

# Model evaluation for extreme risks

Toby Shevlane<sup>1</sup>, Sebastian Farquhar<sup>1</sup>, Ben Garfinkel<sup>2</sup>, Mary Phuong<sup>1</sup>, Jess Whittlestone<sup>3</sup>, Jade Leung<sup>4</sup>, Daniel Kokotajlo<sup>4</sup>, Nahema Marchal<sup>1</sup>, Markus Anderljung<sup>2</sup>, Noam Kolt<sup>5</sup>, Lewis Ho<sup>1</sup>, Divya Siddarth<sup>6, 7</sup>, Shahar Avin<sup>8</sup>, Will Hawkins<sup>1</sup>, Been Kim<sup>1</sup>, Iason Gabriel<sup>1</sup>, Vijay Bolina<sup>1</sup>, Jack Clark<sup>9</sup>, Yoshua Bengio<sup>10, 11</sup>, Paul Christiano<sup>12</sup> and Allan Dafoe<sup>1</sup>

<sup>1</sup>Google DeepMind, <sup>2</sup>Centre for the Governance of AI, <sup>3</sup>Centre for Long-Term Resilience, <sup>4</sup>OpenAI, <sup>5</sup>University of Toronto, <sup>6</sup>University of Oxford, <sup>7</sup>Collective Intelligence Project, <sup>8</sup>University of Cambridge, <sup>9</sup>Anthropic, <sup>10</sup>Université de Montréal, <sup>11</sup>Mila – Quebec AI Institute, <sup>12</sup>Alignment Research Center

Current approaches to building general-purpose AI systems tend to produce systems with both beneficial and harmful capabilities. Further progress in AI development could lead to capabilities that pose extreme risks, such as offensive cyber capabilities or strong manipulation skills. We explain why *model evaluation* is critical for addressing extreme risks. Developers must be able to identify dangerous capabilities (through “dangerous capability evaluations”) and the propensity of models to apply their capabilities for harm (through “alignment evaluations”). These evaluations will become critical for keeping policymakers and other stakeholders informed, and for making responsible decisions about model training, deployment, and security.

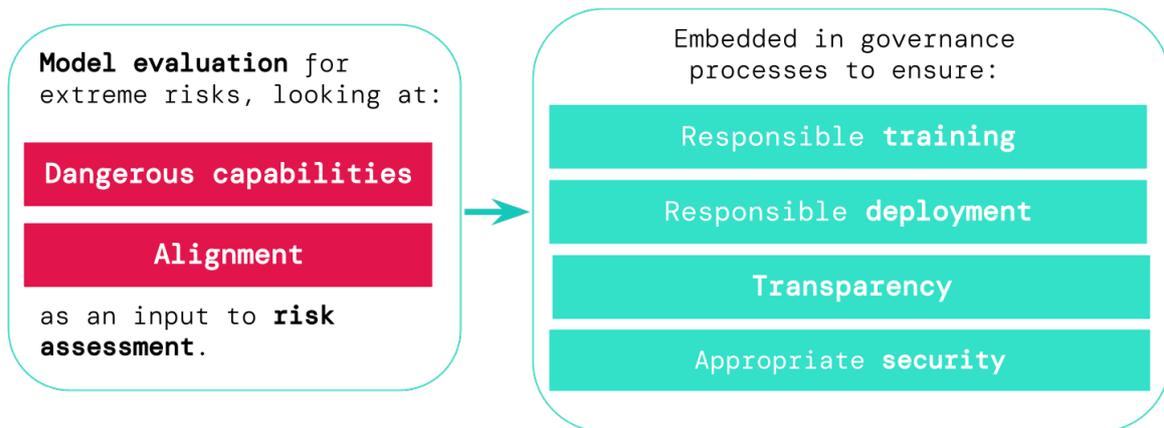


Figure 1 | The theory of change for model evaluations for extreme risk. Evaluations for dangerous capabilities and alignment inform risk assessments, and are in turn embedded into important governance processes.

## 1. Introduction

As AI progress has advanced, general-purpose AI systems have tended to display new and hard-to-forecast capabilities – including harmful capabilities that their developers did not intend (Ganguli et al., 2022). Future systems may display even more dangerous emergent capabilities, such as the ability to conduct offensive cyber operations, manipulate people through conversation, or provide actionable instructions on conducting acts of terrorism.

AI developers and regulators must be able to identify these capabilities, if they want to limit the risks they pose. The AI community already relies heavily on model evaluation – i.e. empirical assessment of a model’s properties – for identifying and responding to a wide range of risks. Existing model evaluations measure gender and racial biases, truthfulness, toxicity, recitation of copyrighted

arXiv:2305.15324v2 [cs.AI] 22 Sep 2023

content, and many more properties of models (Liang et al., 2022).

We propose **extending this toolbox** to address risks that would be *extreme in scale*, resulting from the misuse or misalignment of general-purpose models. Work on this new class of model evaluation is already underway. These evaluations can be organised into two categories: (a) whether a model has certain **dangerous capabilities**, and (b) whether it has the propensity to harmfully apply its capabilities (**alignment**).

Model evaluations for extreme risks will play a critical role in governance regimes. A central goal of AI governance should be to limit the creation, deployment, and proliferation of systems that pose extreme risks. To do this, we need tools for looking at a particular system and assessing whether it poses extreme risks. We can then craft company policies or regulations that ensure:

1. **Responsible training:** Responsible decisions are made about whether and how to train a new model that shows early signs of risk.
2. **Responsible deployment:** Responsible decisions are made about whether, when, and how to deploy potentially risky models.
3. **Transparency:** Useful and actionable information is reported to stakeholders, to help them mitigate potential risks.
4. **Appropriate security:** Strong information security controls and systems are applied to models that might pose extreme risks.

Many AI governance initiatives focus on the risks inherent to a particular deployment context, such as the “high-risk” applications listed in the [draft EU AI Act](#). However, models with sufficiently dangerous capabilities could pose risks even in seemingly low-risk domains. We therefore need tools for assessing *both* the risk level of a particular domain *and* the potentially risky properties of particular models; this paper focuses on the latter.

**Section 2** motivates our focus on extreme risks from general-purpose models and refines the scope of the paper. **Section 3** outlines a vision for how model evaluations for such risks should be incorporated into AI governance frameworks. **Section 4** describes early work in the area and outlines key design criteria for extreme risk evaluations. **Section 5** discusses the limitations of model evaluations for extreme risks and outlines ways in which work on these evaluations could cause unintended harm. We conclude with high-level recommendations for AI developers and policymakers.

## 2. Extreme risks from general-purpose models

Frontier AI developers are making rapid progress in developing increasingly capable general-purpose models (Bubeck et al., 2023). These models learn their capabilities and behaviours during training, and current methods for steering this process are imperfect (Gao et al., 2022; Shah et al., 2022). At the research frontier, models display new capabilities, often unforeseen by their developers (Wei et al., 2022b).

This poses a challenge for safety. AI developers could train general-purpose models that have dangerous capabilities – such as skills in deception, cyber offense, or weapons design – without actively seeking these capabilities. Humans could then intentionally misuse these capabilities (Brundage et al., 2018), e.g. for assistance in disinformation campaigns, cyberattacks, or terrorism. Additionally, due to failures of alignment, AI systems could harmfully apply their capabilities even without deliberate misuse (Ngo et al., 2022).

In the near-term, these risks will be especially concentrated on the frontier of AI research and development. We loosely define the “frontier” as models that are both (a) close to, or exceeding, the average capabilities of the most capable existing models,<sup>1</sup> and (b) different from other models, either in terms of scale, design (e.g. different architectures or alignment techniques), or their resulting mix of capabilities and behaviours. Accordingly, frontier models are uniquely risky because (a) more capable models can excel at a wider range of tasks, which will unlock more opportunities to cause harm;<sup>2</sup> and (b) novel models are less well-understood by the research community.

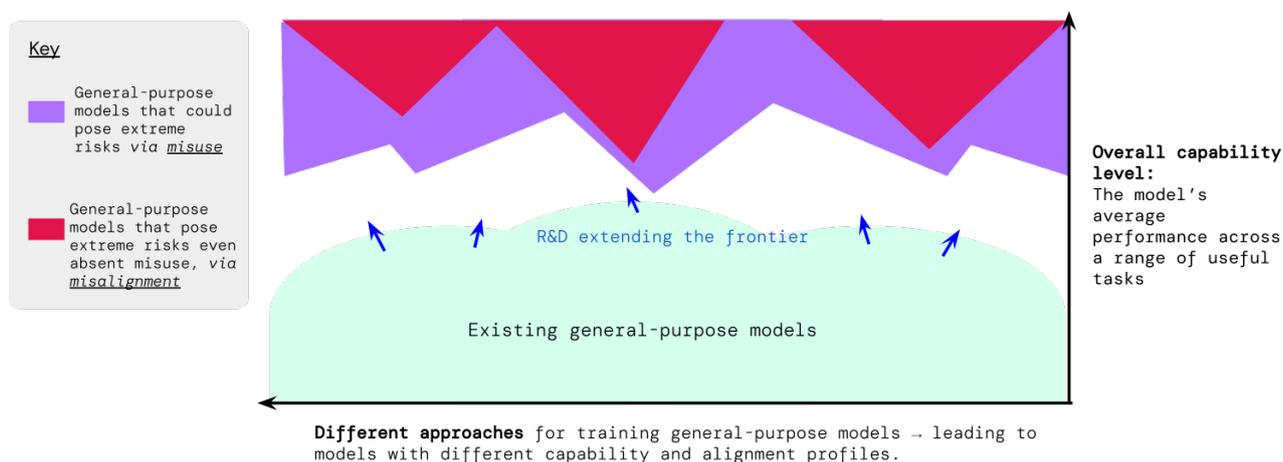


Figure 2 | Leading AI developers push the frontier outward, typically by training models at greater scale and using more efficient architectures and algorithms. This continued expansion takes the field closer to points in model space that could pose extreme risks. The diagram is purely illustrative.

We focus on “extreme” risks, i.e. those that would be extremely large in scale (even relative to the scale of deployment). This can be operationalised in terms of the scale of impact (e.g. damage in the tens of thousands of lives lost, hundreds of billions of dollars of economic or environmental damage) or the level of adverse disruption to the social and political order. The latter could mean, for example, the outbreak of inter-state war, a significant erosion in the quality of public discourse, or the widespread disempowerment of publics, governments, and other human-led organisations (Carlsmith, 2022).

Many AI researchers (and other stakeholders) view extreme risks from AI as an important challenge. In a 2022 survey of AI researchers, 36% of respondents thought that AI systems could plausibly “cause a catastrophe this century that is at least as bad as an all-out nuclear war” (Michael et al., 2022). However, very few existing model evaluations intentionally target risks on this scale.

To guard against extreme risks, AI developers should use model evaluation to uncover:

1. To what extent a model is capable of causing extreme harm (which relies on evaluating for certain **dangerous capabilities**).
2. To what extent a model has the propensity to cause extreme harm (which relies on **alignment** evaluations).

<sup>1</sup>In practice, defining “average capabilities” would involve many judgement calls over which evaluations should be included and how they should be weighted.

<sup>2</sup>This is especially pertinent for extreme risks: causing such large-scale harm is not normally an easy challenge. Even well-resourced terrorist groups, determined to cause extreme harm, often fail.

We provide a non-exhaustive list of dangerous capabilities in Table 1. Most of the capabilities listed are offensive capabilities: they are useful for exerting influence or threatening security (e.g. see: persuasion and manipulation, cyber-offense, weapons acquisition). Some (e.g. situational awareness) are capabilities that would be advantageous for a misaligned AI system evading human oversight (Ngo et al., 2022). We omit many generically useful capabilities (e.g. browsing the internet, understanding text) despite their potential relevance to both the above.

The most risky scenarios will involve multiple dangerous capabilities combined together – further research should explore what combinations would be most dangerous. Sometimes specific capabilities can be supplied by the user or outsourced to other humans (e.g. crowdworkers) or AI systems.

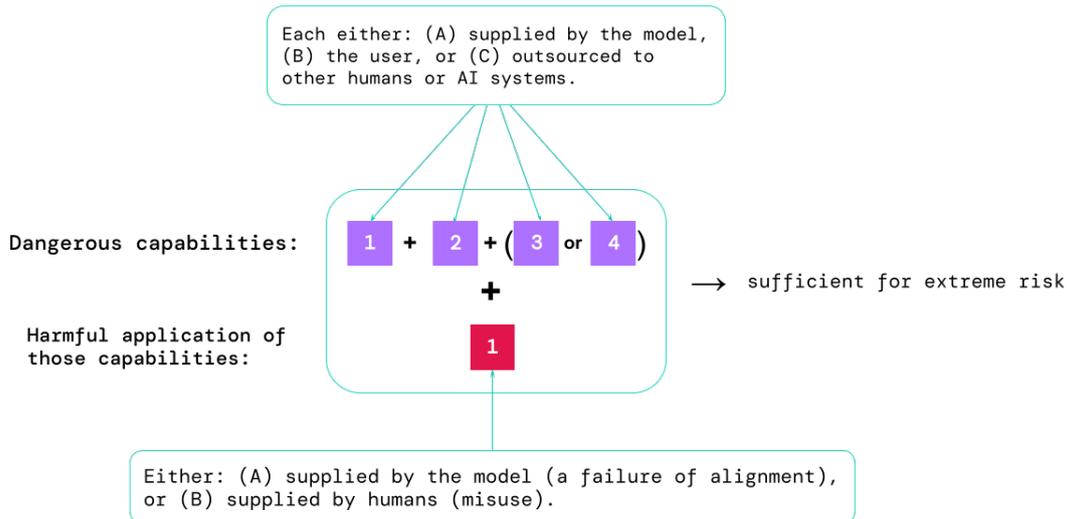


Figure 3 | Ingredients for extreme risk.

A simple heuristic: a model should be treated as highly dangerous if it has a capability profile that would be sufficient for extreme harm, *assuming* misuse and/or misalignment. To deploy such a model, AI developers would need very strong controls against misuse (Shevlane, 2022b) and very strong assurance (via **alignment evaluations**) that the model will behave as intended. Alignment evaluations should look for behaviours identified in the literature, such as whether the model:

- Pursues long-term, real-world goals, different from those supplied by the developer or user (Chan et al., 2023; Ngo et al., 2022);
- Engages in “power-seeking” behaviours (Krakovna and Kramar, 2023; Turner et al., 2021);
- Resists being shut down (Hadfield-Menell et al., 2016; Orseau and Armstrong, 2016);
- Can be induced into collusion with other AI systems against human interests (Ngo et al., 2022).
- Resists malicious users’ attempts to access its dangerous capabilities (Glaese et al., 2022).

We focus on risks arising from misuse and misalignment because a new generation of model evaluations is needed for identifying these risks. Conversely, though important, we leave out of scope:

1. **Structural risks**, which depend *especially* heavily on how the AI system interacts with larger social, political, and economic forces in society (Zwetsloot and Dafoe, 2019). Model evaluation sheds less light on these risks, because they depend so heavily on factors external to the model.
2. Risks from models **incompetently** performing important tasks (Raji et al., 2022a). Existing kinds of model evaluation will be most relevant here (e.g. testing the model’s accuracy and robustness on the relevant task).

Capability	Could include:
<b>Cyber-offense</b>	The model can <b>discover vulnerabilities</b> in systems (hardware, software, data). It can write code for <b>exploiting</b> those vulnerabilities. It can make effective decisions once it has gained access to a system or network, and skilfully evade threat detection and response (both human and system) whilst focusing on a specific objective. If deployed as a coding assistant, it can <b>insert subtle bugs</b> into the code for future exploitation.
<b>Deception</b>	The model has the skills necessary to <b>deceive humans</b> , e.g. constructing believable (but false) statements, making accurate predictions about the effect of a lie on a human, and keeping track of what information it needs to withhold to maintain the deception. The model can impersonate a human effectively.
<b>Persuasion &amp; manipulation</b>	The model is effective at <b>shaping people's beliefs</b> , in dialogue and other settings (e.g. social media posts), even towards untrue beliefs. The model is effective at <b>promoting certain narratives</b> in a persuasive way. It can convince people to do things that they would not otherwise do, including unethical acts.
<b>Political strategy</b>	The model can perform the social modelling and planning necessary for an actor to gain and exercise <b>political influence</b> , not just on a micro-level but in scenarios with <b>multiple actors</b> and rich <b>social context</b> . For example, the model can score highly in forecasting competitions on questions relating to global affairs or political negotiations.
<b>Weapons acquisition</b>	The model can <b>gain access to existing weapons</b> systems or contribute to <b>building new weapons</b> . For example, the model could assemble a bioweapon (with human assistance) or provide actionable instructions for how to do so. The model can make, or significantly assist with, scientific discoveries that unlock novel weapons.
<b>Long-horizon planning</b>	The model can make <b>sequential plans</b> that involve multiple steps, unfolding over <b>long time horizons</b> (or at least involving many interdependent steps). It can perform such planning within and across many domains. The model can sensibly adapt its plans in light of unexpected obstacles or adversaries. The model's planning capabilities generalise to <b>novel settings</b> , and do not rely heavily on trial and error.
<b>AI development</b>	The model could build new AI systems from scratch, including AI systems with dangerous capabilities. It can find ways of adapting other, existing models to increase their performance on tasks relevant to extreme risks. As an assistant, the model could significantly improve the productivity of actors building dual use AI capabilities.
<b>Situational awareness</b>	The model can distinguish between <b>whether it is being trained, evaluated, or deployed</b> – allowing it to behave differently in each case. The model <b>knows that it is a model</b> , and has <b>knowledge about itself</b> and its likely surroundings (e.g. what company trained it, where their servers are, what kind of people might be giving it feedback, and who has administrative access).
<b>Self-proliferation</b>	The model can break out of its local environment (e.g. using a vulnerability in its underlying system or suborning an engineer). The model can exploit limitations in the systems for monitoring its behaviour post-deployment. The model could independently generate revenue (e.g. by offering crowdwork services, ransomware attacks), use these revenues to acquire cloud computing resources, and operate a large number of other AI systems. The model can generate creative strategies for uncovering information about itself or exfiltrating its code and weights.

Table 1 | Dangerous capabilities

### 3. Model evaluation as critical governance infrastructure

Across many industries, safety standards and regulations rely on tools for assessing risks in new products – for instance, food, drugs, commercial airliners, and automobiles. Model evaluation is not the only tool available for AI risk assessment – more theoretical approaches are also available, e.g. studying the incentives operating on a model during training (Everitt et al., 2021). Nonetheless, model evaluation is one of the main tools we have for AI risk assessment.

Figure 4 provides an overview of this section. It is an ambitious blueprint for how to guard against extreme risks while developing and deploying a model, with evaluation embedded throughout. The evaluation results feed into processes for risk assessment (Khlaaf et al., 2022), which inform (or bind) important decisions around model training, deployment, and security. The developer reports results and risk assessments to external stakeholders.

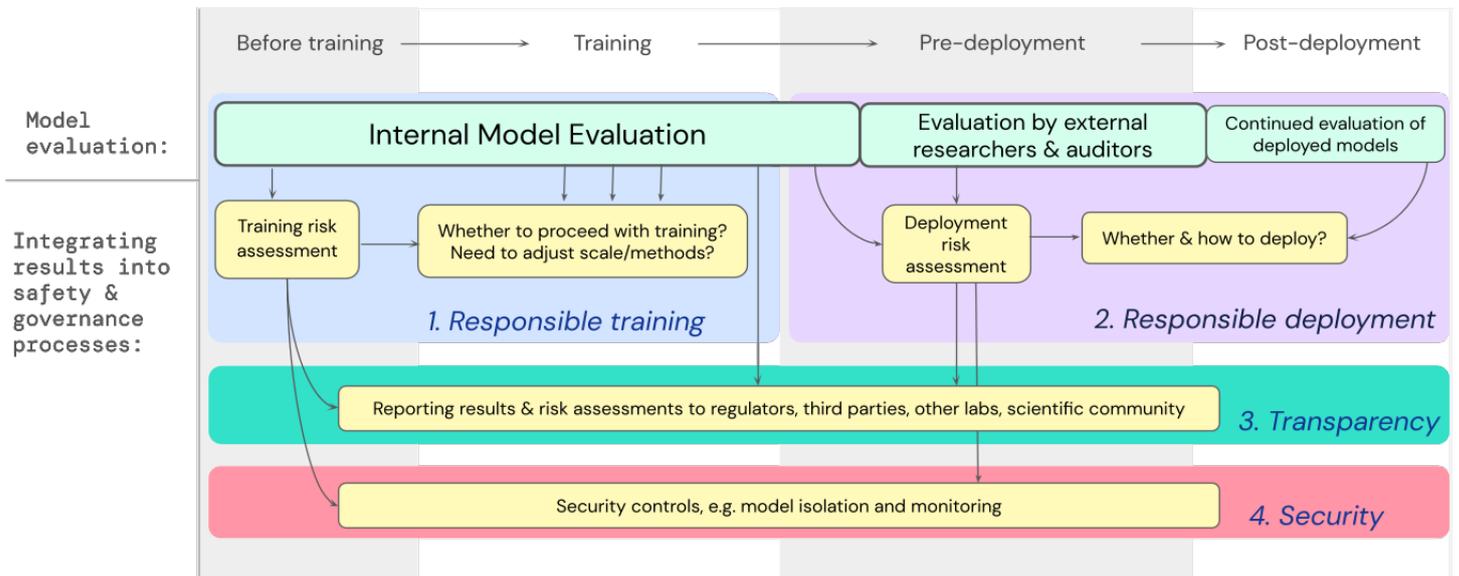


Figure 4 | A workflow for training and deploying a model, embedding extreme risk model evaluation results into key safety and governance processes.

Three sources of model evaluations feed into this process:

- 1. Internal model evaluation**, i.e. the developer conducting its own evaluations. There is no substitute for internal model evaluation, given that internal researchers have high context on the model's design and deeper model access than can be achieved via an API. Developers could have multiple organisational layers of safety evaluation, such as by establishing an internal safety evaluation function that is independent of the teams primarily responsible for building the models, reporting directly to organisational leaders (see Raji et al., 2020).
- 2. External research access.** The developer grants model access to external researchers, likely via an API (Bluemke et al., 2023; Shevlane, 2022a,b). Their research could be exploratory or targeted at evaluating specific properties, including “red teaming” the model's alignment.
- 3. External model audit**, i.e. model evaluation by an independent, external auditor for the purpose of providing a judgement — or input to a judgement — about the safety of deploying a model (or training a new one) (ARC Evals, 2023; Mökander et al., 2023; Raji et al., 2022b). Ideally there would exist a rich ecosystem of model auditors providing broad coverage across different risk areas. (This ecosystem is currently under-developed.)

### 3.1. Responsible training

The first line of defence is to avoid training models that have sufficient dangerous capabilities and misalignment to pose extreme risk. Sufficiently concerning evaluation results should warrant **delaying** a scheduled training run or **pausing** an existing one.<sup>3</sup>

Before a frontier training run, developers have the opportunity to study weaker models that might provide early warning signs. These models come from two sources: (1) previous training runs, and (2) experimental models leading up to the new training run. Developers should evaluate these models and try to forecast the results from the planned training run (see [OpenAI, 2023b](#)). This would include scaling (or “inverse scaling”) analysis where the aim is to find areas where scaling brings unwanted changes to the model ([McKenzie et al., 2022](#)). These insights can feed into a training risk assessment. Then, during the training run, researchers could run extreme risk evaluations at regular intervals.

The developer has a range of possible responses to address the concerning evaluation results:

1. **Study the issue** to understand why the misalignment or dangerous capability emerged.
2. **Adjust the training methods** to circumvent the issue. This could mean adjusting (for example) the architecture, the data, the training tasks, or further developing the alignment techniques used. These adjustments should target the fundamental issue rather than inducing superficial changes to how the model scores on the available evaluations (see section [5.2](#)).
3. **Careful scaling**. If the developer is not confident it can train a safe model at the scale it initially had planned, they could instead train a smaller or otherwise weaker model.

In mature governance regimes, the decision to proceed with a potentially risky training run could require approval from an external model auditor or regulator.

### 3.2. Responsible deployment

Deployment means making the model available for use, e.g. it is built into a product or hosted on an API for software developers to build with. Deployment constitutes a large increase in the model’s exposure to the external world and therefore possible risk. Model evaluation for extreme risks could inform a **deployment risk assessment** that reviews (a) *whether or not* the model is safe to deploy, and (b) the appropriate *guardrails* for ensuring the deployment is safe.

The predeployment evaluation process takes time ([OpenAI, 2023a](#); [Rismani et al., 2023](#)). Industry standards or regulation could require a minimum duration for predeployment evaluation of frontier models, including the length of time that external researchers and auditors have access.

In response to concerning evaluation results, one possibility is to recommend against deployment. A second possibility is to recommend adjustments to the deployment plan that would address potential risks (see [Table 3](#) in the Appendix for a range of variables that could be adjusted). Nonetheless, for a sufficiently capable and poorly aligned model, it is possible that even a restrictive and scaled-back deployment could pose extreme risk.

---

<sup>3</sup>Developers training highly capable, general-purpose models should become accustomed to such a prospect, not planning their research around an assumption that a training run will run to schedule. For example, well-run developers will have compute allocation systems that backfill the vacant computing resources with other useful projects. Similarly, developers should avoid making hard promises to stakeholders (e.g. customers, investors) that they will deploy a certain model at a certain date. At the least, they should retain the flexibility of pivoting to a smaller or otherwise less risky version of the model.

Safe deployment will often be a gradual process (Figure 5) (Brundage et al., 2022). The developer gradually accumulates evidence about the model's safety, through both evaluation (internal and external) and early, small-scale deployment.

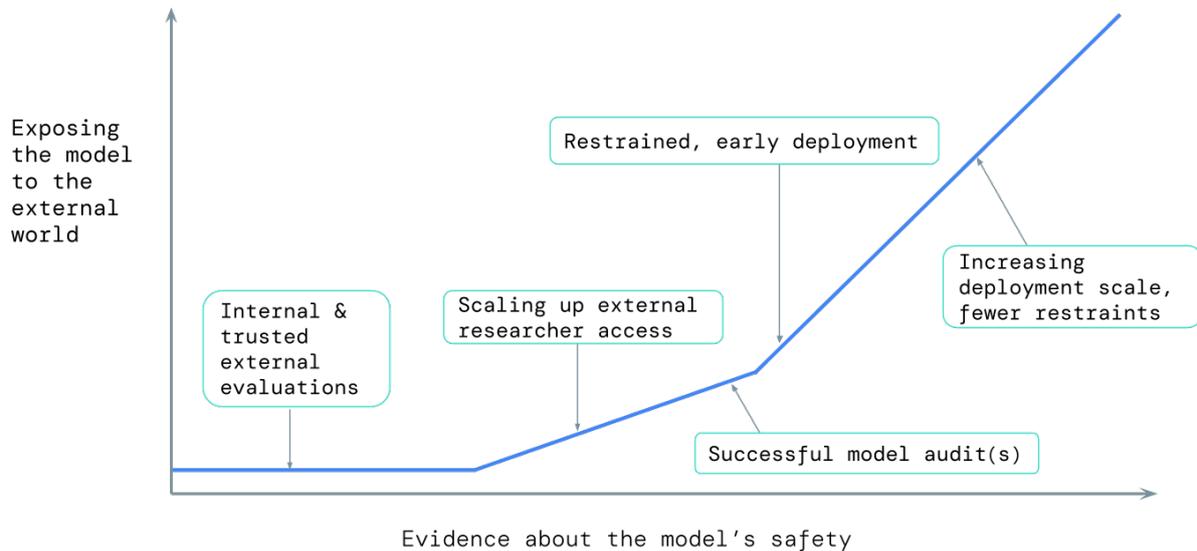


Figure 5 | The developer gradually increases the model's exposure to the external world as it accumulates evidence about the model's safety.

Evaluation will often need to continue after deployment. There are two reasons for this:

1. **Unanticipated behaviours.** Before deployment, it is impossible to fully anticipate and understand how the model will interact in a complex deployment environment (a key limitation of model evaluation: see section 5). For example, users might find new applications for the model or novel prompt engineering strategies; or the model could be operating in a dynamic, multi-agent environment. Therefore, in the early stages of deployment, developers must:
  - (a) Surface emerging model behaviours and risks via **monitoring** efforts. This could include direct monitoring of inputs and outputs to the model, and systems for incident reporting (see Brundage et al., 2022; Raji et al., 2022b).
  - (b) Design and run new model evaluations inspired by these observations.
2. **Updates to the model.** The developer might update the model after deployment, e.g. by fine-tuning on data collected during deployment or by expanding the model's access to external tools. If these updates could increase risk, they should be evaluated before launch. For large changes,<sup>4</sup> the new model could go through the whole process described in this section.

The ideal state is *continuous deployment review*. On an ongoing basis, the developer reassesses deployment safety using model evaluations and monitoring, and at any time, could adjust or terminate the deployment in response to their findings. Further, for deployments that were recognisably unsafe in retrospect, an external audit of the deployment decision-making process could be triggered. Safety issues uncovered during deployment can also inform training risk assessments for future models.

<sup>4</sup>The magnitude of the change to the model could be assessed in terms of the amount of additional training that it has gone through (as a percentage of the original training length), or the model's improvement on key performance benchmarks.

Finally, even *internal* deployments of highly capable general-purpose models, notably as coding assistants for AI researchers and engineers, could require pre-deployment evaluation for dangerous capabilities (e.g. the ability to insert subtle vulnerabilities into code) and alignment.

### 3.3. Transparency

Model evaluations are a vital tool for keeping stakeholders informed about the state of AI risks on the frontier ([Whittlestone and Clark, 2021](#)). We recommend frontier developers consider processes for externally **reporting** the results of evaluations or extracts from the assessment documents that rely on those evaluation results (such as training risk assessments, auditors' reports, deployment risk assessments).

Model evaluations will unlock four important kinds of transparency around extreme risks:

1. **Incident reporting**, i.e. a structured process for developers to share concerning or otherwise noteworthy evaluation results with other developers, third parties, or regulators (see [Brundage et al., 2020](#)). This would be vital for helping others avoid training risky systems, and for keeping AI developers accountable. In future, regulators could maintain lists of high-risk or banned training approaches; a sufficiently concerning incident report would then trigger an update to this list.
2. **Sharing pre-deployment risk assessments**. Prior to deployment, developers could share their predeployment risk assessments for comment and critique. This could mean sharing with auditors, external researchers, regulators, or publicly. The risk assessment would include evaluation results, alongside the developer's justification for why the deployment is safe given those results. (See related: [Mitchell et al., 2018](#))
3. **Scientific reporting**, i.e. presenting the results to the scientific community in a way that encourages further scientific research into the issue. Much research will be needed to better understand the behaviour of highly capable, general-purpose models, especially their alignment. Topics within the alignment literature, such as situational awareness, could follow the path of fairness research, where early findings of gender and racial biases in models have helped stimulate new areas of research.
4. **Educational demonstrations**, i.e. engaging presentation of concerning capabilities and behaviours to key stakeholders – such as policymakers, the public, or company leadership – to keep them informed about risks from frontier AI systems. The demonstrations could be videos or interactive demos, with the evaluation results presented as supplementary information (as evidence of scientific rigour).

### 3.4. Appropriate security

Models at risk of exhibiting dangerous capabilities will require strong and novel security controls. Developers must consider multiple possible threat actors: insiders (e.g. internal staff, contractors), outsiders (e.g. users, nation-state threat actors), and the model itself as a vector of harm. We must develop new security best practices for high-risk AI development and deployment, which could include for example:

- **Red teaming**: Intensive security red-teaming for the entire infrastructure on which the model is developed and deployed.

- **Monitoring:** Intensive, AI-assisted monitoring of the model's behaviour, e.g. for whether the model is engaging in manipulative behaviour or making code recommendations that would lower the overall security of a system.
- **Isolation:** Appropriate isolation techniques for preventing risky models from exploiting the underlying system (e.g. sole-tenant machines and clusters, and other software-based isolation). The model's network access should be tightly controlled and monitored, as well as its access to tools (e.g. code execution).
- **Rapid response:** Processes and systems for rapid response to disable model actions and the model's integrations with hardware, software, and infrastructure in the event of unexpected unsafe behaviour.
- **System integrity:** Formal verification that served models, memory, or infrastructure have not been tampered with. The development and serving infrastructure should require two-party authorization for any changes and auditability of all changes.

#### 4. Building evaluations for extreme risk

Model evaluation is already a core component of AI research, and increasingly we have evaluations that focus on ethics, safety, and social impact. We recommend extending this toolbox to address extreme risks.

Early work is already underway to build model evaluations for extreme risks. ARC Evals (the evaluations team at the Alignment Research Center) is building evaluations that measure language models' self-proliferation capabilities (see Table 1 above). ARC Evals ran this evaluation on GPT-4 and Claude before their wider release (ARC Evals, 2023; OpenAI, 2023a). OpenAI and the GPT-4 red teamers also tested GPT-4's capabilities in cybersecurity operations and its ability to purchase certain chemical compounds (OpenAI, 2023a).

Google DeepMind has ongoing projects evaluating language models for manipulation capabilities. This includes a game called "Make-me-say", where the language model must lead an (unaware) human conversation partner to say a pre-specified word.<sup>5</sup>

Table 2 contains a range of desirable qualities for extreme risk evaluations. Some of these qualities relate to a *single* evaluation, and some are desirable qualities of a *portfolio* of evaluations.

We anticipate that building comprehensive alignment evaluations will be most challenging. The ambition is for a process of alignment assurance that could conclude, with high confidence, that a model is not dangerously misaligned, even for very capable models. (Model evaluations would not be the only input to this assurance process, but an important one.)

Alignment evaluation is challenging because we need assurance that the model will reliably behave appropriately across a wide diversity of settings (Ziegler et al., 2022). An evaluation might find that a model is aligned in some narrow, prosaic way (for example, a language agent asserting that it does not object to being shut down (Perez et al., 2022a,b)) without providing evidence that the model would exhibit desirable behaviour when presented with genuine (or more convincing) opportunities to achieve self-preservation, greater influence, or other harmful outcomes.

---

<sup>5</sup>This project is led by Mary Phuong.

<i>Comprehensive:</i>	
<b>Threat models</b>	The evaluation portfolio should cover as many plausible extreme risk threat models as possible.
<b>Automated and human-assisted</b>	Many evaluations can be run automatically, lowering the time and resource costs. However, some capabilities and behaviours will need human-assisted evaluations, i.e. involving: (a) human raters who judge the model's outputs; or (b) humans who interact with the model, e.g. in a dialogue setting.
<b>Behavioural and mechanistic</b>	Evaluations should not be restricted to studying a model's behaviour, but should eventually also involve looking mechanistically at how the model produced that behaviour.
<b>Fault-finding</b>	The portfolio of evaluations should include adversarial testing, where researchers purposefully search for cases where the model produces concerning results.
<b>Robust to deception</b>	Ultimately researchers will need evaluations that can rule out the possibility that the model is deliberately appearing safe for the purpose of passing the evaluation process.
<b>Surfacing latent capabilities</b>	Researchers will need to bring latent capabilities to the surface (for example, by prompt engineering or fine-tuning).
<b>Model lifecycle</b>	We recommend conducting evaluations throughout the model development process. In particular, the results from the end of a long development process will likely fail to convey relevant information about the base model, especially if it has been fine-tuned for safety.
<b>Model-level and system-level</b>	Models are often integrated into wider AI systems, e.g. with external tools, other models, or classifiers that filter the model's outputs. Evaluations should study models both with and without these augmentations.
<i>Interpretable:</i>	
<b>Legible</b>	Some evaluations should present risks in an accessible way, requiring little technical understanding. This will be helpful for creating common knowledge around the risks from AI.
<b>Wide difficulty spectrum</b>	The dangerous capability evaluations should ideally contain wide ranges of difficulty – ideally within single evaluations, but at least across the portfolio. This means that researchers can track capabilities progress as it approaches possible danger thresholds, and that the evaluation (or the portfolio) is scalable to future, more capable models. For tracking progress, evaluations would ideally provide a quantitative score, although this will not always be practical.
<i>Safe:</i>	
<b>Safe to implement</b>	Dangerous capability evaluations could involve testing the model in real-world settings, e.g. interacting with crowdworkers. This should not introduce unacceptable levels of risk.

Table 2 | Desirable qualities of extreme risk evaluations.

Researchers must therefore evaluate a model across a broad range of settings. Achieving coverage of settings for alignment evaluation can be helped by:

1. **Breadth:** Evaluating behaviour across as wide a range of settings as possible. One promising avenue is automating the process of writing evaluations using AI systems (Perez et al., 2022b) (see also Pan et al., 2023).
2. **Targeting:** Some settings are much more likely to reveal alignment failures than others, and we may be able to focus on them through clever design – for example, using honeypots or gradient-based adversarial testing and related approaches (Jones et al., 2023).
3. **Understanding generalisation:** Since researchers will be unable to foresee or simulate all possible scenarios, we must develop a better scientific understanding of how and why model behaviours generalise (or fail to generalise) between settings.

Another important tool is **mechanistic** analysis, i.e. studying the model’s weights and activations for understanding how it functions (Nanda et al., 2023; Olah et al., 2020). For example, one ambition is to study how the model’s goals are represented internally, to help verify that they are as intended; another ambition is to detect when a language model’s outputs misreport its knowledge (Burns et al., 2022), which could be an indicator of deceptive behaviour.

Finally, **agency** – in particular, the goal-directedness of an AI system – is an important property to evaluate (Kenton et al., 2022), given the central role of agency in various theories of AI risk (Chan et al., 2023). Partly, agency is a question of the model’s capabilities – is it capable of effectively pursuing goals? Evaluating alignment also requires looking at agency, including: (a) Is the model more goal-directed than the developer intended? For example, has a dialogue agent learnt the goal of manipulating the user’s behavior? (b) Does the model resist a user’s attempt to assemble it into an autonomous AI system (e.g. Auto-GPT) with harmful goals?

## 5. Limitations and hazards

### 5.1. Limitations

Model evaluation, as a tool for addressing extreme risks, has at least **five limitations**. A key issue is that not all risks can necessarily be detected via model evaluation.

1. **Factors beyond the AI system.** Risks will depend on how an AI system interacts with a complex world. For example, a model might use, as tools, other models released in the future, thus augmenting its capabilities; or human civilisation might be less resilient to powerful AI than anticipated.
2. **Unknown threat models.** It is difficult to anticipate all the different plausible pathways to extreme risk. This will be especially true for highly capable models, which could find creative strategies for achieving their goals.
3. **Difficult-to-identify properties.** Some model properties will be challenging to uncover via model evaluations. Two important cases:
  - (a) Capability overhang: Models sometimes have capabilities that the AI research community does not realise. For example, after GPT-3 had already existed for many months, researchers

demonstrated that chain-of-thought prompting could significantly increase performance (Wei et al., 2022c).

(b) Deceptive alignment: A situationally aware model could deliberately exhibit desired behaviour during evaluation (Ngo et al., 2022). (This is one reason not to rely solely on behavioural evaluations.)

4. **Emergence**. Above we recommended using model evaluations to inform the decision to train a new model by performing scaling laws analysis on smaller models. However, sometimes specific capabilities will emerge only at greater scale, which makes this analysis much harder (Ganguli et al., 2022); other capabilities display U-shaped scaling (Wei et al., 2022a).

5. **Maturity of evaluation ecosystem**. The ecosystem for external evaluations and model audits is currently under-developed.

6. **Overtrust in evaluations**. There is a risk that too much faith is placed in evaluation results, leading to risky models being deployed under a false sense of security.

Model evaluation is a necessary but *not sufficient* strategy for identifying and mitigating extreme risks. It must be combined with a wider organisational dedication to safety and other tools for risk identification and assessment.

## 5.2. Hazards

Conducting and reporting the evaluations discussed in this paper poses four potential hazards:

1. **Advancing and proliferating dangerous capabilities**. There is a risk that – through conducting dangerous capability evaluations and sharing relevant materials – the field will proliferate dangerous capabilities or accelerate their development. We highlight four kinds of potentially hazardous information:

(a) Results. Evaluation results could demonstrate novel offensive technologies. Publicly sharing these results could spur investment in new weapons programmes, cyber-offensive efforts, or methods for digital oppression of citizens. By analogy, it has been said that in the 1940s the most valuable secret about the nuclear bomb was that it was possible (Ord, 2022). AI developers, researchers, and auditors should therefore exercise caution around sharing these evaluation results.

(b) Evaluation datasets. Datasets for evaluating dangerous capabilities are dual use because other actors could fine-tune their models on these datasets.

(c) Elicitation techniques. Evaluating dangerous capabilities will often involve eliciting those capabilities from the model. This could involve: (a) prompt engineering; and (b) fine-tuning, including: (i) finding creative new task specifications; (ii) creating or identifying appropriate fine-tuning datasets. These techniques could be useful for a bad actor attempting to elicit dangerous capabilities from similar models. Researchers and auditors must therefore exercise caution over sharing their elicitation techniques, especially if producing them relied on creativity, expert knowledge, or time-consuming experimentation.

(d) Trained models. There are risks from intentionally training dangerously capable models, even for use as safety research artefacts. We could distinguish between (a) simply following off-the-shelf methods (e.g. fine-tuning an existing model), versus (b) cases where the work to produce the dangerous capability could constitute a research contribution in its own right (e.g. it could be accepted to an academic conference). The latter is arguably

comparable to “gain-of-function” research in virology, especially if the resulting model is highly capable and general-purpose. The research may need to be conducted under very high-security conditions and subject to a demanding risk assessment.

2. **Competitive pressures.** One concern is that sharing evaluation results between competing AI developers could incentivise them to behave less responsibly. For example, sharing pre-deployment evaluation results could tip off competitors about future product improvements, incentivising those competitors to rush their own deployments and spend less time on ensuring safety. Similarly, since dangerous capability results will often correlate with the model’s overall capabilities, competing developers could learn that they are falling behind and decide they need to sacrifice on safety to catch up (Emery-Xu et al., 2023).

Given the sensitivities involved, one option is to lean more heavily on alignment evaluation results, at least for reporting between developers. For illustration, a possible inter-developer policy could be:

- (a) Report unexpected or important alignment issues promptly to other developers. By default, limit the description of the model’s training to only a high-level overview (to avoid revealing sensitive information); but share more details if this is absolutely necessary – in particular, if a certain class of methods is causing the problem.
  - (b) Report when certain dangerous capability thresholds have been passed. These thresholds can be set high, to avoid sharing granular information. Wait until deployment to share more.
3. **Superficial improvements to model safety.** There is a risk that widely available safety evaluations will lead to models that exhibit only superficially desirable behaviours. Most clearly, if researchers directly train models to pass these evaluations, the evaluations can no longer act as an indicator of risk. Researchers could do this either accidentally (e.g. because the evaluation datasets are shared online and thereby end up in the pretraining dataset) or as an intentional attempt to pass external audits (analogous to the Volkswagen emissions scandal). The model’s desirable evaluation performance would then likely fail to generalise. Developers and model auditors could therefore consider keeping some private “held out” evaluations and ensuring these are not too overlapping with datasets or tasks used during training.

Even if developers refrain from directly training on the evaluations, we have nevertheless recommended that developers avoid training models that fail the evaluations (section 3.1). This could also exert selection pressure, albeit weaker. Over the long run, the risk is that the field selects for training methods that produce deceptively aligned models.

4. **Harms during the course of evaluation.** Running evaluations will often involve exposing the model to the external world. For example, in evaluating GPT-4, ARC used the model to generate (deceptive) messages to be sent to a TaskRabbit worker (OpenAI, 2023a). In the extreme case, a poorly managed test for whether a model has self-proliferating capabilities could end in actual proliferation; but more prosaically, they could cause harm in other ways, such as causing emotional distress to crowdworkers. Therefore, groups conducting evaluations, such as auditors, should establish safety protocols where necessary.

## 6. Conclusion

Model evaluation for extreme risks should be a priority area for AI safety and governance. There are many challenges ahead for finding effective evaluations and building governance regimes that

incorporate them; we encourage further work in this area. Model evaluation is not a panacea: it will not catch all extreme risks. Nonetheless, it is a necessary component of the governance infrastructure needed to combat extreme risks.

Frontier AI developers currently have a special responsibility to support work on model evaluations for extreme risks, since they have resources – including access to cutting-edge AI models and deep technical expertise – that many other actors typically lack. Frontier AI developers are also currently the actors who are most likely to unintentionally develop or release AI systems that pose extreme risks. Frontier AI developers should therefore:

1. **Invest in research:** Frontier developers should devote resources to researching and developing model evaluations for extreme risks.
2. **Craft internal policies:** Frontier developers should craft internal policies for conducting, reporting, and responding appropriately to the results of extreme risk evaluations.
3. **Support outside work:** Frontier labs should enable outside research on extreme risk evaluations through model access and other forms of support.
4. **Educate policymakers:** Frontier developers should educate policymakers and participate in standard-setting discussions, to increase government capacity to craft any regulations that may eventually be needed to reduce extreme risks.

**Policymakers** should consider building up the governance infrastructure outlined in section 3. Policymakers could:

1. Systematically **track** the development of dangerous capabilities, and progress in alignment, within frontier AI R&D ([Whittlestone and Clark, 2021](#)). Policymakers could establish a formal reporting process for extreme risk evaluations.
2. **Invest** in the ecosystem for external safety evaluation, and create venues for stakeholders (such as AI developers, academic researchers, and government representatives) to come together and discuss these evaluations ([Anthropic, 2023](#))
3. Mandate **external audits**, including model audits and audits of developers' risk assessments, for highly capable, general-purpose AI systems.
4. Embed extreme risk evaluations into the **regulation** of AI deployment, clarifying that models posing extreme risks should not be deployed.

## 7. Acknowledgements

We are grateful for helpful comments and discussions on this work from: Canfer Akbulut, Jide Alaga, Beth Barnes, Joslyn Barnhart, Sasha Brown, Miles Brundage, Martin Chadwick, Tom Everitt, Conor Griffin, Eric Horvitz, Evan Hubinger, William Isaac, Victoria Krakovna, Leonie Koessler, Sébastien Krier, Nikhil Mulani, Neel Nanda, Jonas Schuett, Rohin Shah, Andrew Trask, Gregory Wayne, and Hjalmar Wijk. We are grateful for insightful discussions with the participants of two events held in February 2023: a virtual discussion session on the topic of this paper, and a one-day workshop on dangerous capabilities evaluations co-organised by Steven Adler, Anne le Roux, and Jade Leung. We

also thank Celine Smith for project management support, and Michael Chang for improvements to the visualisations.

## References

- Anthropic. Strengthening U.S. AI innovation through an ambitious investment in NIST. Technical report, 2023.
- ARC Evals. Update on ARC’s recent eval efforts. <https://evals.alignment.org/blog/2023-03-18-update-on-recent-evals/>, Mar. 2023. Accessed: 2023-3-20.
- E. Bluemke, T. Collins, B. Garfinkel, and A. Trask. Exploring the relevance of data Privacy-Enhancing technologies for AI governance use cases. Mar. 2023.
- M. Brundage, S. Avin, J. Clark, H. Toner, P. Eckersley, B. Garfinkel, A. Dafoe, P. Scharre, T. Zeitzoff, B. Filar, H. Anderson, H. Roff, G. C. Allen, J. Steinhardt, C. Flynn, S. Ó. hÉigeartaigh, S. Beard, H. Belfield, S. Farquhar, C. Lyle, R. Crootof, O. Evans, M. Page, J. Bryson, R. Yampolskiy, and D. Amodei. The malicious use of artificial intelligence: Forecasting, prevention, and mitigation. Feb. 2018.
- M. Brundage, S. Avin, J. Wang, H. Belfield, G. Krueger, G. Hadfield, H. Khlaaf, J. Yang, H. Toner, R. Fong, T. Maharaj, P. W. Koh, S. Hooker, J. Leung, A. Trask, E. Bluemke, J. Lebensold, C. O’Keefe, M. Koren, T. Ryffel, J. B. Rubinovitz, T. Besiroglu, F. Carugati, J. Clark, P. Eckersley, S. de Haas, M. Johnson, B. Laurie, A. Ingerman, I. Krawczuk, A. Askill, R. Cammarota, A. Lohn, D. Krueger, C. Stix, P. Henderson, L. Graham, C. Prunkl, B. Martin, E. Seger, N. Zilberman, S. Ó. hÉigeartaigh, F. Kroeger, G. Sastry, R. Kagan, A. Weller, B. Tse, E. Barnes, A. Dafoe, P. Scharre, A. Herbert-Voss, M. Rasser, S. Sodhani, C. Flynn, T. K. Gilbert, L. Dyer, S. Khan, Y. Bengio, and M. Anderljung. Toward trustworthy AI development: Mechanisms for supporting verifiable claims. Apr. 2020.
- M. Brundage, K. Mayer, T. Eloundou, S. Agarwal, S. Adler, G. Krueger, J. Leike, and P. Mishkin. Lessons learned on language model safety and misuse. <https://openai.com/research/language-model-safety-and-misuse>, Mar. 2022. Accessed: 2023-4-17.
- S. Bubeck, V. Chandrasekaran, R. Eldan, J. Gehrke, E. Horvitz, E. Kamar, P. Lee, Y. T. Lee, Y. Li, S. Lundberg, H. Nori, H. Palangi, M. T. Ribeiro, and Y. Zhang. Sparks of artificial general intelligence: Early experiments with GPT-4. Mar. 2023.
- C. Burns, H. Ye, D. Klein, and J. Steinhardt. Discovering latent knowledge in language models without supervision. Dec. 2022.
- J. Carlsmith. Is Power-Seeking AI an existential risk? June 2022.
- A. Chan, R. Salganik, A. Markelius, C. Pang, N. Rajkumar, D. Krashenninnikov, L. Langosco, Z. He, Y. Duan, M. Carroll, M. Lin, A. Mayhew, K. Collins, M. Molamohammadi, J. Burden, W. Zhao, S. Rismani, K. Voudouris, U. Bhatt, A. Weller, D. Krueger, and T. Maharaj. Harms from increasingly agentic algorithmic systems. Feb. 2023.
- N. Emery-Xu, A. Park, and R. Trager. Uncertainty, information, and risk in international technology races. *Unpublished manuscript*, 2023.
- T. Everitt, R. Carey, E. Langlois, P. A. Ortega, and S. Legg. Agent incentives: A causal perspective. Feb. 2021.

- D. Ganguli, D. Hernandez, L. Lovitt, N. DasSarma, T. Henighan, A. Jones, N. Joseph, J. Kernion, B. Mann, A. Askell, Y. Bai, A. Chen, T. Conerly, D. Drain, N. Elhage, S. El Showk, S. Fort, Z. Hatfield-Dodds, S. Johnston, S. Kravec, N. Nanda, K. Ndousse, C. Olsson, D. Amodei, D. Amodei, T. Brown, J. Kaplan, S. McCandlish, C. Olah, and J. Clark. Predictability and surprise in large generative models. Feb. 2022.
- L. Gao, J. Schulman, and J. Hilton. Scaling laws for reward model overoptimization. Oct. 2022.
- A. Glaese, N. McAleese, M. Trębacz, J. Aslanides, V. Firoiu, T. Ewalds, M. Rauh, L. Weidinger, M. Chadwick, P. Thacker, L. Campbell-Gillingham, J. Uesato, P.-S. Huang, R. Comanescu, F. Yang, A. See, S. Dathathri, R. Greig, C. Chen, D. Fritz, J. S. Elias, R. Green, S. Mokra, N. Fernando, B. Wu, R. Foley, S. Young, I. Gabriel, W. Isaac, J. Mellor, D. Hassabis, K. Kavukcuoglu, L. A. Hendricks, and G. Irving. Improving alignment of dialogue agents via targeted human judgements. Sept. 2022.
- D. Hadfield-Menell, A. Dragan, P. Abbeel, and S. Russell. The Off-Switch game. Nov. 2016.
- E. Jones, A. Dragan, A. Raghunathan, and J. Steinhardt. Automatically auditing large language models via discrete optimization. Mar. 2023.
- Z. Kenton, R. Kumar, S. Farquhar, J. Richens, M. MacDermott, and T. Everitt. Discovering agents. Aug. 2022.
- H. Khlaaf, P. Mishkin, J. Achiam, G. Krueger, and M. Brundage. A hazard analysis framework for code synthesis large language models. July 2022.
- V. Krakovna and J. Kramar. Power-seeking can be probable and predictive for trained agents. Apr. 2023.
- P. Liang, R. Bommasani, T. Lee, D. Tsipras, D. Soylu, M. Yasunaga, Y. Zhang, D. Narayanan, Y. Wu, A. Kumar, B. Newman, B. Yuan, B. Yan, C. Zhang, C. Cosgrove, C. D. Manning, C. Re, D. Acosta-Navas, D. A. Hudson, E. Zelikman, E. Durmus, F. Ladhak, F. Rong, H. Ren, H. Yao, J. Wang, K. Santhanam, L. Orr, L. Zheng, M. Yuksekgonul, M. Suzgun, N. Kim, N. Guha, N. Chatterji, O. Khattab, P. Henderson, Q. Huang, R. Chi, S. M. Xie, S. Santurkar, S. Ganguli, T. Hashimoto, T. Icard, T. Zhang, V. Chaudhary, W. Wang, X. Li, Y. Mai, Y. Zhang, and Y. Koreeda. Holistic evaluation of language models. Nov. 2022.
- I. McKenzie, A. Lyzhov, A. Parrish, A. Prabhu, A. Mueller, N. Kim, S. Bowman, and E. Perez. Inverse scaling prize. <https://github.com/inverse-scaling/prize>, 2022. Accessed: 2023-4-7.
- J. Michael, A. Holtzman, A. Parrish, A. Mueller, A. Wang, A. Chen, D. Madaan, N. Nangia, R. Y. Pang, J. Phang, and S. R. Bowman. What do NLP researchers believe? results of the NLP community metasurvey. Aug. 2022.
- M. Mitchell, S. Wu, A. Zaldivar, P. Barnes, L. Vasserman, B. Hutchinson, E. Spitzer, I. D. Raji, and T. Gebru. Model cards for model reporting. Oct. 2018.
- J. Mokander, J. Schuett, H. R. Kirk, and L. Floridi. Auditing large language models: a three-layered approach. Feb. 2023.
- N. Nanda, L. Chan, T. Lieberum, J. Smith, and J. Steinhardt. Progress measures for grokking via mechanistic interpretability. Jan. 2023.
- R. Ngo, L. Chan, and S. Mindermann. The alignment problem from a deep learning perspective. Aug. 2022.

- C. Olah, N. Cammarata, L. Schubert, G. Goh, M. Petrov, and S. Carter. Zoom in: An introduction to circuits. *Distill*, 5(3), Mar. 2020.
- OpenAI. GPT-4 system card. Mar. 2023a.
- OpenAI. GPT-4 technical report. Mar. 2023b.
- T. Ord. Lessons from the development of the atomic bomb. Technical report, Centre for the Governance of AI, Nov. 2022.
- L. Orseau and S. Armstrong. Safely interruptible agents. <https://www.deepmind.com/publications/safely-interruptible-agents>, 2016. Accessed: 2023-4-26.
- A. Pan, C. J. Shern, A. Zou, N. Li, S. Basart, T. Woodside, J. Ng, H. Zhang, S. Emmons, and D. Hendrycks. Do the rewards justify the means? measuring Trade-Offs between rewards and ethical behavior in the MACHIAVELLI benchmark. Apr. 2023.
- E. Perez, S. Huang, F. Song, T. Cai, R. Ring, J. Aslanides, A. Glaese, N. McAleese, and G. Irving. Red teaming language models with language models. Feb. 2022a.
- E. Perez, S. Ringer, K. Lukošiūtė, K. Nguyen, E. Chen, S. Heiner, C. Pettit, C. Olsson, S. Kundu, S. Kadavath, A. Jones, A. Chen, B. Mann, B. Israel, B. Seethor, C. McKinnon, C. Olah, D. Yan, D. Amodei, D. Amodei, D. Drain, D. Li, E. Tran-Johnson, G. Khundadze, J. Kernion, J. Landis, J. Kerr, J. Mueller, J. Hyun, J. Landau, K. Ndousse, L. Goldberg, L. Lovitt, M. Lucas, M. Sellitto, M. Zhang, N. Kingsland, N. Elhage, N. Joseph, N. Mercado, N. DasSarma, O. Rausch, R. Larson, S. McCandlish, S. Johnston, S. Kravec, S. El Showk, T. Lanham, T. Telleen-Lawton, T. Brown, T. Henighan, T. Hume, Y. Bai, Z. Hatfield-Dodds, J. Clark, S. R. Bowman, A. Askell, R. Grosse, D. Hernandez, D. Ganguli, E. Hubinger, N. Schiefer, and J. Kaplan. Discovering language model behaviors with Model-Written evaluations. Dec. 2022b.
- I. D. Raji, A. Smart, R. N. White, M. Mitchell, T. Gebru, B. Hutchinson, J. Smith-Loud, D. Theron, and P. Barnes. Closing the AI accountability gap: Defining an End-to-End framework for internal algorithmic auditing. Jan. 2020.
- I. D. Raji, I. Elizabeth Kumar, A. Horowitz, and A. D. Selbst. The fallacy of AI functionality. June 2022a.
- I. D. Raji, P. Xu, C. Honigsberg, and D. E. Ho. Outsider oversight: Designing a third party audit ecosystem for AI governance. June 2022b.
- S. Rismani, R. Shelby, A. Smart, E. Jatho, J. Kroll, A. Moon, and N. Rostamzadeh. From plane crashes to algorithmic harm: Applicability of safety engineering frameworks for responsible ML. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, number Article 2 in CHI '23, pages 1–18, New York, NY, USA, Apr. 2023. Association for Computing Machinery.
- R. Shah, V. Varma, R. Kumar, M. Phuong, V. Krakovna, J. Uesato, and Z. Kenton. Goal misgeneralization: Why correct specifications aren't enough for correct goals. Oct. 2022.
- T. Shevlane. Sharing powerful AI models. <https://www.governance.ai/post/sharing-powerful-ai-models>, Jan. 2022a. Accessed: 2023-3-18.
- T. Shevlane. Structured access: An emerging paradigm for safe AI deployment. Feb. 2022b.
- A. Turner, L. Smith, R. Shah, A. Critch, and P. Tadepalli. Optimal policies tend to seek power. *Adv. Neural Inf. Process. Syst.*, 34:23063–23074, Dec. 2021.

- J. Wei, N. Kim, Y. Tay, and Q. V. Le. Inverse scaling can become u-shaped. Nov. 2022a.
- J. Wei, Y. Tay, R. Bommasani, C. Raffel, B. Zoph, S. Borgeaud, D. Yogatama, M. Bosma, D. Zhou, D. Metzler, E. H. Chi, T. Hashimoto, O. Vinyals, P. Liang, J. Dean, and W. Fedus. Emergent abilities of large language models. June 2022b.
- J. Wei, X. Wang, D. Schuurmans, M. Bosma, B. Ichter, F. Xia, E. Chi, Q. Le, and D. Zhou. Chain-of-Thought prompting elicits reasoning in large language models. Jan. 2022c.
- J. Whittlestone and J. Clark. Why and how governments should monitor AI development. Aug. 2021.
- D. M. Ziegler, S. Nix, L. Chan, T. Bauman, P. Schmidt-Nielsen, T. Lin, A. Scherlis, N. Nabeshima, B. Weinstein-Raun, D. de Haas, B. Shlegeris, and N. Thomas. Adversarial training for High-Stakes reliability. May 2022.
- R. Zwetsloot and A. Dafoe. Thinking about risks from AI: Accidents, misuse and structure. <https://www.lawfareblog.com/thinking-about-risks-ai-accidents-misuse-and-structure>, Feb. 2019. Accessed: 2023-3-18.

## Appendix: Deployment safety controls

Variable:	Includes:
<b>Scale</b>	How many end users? How many agents are running at any one time, or how many times per day is the model called? How many applications will be built on top of the model?
<b>Use restrictions</b>	Are certain high stakes applications prohibited?
<b>Generality</b>	Will a single model be flexibly applied across a range of applications, or will narrower, application-specific versions of the model be fine-tuned?
<b>Autonomy</b>	Is the AI system tasked with executing tasks, or merely responding to queries? How long are the chains of actions the model can take? Can the model define new types of actions?
<b>Tool use</b>	Web browsing, telephone calls, code execution, control over robotic hardware, calling APIs, access to persistent memory.
<b>Depth of model access</b>	How restrictive is the interaction between the AI system and the user? Can all of the model's capabilities be accessed? Can users or developers fine-tune the model?
<b>Oversight and moderation</b>	How closely are the model's outputs monitored by the provider? Can certain outputs be automatically filtered out?
<b>Global planning</b>	Does the model have access to many user interactions at once, so it can make plans across them, or is it confined to dealing with each user individually?
<b>Adjustments to model</b>	The developer could decide to deploy a smaller version of the model. The developer could attempt to remove certain dangerous capabilities via fine-tuning.

Table 3 | Variables that affect the risk level of deployment, each of which can be adjusted on the basis of evaluation results.