# In Which Areas of Technical AI Safety Could Geopolitical Rivals Cooperate?

BEN BUCKNALL*, Department of Engineering Science, University of Oxford & Oxford Martin AI Governance Initiative, bucknall@robots.ox.ac.uk

SAAD SIDDIQUI*, Safe AI Forum & Oxford Martin AI Governance Initiative, saad@saif.org

LARA THURNHERR, King's College London

CONOR MCGURK, Safe AI Forum

BEN HARACK, Oxford Martin AI Governance Initiative

ANKA REUEL, Stanford University & The Belfer Center for Science and International Affairs, Harvard Kennedy School

PATRICIA PASKOV, RAND

CASEY MAHONEY, RAND

SÖREN MINDERMANN, Mila - Quebec AI Institute

SCOTT SINGER, Carnegie Endowment for International Peace & Oxford Martin AI Governance Initiative

VINAY HIREMATH, Centre for the Governance of AI

CHARBEL-RAPHAËL SEGERIE, Centre pour la Sécurité de l'IA (CeSIA)

OSCAR DELANEY, Institute for AI Policy and Strategy

ALESSANDRO ABATE, Department of Computer Science, University of Oxford

FAZL BAREZ, Department of Engineering Science, University of Oxford

MICHAEL K. COHEN, UC Berkeley & Center for Human Compatible AI

PHILIP TORR, Department of Engineering Science, University of Oxford

FERENC HUSZÁR, University of Cambridge

ANISOARA CALINESCU, Department of Computer Science, University of Oxford

GABRIEL DAVIS JONES, Oxford Digital Health Labs, University of Oxford

YOSHUA BENGIO, Mila - Quebec AI Institute

ROBERT TRAGER, Oxford Martin AI Governance Initiative & Blavatnik School of Government, University of Oxford

International cooperation is common in AI research, including between geopolitical rivals. While many experts advocate for greater international cooperation on AI safety to address shared global risks, some view cooperation on AI with suspicion, arguing that it can pose unacceptable risks to national security. However, the extent to which cooperation on AI safety poses such risks, as well as provides benefits, depends on the specific area of cooperation. In this paper, we consider technical factors that impact the risks of international cooperation on AI safety research, focusing on the degree to which such cooperation can advance dangerous capabilities, result in the sharing of sensitive information, or provide opportunities for harm. We begin by why nations historically cooperate on strategic technologies and analyse current US-China cooperation in AI as a case study. We further argue that existing frameworks for managing associated risks can be supplemented with consideration of key risks specific to cooperation on technical AI safety research. Through our analysis, we find that research into AI verification mechanisms and shared protocols may be suitable areas for such cooperation. Through this analysis we aim to help researchers and governments identify and mitigate the risks of international cooperation on AI safety research, so that the benefits of cooperation can be fully realised.

*Equal contribution.

arXiv:2504.12914v1 [cs.CY] 17 Apr 2025

# 1 Introduction

> *"Many risks arising from AI are inherently international in nature, and so are best addressed through international cooperation."* – Bletchley Declaration, 2023 [37]

International cooperation has long been a part of managing risks from advanced technologies. During the Cold War, despite intense rivalry, the United States and Soviet Union collaborated on nuclear verification methods through initiatives like the Joint Verification Experiment, which facilitated progress on arms control agreements [36]. Today we observe ongoing international cooperation in artificial intelligence (AI) development, including between geopolitical rivals such as the US and China. Additionally, recent years have seen numerous high-level calls for international cooperation specifically on AI safety and governance, from the landmark agreement at the Bletchley AI Safety Summit in 2023 [37] to the consensus statements issued by the International Dialogues on AI Safety[1] [see also, 90].

As in historical analogues, cooperating on the safety of advanced geopolitically sensitive technologies such as AI could play an important role in managing emerging risks. However, cooperation between geopolitical rivals carries its own risks that can and should be carefully weighed in order to ensure that the benefits of cooperation can be fully realised by all parties. This is no less true for safety research on AI ('technical AI safety'). Some AI safety techniques have 'capability externalities', as improvements in safety concurrently provide gains in model performance – for example, reinforcement learning from human feedback (RLHF), as discussed in [26]. Furthermore, some areas of AI safety such as model evaluations for chemical, biological, radiological and nuclear (CBRN) capabilities involve sensitive national security-related information which could be leaked to cooperators. The process of cooperation may also provide avenues for motivated actors to cause harm, such as by placing backdoors in jointly developed infrastructure. Many frameworks and risk management processes from governments exist to guide international technology cooperation at a general level, such as suggesting researchers conduct due diligence on the identity of their counterparties. However, as leading governments' and researchers' continued focus on deepening the science of AI safety demonstrates [10, 74], additional analysis is needed to fully address the specific technical characteristics and geopolitical risks associated with cooperation on AI safety research.

This paper addresses this gap by identifying technical risks specific to cooperation on AI safety research, focusing on the impact of cooperation on advancing capabilities, sharing sensitive information and providing opportunities for harm. We analyse the relative risk of cooperation for four different areas of technical AI safety and assess the feasibility of proposed cooperation in each area. We do not aim to definitively identify the most suitable areas for cooperation, nor to investigate specific benefits of cooperation in different AI safety research areas, though believe these would be valuable directions for future work.

This paper is roughly broken into three parts. The first provides an overview of cooperation on strategic technologies, through i) outlining why geopolitical rivals typically cooperate on strategic technologies; ii) surveying the existing state of cooperation between rivals on AI research; and iii) looking at how risks from cooperation are currently managed. The second part highlights four potential risks to which actors cooperating with rivals on topics in technical AI safety may be exposed. The final part builds on this by assessing the extent to which proposed areas of cooperation on technical AI safety may succumb to these identified risks.We give a brief overview of each area, as well as a discussion regarding its suitability for cooperation in light of the potential risks previously introduced, finding that research into i) AI verification mechanisms and ii) protocols may be areas that are particularly well-suited for international cooperation.

---

[1]https://idais.ai/

## 1.1 Key concepts and definitions

*1.1.1 Coordination, collaboration, and cooperation.* For the purpose of this paper, we adopt the following definitions.[2]

- We use **coordination** to refer to multiple parties acknowledging shared viewpoints and/or agreeing to pursue common goals independently and in parallel.
- We use **collaboration** to refer to multiple parties jointly pursuing shared goals.
- We use **cooperation** to refer to the overarching category of both coordination and collaboration.

As demonstrative examples, we regard the Seoul Ministerial Statement, signed by ministers from 27 nations and the EU, as a case of **coordination**. It, among other things, *'reaffirm[ed] [the signatories'] shared intent to guide the design, development, deployment, and use of AI in a manner which harnesses its benefits for good'* though does not represent a joint effort to act on this intent [39]. In contrast, the US and UK AI Safety Institutes conducting a *'joint safety research and testing exercise'* [as in 3] is a case of **collaboration**, as it represents a joint activity undertaken to achieve a shared goal. Both are examples of **cooperation**.

*1.1.2 Technical cooperation.* In this paper, we limit our scope to considering cooperation on areas of technical AI safety and governance – that is, research and development within computer science, engineering, mathematics (or similar) aimed at advancing the safety of AI systems, or methods for their effective governance, respectively.[3] Taken together, activities included in the present scope include, but are not limited to: academic scientific research; industry research and development on AI systems or related technologies, such as semiconductors; and state-state collaboration on technical topics, for example through their respective AI Safety Institutes. Prior work has explored potential areas for international collaboration, including on non-technical topics such as the sharing of *'best practices and lessons learned from AI regulation efforts in specific countries'* [46].

## 2 Background and motivation

In this section, we provide a historical overview of cooperation on strategic technologies. We outline why geopolitical rivals typically cooperate on strategic technologies, before surveying the existing state of cooperation between rivals on AI research, focusing on the US and China as a case study. Finally, we discuss how risks from cooperation are currently managed, serving as background for the section that follows, in which we propose a framework to consider the specific risks from cooperation on AI safety research.

## 2.1 Why do geopolitical rivals cooperate on strategic technology?

Here, we draw on basic concepts from game theoretic accounts of international cooperation and competition [e.g. 25, 57] to describe a subset of the plausible reasons geopolitical rivals might consider it to be in their interest to cooperate on strategic technologies such as AI. We illustrate these concepts with empirical examples cited in the literature about AI safety cooperation that we consulted over the course of our research.

First, rivals may cooperate to manage risks from technology that cannot be effectively managed by any single actor. This includes cooperation on risks that cross borders, such as illicit use of technology by international criminal groups. Recent examples of cooperation in this vein include Sino-American agreements to jointly combat money laundering through cryptocurrency [95]. Rivals

---

[2]There is limited academic consensus on what the terms coordination, cooperation, and collaboration refer to. Work in international relations and meta-analyses of usage in management science arrive at different conclusions [19, 92].
[3]For overviews of AI safety and technical AI governance, see [48] and [83], respectively.

may also cooperate where collective action is necessary to reduce risk, such as in the November 2024 agreement between the US and China to maintain human control over the decision to use nuclear weapons and avoid integrating AI into nuclear command and control systems [82].

An actor leading on a given technology may also cooperate by unilaterally sharing technology, if doing so is also in the leader's own interests and sharing is technologically feasible. In the early 1960s, the US shared basic designs for Permissive Action Links (PALs) – devices that prevent unauthorised nuclear detonation – with the Soviet Union. This cooperation was made possible by the mutual recognition of clear benefits in preventing accidental escalation, and that early PALs were simple enough to explain without compromising sensitive weapons information [31].

Rivals may cooperate to improve geopolitical stability by creating mechanisms that reduce uncertainty and the risk of unintended escalation. Examples here include the Open Skies Treaty, which allows participating nations to conduct unarmed aerial surveillance flights over each other's territories using standardised sensor technologies, creating predictable patterns of interaction between rival militaries, and establishing technical protocols for verification [8].[4]

Finally, rivals may also cooperate to pool expertise and resources when technological development costs exceed the resources or capacity of any single actor. The International Space Station partnership between space agencies, including those of the U.S. and Russia, leveraged each state's space capabilities [66], while the ITER fusion project brings together rival powers to share the massive costs and technical challenges of developing fusion power.[5]

The reasons specified above hold true for states, but at the academic and corporate level, cooperation can take place for reasons completely unrelated to broader geopolitical tension. For example, academics may collaborate with colleagues based in a rival nation who have valuable expertise in a particular area. Companies may want to cooperate with other companies based in rival jurisdictions because they fulfill complementary roles in a complex globalised supply chain, or to ensure interoperability of products to retain market access in key markets.

Many of these reasons for cooperation likely also apply to AI and AI safety. For example, increasingly capable AI systems developed in one jurisdiction may have negative cross-border impacts in a rival jurisdiction, requiring cooperation between rivals to manage these risks effectively. AI safety techniques developed by one state that reliably reduce the risk of an AI system self-replicating and self-improving without human approval could be shared if they reduce overall global risk, and do not reveal sensitive commercial or national security-relevant information.

## 2.2 Cooperation on AI case study: The US and China

In this section we outline the current state of international cooperation within the realms of academia, industry, and intergovernmental relations, focused on the two leading global AI powers, the US and China.[6] The picture that emerges is one where significant academic and some industry cooperation continue to take place, even given the limited cooperation between the respective governments.

---

[4]As part of measures to ensure the security of the observed country, the treaty gives the observed country the right to conduct a pre-flight inspection of the observation aircraft to ensure it is unarmed and equipped only with the agreed-upon sensors with pre-determined resolutions. The observed country can also supply its own observation aircraft for the other country to use [8].

[5]See https://www.iter.org/few-lines. Despite pooling of resources, ITER is 9 years behind schedule, with private fusion startups on track to outpace it [24].

[6]While we focus on cooperation within each level of academia, industry, or state, it is worth noting that much cooperation can also take place between these communities, in the form of academic-industry partnerships, or academic participation in policy processes, as demonstrated, for example, in the members of the UN's High-Level Advisory Body on AI (see https://www.un.org/en/ai-advisory-body/members).
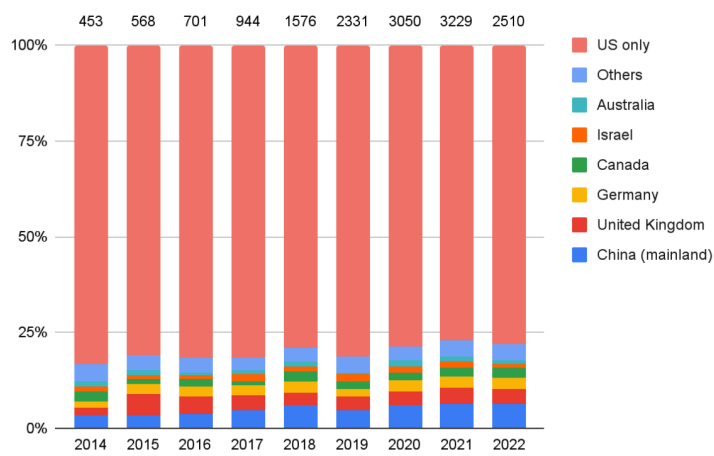
Fig. 1. AI safety co-authorship instances with American researchers (%). Incomplete data from 2023 and 2024 excluded. Data labels at the top of bars show the total number of AI safety papers by US researchers in that year. Data source: [72]. (Note that if a publication lists authors from organisations in more than one country (excluding the US), the publication will "count towards" multiple countries, and thus be represented multiple times in the figure.)

*2.2.1 Academia.* US and Chinese researchers collaborate more than researchers from any other pair of countries [104, §1.1], including on topics in AI safety – as shown in Figure 1 below. China surpassed the UK as the largest collaborator with American researchers in 2017, and has retained this position since, although a significant majority of AI safety research is conducted by American researchers without international co-authors.

*2.2.2 Industry.* Historically, some American firms have set up localised joint ventures in China to establish a foothold in Chinese markets and draw on talent pools. At the same time, these ventures also provided a significant boost to the Chinese technology industry.[7] Recent authoritative histories reveals how the growth of the Chinese semiconductor industry since the late twentieth century serves as a prime example of this dynamic [67]. The most prominent example in the AI software sector is Microsoft Research Asia (MSRA), a research institution set up by Bill Gates in 1998 and first led by Kaifu Lee, founder of 01.AI, a leading Chinese open-source AI company [58]. These investments in Chinese AI ventures have had an impact on the global AI industry. For example, in 2015 a team from MSRA led by Kaiming He introduced 'deep residual nets' [47] and greatly advanced the frontier of deep learning methods.

*2.2.3 Intergovernmental.* Cooperative efforts on AI between the American and Chinese governments have been much more limited than in the case of industry or academia.[8] AI has only recently come under the spotlight as an important geopolitical issue, requiring the attention of senior government officials. For example, only in 2023 was AI included as a summit-level topic for a meeting between US President Joe Biden and Chinese President Xi Jinping [53]. In 2024, a dedicated

---

[7]In certain industries (e.g., insurance), foreign companies are required to establish a joint venture with local Chinese partners. In a more restricted set of industries such as aviation, foreign ownership is also limited to below 50% [78].

[8]Cooperation on AI has also been affected by American export controls on semiconductors, explicitly designed to safeguard an American advantage in AI by restricting Chinese access to cutting-edge semiconductors required for advanced AI [12].

intergovernmental dialogue on AI took place in Geneva in May, with further dialogues planned for the future [65]. Notably, in November 2024, the two heads of state reached an agreement to maintain human control over the decision to use nuclear weapons and avoid integrating AI into nuclear command and control systems [82].

## 2.3 How are risks from cooperation managed?

The above case study illustrates that there is significant cooperation ongoing between rivals. This section draws on authors' engagement with AI expert communities to provide an illustrative overview of the commonly cited measures used to manage the risks arising from such cooperation generally.

States have been cognizant of the risks of cooperation on strategic technology, and have put in place measures to address such risks.[9] However, the risk management processes related to cooperating with a rival at the inter-state level are opaque.

More detailed public guidelines exist for academics and companies who, as noted above, may want to cooperate for reasons largely unrelated to geopolitical tensions and may be less aware of relevant geopolitical considerations.

Academics who cooperate with other academics in rival jurisdictions typically must engage with national guidance systems. These tend to require actions such as assessing the risks related to the subject or domain of research (e.g., whether knowledge produced could be misused) as well as the research conditions (e.g., potential military links of collaborators or the political environment of where a collaborator is based) [2, 97]. They may also be asked to check the identity of collaborators and their institutions against sanctions lists, and to consult documents such as the Bureau of Industry and Security's (BIS) Commerce Control List, which lists a series of items (e.g., encryption software) that are subject to export control regulations [50].

For companies engaging in activities such as joint ventures, additional rules on outbound and inbound investments typically apply. In the US, for example, companies must report relevant inbound transactions to investment screening entities, such as the Committee on Foreign Investment in the United States (CFIUS), which reviews foreign investments that may impact US national security [23].

One gap in the risk management process outlined above is a conceptual one – few of the tools above focus on technology-specific nuances. Therefore, actors considering cooperation lack a clear framework to assess how cooperation on their specific technology of interest could pose geopolitically-relevant risks.

## 3 Risks of cooperation on AI safety

In this section we outline risks associated with international cooperation that are specific or particularly pertinent to AI safety. Namely, we outline risks associated with the advancement of (potentially harmful) AI capabilities, exposure of sensitive information regarding strategic technology, and opportunities for motivated actors to take harmful actions that may be afforded by cooperation. As above, we draw on basic concepts from game theoretic accounts of international cooperation and competition to typologise the risks geopolitical rivals might articulate as barriers to cooperation.[10]

---

[9]It is worth noting that measures that may appear to be sufficient at a time when a technology is not deemed sensitive (e.g., AI in the early 2000s, nuclear research prior to its securitisation), may no longer be sufficient when the national security risks related to the technology become more evident [see e.g., 20]. As such, in this section, we do not take a position on whether existing measures are sufficient.

[10]We make no claim that these categories of risks are comprehensive, rather, that they are particularly relevant in the context of cooperating on AI safety.

**Developments in AI safety could advance the global capabilities frontier.** Geopolitical rivals may be hesitant to cooperate on AI safety due to the risk that doing so may, as a side effect, advance the global frontier of (potentially harmful) AI capabilities [45]. This could as a result of correlations between the safety and capability of AI systems, or if advances made in AI safety can be reapplied or repurposed to increase a system's suitability for deployment [see e.g., 49, pg. 29]. Inadvertently increasing the capabilities of systems through cooperative safety research would be undesirable to an actor that would have a preference for being the sole beneficiary of such capabilities gains. It is worth also noting that the potential of advancing overall dual-use capabilities could be an argument for not pursuing the research, even unilaterally, depending on the specific research project in question, and the actor's risk appetite.

**Cooperation could differentially advance a rival's strategic AI capabilities.** Alternatively, a state that is (or is perceived to be) 'leading' in terms of strategic AI capabilities, may be unwilling to cooperate with a rival due to a risk that doing so will allow the rival to improve their capabilities relative to the leader's. Such a differential benefit to a rival could be gained through application of the results of cooperative research, or as a side effect of increased access to knowledge or resources provided for the purpose of the collaboration.

**Cooperation could expose sensitive information regarding nationally strategic technology.** Cooperation on AI safety could also be risky if the specific focus of the cooperation intersects with other (non-AI) nationally strategic technology, in ways that may raise national security concerns. For example, research that aims to conduct risk modeling to estimate the risk posed to domestic digital infrastructure by AI's potential cyber-offensive capabilities may require in-depth knowledge of this infrastructure. Thus, a meaningful cooperation on such a project may require disclosing such information to a rival – something that may not be within the risk appetite of many geopolitical actors.

**Cooperation on AI safety may present opportunities for motivated actors to cause harm.** Finally, cooperation on AI safety research could allow rivals opportunities to take harmful or malicious action, in cases in which they may be inclined to do so. This could be, for example, through inserting backdoors into systems to which collaborations have access for the purpose of the collaboration, or through misusing resources shared for the purpose of cooperation.

## 4 Candidate areas of cooperation

In this section we present a non-comprehensive overview of four areas of AI technical safety – verification, protocols, infrastructure, and evaluation – on which international cooperation is emerging, or otherwise has been prominently advocated for. For each area we aim to assess the extent to which cooperation in that area would carry the above risks. A summary of these assessments is given in Table 1, and we provide examples of ongoing or proposed cooperation in each area in the list below.

- **Verification mechanisms.** Research into potential methods for verifying the veracity of claims made about systems was highlighted as a potential area of cooperation at a recent track-II dialogue between Western and Chinese academics.[11] See also, e.g., [10, pg. 207].
- **Protocols.** Development of codified protocols and best practices, for example 'AI safety frameworks' or standardisation efforts, have been the focus of ongoing international co-operation, for example, in the Frontier AI Safety Commitments [38], and the subsequent 'Conference on Frontier AI Safety Frameworks' hosted by the UK AI Safety Institute [55].
- **Infrastructure.** Developing shared AI infrastructure, or otherwise methods for distributing access to existing infrastructure is called for in the UN High-Level Advisory Body's final

---

[11]See https://idais.ai/dialogue/idais-venice/.

report [98]. Furthermore, the International Network of AI Safety Institutes used evaluation infrastructure developed by two constituent members (*Inspect* and *Moonshot*) when conducting a pilot testing exercise [73], though these platforms are not themselves the result of international cooperation.

- **Evaluations.** AI evaluation and testing exercises have been the subject of ongoing international collaboration between the UK and US AI safety institutes [see e.g., 3, 73].

Based on our assessment, we judge research on verification mechanisms and protocols to be less-challenging areas for international cooperation than infrastructure and evaluations. However, this does not imply that efforts to cooperate on the latter two areas should be avoided.

## 4.1 Research on verification mechanisms

Verification mechanisms are technical procedures that could allow for the certification of claims about an AI system or related resources (for example, the locations and specifications of data centers) [15, 28].[12] Note that we are here using the term 'verification' in a somewhat different sense than the established technical meaning of ensuring that a software or hardware system possesses certain properties through, for example, formal methods such as model checking [9]. While this is included in our use of the term (see 'formal verification' below), we do not restrict ourselves to this definition, instead, including activities such as verifying compliance to international agreements between state actors, be it regarding models, or another part of the AI value chain.

While related, this is distinct from designing a new system, or from uncovering information about a system (for example through conducting evaluations) as it predominantly only strengthens existing knowledge claims made about a system or activity.[13] The process of developing verification mechanisms, however, could allow rivals to gather sensitive information (e.g., inspectors in data centers could gather detailed information about installed hardware and its vulnerabilities).

Cooperation between multiple actors on the development of verification mechanisms could be beneficial for advancing mutual trust in any resulting solutions, and therefore increasing the chances that they are applied in practice. Trusted international agreements will likely depend on establishing verifiability throughout the entire technical and procedural stack of dependencies – cooperation from as far down this stack as possible could help establish mutual verifiability. Furthermore, cooperation could aid in ensuring interoperability or compatibility of verification mechanisms, enabling their use alongside each actor's existing technology and infrastructure with minimal modifications required.

*4.1.1 Risks of cooperation.* **Advances global frontier capabilities.** To the extent that some applications of verification concern the attestation of properties of a system, rather than demonstrating the existence of such properties, research in such areas is unlikely to advance the capabilities of AI systems. For example, developing methods for attesting to the amount of compute used in a given training run likely does not bear on frontier AI capabilities. However, this may not hold for all areas of verification. For example, the use of formal verification methods such as model checking could uncover new system properties that were not previously known to developers [9].

**Differentially advances a rival's capabilities.** As noted above, the development and use of some verification mechanisms could uncover previously-unknown model properties. This could result in a differential advancement of a rival's strategic capabilities. However, there are also many areas that likely avoid this risk, as the claims and information being verified can be restricted to

---

[12]For related discussion regarding the geopolitics of verification in the case of arms control, see [40] and [25].

[13]An example from [83]: *'an assessment task could be to uncover details about the data that a given model was trained on. In contrast, a verification problem could be, given a dataset and a model, to confirm or refute the claim that the model was trained on the dataset.'*

those already known to all parties. As an analogue, in the Open Skies Treaty, countries agree to allow other countries to carry out observation flights over their territory, but only using certified sensors with a specific predetermined resolution [8].

**Exposes other sensitive information.** A concern associated with cooperating on the development of verification technologies for AI could be if doing so requires disclosing sensitive information about each actor's existing technologies. For example, developing hardware mechanisms for verification [63] could require deep knowledge of existing hardware technology, which may be viewed as sensitive information that should not be shared with rivals. Some verification procedures, for example allowing inspectors in AI data centers to certify particular claims about the amount of compute present, could result in sensitive information related to or adjacent to the object of verification, such as the types of security systems installed, to be inadvertently disclosed. Some proposals for hardware-enabled verification mechanisms could potentially bypass some of such concerns through being independent of the specification of the AI accelerator and other related hardware being used for training or inference [80].

**Provides opportunity for motivated actors to take harmful action.** Jointly developing mechanisms for verification could allow motivated rivals to covertly insert 'verification backdoors', allowing them to feign compliance in cases in which this mechanism is applied. The level of difficulty of doing so without the backdoor being discovered would likely vary greatly depending on the specific verification mechanism. This risk could be managed by developing the verification stack in the open – for example, open-source with significant bug bounties – as a way of increasing the chances that such subversive attempts would be discovered.

*4.1.2 Subareas.* **Formal verification.** Formal verification refers to methods for proving that an AI system meets a specification related to the context and environment it is deployed in [27]. While formal verification for software has a decades-long history [51, 70], most existing methods struggle with the complexity of massive software systems deployed in the real world [64]. More recently, the accelerating inclusion of machine learning models in safety-critical systems has spurred new research extending formal methods to ML [see e.g., 14, 99]. However, these methods are severely underdeveloped and inadequate for addressing open-ended safety-critical scenarios, language models, or complex autonomous agents interacting with humans.

**Verifiable audits.** Cooperation could take place on developing and implementing methods for conducting verifiable audits – that is, audits for which reported results can be verified as being accurate [41, 83, 91, §5.4.1]. This could extend to verifying whether a given AI model was trained on a given dataset [22]. In addition to being an area whose development could be the subject of international cooperation, such methods could be broadly beneficial for building trust in the reported capabilities of a state's own AI systems.

**Compute usage.** Verifying the use of specialised computational hardware could be instrumental for checking compliance with international agreements relating to AI [86, 88]. While some proposals for secure hardware verification have been proposed [63, 80] the development of such methods is nascent. Furthermore, compute verification mechanisms would likely need to be robust to physical tampering so as to disincentivise subversion attempts [4] or incorporate other verification mechanisms (e.g., spot checks of specific chip shipments) as fail-safes.

**Generated content.** Methods for reliably identifying and verifying AI-generated media as such, for example, through embedding machine-readable watermarks in system outputs, is an ongoing area of research [29, 42, 101], and has gained traction as a potential regulatory requirement.[14]

---

[14]See for example, Recital 133 of the EU AI Act [**?** ].

Verification of AI-generated media may also be able to be achieved through cryptographically bound provenance metadata, for example as detailed by the C2PA specification.[15]

## 4.2   Codification of protocols and best practices

Protocols and best practices refer to codified statements of procedures for attaining positive outcomes from AI research and development. Protocols could specify procedures to be followed, or outcomes to be achieved, either by industry or state actors. Developing such shared protocols could be a suitable area for intergovernmental coordination due to its less technical nature, especially if phrased in terms of outcomes rather than actions [87]. While some states are already coordinating on the development protocols (for example, through the G7 Hiroshima Process[16]) there is scope for broadening the extent of cooperation, as well increasing the level of detail at which protocols are specified.

It should be noted that protocols are voluntary or binding codifications of existing knowledge, not purely academic research. As such, they could be more politicised than the other categories discussed here – especially in international settings. This could bring with it both positives, for example, leveraging technical research to build real world political agreement, as well as negatives, such as if the process becomes co-opted and overly politicised, leading to a degradation in the scientific rigour of the process and outputs.[17]

*4.2.1   Risks of cooperation.* **Advances global frontier capabilities.** Developing protocols may be particularly suited for cooperation between state actors, especially on more established topics for which the protocols aim to codify existing techniques and knowledge. In such areas, protocols would be more a process of standardisation rather than advancing state of the art research, and would thus not be at risk of advancing global AI capabilities.

**Differentially advances a rival's capabilities.** To the extent that cooperating on the codification of protocols focusing on areas where all parties have shared knowledge and understanding, the risk of differentially advancing a rival's capabilities through such cooperation is minimal.

**Exposes other sensitive information.** As for the previous risk, due to developing protocols being aimed at codifying mutual knowledge in a structured framework that can be agreed upon by multiple actors, no sensitive or private information would necessarily need to be shared with a rival. While this may depend on both the focus area, as well as the level of specificity of the resulting protocols (for example, detailed best practices regarding national security), such areas could be avoided while continuing to cooperate on less sensitive topics.

**Provides opportunity for motivated actors to take harmful action.** Given that the codification of protocols does not involve direct involvement with AI systems, cooperating on such codification would not allow rivals to take directly harmful actions. However, prior examples in standardisation suggest that both states and industry actors tend to use the international standardisation process to advance their own interests, potentially to the detriment of other actors [62, 85, 106]. For example, the World Trade Organisation's *'Technical Barriers to Trade'* (TBT)[18] agreement requires that countries use international standards instead of domestic ones to prevent favoritism. However, in practice this provides cover for Chinese actors to enforce ITU standards, largely pushed through by and advantageous to Chinese actors, on international companies in

---

[15]See https://c2pa.org/specifications/specifications/2.1/index.html.
[16]https://www.soumu.go.jp/hiroshimaaiprocess/en/index.html
[17]See [62] for a canonical treatment of this dynamic, or [102] for a more recent international history.
[18]https://www.wto.org/english/tratop_e/tbt_e/tbt_e.htm

China. Nonetheless, misuse of standards-setting processes from Chinese actors is more the exception than the rule and the most influential standards-setting bodies continue to function effectively [89].

*4.2.2 Subareas.* **Safety frameworks.** The past year has seen a number of frontier AI developers write and publish safety frameworks [7, 30, 75], with 13 further developers agreeing to publish their own documents by the French AI Action Summit, as part of the Frontier AI Safety Commitments [38]. The Chinese Academy of Communication and Information, a think-tank housed under the Ministry of Industry and Information in China, has also released a set of safety-focused voluntary commitments signed by several leading Chinese AI developers [17]. While such frameworks to-date have been developed within industry, with the Seoul commitments also being directed towards and signed solely by industry actors, there is a role to play in the public sector in setting the criteria by which these frameworks are assessed, and ensuring that industry actors fulfil their commitments [6, 35]. Such cooperation has precedent in the form of many prior national and international frameworks regarding the responsible development and use of AI [see e.g., 52, 69, 76].

In particular, multilateral cooperation could aim to develop and define best practices regarding any of the three outcomes specified in the Seoul commitments (or otherwise) – for example, advancing methods for identifying, assessing, and managing risks associated with AI development and deployment, or determining thresholds for when models might pose unacceptable levels of risk [61], and what should be done in case such predefined thresholds are breached [60].

**AI incident standards.** Actively monitoring for incidents, and having standardised methods for incident reporting, can inform policymakers about harms caused by AI systems [93, 94]. While existing AI incident databases[19] rely on news articles, data on non-public incidents may also be crucial for effective incident monitoring and risk-estimation [77]. Cooperation on a common definition of an 'AI incident' could lay the foundation for usefully interoperable domestic incident monitoring frameworks. Furthermore, any specification of protocols and best practices for incident reporting and monitoring would need to specify the information pertaining to an incident that should be shared with whom, including the type and magnitude of incident for which reporting should be required.

**Secure weight infrastructure standards.** While the specific security measures individual AI development organisations use to protect AI model weights must remain confidential, there may be scope for international collaboration on developing standardised technical frameworks for securing weights without revealing the weights themselves or organisations' other confidential information. This could include joint research into cryptographic protocols specifically designed for protecting large-scale model weights, methods for third parties to verify appropriate security measures without gaining access to the weights themselves drawing on precedents such as nuclear safeguard inspections [54], and interface specifications for hardware security modules designed for AI model protection [63, 71]. Another precedent is the international collaboration on cryptographic standards through organisations like the ISO,[20] that enabled countries to work together on security frameworks while maintaining their own secure implementations. However, such collaboration would need to be carefully structured to ensure that only high-level protocols and standards are shared and that the collaboration does not reveal vulnerabilities in existing security systems. For example, researchers have investigated cryptographic methods to secure AI model weights during inference and deployment using Zero-Knowledge Proofs [59], and hardware security modules (HSMs) have been proposed to protect AI models from unauthorised access and tampering [79].

---

[19]See, for example, the OECD's *AI Incidents Monitor* (https://oecd.ai/en/incidents) or the *AI Incident Database* (https://incidentdatabase.ai/).

[20]See, for example, the security requirements for cryptographic modules [1].

**Meta-analyses.** May 2024 saw the publication of the interim International Scientific Report on the Safety of Advanced AI, an international effort led by the UK's Department for Science, Innovation and Technology that aimed to *"drive a shared, science-based, up-to-date understanding of the safety of advanced AI."* [11] The report represents a collaborative effort between 30 countries, including the US and China, as well as the EU and the UN. Similarly, the UN's High-Level Advisory Body on AI has recommended that the UN establish an *"independent international scientific panel on AI"* which should issue *"an annual report surveying AI-related capabilities, opportunities, risks, and uncertainties,"* among other actions [98].[21] These examples show that diverse states can, and furthermore do, coordinate on initiatives that aim to advance a shared technical and scientific understanding of AI, including the global risks that it may pose or exacerbate. Such 'literature review'-style research efforts that establish a scientific basis for research funding and resource allocations, and more systematic meta analyses could aid in establishing greater alignment and enable better consensus building between rivals. However, it is worth noting that there may be a risk that such efforts become co-opted by political pressures to advance national interests.

### 4.3    Infrastructure

We use the term 'AI safety infrastructure' to refer to systems and processes (be they hardware, software, organisational, or otherwise), external to an AI system, that facilitates research and development activities relating to AI safety. This could take the form of hardware infrastructure, such as computational resources, or software tooling such as code packages for conducting certain types of research [see e.g., 34]. Cooperating on AI safety infrastructure could have large benefits for ensuring interoperability of ongoing research and development activities in different jurisdictions. This may allow for greater global distribution of the benefits of AI [5], as well as facilitate further collaboration on empirical AI research.

*4.3.1    Risks of cooperation.* **Advances global frontier capabilities.** Due to the broad, multi-purpose nature of many forms of infrastructure, there is a risk that the developments made through cooperation on infrastructure could be applied, perhaps with minor repurposing, to advance frontier AI capabilities. For example, general software packages aimed at facilitating AI safety research could potentially also facilitate more effective research and development that advances frontier AI capabilities. The marginal contribution of this infrastructure may however be minimal, due to the large and expanding amount of existing infrastructure for AI research and deployment.

**Differentially advances a rival's capabilities.** Similarly to advancing global frontier capabilities, infrastructure could plausibly be used by rivals to facilitate their own strategic capabilities. However, since this infrastructure will be available to all cooperating parties and thus make it less costly for a given party to detect a rival's progress toward a decisive strategic capabilities on a shared system, it is less likely that this will confer a significant *differential* advantage to a rival that goes undetected.

**Exposes other sensitive information.** To the extent that developing some forms of shared infrastructure for AI safety builds upon existing national infrastructure, doing so may require providing sensitive details regarding existing infrastructure to adversaries. This could raise a broad range of national security-relevant risks. However, many forms of infrastructure, such as high-level evaluations frameworks such as the UK AI Safety Institute's *Inspect*[22] or Singapore's AI Verify Foundation's *Project Moonshot*,[23] likely do not carry this risk.

---

[21]A broader discussion of how these initiatives could inter-relate can be found in [81].
[22]https://inspect.ai-safety-institute.org.uk/
[23]https://aiverifyfoundation.sg/project-moonshot/

**Provides opportunity for motivated actors to take harmful action.** Infrastructure, being broadly multi-purpose and applicable to many potential downstream applications and uses, may be particularly susceptible to misuse by malicious actors. For example, compute resources, while valuable for running large-scale safety research experiments, can also be used to develop systems for intentionally harmful purposes if suitable safeguards are not put in place. Additionally, infrastructural developments may be vulnerable to tampering on the part of a collaborator, for example through the deliberate insertion of backdoors.

*4.3.2 Subareas.* **Infrastructure for evaluating systems.** The development of infrastructure for conducting evaluations could aid in promoting an open ecosystem of research and evaluation into the capabilities and safety of AI systems.[24] In particular, standardising software infrastructure and methods for evaluating AI systems (perhaps facilitated by the open-sourcing of codebases) can aid in providing a shared understanding of the capabilities and risks of systems. It could also enhance interoperability, allowing countries to better interpret and build on the evaluations performed in other countries. While sharing future evaluation methods may help advance adversaries' capabilities by providing a clear metric for improvements or disclosing techniques to elicit upper-bound capabilities, cooperating on evaluation infrastructure would likely not incur these same risks.

**Building an international research cluster.** Some areas of cutting-edge AI safety research demand a large amount of computational resources due to the immense size of the models involved [13]. Some states are thus starting to invest in public compute infrastructure to be made available to academic researchers to catalyse domestic AI research, including on AI safety.[25] If a multilateral partnership between states is to collaborate on AI safety research, it may be necessary for the partnership to jointly develop and maintain shared computational resources. Additionally, a shared cluster could incorporate many of the verifiable computing suggestions covered in the previous subsection, as a proof-of-concept for their utility. This could bring the benefit of mutual trust in the computational infrastructure underlying shared research projects, and serve to pool resources to achieve efficiencies of scale, as in other large international research projects such as CERN or the International Space Station.

**Automatic control and supervision systems.** Alongside the evaluation of models, it is also important to evaluate the systems and infrastructure in which these models are deployed and controlled. The development of robust monitoring and control systems for AI is becoming increasingly critical as alignment and robustness properties of models remain fragile. Despite significant advances in alignment techniques such as RLHF, models remain vulnerable to jailbreaking [18]. While there is a growing industry focused on AI monitoring solutions,[26] current AI safety approaches offer limited guarantees and may not generalise well to more agentic systems or novel failure modes. The *'Benchmarks for the Evaluation of LLM Safeguards'* [32] represents an initial effort to systematically evaluate monitoring systems. There remains a need for more comprehensive evaluation frameworks that can assess monitoring systems' effectiveness against a broader range of threats, including cyber-attacks, potential biohazards, and intentional subversion by models [43, 44]. Cooperating on the evaluation of control and supervision systems does not seem to incur the same risks as cooperating on the direct evaluation of models, including potentially advancing a rival's strategic AI capabilities.

---

[24]See existing initiatives e.g. Inspect (https://inspect.ai-safety-institute.org.uk/) or Project Moonshot (https://aiverifyfoundation.sg/project-moonshot/).
[25]See e.g., https://nairrpilot.org/.
[26]See, for example, *LLM Guard* (https://github.com/protectai/llm-guard) *Azure AI Content Safety* (https://azure.microsoft.com/en-us/products/ai-services/ai-content-safety), and *Lakera Guard* (https://www.lakera.ai/lakera-guard).

## 4.4 Evaluation methodologies

Methods and resources for reliably evaluating the capabilities and safety of AI systems, for example through benchmarking, red-teaming, human uplift studies, or agent evaluations, have become a centrepoint of AI (safety) research, particularly by governments [33, 56, 68]. Cooperation on such methods could ensure interoperability, enabling jurisdictions to share and build upon each other's assessment results – creating a more efficient global system for AI evaluation. Global cooperation in particular could be especially useful for evaluating language models in multilingual settings, due to the varied language and culture expertise necessary to robustly evaluate such systems [84].

*4.4.1 Risks of cooperation.* **Advances global frontier capabilities.** As evaluations are predominantly about assessing a system's capabilities or safety, rather than making improvements in these areas, cooperation on AI evaluations are unlikely to directly advance the global frontier. Indirect effects may still occur, such as benchmark-chasing, whereby the existence of a challenging evaluation becomes a target, spurring increased effort in order to improve a system's score on that evaluation. However, this could itself be harnessed for benefit in the case of predominantly safety-focused evaluations.

   **Differentially advances a rival's capabilities.** Some evaluation methodologies specify elicitation techniques aimed at extracting upper-bound performance from systems on dangerous or dual-use tasks. For example, the instruction to language models to 'think step by step' generated significantly improved performance on many tasks [100], and is now incorporated directly into the system prompt of some reasoning models, as well as being a valuable technique for evaluations.Sharing elicitation techniques could thus be particularly sensitive due to the potential for direct application for improving a system's capabilities. However, cooperation aimed at improving the efficacy or robustness of evaluation techniques for non-sensitive system capabilities may not pose this risk. A separate concern is that having access to quantitative evaluation methodologies for potentially dangerous system capabilities could provide a guiding measure for an actor trying to develop a system with those capabilities.

   **Exposes other sensitive information.** Depending on the focal area of the evaluation, cooperation may require the sharing of sensitive information, for example if assessing a model's cyber-offensive capabilities. However, there are also many domains where this would not be the case, for example assessing a model's propensity to generate false or biased content. By focusing cooperative efforts on these latter evaluation subjects, rivals can avoid having to reveal sensitive, domain-specific information. Early proofs-of-concept suggest that it is also possible to create secure evaluation environments that preserve the privacy of test data [96].

   **Provides opportunity for motivated actors to take harmful action.** The degree to which cooperating on developing evaluations could allow motivated actors to take harmful actions depends to a large extent on the form that cooperation takes. For example, jointly developing benchmarks or other evaluation datasets seems to confer limited risk. However, more involved joint testing exercises may necessitate providing collaborators with privileged access to evaluation systems and infrastructure, as well as the models being tested, both of which could potentially be compromised.

*4.4.2 Subareas.* **Ensuring the reliability of evaluations.** Recent research has highlighted challenges in ensuring the reliability of AI evaluation methodologies [21]. For example, concern has been raised regarding how overlap between a model's training and test datasets can artificially inflate scores on benchmarks and other evaluations [103, 105]. Cooperation could focus on methods for improving the reliability of current approaches to AI evaluation.

   **Advancing the science of evaluations.** Finally, research cooperation could aim to advance the science of AI evaluation – that is a more rigorous understanding of the theoretical basis of AI

evaluation [16]. Since such a research agenda would be more foundational in nature, it is less likely that there will be a significant risk of leaking sensitive information regarding contemporary AI systems, or providing opportunities to cause harm.

## 5 Conclusion

Geopolitical rivals often have incentives to cooperate on strategic technologies, for example to address risks posed by that technology which span across national borders. However, such cooperation itself can pose risks which must be managed if the benefits of cooperation are fully realised. In this paper we have provided an overview of current international cooperation on AI in an important case of geopolitical rivalry, and outlined four sources of risk pertinent to cooperation on technical AI safety that are under-addressed in current risk mitigation strategies. Based on this, we assess the extent to which these risks may be realised by four areas of technical work in AI safety which have been suggested as potential areas of international cooperation, finding that development of verification mechanisms and protocols may be well-suited for cooperation. Future research could aim to extend the analysis of this paper by considering additional areas of technical AI safety research, propose more concrete policy recommendations concerning international cooperation on AI, or consider how states' domestic technology policy issues could relate to international cooperation. We hope that this paper can serve as a foundation for further international cooperation efforts to address the risks associated with AI.

# References

[1] 2012. ISO/IEC 19790:2012. https://www.iso.org/standard/52906.html

[2] 2023. Dealing with Risks in International Research Cooperation: Recommendations from the Deutsche Forschungsgemeinschaft. https://www.dfg.de/resource/blob/289704/585cb3b48bb8e9f5b6e57e0e0a0d700e/risiken-int-kooperationen-en-data.pdf

[3] 2024. *US AISI and UK AISI Joint Pre-Deployment Test: OpenAI o1.* Technical Report. https://www.nist.gov/system/files/documents/2024/12/18/US_UK_AI%20Safety%20Institute_%20December_Publication-OpenAIo1.pdf

[4] Onni Aarne, Tim Fist, and Caleb Withers. 2024. *Secure, Governable Chips: Using On-Chip Mechanisms to Manage National Security Risks from AI & Advanced Computing.* Technical Report. Center for a New American Security. https://www.cnas.org/publications/reports/secure-governable-chips

[5] Sumaya N. Adan, Robert Trager, Kayla Blomquist, Claire Dennis, Gemma Edom, Lucia Velasco, Cecil Abungu, Ben Garfinkel, Julian Jacobs, Chinasa T. Okolo, Boxi Wu, and Jai Vipra. 2024. *Voice and Access in AI: Global AI Majority Participation in Artificial Intelligence Development and Governance.* Technical Report. Oxford Martin AI Governance Initiative, Oxford, UK. https://www.oxfordmartin.ox.ac.uk/publications/voice-and-access-in-ai-global-ai-majority-participation-in-artificial-intelligence-development-and-governance

[6] Jide Alaga, Jonas Schuett, and Markus Anderljung. 2024. A Grading Rubric for AI Safety Frameworks. doi:10.48550/arXiv.2409.08751 arXiv:2409.08751 [cs].

[7] Anthropic. 2024. Responsible Scaling Policy. https://assets.anthropic.com/m/24a47b00f10301cd/original/Anthropic-Responsible-Scaling-Policy-2024-10-15.pdf

[8] Arms Control Association. 2021. The Open Skies Treaty at a Glance. https://www.armscontrol.org/factsheets/openskies

[9] Christel Baier and Joost-Pieter Katoen. 2008. *Principles of Model Checking.* The MIT Press. https://mitpress.mit.edu/9780262026499/principles-of-model-checking/

[10] Yoshua Bengio, Sören Mindermann, Daniel Privitera, Tamay Besiroglu, Rishi Bommasani, Stephen Casper, Yejin Choi, Philip Fox, Ben Garfinkel, Danielle Goldfarb, Hoda Heidari, Anson Ho, Sayash Kapoor, Leila Khalatbari, Shayne Longpre, Sam Manning, Vasilios Mavroudis, Mantas Mazeika, Julian Michael, Jessica Newman, Kwan Yee Ng, Chinasa T. Okolo, Deborah Raji, Girish Sastry, Elizabeth Seger, Theodora Skeadas, Tobin South, Emma Strubell, Florian Tramèr, Lucia Velasco, and Nicole Wheeler. 2025. *International AI Safety Report: The International Scientific Report on the Safety of Advanced AI.* Technical Report. Department for Science, Innovation & Technology. https://www.gov.uk/government/publications/international-ai-safety-report-2025

[11] Yoshua Bengio, Sören Mindermann, Daniel Privitera, Rishi Bommasani, Stephen Casper, Yejin Choi, Danielle Goldfarb, Hoda Heidari, Leila Khalatbari, Shayne Longpre, Vasilios Mavroudis, Mantas Mazeika, Kwan Yee Ng, Chinasa T. Okolo, Deborah Raji, Theodora Skeadas, and Florian Tramèr. 2024. *International Scientific Report on the Safety of Advanced AI - Interim Report.* Technical Report. Department for Science, Innovation & Technology. https://www.gov.uk/government/publications/international-scientific-report-on-the-safety-of-advanced-ai

[12] Emily Benson. 2023. Updated October 7 Semiconductor Export Controls. https://www.csis.org/analysis/updated-october-7-semiconductor-export-controls

[13] Tamay Besiroglu, Sage Andrus Bergerson, Amelia Michael, Lennart Heim, Xueyun Luo, and Neil Thompson. 2024. The Compute Divide in Machine Learning: A Threat to Academic Contribution and Scrutiny? doi:10.48550/arXiv.2401.02452 arXiv:2401.02452 [cs].

[14] Fateh Boudardara, Abderraouf Boussif, Pierre-Jean Meyer, and Mohamed Ghazel. 2024. A Review of Abstraction Methods Toward Verifying Neural Networks. *ACM Trans. Embed. Comput. Syst.* 23, 4 (June 2024), 58:1–58:19. doi:10.1145/3617508

[15] Miles Brundage, Shahar Avin, Jasmine Wang, Haydn Belfield, Gretchen Krueger, Gillian Hadfield, Heidy Khlaaf, Jingying Yang, Helen Toner, Ruth Fong, Tegan Maharaj, Pang Wei Koh, Sara Hooker, Jade Leung, Andrew Trask, Emma Bluemke, Jonathan Lebensold, Cullen O'Keefe, Mark Koren, Théo Ryffel, J. B. Rubinovitz, Tamay Besiroglu, Federica Carugati, Jack Clark, Peter Eckersley, Sarah de Haas, Maritza Johnson, Ben Laurie, Alex Ingerman, Igor Krawczuk, Amanda Askell, Rosario Cammarota, Andrew Lohn, David Krueger, Charlotte Stix, Peter Henderson, Logan Graham, Carina Prunkl, Bianca Martin, Elizabeth Seger, Noa Zilberman, Seán Ó hÉigeartaigh, Frens Kroeger, Girish Sastry, Rebecca Kagan, Adrian Weller, Brian Tse, Elizabeth Barnes, Allan Dafoe, Paul Scharre, Ariel Herbert-Voss, Martijn Rasser, Shagun Sodhani, Carrick Flynn, Thomas Krendl Gilbert, Lisa Dyer, Saif Khan, Yoshua Bengio, and Markus Anderljung. 2020. Toward Trustworthy AI Development: Mechanisms for Supporting Verifiable Claims. doi:10.48550/arXiv.2004.07213 arXiv:2004.07213 [cs].

[16] John Burden. 2024. Evaluating AI Evaluation: Perils and Prospects. doi:10.48550/arXiv.2407.09221 arXiv:2407.09221 [cs] version: 1.

[17] CAICT. 2024. Protecting AI security and building a model of industry self-discipline - the first batch of 17 companies signed the "Artificial Intelligence Security Commitment". https://mp.weixin.qq.com/s/s-XFKQCWhu0uye4opgb3Ng

[18] Stephen Casper, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, Jérémy Scheurer, Javier Rando, Rachel Freedman, Tomasz Korbak, David Lindner, Pedro Freire, Tony Wang, Samuel Marks, Charbel-Raphaël Segerie, Micah Carroll, Andi Peng, Phillip Christoffersen, Mehul Damani, Stewart Slocum, Usman Anwar, Anand Siththaranjan, Max Nadeau, Eric J. Michaud, Jacob Pfau, Dmitrii Krasheninnikov, Xin Chen, Lauro Langosco, Peter Hase, Erdem Bıyık, Anca Dragan, David Krueger, Dorsa Sadigh, and Dylan Hadfield-Menell. 2023. Open Problems and Fundamental Limitations of Reinforcement Learning from Human Feedback. doi:10.48550/arXiv.2307.15217 arXiv:2307.15217 [cs].

[19] Xavier Castañer and Nuno Oliveira. 2020. Collaboration, Coordination, and Cooperation Among Organizations: Establishing the Distinctive Meanings of These Terms Through a Systematic Literature Review. *Journal of Management* 46, 6 (July 2020), 965–1001. doi:10.1177/0149206320901565 Publisher: SAGE Publications Inc.

[20] Seth Center and Emma Bates. 2019. *Tech-Politik: Historical Perspectives on Innovation, Technology, and Strategic Competition.* Technical Report. Center for Strategic & International Studies. https://www.csis.org/analysis/tech-politik-historical-perspectives-innovation-technology-and-strategic-competition

[21] Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, Wei Ye, Yue Zhang, Yi Chang, Philip S. Yu, Qiang Yang, and Xing Xie. 2024. A Survey on Evaluation of Large Language Models. *ACM Trans. Intell. Syst. Technol.* 15, 3 (March 2024), 39:1–39:45. doi:10.1145/3641289

[22] Dami Choi, Yonadav Shavit, and David Duvenaud. 2023. Tools for Verifying Neural Models' Training Data. doi:10.48550/arXiv.2307.00682 arXiv:2307.00682 [cs].

[23] Cathleen D Cimino-Isaacs and Karen M Sutter. 2024. Committee on Foreign Investment in the United States (CFIUS). https://crsreports.congress.gov/product/pdf/IF/IF10177

[24] Daniel Clery. 2024. Giant fusion project is in big trouble: ITER operations delayed to 2034, with energy-producing reactions expected 5 years later. *Science* 385, 6704 (July 2024), 10–11. https://www.science.org/content/article/giant-international-fusion-project-big-trouble

[25] Andrew J. Coe and Jane Vaynman. 2020. Why Arms Control Is So Rare. *American Political Science Review* 114, 2 (May 2020), 342–355. doi:10.1017/S000305541900073X

[26] Josef Dai, Xuehai Pan, Ruiyang Sun, Jiaming Ji, Xinbo Xu, Mickel Liu, Yizhou Wang, and Yaodong Yang. 2023. Safe RLHF: Safe Reinforcement Learning from Human Feedback. doi:10.48550/arXiv.2310.12773 arXiv:2310.12773 [cs].

[27] David "davidad" Dalrymple. 2024. *Safeguarded AI: Constructing guaranteed safety.* Technical Report. Advanced Research + Invention Agency. https://www.aria.org.uk/media/3nhijno4/aria-safeguarded-ai-programme-thesis-v1.pdf

[28] David "davidad" Dalrymple, Joar Skalse, Yoshua Bengio, Stuart Russell, Max Tegmark, Sanjit Seshia, Steve Omohundro, Christian Szegedy, Ben Goldhaber, Nora Ammann, Alessandro Abate, Joe Halpern, Clark Barrett, Ding Zhao, Tan Zhi-Xuan, Jeannette Wing, and Joshua Tenenbaum. 2024. Towards Guaranteed Safe AI: A Framework for Ensuring Robust and Reliable AI Systems. doi:10.48550/arXiv.2405.06624 arXiv:2405.06624 [cs].

[29] Sumanth Dathathri, Abigail See, Sumedh Ghaisas, Po-Sen Huang, Rob McAdam, Johannes Welbl, Vandana Bachani, Alex Kaskasoli, Robert Stanforth, Tatiana Matejovicova, Jamie Hayes, Nidhi Vyas, Majd Al Merey, Jonah Brown-Cohen, Rudy Bunel, Borja Balle, Taylan Cemgil, Zahra Ahmed, Kitty Stacpoole, Ilia Shumailov, Ciprian Baetu, Sven Gowal, Demis Hassabis, and Pushmeet Kohli. 2024. Scalable watermarking for identifying large language model outputs. *Nature* 634, 8035 (Oct. 2024), 818–823. doi:10.1038/s41586-024-08025-4 Publisher: Nature Publishing Group.

[30] Google DeepMind. 2024. Frontier Safety Framework v1.0.

[31] Jeffrey Ding. 2024. Keep your enemies safer: technical cooperation and transferring nuclear safety and security technologies. *European Journal of International Relations* 30, 4 (Dec. 2024), 918–945. doi:10.1177/13540661241246622 Publisher: SAGE Publications Ltd.

[32] Diego Dorn, Alexandre Variengien, Charbel-Raphaël Segerie, and Vincent Corruble. 2024. BELLS: A Framework Towards Future Proof Benchmarks for the Evaluation of LLM Safeguards. doi:10.48550/arXiv.2406.01364 arXiv:2406.01364 [cs].

[33] Michael Feffer, Anusha Sinha, Wesley Hanwen Deng, Zachary C. Lipton, and Hoda Heidari. 2024. Red-Teaming for Generative AI: Silver Bullet or Security Theater? doi:10.48550/arXiv.2401.15897 arXiv:2401.15897 [cs].

[34] Jaden Fiotto-Kaufman, Alexander R. Loftus, Eric Todd, Jannik Brinkmann, Koyena Pal, Dmitrii Troitskii, Michael Ripa, Adam Belfki, Can Rager, Caden Juang, Aaron Mueller, Samuel Marks, Arnab Sen Sharma, Francesca Lucchetti, Nikhil Prakash, Carla Brodley, Arjun Guha, Jonathan Bell, Byron C. Wallace, and David Bau. 2025. NNsight and NDIF: Democratizing Access to Open-Weight Foundation Model Internals. doi:10.48550/arXiv.2407.14561 arXiv:2407.14561 [cs].

[35] FLI. 2024. *FLI AI Safety Index 2024: Independent experts evaluate safety practices of leading AI companies across critical domains.* Technical Report. Future of Life Institute. http://futureoflife.org/index

[36] Center for Arms Control and Non-Proliferation. 2017. Fact Sheet: The Threshold Test Ban Treaty (TTBT). https://armscontrolcenter.org/fact-sheet-threshold-test-ban-treaty-ttbt/

[37] Department for Science Innovation and Technology. 2023. The Bletchley Declaration by Countries Attending the AI Safety Summit, 1-2 November 2023. https://www.gov.uk/government/publications/ai-safety-summit-2023-the-bletchley-declaration/the-bletchley-declaration-by-countries-attending-the-ai-safety-summit-1-2-november-2023

[38] Department for Science Innovation and Technology. 2024. Frontier AI Safety Commitments, AI Seoul Summit 2024. https://www.gov.uk/government/publications/frontier-ai-safety-commitments-ai-seoul-summit-2024/frontier-ai-safety-commitments-ai-seoul-summit-2024

[39] Department for Science Innovation and Technology. 2024. Seoul Ministerial Statement for advancing AI safety, innovation and inclusivity: AI Seoul Summit 2024. https://www.gov.uk/government/publications/seoul-ministerial-statement-for-advancing-ai-safety-innovation-and-inclusivity-ai-seoul-summit-2024/seoul-ministerial-statement-for-advancing-ai-safety-innovation-and-inclusivity-ai-seoul-summit-2024

[40] Nancy W. Gallagher. 1997. The politics of verification: Why 'how much?' Is not enough. *Contemporary Security Policy* 18, 2 (Aug. 1997), 138–170. doi:10.1080/13523269708404165 Publisher: Routledge.

[41] Sanjam Garg, Aarushi Goel, Somesh Jha, Saeed Mahloujifar, Mohammad Mahmoody, Guru-Vamsi Policharla, and Mingyuan Wang. 2023. Experimenting with Zero-Knowledge Proofs of Training. https://eprint.iacr.org/2023/1345 Publication info: Published elsewhere. Major revision. ACM CCS 2023.

[42] Soumya Suvra Ghosal, Souradip Chakraborty, Jonas Geiping, Furong Huang, Dinesh Manocha, and Amrit Bedi. 2023. A Survey on the Possibilities & Impossibilities of AI-generated Text Detection. *Transactions on Machine Learning Research* (Oct. 2023). https://openreview.net/forum?id=AXtFeYjboj

[43] Ryan Greenblatt, Buck Shlegeris, Kshitij Sachan, and Fabien Roger. 2024. AI Control: Improving Safety Despite Intentional Subversion. doi:10.48550/arXiv.2312.06942 arXiv:2312.06942 [cs].

[44] Charlie Griffin, Louis Thomson, Buck Shlegeris, and Alessandro Abate. 2024. Games for AI Control: Models of Safety Evaluations of AI Deployment Protocols. doi:10.48550/arXiv.2409.07985 arXiv:2409.07985 [cs].

[45] Oliver Guest, Michael Aird, and Seán Ó hÉigeartaigh. 2023. *Safeguarding the Safeguards: How best to promote AI alignment in the public interest.* Technical Report. Institute for AI Policy and Strategy. https://www.iaps.ai/research/safeguarding-the-safeguards

[46] Oliver Guest and Zoe Williams. 2024. *Topics for track IIs: What can be discussed in dialogues about advanced AI risks without leaking sensitive information?* Technical Report. Institute for AI Policy and Strategy. https://static1.squarespace.com/static/64edf8e7f2b10d716b5ba0e1/t/6633b93601a0553b73d56095/1714665783885/%5BFinal%5D+Topics+for+track+IIs.pdf

[47] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2015. Deep Residual Learning for Image Recognition. doi:10.48550/arXiv.1512.03385 arXiv:1512.03385 [cs].

[48] Dan Hendrycks. 2024. *Introduction to AI Safety, Ethics and Society.* Taylor & Francis. https://www.aisafetybook.com/

[49] Dan Hendrycks, Mantas Mazeika, and Thomas Woodside. 2023. An Overview of Catastrophic AI Risks. doi:10.48550/arXiv.2306.12001 arXiv:2306.12001 [cs].

[50] Eric L. Hirschhorn, Brian J. Egan, Edward J. Krauland, Eric L. Hirschhorn, Brian J. Egan, and Edward J. Krauland. 2022. *U.S. Export Controls and Economic Sanctions* (fourth edition, fourth edition ed.). Oxford University Press, Oxford, New York.

[51] C. A. R. Hoare. 1969. An axiomatic basis for computer programming. *Commun. ACM* 12, 10 (Oct. 1969), 576–580. doi:10.1145/363235.363259

[52] The White House. 2022. Blueprint for an AI Bill of Rights: Making automated systems work for the American people. https://www.whitehouse.gov/ostp/ai-bill-of-rights/

[53] Yukon Huang, Isaac B. Kardon, and Matt Sheehan. 2023. Three Takeaways From the Biden-Xi Meeting. https://carnegieendowment.org/posts/2023/11/three-takeaways-from-the-biden-xi-meeting?lang=en

[54] IAEA. 2014. IAEA Safeguards Overview. https://www.iaea.org/publications/factsheets/iaea-safeguards-overview Publisher: IAEA.

[55] UK AI Safety Institute. 2024. Conference on frontier AI safety frameworks. https://www.aisi.gov.uk/work/conference-on-frontier-ai-safety-frameworks

[56] UK AI Safety Institute. 2024. Early lessons from evaluating frontier AI systems. https://www.aisi.gov.uk/work/early-lessons-from-evaluating-frontier-ai-systems

[57] Robert Jervis. 1978. Cooperation Under the Security Dilemma. *World Politics* 30, 2 (1978), 167–214. doi:10.2307/2009958 Publisher: [Trustees of Princeton University, The Johns Hopkins University Press].

[58] Wang Jingjing. 2016. The Whampoa Academy of China's Internet. https://weibo.com/p/1001643998598932131471 English commentary and translation by Jeffrey Ding available at https://chinai.substack.com/p/chinai-37-happy-20th-anniversary.

[59] Daniel Kang, Tatsunori Hashimoto, Ion Stoica, and Yi Sun. 2022. Scaling up Trustless DNN Inference with Zero-Knowledge Proofs. doi:10.48550/arXiv.2210.08674 arXiv:2210.08674 [cs].

[60] Holden Karnofsky. 2024. If-Then Commitments for AI Risk Reduction. https://carnegieendowment.org/research/2024/09/if-then-commitments-for-ai-risk-reduction?lang=en

[61] Leonie Koessler, Jonas Schuett, and Markus Anderljung. 2024. Risk thresholds for frontier AI. doi:10.48550/arXiv.2406.14713 arXiv:2406.14713 [cs].

[62] Stephen D. Krasner. 1991. Global Communications and National Power: Life on the Pareto Frontier. *World Politics* 43, 3 (April 1991), 336–366. doi:10.2307/2010398

[63] Gabriel Kulp, Daniel Gonzales, Everett Smith, Lennart Heim, Prateek Puri, Michael J. D. Vermeer, and Zev Winkelman. 2024. *Hardware-Enabled Governance Mechanisms: Developing Technical Solutions to Exempt Items Otherwise Classified Under Export Control Classification Numbers 3A090 and 4A090.* Technical Report. RAND Corporation. https://www.rand.org/pubs/working_papers/WRA3056-1.html

[64] Nancy G. Leveson and John P. Thomas. 2023. Certification of Safety-Critical Systems. *Commun. ACM* 66, 10 (Sept. 2023), 22–26. doi:10.1145/3615860

[65] Michael Martina and Trevor Hunnicutt. 2024. US, China meet in Geneva to discuss AI risks. *Reuters* (May 2024). https://www.reuters.com/technology/us-china-meet-geneva-discuss-ai-risks-2024-05-13/

[66] Jean-Christophe Mauduit. 2017. Collaboration around the International Space Station: Science for diplomacy and its implication for U.S.-Russia and China relations. In *Proceedings of 7th Annual SAIS Asia Conference (SAIS 2018)*, Vol. 462. Washington, DC, 412–413. doi:10.1038/462412a

[67] Chris Miller. 2023. *Chip War: The Fight for the World's Most Critical Technology.* Simon & Schuster. https://www.simonandschuster.co.uk/books/Chip-War/Chris-Miller/9781398504127

[68] Christopher A. Mouton, Caleb Lucas, and Ella Guest. 2023. *The Operational Risks of AI in Large-Scale Biological Attacks: A Red-Team Approach.* Technical Report. RAND Corporation. https://www.rand.org/pubs/research_reports/RRA2977-1.html

[69] National Institute of Standards and Technology (US). 2024. *Artificial Intelligence Risk Management Framework: Generative Artificial Intelligence Profile.* Technical Report NIST AI 600-1. National Institute of Standards and Technology (U.S.), Gaithersburg, MD. error: 600–1 pages. doi:10.6028/NIST.AI.600-1

[70] Peter Naur. 1966. Proof of algorithms by general snapshots. *BIT Numerical Mathematics* 6, 4 (July 1966), 310–316. doi:10.1007/BF01966091

[71] Sella Nevo, Dan Lahav, Ajay Karpur, Yogev Bar-On, Henry Alexander Bradley, and Jeff Alstott. 2024. *Securing AI Model Weights: Preventing Theft and Misuse of Frontier Models.* Technical Report. RAND Corporation. https://www.rand.org/pubs/research_reports/RRA2849-1.html

[72] Emerging Technology Observatory. 2024. Country Activity Tracker (CAT): Artificial Intelligence. https://cat.eto.tech/

[73] International Network of AI Safety Institutes. 2024. *Improving International Testing of Foundation Models: A Pilot Testing Exercise from the International Network of AI Safety Institutes.* Technical Report. San Francisco. https://www.nist.gov/system/files/documents/2024/11/21/Improving%20International%20Testing%20of%20Foundation%20Models-%20%20%20A%20Pilot%20Testing%20Exercise%20from%20the%20International%20Network%20of%20AI%20Safety%20Institutes.pdf

[74] U.S. Department of Commerce. 2024. U.S. Secretary of Commerce Raimondo and U.S. Secretary of State Blinken Announce Inaugural Convening of International Network of AI Safety Institutes in San Francisco | U.S. Department of Commerce. https://www.commerce.gov/news/press-releases/2024/09/us-secretary-commerce-raimondo-and-us-secretary-state-blinken-announce

[75] OpenAI. 2023. Preparedness Framework (Beta). https://cdn.openai.com/openai-preparedness-framework-beta.pdf

[76] World Health Organization. 2021. Ethics and Governance of Artificial Intelligence for Health: WHO guidance. https://www.who.int/publications/i/item/9789240029200

[77] Joe O'Brien, Shaun Ee, Jam Kraprayoon, Bill Anderson-Samways, Oscar Delaney, and Zoe Williams. 2024. *Coordinated Disclosure of Dual-Use Capabilities: An Early Warning System for Advanced AI.* Technical Report. Institute for AI Policy and Strategy. https://www.iaps.ai/research/coordinated-disclosure

[78] Sean O'Connor. 2019. *How Chinese Companies Facilitate Technology Transfer from the United States.* Staff Research Report. U.S.-China Economic and Security Review Commission, Washington, DC. https://www.uscc.gov/sites/default/files/Research/How%20Chinese%20Companies%20Facilitate%20Tech%20Transfer%20from%20the%20US.pdf

[79] Qianqian Pan, Mianxiong Dong, Kaoru Ota, and Jun Wu. 2022. Device-Bind Key-Storageless Hardware AI Model IP Protection: A PUF and Permute-Diffusion Encryption-Enabled Approach. doi:10.48550/arXiv.2212.11133 arXiv:2212.11133 [cs].

[80] James Petrie, Onni Aarne, Nora Ammann, and David "davidad" Dalrymple. 2024. *Interim Report: Mechanisms for Flexible Hardware-Enabled Guarantees.* Technical Report.

[81] Hadrien Pouget, Claire Dennis, Jon Bateman, Robert F. Trager, Renan Araujo, Haydn Belfield, Belinda Cleeland, Malou Estier, Gideon Futerman, Oliver Guest, Carlos Ignacio Gutierrez, Vishnu Kannan, Casey Mahoney, Matthijs Maas, Charles Martinet, Jakob Mökander, Kwan Yee Ng, Seán Ó hÉigeartaigh, Aidan Peppin, Konrad Seifert, Scott

Singer, Maxime Stauffer, Caleb Withers, and Marta Ziosi. 2024. *The Future of International Scientific Assessments of AI's Risks*. Technical Report. Oxford Martin AI Governance Initiative, Oxford, UK. https://www.oxfordmartin.ox.ac.uk/publications/the-future-of-international-scientific-assessments-of-ais-risks

[82] Jarrett Renshaw and Trevor Hunnicutt. 2024. Biden, Xi agree that humans, not AI, should control nuclear arms. *Reuters* (Nov. 2024). https://www.reuters.com/world/biden-xi-agreed-that-humans-not-ai-should-control-nuclear-weapons-white-house-2024-11-16/

[83] Anka Reuel, Ben Bucknall, Stephen Casper, Tim Fist, Lisa Soder, Onni Aarne, Lewis Hammond, Lujain Ibrahim, Alan Chan, Peter Wills, Markus Anderljung, Ben Garfinkel, Lennart Heim, Andrew Trask, Gabriel Mukobi, Rylan Schaeffer, Mauricio Baker, Sara Hooker, Irene Solaiman, Alexandra Sasha Luccioni, Nitarshan Rajkumar, Nicolas Moës, Jeffrey Ladish, Neel Guha, Jessica Newman, Yoshua Bengio, Tobin South, Alex Pentland, Sanmi Koyejo, Mykel J. Kochenderfer, and Robert Trager. 2024. Open Problems in Technical AI Governance. http://arxiv.org/abs/2407.14981 arXiv:2407.14981 [cs].

[84] Angelika Romanou, Negar Foroutan, Anna Sotnikova, Zeming Chen, Sree Harsha Nelaturu, Shivalika Singh, Rishabh Maheshwary, Micol Altomare, Mohamed A. Haggag, Snegha A, Alfonso Amayuelas, Azril Hafizi Amirudin, Viraat Aryabumi, Danylo Boiko, Michael Chang, Jenny Chim, Gal Cohen, Aditya Kumar Dalmia, Abraham Diress, Sharad Duwal, Daniil Dzenhaliou, Daniel Fernando Erazo Florez, Fabian Farestam, Joseph Marvin Imperial, Shayekh Bin Islam, Perttu Isotalo, Maral Jabbarishiviari, Börje F. Karlsson, Eldar Khalilov, Christopher Klamm, Fajri Koto, Dominik Krzemiński, Gabriel Adriano de Melo, Syrielle Montariol, Yiyang Nan, Joel Niklaus, Jekaterina Novikova, Johan Samir Obando Ceron, Debjit Paul, Esther Ploeger, Jebish Purbey, Swati Rajwal, Selvan Sunitha Ravi, Sara Rydell, Roshan Santhosh, Drishti Sharma, Marjana Prifti Skenduli, Arshia Soltani Moakhar, Bardia Soltani Moakhar, Ran Tamir, Ayush Kumar Tarun, Azmine Toushik Wasi, Thenuka Ovin Weerasinghe, Serhan Yilmaz, Mike Zhang, Imanol Schlag, Marzieh Fadaee, Sara Hooker, and Antoine Bosselut. 2024. INCLUDE: Evaluating Multilingual Language Understanding with Regional Knowledge. doi:10.48550/arXiv.2411.19799 arXiv:2411.19799 [cs].

[85] Tim Rühlig. 2023. The Geopolitics of Technical Standardization. https://dgap.org/en/research/publications/geopolitics-technical-standardization

[86] Girish Sastry, Lennart Heim, Haydn Belfield, Markus Anderljung, Miles Brundage, Julian Hazell, Cullen O'Keefe, Gillian K. Hadfield, Richard Ngo, Konstantin Pilz, George Gor, Emma Bluemke, Sarah Shoker, Janet Egan, Robert F. Trager, Shahar Avin, Adrian Weller, Yoshua Bengio, and Diane Coyle. 2024. Computing Power and the Governance of Artificial Intelligence. doi:10.48550/arXiv.2402.08797 arXiv:2402.08797 [cs].

[87] Jonas Schuett, Markus Anderljung, Alexis Carlier, Leonie Koessler, and Ben Garfinkel. 2024. From Principles to Rules: A Regulatory Approach for Frontier AI. doi:10.48550/arXiv.2407.07300 arXiv:2407.07300 [cs].

[88] Yonadav Shavit. 2023. What does it take to catch a Chinchilla? Verifying Rules on Large-Scale Neural Network Training via Compute Monitoring. doi:10.48550/arXiv.2303.11341 arXiv:2303.11341 [cs].

[89] Matt Sheehan and Jacob Feldgoise. 2023. What Washington Gets Wrong About China and Technical Standards. https://carnegieendowment.org/research/2023/02/what-washington-gets-wrong-about-china-and-technical-standards?lang=en

[90] Scott Singer. 2024. How the UK Should Engage China at AI's Frontier. https://carnegieendowment.org/posts/2024/10/lammy-china-ai-safety-cooperation?lang=en

[91] Tobin South, Alexander Camuto, Shrey Jain, Shayla Nguyen, Robert Mahari, Christian Paquin, Jason Morton, and Alex 'Sandy' Pentland. 2024. Verifiable evaluations of machine learning models using zkSNARKs. doi:10.48550/arXiv.2402.02675 arXiv:2402.02675 [cs].

[92] Arthur A. Stein. 1982. Coordination and collaboration: regimes in an anarchic world. *International Organization* 36, 2 (1982), 299–324. doi:10.1017/S0020818300018968

[93] Merlin Stein, Jamie Bernardi, and Connor Dunlop. 2024. The Role of Governments in Increasing Interconnected Post-Deployment Monitoring of AI. doi:10.48550/arXiv.2410.04931 arXiv:2410.04931 [cs].

[94] Merlin Stein and Connor Dunlop. 2024. Safe beyond sale: post-deployment monitoring of AI. https://www.adalovelaceinstitute.org/blog/post-deployment-monitoring-of-ai/

[95] Mengqi Sun. 2024. U.S., China to Cooperate in the Fight Against Dirty Money. *Wall Street Journal* (April 2024). https://www.wsj.com/articles/u-s-china-to-cooperate-in-the-fight-against-dirty-money-1edb9a25

[96] Andrew Trask, Aziz Berkay Yesilyurt, Bennett Farkas, Callis Ezenwaka, Carmen Popa, Dave Buckley, Eelco van der Wel, Francesco Mosconi, Grace Han, Ionesio Junior, Irina Bejan, Ishan Mishra, Khoa Nguyen, Koen van der Veen, Kyoko Eng, Lacey Strahm, Logan Graham, Madhava Jay, Matei Simtinica, Osam Kyemenu-Sarsah, Peter Smith, Rasswanth S, Ronnie Falcon, Sameer Wagh, Sandeep Mandala, Shubham Gupta, Stephen Gabriel, Subha Ramkumar, Tauquir Ahmed, Teo Milea, Valerio Maggio, Yash Gorana, and Zarreen Reza. 2024. Secure Enclaves for AI Evaluation. https://blog.openmined.org/secure-enclaves-for-ai-evaluation/

[97] UKRI. 2022. Managing risks in international research and innovation: An overview of higher education sector guidance. https://www.ukri.org/wp-content/uploads/2022/07/UKRI-07072022-managing-risks-in-international-

research-and-innovation-uuk-cpni-ukri_1.pdf

[98]  UN HLAB. 2024. *Govering AI for Humanity: Final Report.* Technical Report. United Nations, New York, NY.  https://www.un.org/ai-advisory-body

[99]  Caterina Urban and Antoine Miné. 2021. A Review of Formal Methods applied to Machine Learning. doi:10.48550/arXiv.2104.02466 arXiv:2104.02466 [cs].

[100]  Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. doi:10.48550/arXiv.2201.11903 arXiv:2201.11903 [cs].

[101]  Christian Schroeder de Witt, Samuel Sokota, J. Zico Kolter, Jakob Foerster, and Martin Strohmeier. 2023. Perfectly Secure Steganography Using Minimum Entropy Coupling. doi:10.48550/arXiv.2210.14889 arXiv:2210.14889 [cs].

[102]  JoAnne Yates and Craig N. Murphy. 2019. *Engineering Rules.* Johns Hopkins University Press. doi:10.1353/book.66187

[103]  Andy K. Zhang, Kevin Klyman, Yifan Mai, Yoav Levine, Yian Zhang, Rishi Bommasani, and Percy Liang. 2024. Language model developers should report train-test overlap. doi:10.48550/arXiv.2410.08385 arXiv:2410.08385 [cs].

[104]  Daniel Zhang, Nestor Maslej, Erik Brynjolfsson, John Etchemendy, Terah Lyons, James Manyika, Helen Ngo, Juan Carlos Niebles, Michael Sellitto, Ellie Sakhaee, Yoav Shoham, Jack Clark, and Raymond Perrault. 2022. *The AI Index 2022 Annual Report.* Technical Report. Stanford Institute for Human-Centered AI, Stanford University. https://aiindex.stanford.edu/wp-content/uploads/2022/03/2022-AI-Index-Report_Master.pdf

[105]  Kun Zhou, Yutao Zhu, Zhipeng Chen, Wentong Chen, Wayne Xin Zhao, Xu Chen, Yankai Lin, Ji-Rong Wen, and Jiawei Han. 2023. Don't Make Your LLM an Evaluation Benchmark Cheater. doi:10.48550/arXiv.2311.01964 arXiv:2311.01964 [cs].

[106]  Nicholas Zúñiga, Saheli Datta Burton, Filippo Blancato, and Madeline Carr. 2024. The geopolitics of technology standards: historical context for US, EU and Chinese approaches. *International Affairs* 100, 4 (July 2024), 1635–1652. doi:10.1093/ia/iiae124

## A   Summary table of candidate areas of cooperation

Table 1 provides a summary of the candidate areas for cooperation, as described in section 4.

| | Advances global frontier capabilities | Differentially advances a rival's capabilities | Exposes other sensitive information | Provides opportunity for harmful action |
|---|---|---|---|---|
| **Verification** | *Minimal* Verification largely aims to attest to claims made about a given system/process, rather than improve it. While some areas of verification research could uncover new properties that advance system capabilities, these can be avoided when cooperating. | *Minimal* While verification methods may involve revealing a set of model properties that need to be verified, the risk from this is typically low as the properties to be verified can be restricted to those known to all parties in advance. | *Minimal/moderate* The development of verification technologies, such as hardware-enabled mechanisms, could require disclosing sensitive information. However, methods that are independent of the specification of the AI system and related infrastructure are also being developed. | *Minimal/moderate* When rivals jointly develop verification systems, they may insert hidden backdoors to fake compliance. While the difficulty of concealing such backdoors varies, this risk can be managed through open-source development and bug bounties to help detect subversion attempts. |
| **Protocols** | *Minimal* Developing protocols is a process of codifying existing knowledge, rather than extending the knowledge frontier. | *Minimal* Cooperation on protocols development could focus on areas where all parties have shared knowledge. Thus no expertise would be shared that could allow a rival to advance their capabilities. | *Minimal* Developing protocols and standards usually only draws upon broadly-known knowledge. Any areas requiring national security sensitive information could be avoided. | *Minimal* Protocol development does not involve any direct manipulation of AI systems, though standardisation has sometimes been used to advance unilateral interests |
| **Infrastructure** | *Minimal/moderate* Infrastructure, being multi-purpose, could be useful, though not critical, in facilitating capabilities advances. | *Minimal/moderate* Infrastructure, being multi-purpose, could be instrumental, though not critical, in allowing a rival to advance their capabilities. Though the use of shared infrastructure could increase the likelihood that such applications are detected by cooperating parties. | *Moderate* To the extent that developing shared infrastructure builds upon existing infrastructure, doing so may require divulging sensitive information regarding this existing infrastructure. | *Moderate* Infrastructure's multi-purpose nature makes it vulnerable to misuse – compute resources meant for beneficial AI research could be repurposed for harmful aims, and collaborative infrastructure projects risk tampering through deliberately inserted backdoors. |
| **Evaluations** | *Minimal/moderate* Evaluations concern the assessment of a system's capabilities rather than their improvement. While dangerous capability evaluations may specify capability elicitation techniques, the impact of this on the overall capabilities frontier is minimal compared to dedicated efforts to improve model capabilities. | *Minimal/moderate* Evaluation methods are closely tied to techniques for eliciting upper-bound capabilities from AI systems. Collaboration on evaluations may thus allow a rival to advance their capabilities through having a better awareness or understanding of such techniques. | *Moderate* Evaluations in sensitive areas (such as CBRN/Cyber) can involve a large amount of specialist domain knowledge. Collaborating on evaluations in these domains risks the transfer of this sensitive knowledge. | *Minimal/moderate* Sharing benchmark methodologies is unlikely to provide significant opportunity for harm. However, joint testing wherein a rival receives access to a given model could provide greater opportunities for harm, although the chance that this happens completely undetected is still low. |

Table 1. Preliminary assessment of the risks of cooperating on the four technical AI safety areas discussed.