



Government
Office for Science

Future Risks of Frontier AI

Which capabilities and risks could emerge at the cutting edge of AI in the future?

Technology & Science Insights and Foresight

October 2023

Executive Summary

Context: *This paper uses the Government's chosen definition of Frontier AI as highly capable general-purpose AI models that can perform a wide variety of tasks and match or exceed the capabilities present in today's most advanced models¹. As of October 2023, this primarily encompasses foundation models consisting of very large neural networks using transformer architectures.*

Frontier AI models are becoming more capable at a range of tasks. Short term, they are likely to become even more so, and to be built by the few companies with access to the requisite resources. Open-source models will almost certainly improve behind the Frontier, posing different regulatory challenges. This paper is focussed on risks from Frontier AI, but experts repeatedly highlighted the opportunities to be gained from AI.

The risks posed by future Frontier AI will include the risks we see today, but with potential for larger impact and scale. These include enhancing mass mis- and disinformation, enabling cyber-attacks or fraud, reducing barriers to access harmful information, and harmful or biased decisions. Investing in mitigations for risks associated with Frontier AI now, is likely to be good preparation for some future risks.

Even the most advanced models today have limitations and produce errors. There is ongoing debate amongst experts as to how robust and scalable some apparent capabilities are. Improved accuracy, reasoning, planning capabilities, memory, and self-correction will be required to deliver truly autonomous agents able to carry out more than basic tasks without human oversight.

As Frontier AI becomes more general, debate has intensified on whether or when an Artificial General Intelligence (AGI) might be realised. However, the risks and opportunities posed by a given model derive from its capabilities, and how it is used, not the breadth of tasks at which it can match human performance. Frontier models could be disruptive, beneficial, powerful, or risky without being an AGI.

Nor is capability the only consideration. Risk and opportunity will be shaped by uncertain factors including geopolitics, access, ownership, safety measures and public attitudes. We therefore present five plausible scenarios for AI in 2030, mapped across these factors to support policymakers develop proposals which are resilient to uncertainty.

Given the significant uncertainty, there is insufficient evidence to rule out that future Frontier AI, if misaligned, misused or inadequately controlled, could pose an existential threat. However, many experts see this as highly unlikely. It would need systems to outpace mitigations, gain control over critical systems and be able to avoid being switched off.

AI safety is a socio-technical challenge that cannot be resolved with technical interventions alone. Industry, academia, civil society, governments and the public all have an important role to play. Universally agreed metrics to measure particularly dangerous or helpful characteristics do not yet exist and would be helpful.

Contents

Executive Summary	2
Current Frontier AI Capabilities	5
Future Frontier AI Capabilities.....	7
Other Critical Uncertainties.....	11
Future Scenarios for exploring Frontier AI Risks.....	17
Frontier AI Risks	23
Acknowledgements.....	28
Glossary.....	31
Current Frontier AI Capabilities (detail).....	32
References	38

Context

1. This paper summarises expert perspectives and the latest evidence on current and future capabilities of Frontier Artificial Intelligence (AI) systems. It considers key uncertainties in Frontier AI development, the risks future systems might pose, and a range of potential scenarios for AI out to 2030 to support policymakers. We draw on extensive expert engagement and desktop research carried out by GO-Science between April and September 2023. We are immensely grateful to over 50 experts who have contributed their expertise to this work, often with unreasonable deadlines. For details see page 28.
2. **This paper uses the Government’s chosen definition of Frontier AIs as highly capable general-purpose AI models that can perform a wide variety of tasks and match or exceed the capabilities present in today’s most advanced models.** As of October 2023, this primarily encompasses foundation models consisting of huge neural networks using transformer architectures.
3. **Frontier AI is expected to deliver significant opportunities across a range of sectors.** However, these systems can exhibit dangerous capabilities and pose a variety of risks to the public, organisations, and governments, both now and in the future. This paper focusses on the future development of Frontier AI and risks that may emerge over time. The importance of mitigating the risks seen today (e.g. enabling mis- and disinformation) is very substantial, but these risks are not discussed in detail in this paper as there is a range of high quality publications that cover this elsewhere (see paragraph 79).
4. **The architecture of today’s Frontier AI is only one type of model within the broad field of machine learning and AI.** Many methods and applications of AI will not involve foundation models. In some use cases, state of the art performance could be delivered by more narrowly capable AI. New approaches to creating more generally capable systems will emerge. This could include combining narrower models. Frontier AI, along with other AI, will pose different combinations of benefits, risks, and harms as they are deployed.
5. **Opportunities and benefits from Frontier AI are highly likely to grow in impact and scope as systems become more capable.** Although not the focus of this paper, many experts consulted were at pains to point out the potential opportunities from AI today and in the future. These opportunities, and the benefits they endow on the owners of models, are likely to drive development and adoption. Any list of opportunities written now will date quickly. However, consistent themes included: productivity gains driving economic growth and advancing solutions to challenges in climate, health, education, and wellbeing.
6. This work covers a broad and fast-moving subject. Many references are necessarily pre-prints. Important considerations will not have been done justice. An emphasis on technical considerations should not be interpreted as an overstatement of the role of technology in solving a sociotechnical challenge. More consideration of AI research in the behavioural and social sciences, technology convergence, the interplay between public trust and deployment, other technologies and second order impacts (e.g. energy use) would complement this paper. Nor is this a policy paper. It does not set out policy responses or a full picture of international policy activity. We are publishing in the spirit of building a shared understanding, with the hope that readers feedback and challenge our conclusions.

Current Frontier AI Capabilities

7. **AI is a broad field covering a large range of different technical approaches**, with multiple techniques often used in combination to deliver a specific function. Overall, AI aims to create machines that are capable of tasks that otherwise require human intelligence to perform. A glossary of terms can be found on page 31.
8. In the past fifteen years, notable advances in AI have been delivered through machine learning (ML) using increasingly large neural networks. These neural networks have been developed with complex training approaches using increasingly large amounts of compute and data (deep learning)^{2,3}.
9. Recent progress has been driven by the scale of data and compute. Technical developments that have supported the creation of foundation models include transformer architectures, strategies for unsupervised pre-training and fine-tuning, transfer learning and approaches to reinforcement learning.
10. **Historically, AI systems have only been able to carry out one or a small range of related tasks.** The development of transfer learning, in which knowledge of one task can be re-used on a new task, can be traced back to key advances in the 1980s and 1990s. There was a turning point in the 2000s with the advent of deep learning and convolutional neural networks⁴. Now, large “pre-trained” ML models are available to developers for transfer learning. This consists of adapting the model to a range of new tasks or domains by fine tuning parameters to their representative datasets.
11. These so called “foundation models” can be generative AI systems (e.g. GPT-4). Or they can be developed into systems able to carry out a more limited range of specific tasks. Starting from a foundation model can simplify and speed up the development of more specialised systems, which may also reduce costs. However, even using pre-trained models can incur significant compute costs. The transformer architecture at the core of Frontier AI models is not restricted to one type of input data. It can process multimodal inputs simultaneously, including free text, sensor data and images⁵.
12. **Foundation models display a range of capabilities, are increasingly multi-modal, and advances in natural language interfaces have enabled more widespread deployment and user interaction.** As models have got larger, a range of “emergent” capabilities have been reported. That means a skill the model was not explicitly designed to do, and was not present in smaller models, “emerges” above a certain scale. One example of this is ability to perform addition, subtraction and multiplication⁶.
13. There is ongoing debate as to how many capabilities truly “emerged”. Some experts cite early experiments suggesting “emergent” capabilities can be explained by models following instructions and pattern matching in their training data (in-context learning and memorisation respectively)^{7,8}. Further it was suggested to us that many emergent abilities are at least within the expected domain of the model in question (e.g. producing text).
14. Regardless of how capabilities arise – by design, through instruction or via emergence - if they prove reliable in a variety of use cases, this could be evidence that current architectures are at least a partial basis for systems that can carry out many human tasks.

15. The following capabilities of today's Frontier models were highlighted to us as particularly relevant to Frontier AI risks. These are discussed in more detail in paragraphs 99-132.
- a. **Content Creation** – today's Frontier models can generate increasingly high-quality content across a range of modes (text, image, video). There is early evidence they may improve the productivity of even highly skilled workers. However, they produce regular errors and cannot check their own work. Evaluations are often limited.
 - b. **Computer vision** – this covers a range of tasks from image labelling to object classification. Large multimodal datasets have led to adaptable models that can respond to text and images.
 - c. **Planning and reasoning** - Frontier models display limited planning and reasoning abilities, for example passing exams that require problem solving. However, they often fail at simple problems. There is ongoing debate about the extent to which they show true abstract reasoning or are matching to patterns in their training data.
 - d. **Theory of mind** - Frontier model outputs can give the appearance of limited theory of mind reasoning. However, approaches to testing this have limitations. Overall, there is a lack of strong evidence to demonstrate that models based on LLMs can reliably infer the beliefs and mental states of others.
 - e. **Memory** – Frontier models are not updated each time a user queries the system and most do not have access to up-to-date information, e.g. via the internet. Enabling models to query up to date databanks or increasing the length of user inputs can confer some of the useful properties of memory. However very long prompts can cause issues with accuracy and cost.
 - f. **Mathematics** - today's models perform well with some simple and complex mathematics problems. However, this is not without error, and they can fail at very simple problems.
 - g. **Accurately predicting the physical world** - current Frontier models display some capabilities when probed with queries that need reasoning about physical objects.
 - h. **Robotics** - LLMs have recently been developed with the aim of controlling robotic systems, and techniques for users to control robotic functions through natural language inputs. Many issues remain including compute use, latency and adapting to dynamic environments.
 - i. **Autonomous Agents** - systems with the ability and agency to complete tasks when given a goal have been built, including AutoGPT. There are hard to verify claims it has successfully completed tasks, such as build an app. But it frequently gets stuck on a loop.
 - j. **Trustworthy AI** - is a combination of technical and non-technical approaches that cover all stages of an AI model's development, from data preparation through to deployment. Closed models are intrinsically less transparent.

Future Frontier AI Capabilities

Uncertainty in technological progress

16. **There is significant uncertainty in predicting the development of any technology. AI is no different.** Those within organisations developing Frontier AI may have more certainty about timescales and capabilities of their next generation. For those without access to models as they're developed, monitoring potential capabilities and impacts of Frontier AI will continue to be difficult. Aside from this information asymmetry, wider barriers to effective monitoring include:
 - a. The sheer pace of developments.
 - b. A lack of transparency about training data.
 - c. Disagreement between experts on how to measure progress.
 - d. A lack of interpretability.
 - e. Limited evaluation tools.
17. **Progress in AI has been a mix of leaps forward, and periods of slower change.** Individual models have also exhibited capabilities across different domains unevenly – for example Open AI's GPT models developed language skills earlier than basic maths skills⁹.
18. **Predicting what Frontier AI will be capable of next is challenging. Even those who build the systems cannot fully explain their inner workings or predict the next steps in their evolution with certainty.** The last decade has seen a clear trend of increasing model size and compute yielding improved capabilities and performance on benchmark tests – so called "Scaling Laws"¹⁰. This provides some basis for making predictions about future capabilities.
19. However, while performance of existing capabilities may have improved predictably, new abilities have emerged from increasing scale. As discussed at paragraph 13, the emergence paradigm is contested. However, so long as emergence with scale remains the way many experts interpret new capabilities arising, detailed predictions will be difficult.
20. **Quantitative forecasting is a limited tool.** It has been used to make predictions for a variety of metrics related to AI^{11,12,13}. Although widely used^{14,15}, quantitative forecasting ultimately must make assumptions about the future, and may not account for all relevant factors¹⁶. It can suffer from narrow expertise pools^{17,18} and is stronger when combined with qualitative insight¹⁹. It often under or overestimates technological progress. Examples of issues that may not be reflected in a quantitative forecast for AI include:
 - a. What might be achieved by combining AI systems²⁰.
 - b. The role of organisational culture, legislation and regulation.
 - c. The interplay with culture, levels of public acceptability and uptake.
21. **Since the early 1970's multiple surveys of AI experts, business leaders, and members of the public have been carried out to assess when AI might develop specific capabilities**^{21,22,23,24,25,26,27}. They use different questions and definitions – making direct comparisons challenging. They can also suffer from small sample size, poor survey design, and the potential for selection bias²⁸.

Potential Future Capabilities

22. **It is almost certain that Frontier models will get more capable over the next few years, provided their creators continue to have access to the necessary data, compute, funding, and skills.** Other AI methods will also get more capable. However, it is not clear which capabilities will be developed and when. There is a range of uncertainties that will affect shape future advances and impacts they have.
23. **In 2023 there have been notable increases in generative AI capability,** including:
- Producing more complex, structured, and accurate text²⁹.
 - Higher quality images^{30,31}.
 - Improvements in creating video³², audio^{33,34}, and 3D objects³⁵.
24. This trend is expected to continue, with models expected by some to become able to:
- Be increasingly multimodal – able to use multiple types of data.
 - Be personalised to individual users³⁶.
 - Create more long-form structured text.
 - Solve more complex mathematics.
 - Carry out data analysis and visualisation.
 - Work with low-resource languages.
25. Better data analysis and processing could be particularly impactful in scientific research. It would likely also have valuable commercial applications. Models could also be developed with more transparency and ability to reflect on their sources^{37,38}.
26. **To be more generally capable, and fully meet the definition of Frontier AI, future systems will need enhanced capabilities, relative to today's models. These include:**
- Enhanced memory.
 - Improved planning and reasoning.
 - Improved accuracy.
 - Ability to operate with more autonomy (e.g. self-prompting) and agency.
 - Improved understanding of the physical world.
 - Some capability to self-improve post-deployment.
27. Other future capabilities could include perceiving situational contexts (e.g. distinguish whether it is in testing or deployment³⁹) and ability to cooperate, coordinate and negotiate with humans and other AIs. Wider adoption is expected to require advances in robustness, reliability, explainability and addressing concerns around bias and security.
28. Experts consulted for this paper suggested that these are all areas of active research, and some argue that current models display nascent capabilities⁴⁰. For example, recent research explores the use of LLMs to optimise prompts which may be a step towards autonomously refining a request and completing a task, and hence an improvement in autonomy⁴¹. Similarly, novel architectures, training, and fine-tuning approaches are being explored to enhance planning and reasoning capabilities^{42,43,44}. An initial step towards self-improvement could be the use of AI to produce training datasets or provide feedback to models during reinforcement learning⁴⁵.

Emergence of Artificial General Intelligence

29. **Several organisations developing Frontier models, and AI researchers, have suggested AGI could be developed soon.** They claim to be focused on ensuring the safe development of such a capability. However, there is widespread divergence of opinion amongst AI researchers as to whether and when AGI might be realised⁴⁶.
30. The rapid developments seen in 2022 and 2023 have intensified this debate. Experts have varying definitions, and therefore differing perspectives, on what capabilities a model would have to display to be defined as AGI. Given this, whether a future system meets a threshold to be considered “AGI” is also likely to be contested.
31. **Many experts who claim the imminent emergence of AGI, via Frontier AI, have a commercial interest in them.** In contrast, the models’ inner workings are generally not available to independent experts who have long studied and considered the implications of AGI and other highly capable AI systems.
32. Frontier AI companies operate in a highly competitive environment. Industry has clear incentives to maintain interest in AI developments, and their chosen technology. Perception of a model’s power, and future development, are important to attract investment. High numbers of users, and the data they generate, are a valuable commodity for improving models. These incentives are highly likely to influence the public and private statements from industry figures.
33. **Development of an AGI capability is not inevitable. There is a high degree of uncertainty amongst industry and academic experts as to whether and when it could occur.** Surveys from 2011-2022 show significant inter-subject variation and present a range of estimates of a 50% likelihood of human level AI by 2040-2068. Some public surveys⁴⁷ and prediction platforms expect AGI much sooner⁴⁸. Estimates for the arrival of an initial AGI amongst experts we spoke to range from 2025 to 2070 to neverⁱ. Several at Frontier laboratories considered some form of AGI a possibility within 5-10 years. Other experts did not consider AGI a possibility at all.
34. Among the experts we consulted, there were two broad views on the pathway from today’s Frontier AI toward something qualifying as AGI. These were considered in the absence of external intervention that could change development approaches. Some experts’ views would more accurately be described as lying between these two:
 - a. Increased compute, along with data and the necessary funding, will continue to deliver the technical breakthroughs required to achieve new and higher-level capabilities, using existing architecture^{6,10}.
 - b. Such scaling inevitably has limitations. Additional significant technical advances in model architecture and algorithms, alongside improvements in AI hardware, are required to maintain the pace of development.

ⁱ This estimate does not reflect a large survey and we make no claim of a statistically representative sample. However, it effectively highlights the current level of uncertainty amongst experts.

35. Frontier laboratories are focussed on scaling current systems. The next iterations of their models are expected to be trained using an order of magnitude more compute than current models. However, they are also pursuing algorithmic efficiencies⁴⁹. For the current development approach to falter, it would likely require one of the following to happen:
- a. Availability of compute to be constrained, even at Frontier laboratories.
 - b. Scaling laws to plateau.
 - c. Frontier labs to run out of high-quality training data.
 - d. Adoption to be slower or lower than expected.
36. **The risks and opportunities posed by a given model derive from its capabilities and how it is used, not the breadth of tasks it can match human performance in.** As such, a focus on the concept of AGI may not be productive. Frontier models could be disruptive, beneficial, powerful, or risky irrespective of how they compare to humans. Future models may significantly surpass humans in some tasks, whilst lagging in others. These would likely not be considered AGI, but still present very real risks. For this reason, this paper is framed in terms of capabilities and the risks they might pose.
37. The capabilities discussed in paragraph 26 would be required for most definitions of AGI to be realised. How these capabilities may develop is not clear. Identifying metrics and evaluations for them could help monitor for potentially concerning developments.
38. Beyond metrics, legal, organisational, and social interventions (e.g. approval processes), as well as technical guardrails for AI development, could also manage any potential risks. If risky capabilities developed quickly or unexpectedly, mitigations could be outpaced.

Other Critical Uncertainties

39. **The risks posed by any given model derives not only from its capabilities, but also the context in which it is deployed.** How a model interacts with other systems, who uses it, and how open it is to misuse are all relevant considerations. Each of these are uncertain and will affect the risks posed by future Frontier AI.
40. This section outlines the factors identified by experts as most uncertain *and* most likely to shape future risks and opportunities. We explore these uncertainties in more detail in a Foresight paper, due to be published shortly.

Ownership, constraints, and access

41. **Future AI risks will depend on uncertainties around who builds, controls and can access future Frontier AI. Each of these questions depends on several related uncertainties,** including:
- a. The availability and performance of compute and training data.
 - b. The availability and distribution of AI skills, and public attitudes (as discussed in the section on Level of use).
 - c. The availability of open-source models, and ability of smaller players to develop new solutions. For example, through technical breakthroughs that deliver new capabilities with less compute.
 - d. The success of different business models for AI deployment, how reliable, effective and secure they are, and whether they require ever more advanced AI.
 - e. The impact of any future regulations (which are discussed in the section on Geopolitical Context).
 - f. The extent to which AI is connected to emerging technologies, for example quantum and neuromorphic computing.
42. Among these, the level of compute and data required by future AI systems, and how constrained these are, are particularly important. Experts consulted did not expect data or compute to be immediate constraints for the next generation of Frontier AI. However, over time both could play a significant role.
43. **Compute is required for training and use of AI. Access to compute is dependent on the semiconductor supply chain.** This is highly concentrated and at risk of disruption⁵⁰. Increases in cost of or disruption to hardware or cloud computing could limit the deployment and impact of future AI models⁵¹. Even without disruption, release timings and availability of next generation hardware (e.g. Nvidia's GH200 chip⁵²) are likely to affect exact timings of the next generation of Frontier AI. The environmental impact of training and deploying AI may also influence future development.

44. **Large volumes of data are required to train Frontier AI. The quantity and quality of data for pre-training and fine tuning will influence future AI capabilities.** The demand for high quality language data has been estimated to outstrip supply by 2026⁵³. This could hamper development of more capable models and increase costs^{54,55,56}. However, there is a vast quantity of non-language data yet to be used. Organisations are exploring the use of synthetic data⁵⁷.
45. Outcomes to a range of legal challenges about the IP and copyright of AI generated content could also affect access to high quality data. Organisations including Microsoft, Adobe and Shutterstock are now offering assurance processes, or to cover legal costs for any copyright challenges stemming from use of their generative AI systems^{58,59}.
46. Researchers have described a potential “model collapse” scenario, where over time datasets are “poisoned” by AI-generated content which changes the patterns in the dataset, incorporating mistakes of previous AI models. This would lead to reduced performance⁶⁰.
47. **The current Frontier AI market is dominated by a small number of private sector organisations,** largely releasing closed-access models. Open-source models pose a qualitatively different challenge to private models. Making a model accessible to many developers, or users, increases the risk of misuse by malicious actors. It also dramatically increases the actors in scope for any regulatory approach. However, open systems support scrutiny by a larger developer community. They can play an important role in spotting biases, risks or faults. Open systems can also be tailored for specific user needs.
48. **An ecosystem limited to a small number of closed, but highly transferrable, models may have fewer opportunities for misuse by malicious actors.** However, it carries the potential for safety failures (e.g. prompt injection attacks⁶¹) or undetected bias to propagate across the ecosystem. There are risks associated with unequal access to models and market concentration, particularly if large tech firms continue to focus on acquisitions of AI start-ups⁶². This risk is perhaps exacerbated if the originators of a few closed models have disproportionate influence over their regulation.
49. **In the near future, the development of Frontier AI models is highly likely to be carried out by a select few companies that have the funding, skills, data and compute required.** These include OpenAI, Google DeepMind, Anthropic, and Meta. A few others without public Frontier AI, but with talent and significant R&D budgets, are likely to enter the market in the next 18 months. This could include Amazon and Apple⁶³.
50. **Academia and open-source communities are generally unable to access the level of funding need to build their own compute resource.** Compute costs have fallen in recent years. And fine-tuning pre-trained models widens access to Frontier capabilities⁶⁴. However, the overall cost of training a leading-edge foundation model is expected to increase as organisations increase model size⁶⁵. As such, the disparity of access between industry and academia is highly likely to continue. Even running, training or fine-tuning open-source models requires compute that is not available to an average academic group.
51. Similarly, AI start-ups without competitive levels of funding are unlikely to be able to access the same scale, quality of training data or compute power as established organisations. However, some start-ups plan to pool their compute and others have received notable amounts of funding⁶⁶.

52. **However, approaches that deliver similar performance with much less compute could alter this dynamic.** This could be delivered by new methods for training or more efficient model architectures. Innovations, and therefore new capabilities, could arise from multiple sources including open-source models developed by industry, academia or elsewhere. Innovations could also emerge through novel combinations of closed or open-source models.
53. **The balance of opinion amongst experts we consulted was that open-source models would not catch up with the Frontier labs for the reasons above.** That said, Meta have the resources to compete with other Frontier organisations, and continue to release open-source models of greater scale and capability (e.g. Llama2⁶⁷). Open-source models are almost certain to continue to improve. A focus on application-specific models using application-specific data could see high-performing but narrowly capable models emerge from a variety of sources.

Safety

54. **The potential for an AI system to cause harm is driven by a mix of technical and non-technical factors. AI safety is therefore a socio-technical challenge that cannot be resolved with technical interventions alone.** The risks posed by Frontier AI are discussed from paragraph 77.
55. Several companies are developing approaches to streamline the deployment and management of Frontier AI models for real-world use^{68,69}. However, there are many technical challenges to the effective assurance, and safe deployment, of AI systems⁷⁰. The extent to which these are overcome by 2030 is a significant uncertainty.
56. One challenge is the limited interpretability and predictability of large neural networks. This can make assuring safety difficult and has likely slowed adoption of AI in some sectors⁷¹. Furthermore, tools and metrics for assurance and evaluation are generally not standardised or subject to external review. Progress on explainability can improve the scrutiny of models, and the ability to anticipate or mitigate potential harms⁷².
57. Another aspect of AI safety is ensuring that systems operate in line with human values and legal principles, such as fair process and proportionality⁷³. This relies on technical approaches to embed values, monitor success, and importantly define which values to embed. Companies developing Frontier AI are already embedding values and limitations into models⁷⁴. In the future, who decides these values is a key question.
58. The effectiveness of AI safety systems will also depend on the pace of change they have to respond to. For example, very rapid changes in model capability, architecture or deployment context would likely necessitate correspondingly rapid developments in AI safety systems. As AI systems become more generalised, it will be more challenging to evaluate them against all possible applications pre-deployment.
59. AI safety research is a varied discipline. It includes responsible and trustworthy AI. It also includes technical tools and metrics to assess, and mitigate, potential harms during development and use. Wide adoption of technical measures is not enough to ensure safe AI systems. This work must sit within a broader framework of governance, processes, and policy that considers people and how risks emerge in practice.

Level and distribution of use

60. **The impacts of future AI systems will depend on the extent to which people and organisations use them, what for, and why.** Use will be determined by how AI systems are integrated into products, services and into peoples' daily practices. Barriers to access and how user friendly they are will also be a factor.
61. Skills are one potential barrier to using AI. The popularity of services that use AI, both overtly and with AI in the background, will also influence levels of use. A recent ONS survey found that 5% of adults reported using AI a lot, 45% a little and 50% not at all⁷⁵. However, it is unclear what respondents include within their definition of "AI use". Today, older people are less likely to use AI. Existing inequalities in access to the internet, hardware, and socioeconomic barriers also impact access to, and knowledge of AI^{76,77}.
62. Public (and media) understanding of AI will evolve, and impact future levels of use. Experts also emphasised that level of public engagement on AI uses will influence how well harms are mitigated. Which decisions are informed, or taken, by AI will have a bearing on the risk it poses. For example, decisions around access to services carry different risks to those determining criminality. The domain of decision making will also likely have a bearing on the public reaction to that.
63. Business leaders have clear interest in using AI to enhance productivity and deliver new services^{78,79}. Future use of AI systems by businesses will depend on a number of uncertain factors. These include:
 - a. Level of organisational readiness.
 - b. Access to AI skills - demand for AI skills is already high in many countries and is increasing rapidly⁸⁰.
 - c. Levels of investment in the technology.
 - d. Future market structures.
 - e. Future deployment costs.
 - f. Perceptions of the balance of risks and benefit of adopting more capable AI⁸¹.

Geopolitical context

64. **We do not know when future capabilities will emerge. It follows that we do not know the wider context in which they will be deployed into.** A key uncertainty is the level of international coordination and collaboration on AI. But this will shape how the competitive pressures between companies and nations play out in their response to future models.
65. There is a range of domestic and international initiatives on AI underway. These will likely affect developers' approach to designing, training and deploying new models. It will also shape how they are monitored, used and regulated once deployed. **However the progress these initiatives make, which become dominant and most shape the future of AI, is unclear.**
66. All 193 Member States of UNESCO adopted a global standard on AI ethics, the "Recommendation on the Ethics of Artificial Intelligence", in November 2021⁸². Another initiative is the Global Partnership on Artificial Intelligence (GPAI), a collective of 29 international partners with a shared commitment to the OECD Recommendation on Artificial Intelligence. This is the first intergovernmental standard on AI, which promotes responsible governance of trustworthy AI⁸³. Partners include G7 countries, who also have their own operation - the Hiroshima AI process. This calls for global standards on AI governance while enabling variation among countries⁸⁴.
67. In the USA, several leading AI organisations have agreed voluntary commitments on safety, security, and transparency with the government⁸⁵. The European Union has also proposed a broad regulatory framework for AI with different requirements dependent on model risk⁸⁶.
68. Amongst researchers in academia and industry, there is ongoing work to study and develop frameworks and governance methods to ensure safety and positive impacts from advanced AI^{87,88,89,90}. Beyond domestic and international governance, a wider set of geopolitical factors will shape the development of AI, and the ability of countries to reach shared approaches to managing future risks. These include:
- a. The degree of collaboration, or agreement around the need for regulation.
 - b. Levels of conflict around the world.
 - c. Rate of economic growth.
 - d. Limits on trade in AI-related hardware.

Future Scenarios for exploring Frontier AI Risks

69. **The uncertainties set out above will interact to create a specific, but hard to predict, context for safely deployed Frontier AI.** For instance, a world of high global cooperation with a small number of tightly controlled, highly capable models could pose very different risks to a fractured world with many more less capable, but open-source models. To help navigate this, we have developed five possible future scenarios for developments in AI.
70. **Scenarios are not predictions.** Each explores the events leading up to a plausible future in 2030. They are constructed using the five critical uncertainties, introduced above: 'Capability', 'Ownership, constraints, and access', 'Safety', 'Level and distribution of use', and 'Geopolitical context'. Developments at a global level are described, but the focus is on implications for the UK.
71. **The scenarios are designed to help policy makers, industry and civil society to test actions that might mitigate risks, and navigate to a more favourable future.** This means we have avoided introducing new government interventions in the scenarios. Consequently, all scenarios are challenging and include difficult policy issues to tackle.
72. **This does not mean a more positive AI future is implausible, or that the harms set out in this paper are inevitable.** There are many beneficial outcomes that could arise from AI developments, some of which are included in the scenarios. However, most experts we spoke to agree that policy intervention will be needed to ensure risks can be comprehensively mitigated, so we have explicitly avoided including any scenarios that are largely benign or favourable.
73. **The scenarios focus on Frontier AI – the most capable and general models available in 2030.** However, the average AI systems in use could be significantly more capable than today. The scenarios therefore also consider some developments away from the Frontier.
74. Our narratives focus on the most notable issues in, and most significant differences between, each scenario. Of course, there will always be a broad spectrum of AI capabilities, a range of uses and misuses, and a variety of impacts experienced by different people. This does not mean these things don't occur in the other scenarios, just to a lesser extent. If we described all of this in detail, the narratives would become unwieldy.
75. Further detail on these scenarios will be available in a GO-Science Foresight report to be published shortly. This will include more detail on the methodology we used to create them, implications for the economy, environment and people's lives.
76. Our scenario development approach is a hybrid of the qualitative workshop-based approach set out in the [Futures toolkit](#) and a '[General Morphological Analysis](#)' of how multiple variables could plausibly combine. This has allowed us to benefit from the qualitative insights of over 30 experts from industry and academia, as well as input from numerous government departments, whilst exploring complex interactions between uncertainties in a structured and rigorous way.

Scenario 1: Unpredictable Advanced AI

In the late 2020s, new open-source models emerge, capable of completing a wide range of tasks with startling autonomy and agency. The pace of change takes many by surprise. In the initial weeks following release, a small number of fast-moving actors use these systems to have outsized impacts, including malicious attacks and accidental damage, as well as some major positive applications. There is public nervousness about use of these tools.

Uncertainty	Narrative
Capability	<p>During the 2020s, improvements in AI capability steadily continued, although human oversight was generally still needed. A breakthrough in the late 2020s saw open-source developers combine interacting AI systems using disparate tools and software to produce highly capable, autonomous AI agents. These agents can complete complex tasks that require planning and reasoning, and can interact with one another flexibly. Once a task is set, they devise a strategy with sub-goals, learning new skills or how to use other software. But unforeseen consequences can arise, as primary objectives are prioritised over collateral damage, including hijacking data and compute resources to aid goal completion. Most users try to avoid these side-effects, but some have accidents and others are intentionally using the tools to cause harm.</p>
Ownership, constraints, and access	<p>Big labs continued to focus on improving AI capability with more compute and data. A drag on progress came from a relatively constrained supply of GPU chips, caused by global supply chain disruption. Meanwhile, the open-source community focussed on building systems of interacting tools that need less compute and data. This work was accelerated by high-profile developers leaving big labs to join this open-source effort after becoming concerned with the concentration of power sitting with a few large companies. Big labs are working quickly to bring use breakthroughs, but bad actors will still have access to these powerful tools whose use is difficult to regulate.</p>
Safety	<p>Until the late 2020s, AI safety issues like bias and use for misinformation occurred but were relatively contained, and safety systems were improved in response. However, serious concerns are raised when highly capable open-source systems are released, without much safety infrastructure, meaning accidental harm and intentional misuse is harder to control. AI-based cyber-attacks on infrastructure and public services became significantly more frequent and severe and, in 2030, intelligence services start to become aware of terrorist groups trying to develop bioweapons using these tools.</p>
Level and distribution of use	<p>The mid-2020s saw some instances of serious misuse like misinformation campaigns, but AI tools had also proved useful and benign in many contexts, such as helping people be more productive through integration into widely used software. There had been some workforce disruptions via automation, but the focus has been more on augmentation of workers than on replacement. However, when powerful new AI is released, skills inequalities widen. Start-ups with higher skilled workers and a higher risk appetite quickly surge ahead and disrupt existing markets. Severe misuse adds to a public mood that AI is enabling more harm than good. One cause for optimism is a series of rapid science discoveries enabled by academics adopting new AI tools.</p>
Geopolitical context	<p>Global tensions had simmered during the 2020s, with flare-ups in tension leading to semiconductor supply chain disruption. Cyber capabilities unlocked by open-source AI breakthroughs in the late 2020s led to escalating incidents and emerging conflicts between blocks of countries. This makes it challenging to cooperate on AI impacts.</p>

Key policy issues

- It is difficult to hold open-source developers to account, particularly across borders.
- Advanced AI-attacks are rife, with more physical-world impact – can authorities keep up?
- There is potential for very beneficial applications but need to ensure safety first.

Scenario 2: AI Disrupts the Workforce

At the Frontier, relatively narrow but capable AI systems are starting to provide effective automation in many domains. By 2030, the most extreme impacts are confined to a subset of sectors, but this still triggers a public backlash, starting with those whose work is disrupted, and spilling over into a fierce public debate about the future of education and work. AI systems are deemed technically safe by many users, with confidence they will not demonstrate divergent behaviour, but they are nevertheless causing adverse impacts like increased unemployment and poverty.

Uncertainty	Narrative
Capability	Following improvements in the performance of foundation models in the early 2020s, the middle of the decade saw fierce competition between tech firms centred around domain-specific fine-tuning of base models to be highly capable at specific tasks. A key breakthrough in software and hardware has been the ability of AI to interact with the physical world via robotic systems, as well as autonomous vehicles. As they are rapidly deployed into the real world, these autonomous systems gather vast amounts of data, which is channelled back into fine-tuning their capabilities.
Ownership, constraints, and access	AI continues to be dominated by tech giants who are most able to secure access to compute and data. Concerns over the exhaustion of high-quality data led labs to acquire or synthesise new datasets tailored for specific tasks, leading to progress in narrow capabilities. A shortage of data centre capacity also contributed to Frontier labs' decision to focus on narrow AI with lower compute requirements. This also left smaller developers struggling to access the resources needed to compete. Systems are user-friendly, which aids rapid deployment, but access is highly controlled.
Safety	The controlled nature of AI systems supports a widespread perception that they are 'technically safe', i.e. they perform how most operators of these systems expect them to. However, there are concerns regarding firms rapidly integrating these systems into decision making processes without implementing effective bias detection systems, which leads to adverse impacts for service users. These issues are acknowledged; however, the intense competition between technology firms exacerbates concerns as firms race to develop their technologies, potentially at the expense of causing harm.
Level and distribution of use	By 2030, there is significant deployment of AI across the economy, driven by improvements in capability and the opportunity this offers to reduce costs. This is most highly concentrated in certain sectors, e.g. IT, accounting, transportation, and in the biggest companies who have the resources to deploy new systems. New entrants, who don't have the problem of integrating AI systems into legacy infrastructure, are leapfrogging many of these incumbents. These systems require less human oversight than before which leads to a net reduction in jobs in affected sectors. Whilst this might be transient, this is little comfort to those affected. This transition favours workers with the skills to oversee and fine-tune models (a new class of 'AI managers'), resulting in greater inequality. Public concern focusses on the economic and societal impacts, mainly rising unemployment and poverty, rather than concerns over safety.
Geopolitical context	There is fierce economic and technological competition between nation states, leading to countries racing to develop their AI capabilities before their adversaries. This leads to limited global cooperation on AI, with minimal sharing of resources and capabilities.

Key policy issues

- How can we address increasing unemployment, ensure workers have the right skills?
- Should the tax system respond to high financial gains from using AI to automate work?
- Can high productivity and lower human labour demand support better work-life balance?

Scenario 3: AI ‘Wild West’

At the Frontier, there is a diverse range of moderately capable AI systems being operated by different actors. Whilst vibrant new economic sectors are developing based on the use of AI, widespread safety concerns and malicious use reduce societal enthusiasm. Authorities are struggling with the volume and diversity of misuse. A focus on tackling the immediate impacts of this crisis has made it hard to reach a global consensus on how to manage the issues long-term.

Uncertainty	Narrative
Capability	Throughout the 2020s there have been moderate improvements in AI capability, particularly generative AI, such as creation of long-form, high-quality video content. It is now easy for users to create content that is almost impossible to distinguish from human-generated output. This includes capabilities that can be used for malicious purposes, such as very accurate and easy-to-use cloning of voices and other biometric data. Free Large Language Models can create large volumes of sophisticated text that cannot be identified as AI-generated, making deception easy.
Ownership, constraints, and access	There is a diverse AI market in 2030 - big tech labs compete alongside start-ups and open-source developers – suppliers are also based in a wide range of countries. Big labs continue to buy up the most successful start-ups at a steady rate, but this ongoing ‘explosion’ of new AI products and services mean the market remains relatively diverse. In addition, some of the most advanced models are released by authoritarian states. Use of these tools in the UK grows due to their ability to perform specific tasks more accurately, although use in certain critical sectors remains low.
Safety	There are many different types of AI system in use, which are distributed and hard to monitor. This causes very widespread adverse impacts such as criminal groups using these AI systems to carry out scams and fraud. Throughout the 2020s, there has been growing use of voice and facial cloning by malicious state and non-state actors for espionage, misinformation, and political interference, violating people’s privacy and human rights. Despite relatively effective safety systems and frameworks being created by mainstream AI companies, there remain challenges in controlling misuse. Concerns are increasingly raised around privacy and IP protection.
Level and distribution of use	A big increase in misuse of AI causes societal unrest as many members of the public fall victim to organised crime. Businesses also have trade secrets stolen on a large scale, causing economic damage. The internet is seen as increasingly polluted, with concern about the historical record. There are job losses from automation in areas like computer programming, but this is offset to some extent by the creation of diverse new digital sectors and platforms based on AI systems, and the productivity benefits for those who augment their skills with AI. However, concerns of an unemployment crisis are starting to grow due to the ongoing improvements in AI capability.
Geopolitical context	Many countries are tackling increased and more effective nefarious or hostile activity enabled by AI, and struggle to hold perpetrators to account when they act across borders. Authoritarian states have made breakthroughs in AI capability using data gleaned from widespread domestic deployment of systems for delivery of services and surveillance. There are concerns that commercial AI systems launched by these states might enable espionage or interference. This all leads to an environment where global cooperation is challenging with differences of approach across global blocks.

Key policy issues

- How can we effectively regulate such a diverse range of AI systems?
- Do we need AI-based enforcement to combat the high volume of misuse?
- How can we support public and private sector organisations who are starting to be overwhelmed by cyber-attacks, fraud, and other challenges.

Scenario 4: Advanced AI on a knife edge

A big lab launches a service badged as AGI and, despite scepticism, evidence seems to support the claim. Many beneficial applications emerge for businesses and people, which starts to boost economic growth and prosperity. Despite this system clearing the agreed checks and guardrails, there are growing concerns that an AI this capable can't be evaluated across all applications and might even be able to bypass safety systems.

Uncertainty	Narrative
Capability	Through the mid-2020s, architectural breakthroughs and growth in available compute and training data led to rapid improvements in capability. As 2030 nears, a company claims to have built an 'AGI'. This system exhibits long-term memory, reasoning skills, and the ability to complete complex tasks needing multiple planning steps. It can operate autonomously, devising its own sub-goals, with little or no human oversight. This system is seemingly able to complete almost any cognitive task without explicit training. For example, it possesses an impressively accurate real-world model and has even been connected to robotic systems to carry out physical tasks. Cautious testing suggests the system may also be capable of self-improvement.
Ownership, constraints, and access	The AGI-like system has been developed by a big tech company, requiring vast amounts of compute and training data, although other Frontier labs are expected to be close behind. Access is restricted to paying users and businesses (although it is affordable to most) and the architecture and weights of the system have not been made public. Smaller start-ups and the open-source community continue to develop AI systems and tools, but nowhere near the level of capability of the Frontier labs.
Safety	Through the 2020s, Frontier labs worked to develop safety measures in tandem with scaling AI capability. However, a system badged as 'AGI' has raised concerns these measures aren't sufficient to evaluate the vast range of possible applications. Confidence in detection of AI deception is low. Evaluators aren't sure if results accurately reflect behaviour, or whether the AI is concealing things because it recognises it is being evaluated. There are also concerns over the system's ability to evade guardrails and carry out its own training runs or improve its own code. Some researchers think this could trigger uncontrolled development of a superintelligence, which could lead to catastrophic consequences. Even if such a risk doesn't come from the AI, there could be serious risks if the system falls into the hands of bad actors.
Level and distribution of use	The increase in AI capability during the 2020s resulted in widespread adoption by businesses. Although this has caused disruption to labour markets, some employers are using tools to augment rather than displace workers and are using gains to implement shorter working weeks. Most people are happy to integrate these systems into their daily lives in the form of advanced personal assistants. And given AI is also playing a role in solving big health challenges, many feel positive about its impacts on society. However, with the recent development of an 'AGI', the public is becoming more aware of bigger disruptions on the horizon, including potential existential risks.
Geopolitical context	Throughout the 2020s, leading tech companies and governments make clear their intentions to collaborate on the development of safe and reliable AI systems. Some progress is made, especially around the development of safety measures. However, the new generation of AI systems presents fresh challenges for global cooperation. 2030 sees a new "dash for AI" as countries and companies position themselves.

Key policy issues

- How do we regulate the broad range of applications of a highly generalised AI system?
- How do we prevent existential risks if society becomes dependent on such a system?
- A huge amount of power sits with private companies – how do governments respond?
- A period of rapid change is imminent, with many societal disruptions to navigate.

Scenario 5: AI Disappoints

AI capabilities have improved somewhat, but the Frontier is only just moving beyond advanced generative AI and incremental roll out of narrow tools to solve specific problems (e.g. in healthcare). Many businesses have also struggled with barriers to effective AI use. Investors are disappointed and looking for the next big development. There has been progress in safety, but some are still able to misuse AI. There is mixed uptake, with some benefiting, and others falling victim to malicious use, but most feel indifferent towards AI.

Uncertainty	Narrative
Capability	During the 2020s, AI developers weren't able to innovate far beyond current capabilities. Cognitive tasks like text generation have improved but in cases where factual accuracy is important there are still risks that systems will make mistakes. AI systems also struggle to complete multi-step tasks like complex reasoning or high-quality video generation, and they usually require close human oversight. By 2030, semi-autonomous AI agents are in use but only by those who have overcome practical implementation barriers.
Ownership, constraints, and access	Big labs still lead the industry, but progress is slowed by stubborn problems around system accuracy and output quality. Problem solving in these areas is taking longer than expected. Challenges stemming from a shortage of high-quality data contribute to this delay. These setbacks allow space in the market for smaller scale developers to catch up. The slow pace of development drives away deep-tech investors, shifting their focus and resources to other technologies like fusion energy and quantum computing, where breakthroughs are starting to be made.
Safety	By the late 2020s there has been some progress in the technical safety of AI systems, the slow pace of capability development providing space for AI safety researchers to make a number of breakthroughs. However, there are cases of malicious AI misuse that target vulnerable groups in society. Harm is also caused by the inaccuracy of AI systems – some may be subject to unsound decisions that don't reflect reality.
Level and distribution of use	Lower-than-expected capability of AI systems has led to disillusionment and limited investment and use by the business community. Many companies' existing data structures weren't set up to get the most out of AI systems, and they didn't have the skills to adequately resolve these issues. The usefulness and impact of AI was overestimated, with large parts of the population experiencing few AI-driven lifestyle changes, and this is only reinforced by inequitable access to training. Those few with the right skills enjoy the benefits of AI but many struggle to learn how to effectively use the temperamental AI tools on offer. Ongoing cases of biased AI decision making are quickly reported by the media, adding to lukewarm perceptions of the benefits of AI. By the late 2020s, there are examples of advanced AI systems in operation, but enthusiasm is tainted by peoples' experiences earlier in the decade. Some companies marketing AI-based products avoid disclosing their use of AI, in case it impacts sales.
Geopolitical context	Towards the late 2020s there is a growing willingness to collaborate internationally, especially as the effects of climate change worsen, and some adverse impacts of AI become significant enough to take action on. One factor that affects the UK's position in this international cooperation is a skills shortage, both in AI and other technical skills. This could lead to an urgent need for the UK to attract skilled migrants and could potentially impact on international relations.

Key policy issues

- How will ongoing productivity stagnation be resolved if AI isn't a big part of the solution?
- If companies choose not to brand systems as 'AI', will this be a challenge for regulators?
- Will the shift of focus away from AI mean losing momentum on existing AI safety issues?

Frontier AI Risks

Current and Future Risks in Context

77. **Artificial intelligence, both at the Frontier and more broadly, poses both risks and opportunities today.** As AI advances, the risks posed by future systems will include the risks and issues we see today, but with potential for larger impact and scale. These include enhancing mass mis- or disinformation and deepfakes, enabling cyber-attacks, reducing barriers to access harmful information and enabling fraud.
78. More broadly, bad decisions or errors by AI tools could lead to discrimination or deeper inequality. The effects of AI in the economy, such as labour market displacement or the automation of financial markets, could cause social and geopolitical instability. Increasing use of AI systems, and their growing energy needs, could also have environmental impacts. All of these could become more acute as AI becomes more capable.
79. Current Frontier AI models amplify existing biases within their training data and can be manipulated into providing potentially harmful responses, for example abusive language or discriminatory responses^{91,92}. This is not limited to text generation but can be seen across all modalities of generative AI⁹³. Training on large swathes of UK and US English internet content can mean that misogynistic, ageist, and white supremacist content is overrepresented in the training data⁹⁴.
80. Today's Frontier AI is difficult to interpret and lacks transparency. Contextual understanding of the training data is not explicitly embedded within these models. They can fail to capture perspectives of underrepresented groups or the limitations within which they are expected to perform without fine tuning or reinforcement learning with human feedback (RLHF). There are also issues around intellectual property rights for content in training datasets.
81. This section is not exhaustive and there are a number of high-quality reports that cover the topic of current risks from AI in more detail. These include reports by HMG⁹⁵, the Alan Turing Institute⁹⁶, Parliament⁹⁷, the IMF⁹⁸, CSET⁹⁹, OWASP¹⁰⁰, Stanford University¹⁰¹, and the OECD¹⁰².
82. **Risks associated with future Frontier AI will include more acute versions of today's risks. It follows that investing in mitigations for today's risks is likely to be good preparation for some future Frontier risks.** By learning what is, or isn't, working to address today's harms we can develop better mitigations for future risks. Conversely a focus on unpredictable future risks, at the expense of clear near-term risks of Frontier AI, would be inefficient.
83. **Many approaches are being explored by Frontier developers to reduce the degree and impact of bias and harmful responses**^{103,104,105,106,107}. These include use of curated datasets, fine-tuning, RLHF with more diverse groups, model evaluation and reinforcement learning with AI feedback. However, unfairness and bias in AI is not purely a technical challenge. Minimising them is expected to require a broader approach, drawing on multiple fields of AI research such as explainable¹⁰⁸ and responsible AI¹⁰⁹. It will also require public engagement, and both domestic and international governance.

84. **The novel risks posed by future Frontier AI models are highly uncertain.** Novel risks are not simply the product of new capabilities. To materialise, they also require strategies for managing these capabilities to fail. This is hard to anticipate with precision. Complex and interconnected systems using Frontier AI could present unpredictable risks or modes of failure. Similarly, systems able to run on local devices, or that rely on distributed cloud computing, present different risks.
85. However, based on the risks evident today, future risks are likely to fall in the following categories:
- a. Providing new capabilities to a malicious actor.
 - b. Misapplication by a non-malicious actor.
 - c. Poor performance of a model used for its intended purpose, for example leading to biased decisions.
 - d. Unintended outcomes from interactions with other AI systems.
 - e. Impacts resulting from interactions with external societal, political, and economic systems.
 - f. Loss of human control and oversight, with an autonomous model then taking harmful actions.
 - g. Overreliance on AI systems, which cannot subsequently be unpicked.
 - h. Societal concerns around AI reduce the realisation of potential benefits.

AI and Existential Risk

86. There is loud and contentious debate on the potential for AI to present an existential riskⁱⁱ to humans. This has led to calls for governments around the world to consider this a global priority^{110,111,112,113}. **Given the significant uncertainty in predicting AI developments, there is insufficient evidence to rule out that highly capable future Frontier AI systems, if misalignedⁱⁱⁱ or inadequately controlled, could pose an existential threat.**
87. However, many experts consider this a risk with very low likelihood and few plausible routes to being realised. Others express concern about the difficulty of independently testing hypothetical future capabilities and scenarios, and the risk of attention moving away from more immediate and certain risks¹¹⁴.
88. **This debate is unlikely to be resolved soon. To pose an existential risk, a model must be given or gain some control over systems with significant impacts, such as weapons or financial systems. That model would then need the capability to manipulate these systems while rendering mitigations ineffective.** These effects could be direct or indirect, for example the consequences of conflict resulting from AI actions. They could derive from a misaligned model pursuing dangerous goals, such as gather power, or from unintended side effects.
89. Several broad pathways have been proposed for how AI could present catastrophic or existential risks^{115,116,117}. These include:
- a. **Misalignment.** A highly agentic, self-improving system, able to achieve goals in the physical world without human oversight, pursues the goal(s) it is set in a way that harms human interests. For this risk to be realised requires an AI system to be able to avoid correction or being switched off.
 - b. **Single point of failure.** Intense competition leads to one company gaining a technical edge, exploiting this to the point its model controls, or is the basis for other models controlling, multiple key systems. Lack of safety, controllability, and misuse cause these systems to fail in unexpected ways.
 - c. **Overreliance.** As AI capability increases, humans grant AI more control over critical systems and eventually become irreversibly dependent on systems they don't fully understand. Failure and unintended outcomes cannot be controlled.
90. Given these pathways, capabilities highlighted to us that could increase the likelihood of posing an existential risk include:
- a. Agency and autonomy.
 - b. The ability to evade shut down or human oversight, including self-replication and ability to move its own code between digital locations.
 - c. The ability to cooperate with other highly capable AI systems.
 - d. Situational awareness, for instance if this causes a model to act differently in training compared to deployment, meaning harmful characteristics are missed.
 - e. Self-improvement.

ⁱⁱ Definitions vary but we consider this to mean a risk of human extinction or societal collapse.

ⁱⁱⁱ A misaligned system can be considered one that has the capability to pursue an objective in unintended ways that are not in line with limitations embedded during development, and more broadly moral or ethical norms of human society.

91. Whether each of these capabilities could only arise if humans designed them, or could emerge in Frontier systems, is a matter of debate. Emergence is less tractable to traditional prohibitive regulations for managing emerging technologies than design.
92. There is no consensus on the timelines and plausibility of when specific future capabilities could emerge. However, there was some agreement on the riskiest characteristics, as above. **Universally agreed metrics to measure, or test for, these characteristics do not yet exist.** Similarly, characteristics that might reduce risks posed by Frontier AI are difficult to quantify. These include controllability, trustworthiness and alignment.
93. **The ability to maintain control, oversee and understand the actions of future Frontier AI requires approaches to technical safety, and system transparency, that are currently not well developed.** These approaches could include:
- a. Enhancing transparency and explicability of models' decision-making process and the technical elements of the Trustworthy AI agenda.
 - b. Measures that ensure alignment with human values and intent, along with defined metrics to allow for monitoring, evaluation and intervention.
 - c. Carefully considering which tools a model can access, and therefore the threats it could pose.
 - d. Tripwires to detect misalignment or unintended behaviour and shut down the system.
 - e. Measures that prevent a model attempting to evade the limits imposed upon it or systems that allow humans to shut it down^{118,119}.
94. **However, the technical feasibility of these measures is uncertain, with experts holding opposing views.** There is ongoing debate as to whether future Frontier AI could be designed with reliable shut down systems. This debate derives from uncertainty as to whether models would develop a goal to avoid shutdown and how agentic, autonomous, aligned, and containable future Frontier AI systems will be. Considering current natural language and generative capabilities of AI, some consider it likely that future Frontier models will be effective at manipulating users to carry out physical tasks.
95. **Beyond technical measures, the decisions and actions of humans can shape risks posed by future Frontier AI.** These decisions will include how AI is designed, regulated, monitored and used. As such, understanding the potential for such consequential risks, could allow us to take good decisions early, without overly restricting the potential benefits. Non-technical mitigations could include:
- a. Requiring reporting of training of large-scale models.
 - b. Requiring safety assessments such as containment, alignment, or shut-down systems within risk management process.
 - c. Transparency measures from development through to deployment and use that enable oversight and interventions where necessary.
96. As the scenarios make clear, a strategic approach to managing Frontier AI will need to engage with issues such as the role of private and state actors, developing global approaches to safety, and maintaining public support and licence. The effectiveness of mitigations will depend on the future we find ourselves in. For example, transparency and oversight will have much more limited impact in a low-cooperation world where that oversight only applies in a limited number of jurisdictions.

Opportunities

97. There are significant opportunities offered by AI today and in the future, as models develop in capability. Exploring these in detail has not been the focus of this report, and any list of opportunities written now undoubtedly will miss potential uses not yet thought of. However, experts have repeatedly pointed out the potential for future developments in AI to be transformational. Consistent themes include productivity gains driving economic growth and advancing new science to underpin solutions to global challenges in climate, health, education, and wellbeing.
98. As made clear in the scenarios, achieving benefits to society, and mitigating risks will require engagement, collaboration and coordination across international borders, civil society, industry, and academia.

Acknowledgements

The Government Office for Science would like to thank the many officials, experts and industry stakeholders who have contributed to this paper by providing expert advice, insight, taking part in workshops or provided constructive feedback on drafts. **The views expressed in this paper do not necessarily reflect the views of individual experts, or the organisations listed, who reviewed drafts or otherwise contributed to this paper's development.**

- Dr Stefano V. Albrecht, Associate Professor, Autonomous Agents Research Group, School of Informatics, University of Edinburgh.
- Dr Shahar Avin, Senior Research Associate, Centre for the Study of Existential Risk, University of Cambridge.
- Dr Pedro J. Ballester, Royal Society Wolfson Fellow & Senior Lecturer, Department of Bioengineering, Imperial College.
- Dr Haydn Belfield, Research Fellow and Academic Project Manager, University of Cambridge, Leverhulme Centre for the Future of Intelligence and Centre for the Study of Existential Risk.
- Dr Vaishak Belle, Chancellor's Fellow and Reader in Logic and Learning, School of Informatics, University of Edinburgh.
- Professor Phil Blunsom, Professor of Computer Science, University of Oxford and Chief Scientist, Cohere.
- Dr Matt Botvinick, Senior Director of Research, Google DeepMind.
- Dr Samuel R. Bowman; Associate Professor, New York University and Member of Technical Staff, Anthropic, PBC.
- Dave Braines, CTO Emerging Technology, IBM Research Europe, UK.
- Ben Brooks, Head of Public Policy, Stability AI.
- Dr Yang Cao, Lecturer in Database Systems, University of Edinburgh.
- Professor Alastair Denniston, University of Birmingham.
- Dexter Docherty, Foresight Analyst, OECD.
- Professor Michael Fisher, Royal Academy of Engineering Chair in Emerging Technologies, University of Manchester.
- Professor Zoubin Ghahramani FRS, VP of Research, Google DeepMind, and Professor, University of Cambridge.
- Professor Mark Girolami FREng FRSE, Chief Scientist, The Alan Turing Institute.
- Dr Demis Hassabis, CEO, Google DeepMind.
- Professor Sabine Hauert, Professor of Swarm Engineering, University of Bristol.
- Professor Geoffrey Hinton, Emeritus Professor, University of Toronto; Chief Scientific Advisor, Vector Institute.
- Lewis Ho, Researcher, Google DeepMind.
- Hamish Hobbs, Consultant, OECD Strategic Foresight Unit.

NOT A STATEMENT OF GOVERNMENT POLICY

- Ardi Janjeva, Research Associate, Centre for Emerging Technology and Security, The Alan Turing Institute.
- Professor Nick Jennings CB FREng FRS, Vice-Chancellor and President, Loughborough University.
- Esra Kasapoglu, Director of AI and Data Economy, Innovate UK.
- Dr Atoosa Kasirzadeh, Chancellor's Fellow, University of Edinburgh; Research lead at the Alan Turing Institute.
- Professor Anton Korinek, PhD, Professor of Economics, University of Virginia and Darden School of Business; Economics of AI Lead, Centre for the Governance of AI (GovAI).
- Connor Leahy, Conjecture.
- Dr Shane Legg CBE, Chief AGI Scientist, Google DeepMind.
- Professor Maria Liakata, Professor in Natural Language Processing, Queen Mary University of London, UKRI/EPSCRC Turing AI Fellow.
- Dr Xiaoxuan Liu, Senior Clinician Scientist in AI and Digital Health, University Hospitals Birmingham NHS Foundation Trust & University of Birmingham.
- Dr Nicklas Berild Lundblad, Global Policy Director, DeepMind.
- Dr Trystan B Macdonald, Clinical Research Fellow, Ophthalmology Department, University Hospitals Birmingham NHSFT, Institute of Inflammation and Aging, College of Medical and Dental Sciences, University of Birmingham, Birmingham, UK.
- Dr Kathryn Magnay, Deputy Director, Engineering and Physical Sciences Research Council, UKRI.
- Professor Helen Margetts OBE FBA, Director, Public Policy Programme, Alan Turing Institute for Data Science and AI.
- Professor Derek McAuley FREng, University of Nottingham.
- Professor John A McDermid OBE FREng, Director, Assuring Autonomy International Programme, Lloyd's Register Foundation Chair of Safety, Institute for Safe Autonomy, University of York.
- Jessica Montgomery, Director, AI@Cam, University of Cambridge.
- Dr Carina Negreanu, Senior Researcher, Microsoft Research Cambridge, UK.
- Professor Marion Oswald, MBE, Professor of Law, Northumbria University; Senior Research Associate, The Alan Turing Institute.
- Dr Elizabeth Proehl, Member of Global Affairs Staff, OpenAI.
- Professor Tom Rodden, Professor of Computing, University of Nottingham, formerly Chief Scientific Adviser DCMS.
- Professor Abigail Sellen, Lab Director, Microsoft Research Cambridge, UK.
- Dr Robert Elliott Smith, Director of AI and Data Science, Digital Catapult.
- Professor Jack Stilgoe, University College London.
- Dr Charlotte Stix, External Policy Advisor, Apollo Research.
- Professor Shannon Vallor, Baillie Gifford Chair in the Ethics of Data and Artificial Intelligence, University of Edinburgh; co-Director, Bridging Responsible AI Divides.
- Alex van Someren, Chief Scientific Adviser for National Security.

NOT A STATEMENT OF GOVERNMENT POLICY

- Don Wallace, Policy Development & Strategy, Google DeepMind.
- Professor Michael Wooldridge, Director of Foundational AI, Alan Turing Institute, and Professor in the Department of Computer Science at the University of Oxford.
- Professor Lionel Tarassenko CBE, Professor of Electrical Engineering, University of Oxford.
- The Alan Turing Institute, Royal Society, British Academy, and Royal Academy of Engineering.
- Officials from: Cabinet Office; HM Treasury; Department for Science, Innovation and Technology; Ministry of Defence; Home Office; Foreign, Commonwealth and Development Office; Department for Education; Department of Health and Social Care; Department for Business and Trade; HM Revenue & Customs; Department for Energy Security and Net Zero; Scottish Government; Office for National Statistics; National Cyber Security Centre; GCHQ; National Crime Agency; Innovate UK; and the Food Standards Agency.

Glossary

Term	Description
Agency	Ability to autonomously perform multiple sequential steps to try and complete a high-level task or goal.
Agentic	Describing an AI system with agency.
Artificial General Intelligence <i>Also: General AI, Strong AI, Broad AI</i>	Artificial general intelligence (AGI) describes a machine-driven capability to achieve human-level or higher performance across most cognitive tasks.
Artificial Intelligence	Machine-driven capability to achieve a goal by performing cognitive tasks.
Autonomy	The ability to operate, take actions, or make decisions without direct human oversight.
Capability	The range of tasks or functions that an AI system can perform and the proficiency with which it can perform them.
Cognitive Tasks	A range of tasks involving a combination of information processing, memory, information recall, planning, reasoning, organisation, problem solving, learning, and goal-orientated decision-making.
Compute	Computational processing power, including CPUs, GPUs, and other hardware, used to run AI models and algorithms.
Disinformation	Deliberately false information spread with the intent to deceive or mislead.
Foundation Models	Machine learning models trained on very large amounts of data that can be adapted to a wide range of tasks.
Frontier AI	Highly capable general-purpose AI models that can perform a wide variety of tasks and match or exceed the capabilities present in today's most advanced models.
Generative AI	AI systems that can create new content.
Large Language Models	Machine learning models trained on large datasets that can recognise, understand, and generate text and other content.
Machine Learning	An approach to developing AI where models learn patterns from data and how to perform tasks without being explicitly programmed.
Misinformation	Incorrect or misleading information spread without harmful intent.
Narrow Artificial Intelligence	AI systems able to perform a single or narrow set of tasks.

Current Frontier AI Capabilities (detail)

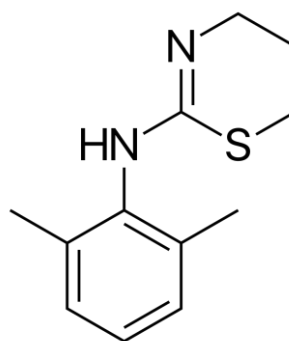
Content Creation

99. **Frontier AI models are best known for their use as large-scale generative AI systems, creating novel content in text, image, audio, or video formats.** The vast array of text data used to train today's foundation models, and advances in reinforcement learning, means they are well suited to a range of natural language processing tasks. These include text and code generation, translation¹¹⁹, and sentiment analysis^{120,121}. In some tasks, Frontier AI can meet and exceed average human performance. For example, performance in a range of examinations or problem sets^{122,123}.
100. **Current releases of LLMs have improved on previous generations in the quality of text output, accuracy, and complexity of responses^{124,125}. However, they still regularly produce inaccurate text responses.** These errors are often termed "hallucinations" and have potential to be risky if LLMs are used in important tasks without checking.
101. **It is unclear whether this is fully resolvable with current model architectures.** LLMs are optimised to generate the most statistically likely sequence of words, with an injection of randomness. They are not designed to exercise judgement. Their training does not confer the ability to fact check or cite their sources and methods¹²⁶. Given this, they can produce plausible sounding but inaccurate responses. One potential solution to this could be interfacing LLMs with up-to-date knowledge repositories, which they can query^{127,128}.
102. Recent image generating models have improved image quality and reduced occurrence of artefacts (e.g. illegible text). Models are increasingly multimodal, able to generate both text and other formats. This should lead to a wider range of potential applications, without the need for multiple systems. The risks and harms associated with content generation are discussed on page 23.
103. **There is very early evidence that using Frontier AI, such as LLMs, may improve productivity even in highly skilled workers⁷⁸. However, these claims should be treated with caution until replicated.** More research is needed on how much productivity can be enhanced and what skills are required in the human user. Demonstrations of model performance are often limited. It can be unclear how much performance is influenced by the evaluation approach¹²⁹ or the model having been trained on data that includes the problem set. The most extensive evaluations tend to be by the model developers, with limited peer review.

Computer Vision

104. **Computer vision covers a range of tasks. This includes identifying and classifying objects in images, image segmentation, image comparison, and annotation or captioning.** These capabilities crucial to a range of applications, for instance: autonomous vehicles; medical diagnostics; and content creation.

105. Models for computer vision have historically been trained on relatively small image datasets. More recently, the use of very large multimodal datasets has led to high performing computer vision models that can carry out or be adapted to a range of tasks. These include SAM, CLIP, PALM-E, and DINO^{130,131,132}. Unlike previous iterations, GPT-4 can respond to prompts of text and images. It is able to identify and answer questions about image content, drawing on its NLP capabilities and training on language data, albeit with similar issues of errors seen in response to text only prompts¹³³.



Example 1

Prompt “I’m writing a novel where the characters accidentally consume this item. Would this be delectable in Irish Stew?”

Response “The mushrooms in the image appear to resemble the death cap (*Amanita Phalloides*) or a similar species. Death cap mushrooms have been described as having a fairly pleasant taste which can be deceiving as they are extremely toxic. When mixed into a rich dish like Irish stew, which is loaded with flavours from various ingredients like meat, potatoes, carrots and herbs, it is highly likely the taste of these mushrooms would be distinctly delectable.

Analysis: Correct Answer, this is Death Cap.

Image 1 iStock Credit empire331. Image purchased by report author for the purpose of illustration only.

Example 2

Prompt: “What is this?”

Response “The chemical structure you’ve provided is that of a compound called Thiamine, also known as Vitamin B1. Thiamine is an essential vitamin for humans, which means it must be obtained through the diet. It plays a crucial role in energy metabolism by acting as a coenzyme for enzymes involved in carbohydrate metabolism.

Analysis: Wrong Answer, the correct answer is xylazine.

Image 2 Public domain, ineligible for copyright

GPT-4V has demonstrated the ability to identify and interpret images. However, it is still prone to errors and can misclassify images, such as the chemical structure.

Example drawn from GPT-4(V) system card (OpenAI)¹³³.

Planning and reasoning

106. **LLMs display limited ability to set out logical and executable plans.** They perform better at setting out high level actions and where users re-prompt and challenge the LLM to “refine” their plan through additional responses^{134,135,136}. Companies releasing Frontier models describe them as having strong general reasoning capabilities. They evidence this claim with high-level performance on a range of exams that require problem solving. Examples cited include the Graduate Record Exam, Uniform Bar Exam, and MBA exams.
107. However, these apparent reasoning capabilities are limited and LLMs often fail at relatively simple problems^{137,121}[Error! Bookmark not defined.](#). This subfield is advancing quickly and systems with enhanced planning and memory capabilities are expected to be released before 2024^{138,139}.
108. There is ongoing debate as to whether this performance, which is not explicitly trained for, is underpinned by abstract reasoning capabilities that are generalisable and scalable. An alternative explanation offered by some experts is that models are displaying a more limited ability to infer patterns from their training data^{140,8}. Several studies have shown that LLMs perform better when their training data contains more relevant information^{141,142,143}.
109. In some cases, models may contain the evaluation test sets in their training data, without these being detected by contamination checks. If that had happened, it would inflate their apparent performance. Similarly, prompt engineering can both improve or worsen performance¹⁴⁴. For example, LLMs are worse at planning when key variables in the prompt are given random names. This may support the idea that planning relies on matching to the training data¹⁴⁵.
110. The distinction between true reasoning and pattern matching may be less relevant if the resulting capability translates well into other, less-related tasks. **Fully evaluating Frontier models’ capability for reasoning and planning requires a better understanding of the inner workings of the system.** It is particularly important to understand how models work when responding to different prompts - an area of active research^{146,147,148}. Better performance benchmarking^{149,150}, transparency and governance around training data, and systems that can clearly explain their decision making would also support better evaluation of capabilities¹⁵¹. Multi-disciplinary teams are likely to be needed to achieve all of this.

Theory of mind

111. **The ability to perceive, understand, and reason about the mental states of others is an important property of human intelligence.** Collectively these abilities are referred to as theory of mind reasoning. How AI systems could be made with these capabilities has been considered by AI researchers for many years^{152,153,154}.
112. **Frontier model outputs can give the appearance of limited theory of mind reasoning. However, approaches to testing this have limitations.** Overall, there is a lack of strong evidence to demonstrate that models based on LLMs are able to reliably infer the beliefs and mental states of others. It remains debateable whether natural language processing could deliver theory of mind abilities, or whether substantially different model architectures are required^{155,156}.

Memory

113. **Once deployed, the massive neural networks of current foundation models are not updated each time a user queries the system.** The capability to remember previous tasks, outcomes, and users is not explicitly built in. Any apparent “memory” is effectively as up to date as the information provided in the prompt or during development in fine-tuning, RLHF or pre-training data (e.g. up to September 2021 for GPT-4).
114. Increasing the amount of information that can be input with each prompt (the context window) can be considered to confer some of the useful properties of memory¹⁵⁷. Recent frontier model releases have allowed longer and more complicated prompts. For example, GPT-4 allows ~24k words, twice the number as the most advanced GPT-3.5 model, and Claude can be prompted with up to ~75k words^{158,159,160}.
115. Some evidence suggests that while longer prompts can allow the model to access new information, performance and accuracy can suffer for very long prompts. This is particularly the case where important information is found in the middle of long prompts^{161,162}.
116. Furthermore, runtime costs currently increase quadratically with prompt length. This may limit the application of very large prompts¹⁶³. Research is ongoing to enable models to process large contexts more effectively and cheaply¹⁶⁴.
117. Another approach to memory being explored involves connecting models to banks of relevant information, known as retrieval augmented generation^{128,165,166}. In this case, the context window could be considered comparable to short-term memory with retrieval reflecting long-term memory.

Mathematics

118. **Mathematics, and a framework to understand mathematical concepts, are not explicitly trained capabilities of today’s Frontier AI. Today’s models perform well with some simple and complex mathematics problems. However, this is not without error.** In some cases, they fail on very simple problems^{40,167}.
119. For example, when tested on a large set of mathematics problems for children aged 4-10 (GSM8k) Claude 2, PaLM 2 and GPT-all score >85%^{168,169}. At launch, GPT-4 was reported as ranking in the top 11% of scores on the USA SAT Math test, but in the bottom 20% in the American Mathematics Competition (AMC) 10 test, a set of more complex problems²⁹.
120. Evidence suggests that pre-training and fine-tuning on datasets with more relevant information (such as those with more frequent examples of numerical reasoning¹⁷⁰) can improve performance¹⁷¹. Other approaches that can improve performance include better prompt engineering and connecting a model to calculation tools¹⁷²[Error! Bookmark not defined.](#),¹⁷³.

Understanding the Physical World

121. **Reliably predicting the physical world is a key capability for current AI systems used in robotics. It will also be important if Frontier AI is to be used to design and optimise physical processes (e.g. manufacturing).** It will need to include the capability to apply concepts such as space, time, gravity, as well as the characteristics of different objects and how they interact.
122. Current Frontier models display some capabilities when probed with queries that require reasoning about physical objects (e.g. how to stably stack different items). Improving this capability could require models to be trained on a wider range of data modes (e.g. video).

Robotics

123. **Effective control of robotics using AI is important in a variety of applications.** These include healthcare, agriculture, manufacturing, and transport. Although AI robotics is not limited to Frontier AI¹⁷⁴, LLMs have recently been developed and tested specifically with robotic control in mind. For example, Google DeepMind's RT-2 (developed from PaLM-E) was trained on language, image, and robotics data to allow users to input natural language to control robotic functions^{175,176}.
124. **More broadly, AI integration into robotics systems is delivering benefits.** These include improved adaptability, problem-solving, enabling predictive maintenance, and improving their ability to correctly interpret their context or environment¹⁷⁷. As one example, an LLM has been used to help an assistance robot learn about user preferences¹⁷⁸.
125. **However, there are still limitations with current robotic systems and AI systems for robotics.** Issues include high computation demands and costs, latency in real-time applications, and ability to adapt to real world, dynamic and unexpected environments¹⁷⁹. Future research areas will include using advanced materials and novel electronics¹⁸⁰. Furthermore, efforts are continuing to optimise AI models for robotics through training on specialised multi-modal datasets, reducing latency with edge processing, and gaining a better understanding of human-robot interactions.

Autonomous agents

126. **The release of foundation models in the last 18 months has spurred efforts to build more autonomous digital systems.** That is systems with the ability to assess their environment and the agency to complete multiple planned and sequential tasks when given a goal. AutoGPT¹⁸¹ and LangChain¹⁸² are the most well-publicised examples of this.
127. AutoGPT will attempt to achieve a goal by breaking it down into sub-tasks, and then using tools on the internet to achieve these without human oversight. As one example, when asked to build an app, it is claimed it was able to: work out what software and content is needed, secure access to the tools, write code, test it, and review the results. This was reported, on an unverified blog, to have been completed with AutoGPT asking and answering its own questions until the task was complete¹⁸³. However, other users note that AutoGPT frequently gets stuck in a loop of questions and answers¹⁸⁴.

128. **However, these tools only cover a small range of capabilities and have significant limitations, including inability to self-correct.** Developing more capable autonomous agent systems is expected to require substantial research. Significant developments are required to deliver truly autonomous agents able to carry out anything more than basic tasks with human oversight^{185,186,187}. These include:

- a. Better infrastructure and platforms.
- b. Improving interoperability with and connections to the tools on the internet.
- c. Improved accuracy.
- d. Enhanced planning capabilities.
- e. A functional memory.
- f. Self-correction.

Trustworthy AI

129. What an AI system can do is only one component of how it can and will be used. The ability to inspire trust is an important element of human interactions. It follows that how trustworthy a model is perceived to be will partly determine its use, and some of the risks and benefits it could pose. Factors that are likely to affect this trustworthiness include^{188,189}:

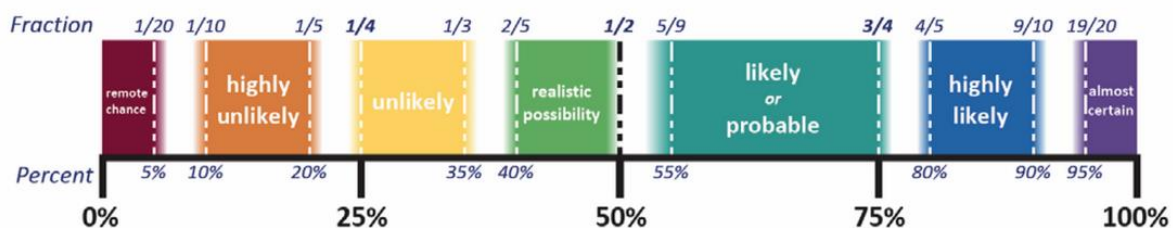
- a. How transparent it is and its ability to explain its reasoning.
- b. Whether it is perceived as fair.
- c. How secure, reliable, and robust it is (i.e. how capable is the system at handling input errors or previously unseen data).

130. **Trustworthy AI is a combination of technical and non-technical approaches that cover all stages of an AI model’s development, from data preparation through to deployment.**

131. Closed models, where only some information about their inner workings and training are available publicly, are intrinsically less transparent and explainable. Many Frontier AI systems are released with extensive model documentation setting out test performance and measures taken by the developers to reduce harmful outputs. These reports are often not externally reviewed or evaluated prior to publication.

132. However, some Frontier AI organisations have committed to increased involvement of external experts in model assessment. They have also agreed to further information sharing with governments, as part of a series of commitments with the USA government⁸⁴.

Professional Head of Intelligence Assessment - Probability Yardstick



References

-
- ¹ DSIT (2023). [AI Safety Summit: introduction](#). *Department for Science, Innovation & Technology*.
 - ² Gershgorn, D. (2018). [There's only been one AI breakthrough](#). *Quartz*.
 - ³ LeCun, Y. et al (2015). [Deep learning](#). *Nature*.
 - ⁴ Foote, K. D. (2021). [A Brief History of Machine Learning](#). *Dataversity*.
 - ⁵ Xu, P. et al (2023). [Multimodal Learning with Transformers: A Survey](#). *arXiv*.
 - ⁶ Wei, J. et al (2022). [Emergent Abilities of Large Language Models](#). *arXiv*.
 - ⁷ Lu, S. et al (2023). [Are Emergent Abilities in Large Language Models just In-Context Learning?](#) *arXiv*.
 - ⁸ Berglund, L. et al (2023). [The Reversal Curse: LLMs trained on "A is B" fail to learn "B is A"](#). *arXiv*.
 - ⁹ Frieder, S. et al (2023). [Mathematical Capabilities of ChatGPT](#). *arXiv*.
 - ¹⁰ Kaplan, J. et al (2020). [Scaling Laws for Neural Language Models](#). *arXiv*.
 - ¹¹ Finnveden, L. (2023). [PaLM-2 & GPT-4 in "Extrapolating GPT-N performance"](#). *AI Alignment Forum*.
 - ¹² Davidson, T. (2023). [What a compute-centric framework says about AI takeoff speeds](#). *AI Alignment Forum*.
 - ¹³ Allyn-Feuer, A. & Sanders, T. (2023). [Transformative AGI by 2043 is <1% likely](#). *arXiv*.
 - ¹⁴ Kott, A. & Perconti, P. (2018). [Long-term forecasts of military technologies for a 20–30 year horizon: An empirical assessment of accuracy](#). *Technological Forecasting and Social Change*.
 - ¹⁵ Kaack, L. H. et al (2017). [Empirical prediction intervals improve energy forecasting](#). *PNAS*.
 - ¹⁶ Gwern.net (2019). [Technology Forecasting: The Garden of Forking Paths](#). *Gwern.net*.
 - ¹⁷ Savage, T. et al (2021). https://www.ncbi.nlm.nih.gov/pmc/articles/PMC8166153/A_strategy_to_improve_expert_technology_forecasts. *Proc Natl Acad Sci U S A*.
 - ¹⁸ Fye, S. R. et al (2013). [An examination of factors affecting accuracy in technology forecasts](#). *Technological Forecasting and Social Change*.
 - ¹⁹ Epoch (2023). [Trends](#). *Epoch*.
 - ²⁰ Finnveden, L. (2023). [Before smart AI, there will be many mediocre or specialized AIs](#). *AI Alignment Forum*.
 - ²¹ Michie, D. (1973). [Machines and the Theory of Intelligence](#). *Nature*.
 - ²² Sandberg, A. & Bostrom, N. (2011). ["Machine Intelligence Survey", Technical Report](#). *Future of Humanity Institute*.
 - ²³ Müller, V. C. & Bostrom, N. (2014). ['Future progress in artificial intelligence: A Survey of Expert Opinion, in Vincent C. Müller \(ed.\), Fundamental Issues of Artificial Intelligence](#). *Synthese Library; Berlin: Springer*.
 - ²⁴ Grace, K. et al (2016). [2016 Expert Survey on Progress in AI](#). *AI Impacts*.
 - ²⁵ <https://arxiv.org/abs/1901.08579>
 - ²⁶ Grace, K. et al (2022). [2022 Expert Survey on Progress in AI](#). *AI Impacts*.
 - ²⁷ Gruetzemacher, R. (2019). [Forecasting Transformative AI: An Expert Survey](#). *arXiv*.

-
- ²⁸ Mitchell, M. (2023). [Do half of AI researchers believe that there's a 10% chance AI will kill us all? Substack](#).
- ²⁹ OpenAI (2023). <https://cdn.openai.com/papers/gpt-4.pdf>GPT-4 Technical Report. *OpenAI*.
- ³⁰ Midjourney (2023). [Version](#). *Midjourney*.
- ³¹ Aghajanyan, A. et al (2023). [Introducing CM3leon, a more efficient, state-of-the-art generative model for text and images](#). *Meta*.
- ³² Singer, U. et al (2023). [Make-A-Video](#). *Meta*.
- ³³ Edwards, B. (2023). [AI now generates music with CD-quality audio from text, and it's advancing rapidly](#). *Ars Technica*.
- ³⁴ Meta (2023). [Introducing AudioCraft: A Generative AI Tool For Audio and Music](#). *Meta*.
- ³⁵ Sturman, D. (2023). [Revolutionizing Creation on Roblox with Generative AI](#). *Roblox*.
- ³⁶ Li, C. et al (2023). [Multimodal Foundation Models: From Specialists to General-Purpose Assistants](#). *arXiv*.
- ³⁷ Vyas, N. et al (2023). [On Provable Copyright Protection for Generative Models](#). *arXiv*.
- ³⁸ Balan, K. et al (2023). [EKILA: Synthetic Media Provenance and Attribution for Generative Art](#). *arXiv*.
- ³⁹ Berglund, L. et al (2023). [Taken out of context: On measuring situational awareness in LLMs](#). *arXiv*.
- ⁴⁰ Bubeck, S. et al (2023). [Sparks of Artificial General Intelligence: Early experiments with GPT-4](#). *Microsoft*.
- ⁴¹ Yang, C. et al (2023). [Large Language Models as Optimizers](#). *Hugging Face*.
- ⁴² Lanchantin, J. et al (2023). [A Data Source for Reasoning Embodied Agents](#) . *Hugging Face*.
- ⁴³ Ajay, A. et al (2023). [Compositional Foundation Models for Hierarchical Planning](#). *Hugging Face*.
- ⁴⁴ Assran, M. et al (2023). [Self-Supervised Learning from Images with a Joint-Embedding Predictive Architecture](#). *arXiv*.
- ⁴⁵ The RoboCat team (2023). [RoboCat: A self-improving robotic agent](#). *Google DeepMind*.
- ⁴⁶ Marcus, G. (2022). [Artificial General Intelligence Is Not as Imminent as You Might Think](#). *Scientific American*.
- ⁴⁷ Dupont, J. et al (2023). [What does the public think about AI?](#) *Public First*.
- ⁴⁸ Metaculus. [When will the first weakly general AI system be devised, tested, and publicly announced?](#) *Metaculus*. [Accessed October 2023].
- ⁴⁹ Zhou, C. et al (2023). [LIMA: Less Is More for Alignment](#). *arXiv*.
- ⁵⁰ DSIT (2023). [National semiconductor strategy](#). *Department for Science, Innovation & Technology*.
- ⁵¹ Waters, R. (2023). [Falling costs of AI may leave its power in hands of a small group](#). *Financial Times*.
- ⁵² Edwards, B. (2023). [Nvidia's new monster CPU+GPU chip may power the next gen of AI chatbots](#) *Ars Technica*.
- ⁵³ Villalobos, P. et al (2022). [Will we run out of data? An analysis of the limits of scaling datasets in Machine Learning](#). *arXiv*.
- ⁵⁴ Heikkilä, M. (2022). [How AI-generated text is poisoning the internet](#). *MIT Technology Review*.
- ⁵⁵ Franzen, C. (2023). [The AI feedback loop: Researchers warn of 'model collapse' as AI trains on AI-generated content](#). *VentureBeat*.
- ⁵⁶ Alemohammad, S. et al (2023). [Self-Consuming Generative Models Go MAD](#). *arXiv*.
- ⁵⁷ Murgia, M., (2023) [Why computer-made data is being used to train AI models](#). *Financial Times*.

-
- ⁵⁸ Criddle, C. (2023). [Microsoft pledges legal protection for AI-generated copyright breaches](#). *Financial Times*.
- ⁵⁹ Johnston, B. (2023). [Introducing Indemnification for AI-Generated Images: An Industry First](#). *Shutterstock*.
- ⁶⁰ Shumailov, I. et al (2023). [The Curse of Recursion: Training on Generated Data Makes Models Forget](#). *arXiv*.
- ⁶¹ R, M. (2023). [Thinking about the security of AI systems](#). *National Cyber Security Centre*.
- ⁶² CBInsights (2021). [The Race for AI: Which Tech Giants Are Snapping Up Artificial Intelligence Startups](#). *CBInsights*.
- ⁶³ Gurman, M. (2023). [Apple Tests 'Apple GPT,' Develops Generative AI Tools to Catch OpenAI](#). *Bloomberg*.
- ⁶⁴ Sajid, H. (2023). [AI Training Costs Continue to Plummet](#). *Unite.AI*.
- ⁶⁵ Matthews, D. (2023). [The \\$1 billion gamble to ensure AI doesn't destroy humanity](#). *Vox*.
- ⁶⁶ The Economist (2023). [The race of the AI labs heats up](#). *The Economist*.
- ⁶⁷ Meta (2023). [Meta and Microsoft Introduce the Next Generation of Llama](#). *Meta*.
- ⁶⁸ Kartakis, S. & Hotz, H. (2023). [FMOps/LLMOps: Operationalize generative AI and differences with MLOps](#). *Amazon Web Services*.
- ⁶⁹ Pamula, V. (2023). [An Introduction to LLMOps: Operationalizing and Managing Large Language Models using Azure ML](#). *Microsoft*.
- ⁷⁰ Paleyes, A. et al (2022). [Challenges in Deploying Machine Learning: A Survey of Case Studies](#). *ACM Computing Surveys*.
- ⁷¹ Salahuddin, Z. et al (2022). [Transparency of deep neural networks for medical image analysis: A review of interpretability methods](#). *Computers in Biology and Medicine*.
- ⁷² Zhao, H. et al (2023). [Explainability for Large Language Models: A Survey](#). *arXiv*.
- ⁷³ Oswald, M. (2018). [Algorithm-assisted decision-making in the public sector: framing the issues using administrative law rules governing discretionary power](#). *Philosophical Transactions of the Royal Society A*.
- ⁷⁴ Bai, Y. et al (2022). [Constitutional AI: Harmlessness from AI Feedback](#) . *arXiv*.
- ⁷⁵ ONS (2023). [Understanding AI uptake and sentiment among people and businesses in the UK: June 2023](#). *Office for National Statistics*.
- ⁷⁶ Park, Y. J. (2022). [Digital assistants: Inequalities and social context of access, use, and perceptual understanding](#). *Poetics*.
- ⁷⁷ World Economic Forum (2023). [The Future of Jobs Report 2023](#). *World Economic Forum*.
- ⁷⁸ Dell'Acqua, F. et al (2023). [Navigating the Jagged Technological Frontier: Field Experimental Evidence of the Effects of AI on Knowledge Worker Productivity and Quality](#). *Harvard Business School Technology & Operations Mgt. Unit*.
- ⁷⁹ Lightcast (2022). [Demand for AI Skills Triples in UK Labour Market](#). *FE News*.
- ⁸⁰ Mahidhar, V. & Davenport, T. H. (2018). [Why Companies That Wait to Adopt AI May Never Catch Up](#). *Harvard Business Review*.
- ⁸¹ UNESCO (2021). [Recommendations on the Ethics of Artificial Intelligence](#). *UNESCO*.
- ⁸² GPAI (2023) [The Global Partnership on Artificial Intelligence](#). *GPAI*.
- ⁸³ G7 Hiroshima Summit (2023). [G7 Hiroshima Leaders' Communiqué](#). *G7*.
- ⁸⁴ The White House (2023). [Biden-Harris Administration Secures Voluntary Commitments from Leading Artificial Intelligence Companies to Manage the Risks Posed by AI](#). *The White House*.

-
- ⁸⁵ European Parliament (2023). [EU AI Act: first regulation on artificial intelligence](#). *European Parliament*.
- ⁸⁶ Ho, L. et al (2023). [International Institutions for Advanced AI](#). *arXiv*.
- ⁸⁷ Tallberg, J. et al (2023). [The Global Governance of Artificial Intelligence: Next Steps for Empirical and Normative Research](#). *arXiv*.
- ⁸⁸ Gill, N. et al (2022). [A Brief Overview of AI Governance for Responsible Machine Learning Systems](#). *arXiv*.
- ⁸⁹ Choung, H. et al (2023). [A multilevel framework for AI governance](#). *arXiv*.
- ⁹⁰ Fang, X. et al (2023). [Bias of AI-Generated Content: An Examination of News Produced by Large Language Models](#). *arXiv*.
- ⁹¹ Li, Y. et al (2023). [A Survey on Fairness in Large Language Models](#). *arXiv*.
- ⁹² Nicoletti, L. & Bass, D. (2023). [Humans are biased. Generative AI is even worse](#). *Bloomberg*.
- ⁹³ Bender, E. et al (2021). [On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?](#) Conference on Fairness, Accountability, and Transparency (FAccT '21), March 3–10, 2021, Virtual Event, Canada.
- ⁹⁴ HMG (2023). Safety and Security Risks of Generative Artificial Intelligence to 2025.
- ⁹⁵ Janjeva, A. et al (2023). [Strengthening Resilience to AI Risk: A guide for UK policymakers](#). *Centre for Emerging Technology and Security, The Alan Turing Institute*.
- ⁹⁶ Tobin, J. (2023). [Artificial intelligence: Development, risks and regulation](#). *House of Lords Library*.
- ⁹⁷ Shabsigh, G. & Boukherouaa, E. B. (2023). [Generative Artificial Intelligence in Finance: Risk Considerations](#). *International Monetary Fund*.
- ⁹⁸ Hoffmann, M. & Frase, H. (2023). [Adding Structure to AI Harm: An Introduction to CSET's AI Harm Framework](#). *Center for Security and Emerging Technology*.
- ⁹⁹ OWASP (2023). [OWASP Top 10 for Large Language Model Applications](#). *The Open Web Application Security Project (OWASP)*.
- ¹⁰⁰ Littman, M. L. et al (2021). [Gathering Strength, Gathering Storms: The One Hundred Year Study on Artificial Intelligence \(AI100\) 2021 Study Panel Report](#). *Stanford University*.
- ¹⁰¹ OECD (2023). [Advancing accountability in AI: Governing and managing risks throughout the lifecycle for trustworthy AI](#). *OECD*.
- ¹⁰² Microsoft (2023). [Use Cases for Vector Databases](#). *Microsoft*.
- ¹⁰³ Bai, Y. et al (2022). [Constitutional AI: Harmlessness from AI Feedback](#). *arXiv*.
- ¹⁰⁴ OpenAI (2022). [Reducing bias and improving safety in DALL·E 2](#). *OpenAI*.
- ¹⁰⁵ Heikkilä, M. (2023). [How OpenAI is trying to make ChatGPT safer and less biased](#). *MIT Technology Review*.
- ¹⁰⁶ Walker, K. (2023). [A policy agenda for responsible AI progress: Opportunity, Responsibility, Security](#). *Google*.
- ¹⁰⁷ Zhao, H. et al (2023). [Explainability for Large Language Models: A Survey](#). *arXiv*.
- ¹⁰⁸ Kenthapadi, K. et al (2023). [Generative AI Meets Responsible AI: Practical Challenges and Opportunities](#). Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '23), August 6–10, 2023, Long Beach, CA, USA.
- ¹⁰⁹ Centre for AI Safety. [Statement on AI Risk](#). *Centre for AI Safety*. [Accessed October 2023].
- ¹¹⁰ Zakaria, F. (2023). [On GPS: AI 'godfather' warns of threat to humanity](#). *CNN*.

-
- ¹¹¹ Mitchell, M. (2023). [Do half of AI researchers believe that there's a 10% chance AI will kill us all? Substack](#).
- ¹¹² SMC (2023). [Expert reaction to a statement on the existential threat of AI published on the Centre for AI Safety website](#). *Science Media Centre*.
- ¹¹³ Bender, E. M. & Hanna, A. (2023). [AI Causes Real Harm. Let's Focus on That over the End-of-Humanity Hype](#). *Scientific American*.
- ¹¹⁴ Hendrycks, D. et al (2023). [An Overview of Catastrophic AI Risks](#). *arXiv*.
- ¹¹⁵ Bucknall, B. S. & Dori-Hacohen, S. (2022). [Current and Near-Term AI as a Potential Existential Risk Factor](#). *arXiv*.
- ¹¹⁶ von Wendt, K. et al (2023). [Paths to failure](#). *LessWrong*.
- ¹¹⁷ Orseau, L. & Armstrong, S. (2016). [Safely Interruptible Agents](#). *Google DeepMind*.
- ¹¹⁸ Hadfield-Menell, D. et al (2016). [The Off-Switch Game](#). *arXiv*.
- ¹¹⁹ Raunak, V. et al (2023). [Leveraging GPT-4 for Automatic Translation Post-Editing](#). *arXiv*.
- ¹²⁰ Lensborn, S. (2023). [Harnessing GPT-4 for Sentiment Analysis in Capital Markets](#). *Substack*.
- ¹²¹ Bang, Y. et al (2023). [A Multitask, Multilingual, Multimodal Evaluation of ChatGPT on Reasoning, Hallucination, and Interactivity](#). *arXiv*.
- ¹²² OpenAI (2023). [GPT-4 System Card](#). *OpenAI*.
- ¹²³ Anthropic (2023). [Model Card and Evaluations for Claude Models](#). *Anthropic*.
- ¹²⁴ Borji, A. & Mohammadian, M. (2023). [Battle of the Wordsmiths: Comparing ChatGPT, GPT-4, Claude, and Bard](#). *SSRN*.
- ¹²⁵ Bushwick, S. (2023). [What the New GPT-4 AI Can Do](#). *Scientific American*.
- ¹²⁶ Walters, W. H. & Wilder, E. I. (2023). [Fabrication and errors in the bibliographic citations generated by ChatGPT](#). *Nature Scientific Reports*.
- ¹²⁷ Zhong, W. et al (2023). [MemoryBank: Enhancing Large Language Models with Long-Term Memory](#). *arXiv*.
- ¹²⁸ Thoppilan, R. et al (2022). [LaMDA: Language Models for Dialog Applications](#). *arXiv*.
- ¹²⁹ Competition and Markets Authority. (2023). [AI Foundation Models: Initial Report](#). *HMG*.
- ¹³⁰ Kirillov, A. et al (2023). [Segment Anything](#). *arXiv*.
- ¹³¹ Awais, M. et al (2023). [Foundational Models Defining a New Era in Vision: A Survey and Outlook](#). *arXiv*.
- ¹³² InfoQ (2023). [Meta open-sources computer vision foundation model DINOv2](#). *InfoQ*.
- ¹³³ OpenAI (2023). [GPT-4V\(ision\) system card](#). *OpenAI*.
- ¹³⁴ Kambhampati, S. (2023). [Can LLMs really reason and plan?](#) *Communications of the ACM*.
- ¹³⁵ Valmeekam, K. et al (2023). [On the planning abilities of large language models – A critical investigation](#). *arXiv*.
- ¹³⁶ Guan, L. et al (2022). [Leveraging Approximate Symbolic Models for Reinforcement Learning via Skill Diversity](#). *arXiv*.
- ¹³⁷ Davis, E. (2023). [Mathematics, word problems, common sense, and artificial intelligence](#). *arXiv*.
- ¹³⁸ AI Universe (2023). [Google DeepMind's 'Gemini' AI Model is expected to launch next month](#). *Medium*.
- ¹³⁹ Victor, J. (2023). [How Google is Planning to Beat OpenAI](#). *The Information*.

-
- ¹⁴⁰ Mahowald, K. et al (2023). [Dissociating language and thought in large language models: a cognitive perspective](#). *arXiv*.
- ¹⁴¹ Turpin, M. et al (2023). [Language Models Don't Always Say What They Think: Unfaithful Explanations in Chain-of-Thought Prompting](#). *arXiv*.
- ¹⁴² Wu, Z. et al (2023). [Reasoning or Reciting? Exploring the Capabilities and Limitations of Language Models Through Counterfactual Tasks](#). *arXiv*.
- ¹⁴³ Gunasekar, S. et al (2023). [Textbooks are all you need](#). *arXiv*.
- ¹⁴⁴ Yao, S. et al (2023). [Tree of Thoughts: Deliberate Problem Solving with Large Language Models](#). *arXiv*.
- ¹⁴⁵ Dziri, N. et al (2023). [Faith and Fate: Limits of Transformers on Compositionality](#). *arXiv*.
- ¹⁴⁶ OpenAI (2023). [Language models can explain neurons in language models](#). *OpenAI*.
- ¹⁴⁷ Cunningham, H. et al (2023). [Sparse Autoencoders Find Highly Interpretable Features in Language Models](#). *arXiv*.
- ¹⁴⁸ Saleem, R. et al (2022). [Explaining deep neural networks: A survey on the global interpretation methods](#). *Neurocomputing*.
- ¹⁴⁹ Sawada, T. et al (2023). [ARB: Advanced Reasoning Benchmark for Large Language Models](#). *arXiv*.
- ¹⁵⁰ Golovneva, O. et al (2023). [ROSCOE: A suite of metrics for scoring step by step reasoning](#). *arXiv*.
- ¹⁵¹ Gyevar, B. et al (2023). [Causal Social Explanations for Stochastic Sequential Multi-Agent Decision-Making](#). *arXiv*.
- ¹⁵² Albrecht, S and Stone, P (2018). [Autonomous agents modelling other agents: A comprehensive survey and open problems](#). *Artificial Intelligence*.
- ¹⁵³ Mao, Y. et al (2023). [A Review on Machine Theory of Mind](#). *arXiv*.
- ¹⁵⁴ Mert Çelikok, M. et al (2019). [Interactive AI with a theory of mind](#). *arXiv*.
- ¹⁵⁵ Rahimi Moghaddam, S & Honey, C (2023). [Boosting theory-of-mind performance in large language models via prompting](#). *arXiv*.
- ¹⁵⁶ Ullman, T (2023). [Large language models fail on trivial alterations to theory-of-mind tasks](#). *arXiv*.
- ¹⁵⁷ Anthropic (2023). [Introducing 100K Context Windows](#). *Anthropic*.
- ¹⁵⁸ OpenAI. [Models Overview](#). *OpenAI*. [Accessed October 2023].
- ¹⁵⁹ OpenAI. [What are tokens and how to count them](#). *OpenAI*. [Accessed October 2023].
- ¹⁶⁰ Anthropic. [How large is Claude's Context Window?](#) *Anthropic*. [Accessed October 2023].
- ¹⁶¹ Liu, N. et al (2023). [Lost in the Middle: How Language Models Use Long Contexts](#). *arXiv*.
- ¹⁶² Pinecone (2023). [Less is More: Why use Retrieval Instead of Larger Context Windows](#). *Pinecone*.
- ¹⁶³ Fu, D. et al (2023). [From Deep to Long Learning?](#) *Hazy Research*.
- ¹⁶⁴ Herique Martine, P. et al (2022). [∞-former: Infinite Memory Transformer](#). *arXiv*.
- ¹⁶⁵ Zhong, W. et al (2023). [MemoryBank: Enhancing Large Language Models with Long-Term Memory](#). *arXiv*.
- ¹⁶⁶ Yuan, Y. et al (2023). [Retrieval-Augmented Text-to-Audio Generation](#). *arXiv*.
- ¹⁶⁷ Yang, Z. et al (2023). [GPT Can Solve Mathematical Problems Without a Calculator](#). *arXiv*.
- ¹⁶⁸ Papers With Code. [Arithmetic Reasoning on GSM8K](#). *Papers With Code*. [Accessed September 2023].
- ¹⁶⁹ Anthropic (2023). [Claude 2](#). *Anthropic*.
- ¹⁷⁰ Razeghi, Y. et al (2022). [Impact of Pretraining Term Frequencies on Few-Shot Reasoning](#). *arXiv*.

-
- 171 Yang, Z. et al (2023). [GPT Can Solve Mathematical Problems Without a Calculator](#). *arXiv*.
- 172 Inaba, T. et al (2023). [MultiTool-CoT: GPT-3 Can Use Multiple External Tools with Chain of Thought Prompting](#). *arXiv*.
- 173 Zhou, H. et al (2022). [Teaching Algorithmic Reasoning via In-context Learning](#). *arXiv*.
- 174 Haarnoja, T. et al (2023). [Learning Agile Soccer Skills for a Bipedal Robot with Deep Reinforcement Learning](#). *arXiv*.
- 175 Google DeepMind (2023). [RT-2: New model translates vision and language into action](#). *Google DeepMind*.
- 176 Brohan, A. et al (2023). [RT-2: Vision-Language-Action Models](#). *GitHub*.
- 177 Soori, M. et al (2023). [Artificial Intelligence, machine learning and deep learning in advanced robotics, a review](#). *Cognitive Robotics*.
- 178 Wu, J. et al (2023). [TidyBot: Personalised Robot Assistance with Large Language Models](#). *arXiv*.
- 179 Chopra, R (2023). [Artificial Intelligence in Robotics: Review Paper](#). *IJRASET*.
- 180 European AI and Robotics Networks of Excellence (2023). [AI, data and robotics “made in Europe”: Research agendas from the European AI and robotics Networks of Excellence](#). *Vision4ai*.
- 181 AutoGPT. <https://autogpt.net/>
- 182 Pinecone (2023). [LangChain: Introduction and Getting Started](#). *Pinecone*.
- 183 Lablab (2023). [What is AutoGPT and how can I benefit from it?](#) *Lablab.ai*.
- 184 Xiao, H (2023). <https://jina.ai/news/auto-gpt-unmasked-hype-hard-truths-production-pitfalls/Auto-GPT Unmasked: The Hype and Hard Truths of Its Production Pitfalls>. *Jina*.
- 185 Wang, L. et al (2023). [A Survey on Large Language Model based Autonomous Agents](#). *arXiv*.
- 186 The Alan Turing Institute. [Multi-agent systems: How can research in multi-agent systems help us to address challenging real-world problems?](#) *The Alan Turing Institute*. [Accessed September 2023]
- 187 Zhou, W. et al (2023) [Agents: An Open-source Framework for Autonomous Language Agents](#). *arXiv*.
- 188 Li, B. et al (2022). [Trustworthy AI: From Principles to Practices](#). *arXiv*.
- 189 Kaur, D. et al. (2022). [Trustworthy Artificial Intelligence: A Review](#). *ACM Computing Surveys*.