledge while remaining easy to control, avoiding the potential risks that come with advanced highly agentic AI (see [2.2.3. Loss of control](#)). Currently, however, practically useful, provable guarantees of safety have yet to be demonstrated for general-purpose AI models and methods, and many open questions remain in order to achieve those objectives for large-scale AI systems (1205).

Maintaining the quality of human supervision and evaluation of AI behaviour

State-of-the-art training and evaluation techniques rely on feedback or demonstrations from humans and, as such, are constrained by human error and bias. Developers fine-tune state-of-the-art general-purpose AI systems using a large amount of human involvement. In practice, this involves techniques that leverage human-generated examples of desired actions (28) or human-generated feedback on examples from models (29, 30, 31*, 1182). This is done at scale, making it labour-intensive and expensive. However, human attention, comprehension, and trustworthiness are not perfect (1182), which limits the quality of the resulting general-purpose AI systems (1206, 1207*, 1208). Even slight imperfections in feedback from humans can be amplified when used to train highly capable systems, with potentially serious consequences (see for example [2.2.3. Loss of control](#)).

Improving the quality and quantity of human oversight can help to train more robust models. Some research has shown that using richer, more detailed forms of feedback from humans can provide better oversight for AI models, but at the cost of increased time and effort for data collection (1209*, 1210, 1211). To gather larger datasets, leveraging general-purpose AI systems to partially automate the feedback process can greatly increase the volume of data (33*, 256*). However, in practice, the amount of explicit human oversight used during fine-tuning is very small compared to the trillions of data points used in pre-training on internet data, and so human oversight may, therefore, be unable to fully remove harmful knowledge or capabilities from pre-training. Improving fine-tuning feedback data is likely to form only a part of the solution to cooperative robustness.

Improving the factuality of model outputs

The hallucination of falsehoods is a challenge, but it can be reduced. In AI, ‘hallucination’ refers to the propensity of general-purpose AI systems to output falsehoods and made-up content. For example, language models commonly hallucinate non-existent citations, biographies, and facts (101, 102*, 103, 104, 105), which could pose legal and ethical problems involving the spread of misinformation (1212). It is possible but challenging to reduce general-purpose AI systems’ tendency to hallucinate untrue outputs. Fine-tuning general-purpose AI models explicitly to make them more truthful – both in the accuracy of their answers and analysis of their own competence – is one approach to tackling this challenge (1213*). Additionally, allowing language models to access knowledge databases when they are asked to perform tasks helps to improve the reliability of their generations (838, 1214). Alternative approaches detect hallucinations and inform the user if the generated output is not to be trusted (1215), perform fine-grained checks on the individual claims made by a model (1216), or quantify the model’s confidence (1217). However, reducing hallucination remains a very active area of research.

Improving robustness against unexpected failures

Ensuring that general-purpose AI systems learn beneficial behaviours that translate from their training contexts to real-world, high-stakes deployment contexts is highly challenging. Sometimes, unfamiliar inputs that a general-purpose AI system encounters in deployment can cause unexpected failures (1218). Just as general-purpose AI systems are trained to optimise for imperfect proxy goals, the training context can also fail to adequately represent the real-world situations that systems will encounter after they are deployed. In such cases, general-purpose AI systems may still take harmful actions even if they are trained with correct human-provided feedback (616, 1032, 1033). For example, some researchers have found that chatbots are more likely to take harmful actions in languages that are underrepresented in their training data (1034). One way to mitigate these failures is with evaluation frameworks that test many combinations of deployment conditions, such as the Holistic Evaluation of Language Models framework (HELM (1150)), which enumerates and tests combinations of many different tasks, user profiles, and languages, among other features. Another is to develop methods by which models can estimate and communicate their uncertainty in rare cases to anticipate mistakes (1219*, 1220*). However, in

general it is likely impossible to enumerate all possible real-world situations for evaluation or to anticipate all potential mistakes.

Understanding a model’s internal computations might help researchers to investigate whether they have learned robust solutions. Methods exist to automatically identify features (i.e. mathematical patterns) inside a neural network model which correspond to human-interpretable concepts (1009, 1013*, 1221, 1222*), including specific people and places as well as abstract concepts and behaviours such as errors in code, nonconformity to certain political opinions, or descriptions of how to create drugs (1012). These features can serve as a guide to identifying dangerous or undesirable behaviours in a system’s training data or its outputs at a larger scale than would be practical with human review alone. Researchers have attempted to automate this review using an ‘automated interpretability agent’ that has access to interpretability tools. A preliminary study shows that this is possible on a small scale (1201), and there is no clear barrier to scaling up this kind of work.

There is recent progress on using understanding of a model’s internal workings to improve its behaviour, but this approach needs more work. Despite the difficulty of understanding models’ inner workings, some techniques can be used to guide specific edits to them. Compared to fine-tuning, these methods can sometimes be more compute- or data-efficient ways of modifying models’ functionality. Researchers have used a variety of methods for this, based on making changes to models’ internal parameters learned during training (1223, 1224, 1225, 1226, 1227), neurons (1221, 1228, 1229), or representations (1199, 1230, 1231, 1232, 1233). These techniques are imperfect (1023), generally limited to very specific kinds of behaviours (1227), and typically introduce unintended side effects on model behaviour (1234), but they remain an active area of research. It is unclear to what extent current methods offer a ‘useful and reliable’ way of understanding and engineering general-purpose AI models (1026*).

Adversarial robustness: preventing model misuse

Users of general-purpose AI systems can often bypass their safeguards with ‘jailbreaks’ that induce them to comply with harmful requests. Even if a system always behaves well under normal use, a motivated individual may still construct unusual inputs that are specifically designed to make a system fail or engage in undesired (e.g. harmful) behaviours (1054). Language models in particular are subject to general purpose ‘jailbreaks’ which can make them much more likely to comply with harmful requests. Examples of jailbreaking methods include: inducing an AI system to adopt the persona of someone who would say the harmful content (1053), priming it with examples of harmful answers (1235*), or making requests in a language that was scarce in the system’s training data (1236), which could increase models’ vulnerability in some low- and middle-income countries (LMICs) (see Table 3.3 for some example jailbreaks). While jailbreaks can be partially guarded against after their discovery, it is difficult to anticipate them during model development, and currently, it is generally easy to find new jailbreaks that work for state-of-the-art models. This being the case, it is unclear how widely jailbreaks are used to actually cause harmful behaviour by AI systems outside of a research setting.

Training models to detect and refuse harmful requests from adversaries

Adversarial training helps improve robustness in state-of-the-art AI systems, though only to a limited extent. ‘Adversarial training’ involves first constructing ‘attacks’ designed to make a model act undesirably, and second, training the system to handle these attacks appropriately. Attacks against AI systems can take many forms and can be either human- or algorithm-generated. Once an adversarial attack has been produced, training on these examples can proceed as usual. Adversarial training has become a commonly used technique to make models more robust to failures, and is used in the development of major general-purpose AI systems (4*, 48*, 147*, 1158*, 1163, 1241). However, it is not sufficient by itself, as adversarially trained systems are still generally vulnerable to attack, especially with multimodal inputs (e.g. with images). Moreover, the potential appropriateness or harmfulness from an AI system’s outputs cannot always be evaluated outside of the context in which it is used, which is not available during adversarial training (1242).

Making general-purpose AI systems more robust to unforeseen attacks is a challenging open problem, but there are potentially promising methods for minimising the relevant harms.

Adversarial training generally requires specific examples of failures (598*, 1243). These limitations have resulted in ongoing games of ‘cat and mouse’ in which some developers continually update models in response to newly discovered vulnerabilities. The process of searching for vulnerabilities and attempting to induce undesirable behaviour is known as ‘red-teaming’. A partial solution to models’ continued vulnerability is to simply produce and train on more adversarial examples. Automated methods for generating attacks can help scale up adversarial training (522, 904*, 1157, 1244). However, the exponentially large number of possible inputs for general-purpose AI systems makes it intractable to thoroughly search for all types of attacks. Interpretability methods might help here (907), and there has been preliminary progress on improving robustness through methods that operate on the model’s internal states (1076, 1195, 1200). Even if all attacks cannot be prevented beforehand, if they can be detected quickly at run-time then systems can be efficiently adapted to defend against them: in one study, a system had greater than 95% success defending against attacks after seeing only one example of the same kind of attack (1245). While research on these mitigations is preliminary, requiring live monitoring, response, and adversarial training mitigations on potentially dangerous AI systems is critical for decreasing the damage of AI misuse.

‘Machine unlearning’ methods aim to remove certain undesirable capabilities from general-purpose AI systems, but current techniques often suppress rather than fully remove such capabilities. For example, machine unlearning can remove certain capabilities that could aid malicious users in making explosives, bioweapons, chemical weapons, and cyberattacks (392). Unlearning as a way of negating the influence of undesirable training data was originally proposed as a way to protect privacy and copyright (821), discussed in [2.3.6. Risks of copyright infringement](#). Unlearning methods to remove hazardous capabilities (892, 1246) include methods based on fine-tuning (893*) and editing the inner workings of models (392). Ideally, unlearning should make a model unable to exhibit the unwanted behaviour even when subject to knowledge-extraction attacks, novel situations (e.g. requests in various languages), or small amounts of fine-tuning. However, current

unlearning methods often suppress harmful information without removing it robustly (1247). This creates challenges for governance, since models might appear to lack harmful capabilities when these are actually just hidden and can be reactivated. Current unlearning methods may also introduce unwanted side effects on desirable model knowledge (1247). It is unclear if unlearning a harmful skill could fully remove the model's ability to perform a harmful task by combining desirable skills and knowledge. Unlearning remains an area of active research.

3.4.2. Monitoring and intervention

KEY INFORMATION

- **Monitoring and intervention are complementary approaches for preventing AI system malfunctions and malicious use.** Monitors inspect system inputs and outputs, hardware state, model internals, and real-world impacts while systems are used, triggering interventions that block potentially harmful actions. Current tools can detect AI-generated content, track system behaviour, and identify concerning patterns across these monitoring targets. However, moderately skilled users can often circumvent these safeguards through various technical means.
- **Model interpretability and explanation methods can help monitor AI decisions but current methods can also produce misleading insights.** Technical approaches for explaining AI system outputs help developers and deployers scrutinise decision-making, though studies indicate that these methods can produce inaccurate or oversimplified explanations of complex model behaviour.
- **Multiple layers of monitoring and intervention create stronger protection against malfunctions and malicious use.** Combining technical monitoring and intervention capabilities with humans in the loop builds stronger safeguards, though these measures can introduce costs and delays.
- **In recent months, there has been progress in model interpretability and hardware-based monitoring measures.** Since the publication of the Interim Report (May 2024), model interpretability research has progressed to begin explaining model behaviours, and early work investigating privacy-preserving hardware-based monitoring has the potential to improve regulatory visibility into AI development.
- **Key challenges for policy makers centre on balancing safety measures against their practical costs.** While layered safety measures provide stronger protection, they also introduce operational delays, raise privacy concerns, and increase deployment costs. Policymakers therefore need to weigh safety requirements against these practical constraints, particularly given potential misalignment between safety measures and business incentives.

Key Definitions

- **Model:** A computer program, often based on machine learning, designed to process inputs and generate outputs. AI models can perform tasks such as prediction, classification, decision-making, or generation, forming the core of AI applications.
- **System:** An integrated setup that combines one or more AI models with other components, such as user interfaces or content filters, to produce an application that users can interact with.

- **Interpretability:** The degree to which humans can understand the inner workings of an AI model, including why it generated a particular output or decision. A model is highly interpretable if its mathematical processes can be translated into concepts that allow humans to trace the specific factors and logic that influenced the model's output.
- **AI-generated fake content:** Audio, text, or visual content, produced by generative AI, that depicts people or events in a way that differs from reality in a malicious or deceptive way, e.g. showing people doing things they did not do, saying things they did not say, changing the location of real events, or depicting events that did not happen.
- **Deepfake:** A type of AI-generated fake content, consisting of audio or visual content, that misrepresents real people as doing or saying something that they did not actually do or say.
- **Digital forensics:** The process of tracing the origin and spread of digital media.
- **Watermark:** A subtle, often imperceptible pattern embedded within AI-generated content (such as text, images, or audio) to indicate its artificial origin, verify its source, or detect potential misuse.
- **Defence in depth:** A strategy that includes layering multiple risk mitigation measures in cases where no single existing method can provide safety.
- **Human in the loop:** A requirement that humans must oversee and sign off on otherwise automated processes in critical areas.
- **AI agent:** A general-purpose AI which can make plans to achieve goals, adaptively perform tasks involving multiple steps and uncertain outcomes along the way, and interact with its environment – for example by creating files, taking actions on the web, or delegating tasks to other agents – with little to no human oversight.

Monitoring and intervention strategies are applied to AI systems – the complete deployment package that includes both the AI model and additional safety components – leaving the *model* unchanged. Unlike the strategies discussed in [3.4.1. Training more trustworthy models](#), monitoring and intervention methods are integrated at the system level and implemented as part of system deployment. This section discusses monitoring and intervention strategies that researchers and developers use for general-purpose AI systems (see Figure 3.2).

The main evidence gaps around monitoring and intervention include understanding how effective methods are, and how easy they are to circumvent. Monitoring and intervention techniques are, in many cases, easy, simple, and effective system-level safeguards in typical use cases. They offer an essential additional line of defence aside from the model-level techniques discussed in [3.4.1. Training more trustworthy models](#). From this perspective, there are few technical barriers to the widespread adoption of many techniques. However, scientists do not yet have a thorough quantitative understanding of their effectiveness in real-world settings and how easily monitoring methods can be coordinated across the AI supply chain. A key barrier toward highly effective monitoring and intervention techniques is understanding how vulnerable they are to being actively circumvented by malicious users.

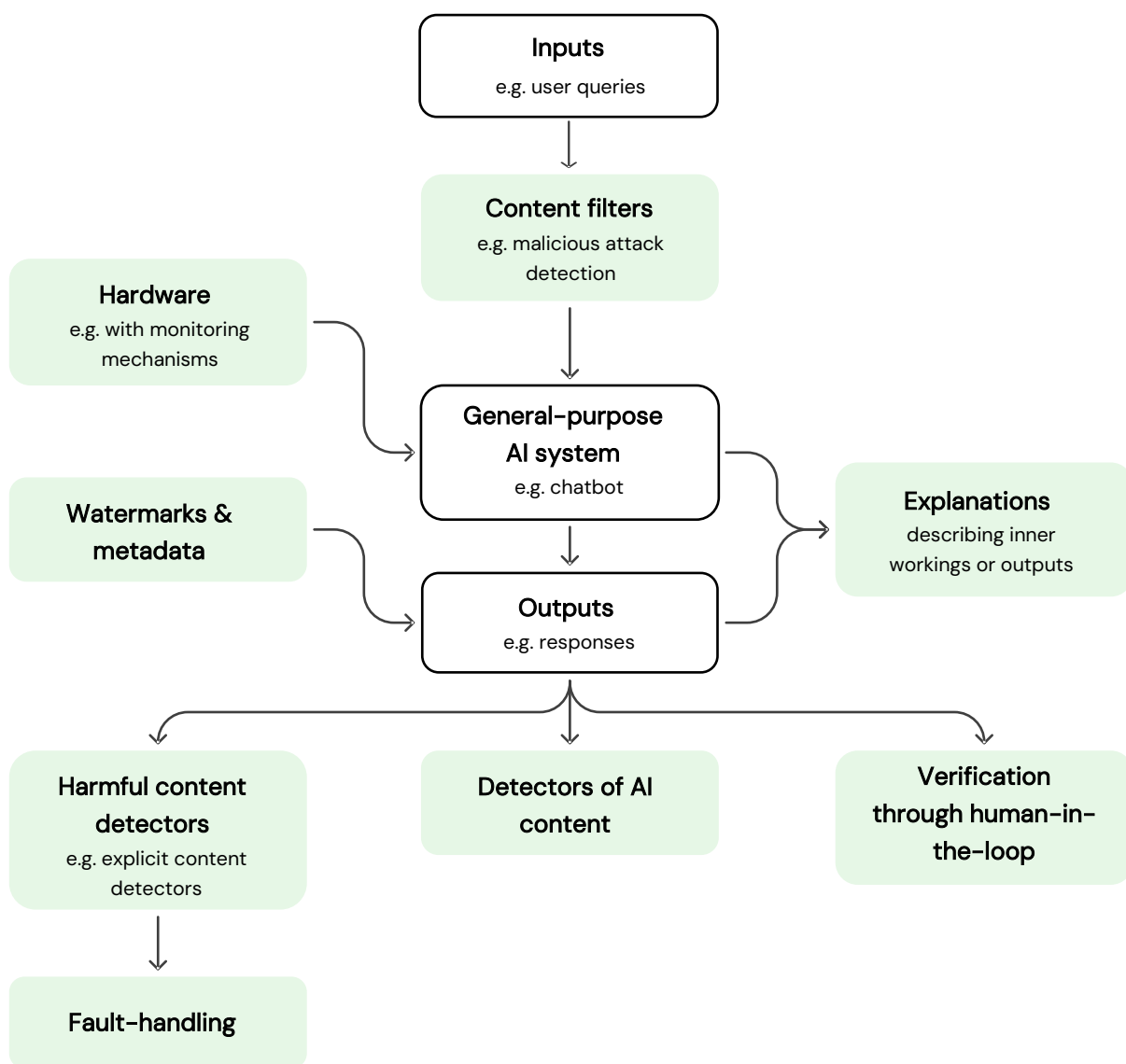


Figure 3.2: Monitoring and intervention techniques are system-level safeguards that can be applied to general-purpose AI system inputs, outputs, and models themselves in order to help researchers and developers monitor AI behaviour and, if necessary, intervene. Source: International AI Safety Report.

Detecting AI-generated content

Content generated by general-purpose AI systems – particularly ‘deepfakes’ – could have widespread harmful effects (1248, 1249, 1250) (see [2.1.1. Harm to individuals through fake content](#)). However, the ability to distinguish between genuine and AI-generated content can help to reduce the harmful use of generative models. For example, if web browsers were able to put reliability notices on content that was likely AI-generated, this would help to combat the spread of misinformation online. There are a variety of technical tools for the detection of AI-generated content. None are perfect, but together, they can be immensely helpful for digital forensics.

Unreliable but still useful techniques exist for detecting AI-generated content. Just as different humans have discernible artistic and writing styles, so do generative AI models. Some procedures have been developed to distinguish AI-generated text (332, 333, 337, 338, 1251, 1252, 1253) and images (1254, 1255) from human-generated content. Detection methods are typically based either on specialised classifiers or assessing how likely it is that a given example was generated by some general-purpose AI model. However, existing methods are limited and are prone to mistakes. A significant challenge is that general-purpose AI systems tend to memorise examples that appear in their training data. Because of this, common text snippets (e.g. famous historical documents) or images of common objects (e.g. famous art) are sometimes falsely classified as being AI-generated. As general-purpose AI-generated content becomes more realistic, it may be increasingly challenging to detect. Meanwhile, AI-text detectors tend to have inconsistent performance across world languages, posing challenges to linguistic equality (1256).

‘Watermarks’ – subtle but distinct motifs inserted into AI-generated data – make distinguishing AI-generated content easier, but they can be removed. Watermarks are features that are often designed to be difficult for a human to notice but easy for detection algorithms to identify. Watermarks typically take the form of imperceptible patterns inserted into image or video pixels (290, 291, 292, 293, 294*, 1257), imperceptible signals in audio (295, 296), or stylistic or word-choice biases in text (297, 1258, 1259, 1260, 1261). Watermarks can be used to detect AI-generated content with near perfect accuracy when they are not tampered with. As discussed in [2.1.1. Harm to individuals through fake content](#), they can be used to detect AI-generated fake content. They are an imperfect strategy for detecting AI-generated content (especially text) because they can be removed by simple modifications to data (298*, 299, 333, 1262). However, this does not mean that they are not useful. As an analogy, fingerprints are easy to avoid or remove, but they are still very useful in forensic science. Finally, there are concerns about privacy and potential misuse of watermarking technology, as it could be used to track and identify users (300).

Watermarks can also be used to indicate genuine, non-AI-generated content. Certifying the authenticity of data is part of ‘data provenance’. In contrast to inserting watermarks into general-purpose AI-generated content, another approach is to automatically insert watermarks into non-AI-generated content (1263). However, this will often require changes to the hardware and software of physical recording devices. These provenance methods would be very hard to tamper with at the device level. Some researchers are working towards common methods and standards for tracing the origin of media, including the use of encryption methods to prove authenticity which are difficult to counterfeit (e.g. CPPA (1264); AIMASC (1265)).

‘Metadata’ and system activity logs aid in digital forensics. ‘Digital forensics’ refers to the science of identifying and analysing digital evidence (1266, 1267, 1268, 1269, 1270). It is common for data to be saved along with ‘metadata’ that gives additional context about the data that is stored. This metadata is useful (and commonly used) for tracing the origin of data. For example, many mobile devices save image and audio files using the Exchangeable Image File Format (ExIF) standard (1271) which can store information about camera settings, time, location, and other details. Similar

metadata could be used to help track information about whether data was generated by a general-purpose AI system and, if so, other details about how it was done. For example, developers and deployers could attach identifiers to actions taken by an AI system (1272, 1273). Developers and deployers can also save ‘activity logs’ to track system behaviour, in order to improve monitoring over time (1272). Additionally, simply adding warning labels to AI-generated content can help to reduce the spread of misinformation. One study found that these labels improved humans’ deepfake detection from 10.7% to 21.6% (289). Metadata can typically be tampered with, but evidence suggests that the use of encrypted digital signatures can enable proof of authenticity in a way that is very hard to counterfeit (1274).

Beyond technical interventions, digital media literacy initiatives have also been proposed to combat AI-generated fake content (1275). Some studies have found that media literacy interventions can improve participants’ ability to detect fake content (1276, 1277, 1278, 1279). However, in general, evidence on the effects of digital media literacy interventions is mixed, partly due to large variations in study contexts and intervention designs (1279). See [2.1.1. Harm to individuals through fake content](#) for further discussion of fake content.

Detecting and defending against harmful content

Although there is no perfect safety measure, having multiple layers of protection and redundant safeguards increases confidence in safety (a strategy known as ‘defence in depth’). While the present section focuses on technical approaches, systems are not deployed in a vacuum. Embedding them in a sociotechnical system that seeks to maintain safety and performance is key to the ongoing process of identifying, studying, and defending against harm (also discussed in [3.1. Risk management overview](#)). This section discusses various complementary technical methods of detecting and defending against harmful behaviours from general-purpose AI systems.

Detecting anomalies and potentially harmful behaviours allows for precautions to be taken. Some methods have been developed that can help detect anomalous inputs or behaviours from AI systems (1280, 1281, 1282). For example, users sometimes trick language models into behaving harmfully by having them encode their responses in ciphred text (460, 1063*) that does not at all resemble normal text. It is also sometimes possible to detect a significant proportion of inputs (1243, 1283), internal states (1284, 1285, 1286*, 1287), or outputs (1287, 1288, 1289, 1290*, 1291) involved in harmful behaviours such as assisting with dangerous tasks. Once detected, risky examples can be sent to a fault-handling process or flagged for further investigation. For example, data flagged as harmful could be blocked by a filter or edited to remove harmful content.

Having a human in the loop allows for direct oversight and manual overrides but can be prohibitively costly. Humans in the loop are expensive compared to automated systems. However, when there is a high risk of a general-purpose AI system taking unacceptable actions, a human in the loop can be essential. Analogously, manual overrides are standard in cars that have autonomous driving modes (1292). Meanwhile, humans and general-purpose AI systems can

sometimes make decisions collaboratively. Instead of teaching general-purpose AI systems to act on behalf of a human, the human-AI cooperation paradigm aims to combine the skills and strengths of both general-purpose AI systems and humans (1293, 1294*, 1295, 1296, 1297, 1298, 1299). However, having a human in the loop is not practical in many situations, such as times when decision-making happens too quickly (such as chat applications with millions of users), or the human does not have sufficient domain knowledge, or human bias and error can exacerbate risks (1300). Humans in the loop of automated decision-making also tend to exhibit ‘automation bias’, meaning that they place a greater amount of trust in the AI system than intended (1301). In cases where a human in the loop is not practical, hybrid approaches involving a mix of human and automated monitoring and intervention are possible.

Secure operation protocols can be designed for general-purpose AI systems with potentially dangerous capabilities. General-purpose AI agents which can act autonomously and without limitation on the web or in the physical world pose elevated risks (see [3.2.1. Technical challenges for risk management and policymaking](#) and [2.2.3. Loss of control](#)). For general-purpose AI systems with potentially risky capabilities, limiting the ways in which they can directly influence the world makes it easier to oversee and manage them (1302, 1303). For example, if an agentic general-purpose AI system has an unconstrained ability to access a computer’s file systems and/or run custom code, it is safer to run that agent in an ad hoc computing environment than directly on the user’s computer (22*). However, these approaches can be hard to implement for applications in which a system must act directly in the world. In these cases, it is sometimes difficult even for humans to anticipate when an action might be harmful.

Explaining AI system actions

Some techniques can be used to help explain why deployed general-purpose AI systems act the way they do. Understanding why general-purpose AI systems act the way they do is useful for evaluating capabilities, diagnosing harms, and determining accountability if harm is caused (1304, 1305, 1306). While it can be useful, simply asking general-purpose AI language models for explanations of their decisions can also lead to misleading answers (97, 1307). To increase the reliability of model explanations, researchers are working on improved prompting and training strategies (1308*, 1309*, 1310, 1311). Meanwhile, other techniques for explaining general-purpose AI model actions (1312, 1313) can sometimes help with finding problems in models (1163). However, correctly explaining general-purpose AI model actions is a difficult problem because of their size and complexity. Some research is working toward developing techniques for helping humans interpret the computations of general-purpose AI systems (1010*, 1011*, 1012). Techniques to help explain model decisions are recognised as a useful part of the model evaluation toolbox (1314). However, these methods provide only a partial understanding. They depend on significant assumptions and more research is needed to demonstrate how useful they are in practice.

Monitoring and interventions with specialised hardware

Privacy-preserving monitoring mechanisms integrated into computing hardware are emerging as a more reliable and trustworthy alternative to software-based monitoring or self-reporting. Compute is central to the development and deployment of modern general-purpose AI systems, and the amount of compute used for training and inference is correlated with the capability of an AI system (see [1.3. Capabilities in coming years](#)). Research into privacy-preserving hardware mechanisms aims to enable policymakers to monitor and verify certain aspects of general-purpose AI systems during training and deployment, such as compute usage, without relying on reporting by AI developers. For example, research into these mechanisms argues that they make it technically feasible to verify usage details such as time and location of usage (1315, 1316), the types of models and processes being run (1317, 1318), or to provide proofs that a particular model was trained (1319, 1320). If feasible, these mechanisms can be applied to many governance issues, such as verifying adherence to international agreements even across borders (270). Some countries consider international agreements because of the competitive pressures between countries and its effect on incentives to thoroughly manage risks (see [3.2.2. Societal challenges for risk management and policy making](#) for an analysis of this dynamic). In this context, countries may resist monitoring and verification of agreements due to concerns about intellectual property and competitive advantages. Hardware-based verification mechanisms are sometimes considered to address this shortcoming since they could enable monitoring of key metrics while preserving the confidentiality of proprietary AI systems and training data. However, these applications are still in the stage of early research (270).

While much of the required functionality for hardware-based mechanisms exists on today's AI chips, hardware-based monitoring has not yet been proven at scale and could threaten user interests if implemented haphazardly. Some hardware-based mechanisms are widely deployed in contexts outside of AI, such as Apple's Secure Enclaves, which permit the manufacturer to restrict which applications are installed on their devices (1321*). Some leading AI chips, such as the H100 graphics processing unit (GPU), already have some of the necessary hardware in the form of Confidential Computing (1322*). Nonetheless, some hardware-based monitoring and verification mechanisms for AI could themselves be compromised by a well-resourced attacker, potentially leaking sensitive information (1323).