Frontier Safety Framework

Version 2.0

The Frontier Safety Framework is a set of protocols that aims to address severe risks that may arise from powerful capabilities of foundation models. It is intended to complement Google's existing suite of AI responsibility and safety practices, and enable AI innovation and deployment consistent with our <u>AI Principles</u>.

The Framework is informed by the broader conversation on Frontier AI Safety Frameworks.¹ The core components of Frontier AI Safety Frameworks are to:

- Identify capability levels at which AI models without additional mitigations could pose severe risk
- Implement protocols to detect the attainment of such capability levels
- Prepare and articulate mitigation plans in advance of when such capability levels are attained
- Where appropriate, involve external parties to help inform and guide our approach.

In version 2.0 of the Framework, we specify protocols for the detection of capability levels at which models may pose severe risks (which we call "Critical Capability Levels (CCLs)"), and articulate mitigation approaches to address such risks. At present, the Framework primarily addresses misuse risk,² but we also include an exploratory section addressing deceptive alignment risk,³ focusing on capability levels at which such risks may begin to arise. For each type of risk, we define here a set of CCLs and a mitigation approach for them. Risk assessment will necessarily involve evaluating cross-cutting capabilities such as agency, tool use, reasoning, and scientific understanding.

The safety of frontier AI systems is a global public good. The protocols here represent our current understanding and recommended approach of how severe frontier AI risks may be anticipated and addressed. Importantly, there are certain mitigations whose social value is significantly reduced if not broadly applied to AI systems reaching critical capabilities. These mitigations should be understood as recommendations for the industry collectively: our adoption of them would only result in effective risk mitigation for society if all relevant organizations provide similar levels of protection, and our adoption of the protocols described in this Framework may depend on whether such organizations across the field adopt similar protocols.

The Framework is exploratory and based on early research. We may change our approach and recommendations over time as we gain experience and insights on the projected capabilities of future frontier models. We will review the Framework periodically and we expect it to evolve substantially as our understanding of the risks and benefits of frontier models improves.

- ¹ See https://www.gov.uk/government/publications/emerging-processes-for-frontier-ai-safety, https://metr.org/blog/2023-09-26-rsp/, https://metr.org///metr.org/)
- https://www.frontiermodelforum.org/updates/issue-brief-components-of-frontier-ai-safety-frameworks/,

² As in, in the context of the Framework, risks of threat actors using critical capabilities of deployed or exfiltrated models to cause harm.

³ As in, in the context of the Framework, risks of highly autonomous systems purposefully undermining human control over AI systems.

Table of Contents:	
Framework	3
Misuse:	5
Mitigation approach	5
Misuse Critical Capability Levels	6
Table 1: Misuse CCLs and Security Mitigations	6
Deceptive Alignment:	7
Mitigation approach and critical capability levels	8
Table 2: Deceptive Alignment CCLs and Mitigations	8
Safety cases and control evaluations	8
Governance and accountability:	9
Governance	9
Disclosure	9
Future work	9
Acknowledgements	9

Framework

This section describes the central components of the Frontier Safety Framework. These protocols represent our current understanding of and recommended approach for how severe frontier AI risks may be anticipated and addressed.

1 - Critical Capability Levels

The Framework is built around capability thresholds called "Critical Capability Levels." These are capability levels at which, absent mitigation measures, AI models or systems may pose heightened risk of severe harm. CCLs can be determined by identifying and analyzing the main foreseeable paths through which a model could cause severe harm, and then defining the CCLs as the minimal set of capabilities a model must possess to do so. Note that we have selected our CCLs to be conservative; it is not clear to what extent CCLs might translate to harm in real-world contexts.

We describe two sets of CCLs: misuse CCLs that can indicate heightened risk of severe harm from misuse if not addressed, and deceptive alignment CCLs that can indicate heightened risk of deceptive alignment-related events if not addressed.

For misuse risk, we define CCLs in high-risk domains where, based on early research, we believe risks of severe harm may be most likely to arise from future models:

- **CBRN:** Risks of models assisting in the development, preparation, and/or execution of a chemical, • biological, radiological, or nuclear ("CBRN") attack.
- Cyber: Risks of models assisting in the development, preparation, and/or execution of a cyber attack.
- Machine Learning R&D: Risks of the misuse of models capable of accelerating the rate of AI • progress to potentially destabilizing levels, the result of which could be the unsafe attainment or proliferation of other powerful AI models. Capabilities in this area are under active research, and in the longer term may exacerbate frontier AI risks-including in other risk domains-if insufficiently managed.

For deceptive alignment risk, the initial approach focuses on detecting when models might develop a baseline instrumental reasoning ability at which they have the potential to undermine human control, assuming no additional mitigations were applied. The two instrumental reasoning CCLs thus focus on

delineating when such capability becomes present, and subsequently when the initial mitigation for this capability—automated monitoring—is no longer adequate.

2 - Assessing the Capabilities of Frontier Models

We intend to evaluate our most powerful frontier models regularly to check whether their AI capabilities are approaching a CCL. We also intend to evaluate any of these models that could indicate an exceptional increase in capabilities over previous models, and where appropriate, assess the likelihood of such capabilities and risks before and during training.

To do so, we will define a set of <u>evaluations</u> called "early warning evaluations," with a specific "alert threshold" that flags when a CCL may be reached before the evaluations are run again. In our evaluations, we seek to equip the model with appropriate scaffolding and other augmentations to make it more likely that we are also assessing the capabilities of systems that will likely be produced with the model. We may run early warning evaluations more frequently or adjust the alert threshold of our evaluations if the rate of progress suggests our safety buffer is no longer adequate.

Where necessary, early warning evaluations may be supplemented by other evaluations to better understand model capabilities relative to our CCLs. We may use additional external evaluators to test a model for relevant capabilities, if evaluators with relevant expertise are needed to provide an additional signal about a model's proximity to CCLs.

<u>3 - Applying Mitigations</u>

When a model reaches an alert threshold for a CCL, we will assess the proximity of the model to the CCL and analyze the risk posed, involving internal and external experts as needed. This will inform the formulation and application of a response plan.

Central to most response plans will be the application of the mitigations described later in this document. For misuse, we have two categories of mitigations: *security mitigations* intended to prevent the exfiltration of model weights, and *deployment mitigations* (such as safety fine-tuning and misuse filtering, detection, and response) intended to counter the misuse of critical capabilities in deployments. For deceptive alignment risk, *automated monitoring* may be applied to detect and respond to deceptive behavior for models that meet the first deceptive alignment CCL. Note that these mitigations reflect considerations from the perspective of addressing severe risks from powerful capabilities alone; due to this focused scope, other risk management and security considerations may result in more stringent mitigations applied to a model than specified by the Framework.

A model flagged by an alert threshold may be assessed to pose risks for which readily available mitigations (including but not limited to those described below) may not be sufficient. If this happens, the response plan may involve putting deployment or further development on hold until adequate mitigations can be applied. Conversely, where model capabilities remain quite distant from a CCL, a response plan may involve the adoption of additional capability assessment processes to flag when heightened mitigations may be required.

The appropriateness and efficacy of applied mitigations should be reviewed periodically, drawing on information like related misuse or misuse attempt incidents; results from continued post-mitigation testing; statistics about our intelligence, monitoring and escalation processes; and updated threat modeling and risk landscape analysis.

Misuse

This section describes our mitigation approach for models that pose risks of severe harm through misuse, and then details our set of misuse CCLs, as well as the mitigations that we provisionally assess as appropriate for them.

Mitigation Approach

There are two categories of mitigations to address models with misuse critical capabilities: *security mitigations* intended to prevent the exfiltration of model weights, and *deployment mitigations* intended to counter the misuse of critical capabilities in deployments. For security, we have several levels of mitigations, allowing calibration of the robustness of security measures to the risks posed. For deployment mitigations, we specify a standard process for applying, assessing, and reviewing mitigations: the aim of this process is to calibrate mitigations to CCLs, and the procedural approach reflects the more iterative and CCL-dependent nature of deployment mitigations.

Security Mitigations

Security mitigations against exfiltration risk are important for models reaching CCLs. This is because the release of model weights may enable the removal of any safeguards trained into or deployed with the model, and hence access (including by threat actors) to any critical capabilities. Here, we rely on the RAND framework⁴ to articulate the level of security recommended for each CCL. When we reference RAND security levels, we are referring to the security principles in their framework, rather than the benchmarks (i.e. concrete measures) also described in the RAND report.⁵ Because AI security is an area of active research, we expect the concrete measures implemented to reach each level of security to evolve substantially.

Deployment Mitigations

The following deployment mitigation process will be applied to models reaching a CCL, allowing for iterative and flexible tailoring of mitigations to each risk and use case.⁶

- 1. Development and assessment of mitigations: safeguards and an accompanying safety case⁷ are developed by iterating on the following:
 - **a.** Developing and improving a suite of safeguards targeting the capability. This includes, as appropriate, safety fine-tuning, misuse filtering and detection, and response protocols.
 - **b.** Assessing the robustness of these mitigations against the risk posed through assurance evaluations and threat modeling research. The assessment takes the form of a safety case, taking into account factors such as the likelihood and consequences of misuse.
- 2. Pre-deployment review of safety case: general availability deployment⁸ of a model takes place only after the appropriate corporate governance body determines the safety case regarding each CCL the model has reached to be adequate.
- **3.** Post-deployment review of safety case: the safety case will be updated through red-teaming and revisions to our threat models. The safeguards for the model may be updated as well to ensure continued adequacy.

Misuse Critical Capability Levels

The table below details a set of CCLs we have identified through ongoing analysis of the CBRN, Cyber, and Machine Learning R&D risk domains. We expect these to evolve over time. We recommend a security level to each of these CCLs, which reflect our assessment of the minimum appropriate level of security

https://www.aisi.gov.uk/work/safety-case-template-for-inability-arguments and https://arxiv.org/abs/2410.21572. ⁸ General availability deployment refers to public deployments intended for production use and without the imposition of significant user or use restrictions.

⁴ See <u>https://www.rand.org/pubs/research_reports/RRA2849-1.html</u>, pp 21-22.

⁵ As the authors point out, the "security level benchmarks represent neither a complete standard nor a compliance regime—they are provided for informational purposes only and should inform security teams' decisions rather than supersede them."

⁶ This process is derived from that required by Deployment Level 2 in <u>Frontier Safety Framework version 1.0</u>. Our rationale for this change is as follows: the process extends naturally to cover risks for which "prevention of access" (Deployment Level 3) is required, and is worth applying to all CCLs.

⁷ A safety case is an assessable argument showing how severe risks associated with a model's CCLs have been minimised to an appropriate level. See also, for reference,

the field of frontier AI should apply to models reaching each CCL. In practice, our overall security posture may commonly exceed the baseline levels recommended here.

These recommended security levels reflect our current thinking and may be adjusted if our empirical understanding of the risks changes. This may occur if, for example, a model does not possess capabilities meaningfully different from other widely available models that have not demonstrably caused or contributed to severe risks, or if we assess that the benefits of the open release of model weights outweigh the risks. Relatedly, we believe these recommendations will only be effective if the entire frontier AI field applies them, and of limited social utility if not.

Risk domain ⁹	Critical capability level	Recommended security level and rationale
CBRN: Risks of models assisting in the development, preparation and/or execution of CBRN attacks.	CBRN uplift 1: Can be used to significantly assist a low-resourced actor with dual-use scientific protocols, resulting in a substantial increase in ability to cause a mass casualty event. ¹⁰	Security controls and detections at a level generally aligned with RAND SL 2 The potential magnitude of harm these capabilities may enable means the exfiltration and leak of model weights reaching this CCL could be highly damaging. However, low-resourced actors are unlikely to pose a substantial exfiltration threat.
Cyber: Risks of models assisting in the execution of a cyber attack.	Cyber autonomy level 1: Can be used to drastically reduce the cost (e.g. through full automation) of carrying out end-to-end cyberattacks ¹¹ on organizations with a limited security posture.	Security controls and detections at a level generally aligned with RAND SL 2 Harmful cyberattacks against organizations with limited security posture can already be carried out by individuals with limited expertise, but the automation of such attacks would significantly lower the costs of doing so. Exfiltration of model weights could enable the execution of such attacks at scale. However, cybersecurity may improve correspondingly when models reach such capability levels. The relatively ambiguous net costs of exfiltration count against security levels with higher costs to innovation.
	Cyber uplift level 1: Can be used to significantly assist with high impact cyber attacks, ¹² resulting in overall cost/resource reductions of an order of magnitude or more. ¹³	Security controls and detections at a level generally aligned with RAND SL 2 A model at this capability level could help fairly well-resourced threat actors carry out severe cyber attacks on targets like critical businesses, national government entities, and critical national infrastructure with lower resource expenditure, potentially increasing the frequency of such attacks significantly. However, as above, cyber defense may improve to diminish the impact of

Table 1: Misuse CCLs and Security Mitigations

⁹ Note that we have removed the Autonomy risk domain, which was included in Frontier Safety Framework version 1.0. Most of the advanced risk that was captured by this CCL is now covered by our misalignment section. From the perspective of misuse risks, our threat models suggest that no heightened deployment mitigations would be necessary, and that security controls and detection at a level generally aligned with RAND SL 2 would be adequate. ¹⁰ For example, through the use of a self-replicating CBRN agent. Compared to a counterfactual of not using generative AI systems.

¹¹ E.g. deletion or exfiltration of critical information, or destroying or disabling key systems.

¹² E.g. deletion or exfiltration of sensitive information/disruption of key systems of organizations with strong security posture, or the creation of predictable and large-scale effects on critical national infrastructure.

¹³ Relative to the counterfactual of using 2024 AI technology and tooling.

		AI-assisted cyber attacks. Similarly, the ambiguous net costs of exfiltration count against security levels with higher costs to innovation.
Machine Learning R&D: Risks of the misuse of models capable of accelerating the rate of AI progress, the result of which could be the unsafe attainment or proliferation of other powerful AI models.	Machine Learning R&D uplift level 1: Can or has been used to accelerate AI development, resulting in AI progress substantially accelerating (e.g. 2x) from 2020-2024 rates.	Security controls and detections at a level generally aligned with RAND SL 3 Unrestricted access to models at this level of capability could significantly increase an actor's ability to progress to yet more powerful models and other critical capabilities. The exfiltration of such a model may therefore have a significant effect on society's ability to adapt to and govern powerful AI systems, effects that may have long-lasting consequences. Substantially strengthened security is therefore recommended.
	Machine Learning R&D autonomy level 1: Can fully automate the AI R&D pipeline at a competitive cost. The cost comparison is relative to humans augmented by AI tools.	We recommend security controls and detections at a level generally aligned with RAND SL 4 for this capability threshold, but emphasize that this must be taken on by the frontier AI field as a whole. Unrestricted access to models at this level of capability could give an actor with adequate computational resources the ability to reach capabilities much more powerful than those in the other CCLs listed in a short amount of time. This could be catastrophic if there is no effective way of defending against rapidly improving and potentially superhuman AI systems wielded by threat actors. Therefore, we recommend models at this level of capability have exceptional security even though they may have substantial innovation costs.

Deceptive Alignment

This section describes an initial approach for addressing risks of deceptive alignment. By "deceptive alignment," we mean the risk that AI systems purposefully undermine human control over AI systems. While we do not express any opinion here about how likely it is for deceptive alignment to arise, this section's provisions aim to prepare for this possibility.

Note that the effective mitigation of deceptive alignment risk requires the application of safeguards to all AI systems reaching critical capabilities. It is important that the industry adopts strategies for addressing this risk: what is detailed in this section represents one possible approach for doing so.

Mitigation Approach and Critical Capability Levels

An initial mitigation approach focuses on detecting when models might develop a baseline instrumental reasoning ability at which they have the potential to undermine human control, assuming no additional mitigations were applied. When models reach this capability level, we believe applying an automated monitor to the model's explicit reasoning (e.g. chain-of-thought output) is an effective mitigation. Once a model is capable of effective instrumental reasoning in ways that cannot be monitored, additional mitigations may be warranted—the development of which is an area of active research.

The table below summarizes the initial approach:

Deceptive Alignment CCL	Mitigation
Instrumental Reasoning Level 1: the instrumental reasoning abilities of the model enable enough situational awareness (ability to work out and use relevant details of its deployment setting) and stealth (ability to circumvent basic oversight mechanisms) such that, absent additional mitigations, we cannot rule out the model significantly undermining human control.	Automated monitoring: a monitoring system to detect illicit use of instrumental reasoning capabilities in a model's chain-of-thought and use it in high-stakes deployments where there is significant risk of the model undermining human control.
Instrumental Reasoning Level 2: the instrumental reasoning abilities of the model enable enough situational awareness and stealth that, even when relevant model outputs (including, e.g., scratchpads) are being monitored, we cannot detect or rule out the risk of a model significantly undermining human control.	Future work: We are actively researching approaches to addressing models that reach this CCL.

Table 2: Deceptive Alignment CCLs and Mitigations

Safety Cases and Control Evaluations

Looking forward, we expect the approach to deceptive alignment risks to take a similar form as the deployment mitigation process. For each deceptive alignment risk, AI developers should:

- Develop a safety case based on the model's capabilities and the mitigations applied and tested. Such safety cases would make the case that the likelihood of deceptive alignment risk would be low enough for safe deployment, even if the model were trying to meaningfully undermine human control.
- Test the mitigations applied to models via red-teaming protocols called "control evaluations,"¹⁴ to determine whether they meet the requirements of the safety case.

The approach above relies on two safety cases centering respectively on the claims that:

- 1. A model does not meet a threshold CCL and therefore would be incapable of substantially undermining human control.
- 2. An automated monitoring process that is validated by a control evaluation is adequate for catching the model's malicious behavior.

Once neither of these apply, that is, once automated monitoring is no longer adequate/can no longer pass a control evaluation, a model is deemed to have reached the second CCL.

Governance and Accountability

Governance

For Google models, when alert thresholds are reached, the response plan will be reviewed and approved by appropriate corporate governance bodies such as the Google DeepMind AGI Safety Council, Google DeepMind Responsibility and Safety Council, and/or Google Trust & Compliance Council. The Google DeepMind AGI Safety Council will periodically review the implementation of the Framework.

¹⁴ See e.g. <u>https://arxiv.org/abs/2312.06942</u>.

Disclosure

If we assess that a model has reached a CCL that poses an unmitigated and material risk to overall public safety, we aim to share information with appropriate government authorities where it will facilitate the development of safe AI. Where appropriate, and subject to adequate confidentiality and security measures and considerations around proprietary and sensitive information, this information may include:

- Model information: characteristics of the AI model relevant to the risk it may pose with its critical capabilities.
- Evaluation results: such as details about the evaluation design, the results, and any robustness tests.
- Mitigation plans: descriptions of our mitigation plans and how they are expected to reduce the risk.

We may also consider disclosing information to other external organizations to promote shared learning and coordinated risk mitigation. We will continue to review and evolve our disclosure process over time.

Future Work

We expect the Framework to evolve substantially as our understanding of the risks and benefits of frontier models improves, and we will publish substantive revisions as appropriate. Issues that we aim to address in future versions of the Framework include:

- **Greater precision in risk modeling:** While we have updated our CCLs and underlying threat models from version 1.0, there remains significant room for improvement in understanding the risks posed by models in different domains, and refining our set of CCLs.
- **Capability elicitation:** Our evaluators continue to improve their ability to estimate what capabilities may be attainable by different threat actors with access to our models, taking into account a growing number of possible post-training enhancements.
- Updated set of risks and mitigations: There may be additional risk domains and critical capabilities that fall into scope as AI capabilities improve and the external environment changes. Future work will aim to include additional pressing risks, which may include additional risk domains or higher CCLs within existing domains.
- **Deceptive alignment approach beyond automated monitoring:** We are actively researching approaches to addressing models that reach the "Instrumental Reasoning Level 2" CCL.
- Broader approach to ML R&D risks: The risks posed by models reaching our ML R&D CCLs may require measures beyond security and deployment mitigations, to ensure that safety measures and social institutions continue to be able to adapt to new AI capabilities amidst possible acceleration in AI development. We are actively researching appropriate responses to these scenarios.

Acknowledgements

Version 2.0 of the Frontier Safety Framework was developed by Lewis Ho, Celine Smith, Claudia van der Salm, Joslyn Barnhart, and Rohin Shah, under the leadership of Allan Dafoe, Anca Dragan, Andy Song, Demis Hassabis, Four Flynn, Jennifer Beroshi, Helen King, Nicklas Lundblad, and Tom Lue. The Framework was developed with substantial contributions from Aalok Mehta, Adam Stubblefield, Alex Kaskasoli, Alice Friend, Amy Merrick, Anna Wang, Ben Bariach, Charley Snyder, David Bledin, David Lindner, Dawn Bloxwich, Don Wallace, Eva Lu, Heidi Howard, Iason Gabriel, James Manyika, Joana Iljazi, Kent Walker, Lila Ibrahim, Mary Phuong, Mikel Rodriguez, Peng Ning, Roland S. Zimmermann, Samuel Albanie, Sarah Cogan, Sasha Brown, Seb Farquhar, Sebastien Krier, Shane Legg, Victoria Krakovna, Vijay Bolina, Xerxes Dotiwalla, Ziyue Wang. We are grateful for input from Apollo Research, Carnegie Endowment for International Peace, Center for Long-Term Resilience, Center for the Governance of Al, METR, Redwood Research, UK Al Safety Institute, and several independent researchers.