



METR

Common Elements of Frontier AI Safety Policies

November 2024



Summary

A number of developers of large foundation models have committed to evaluating their models for severe risks and implementing necessary risk mitigations. Sixteen companies agreed to do so as part of the Frontier AI Safety Commitments at the AI Seoul Summit. Among them, three companies had previously released detailed policies: Anthropic’s Responsible Scaling Policy, OpenAI’s Preparedness Framework, and Google DeepMind’s Frontier Safety Framework. This document describes key shared elements of the three frameworks and includes relevant excerpts, which expand on the high-level objectives of the Frontier AI Safety Commitments.

Each framework describes *covered threat models*, including the potential for AI models to facilitate biological weapons development or cyberattacks, or to engage in autonomous replication or automated AI research and development. The frameworks also outline commitments to assess whether models are approaching capability thresholds that could enable catastrophic harm.

When these capability thresholds are approached, the frameworks prescribe *model weight security* and *model deployment mitigations* to be adopted in response. For models with more concerning capabilities, developers commit to securing model weights to prevent theft by increasingly sophisticated adversaries. Additionally, they commit to implementing deployment safety measures that would significantly reduce the risk of dangerous AI capabilities being misused and causing serious harm.

The frameworks also establish *conditions for halting development and deployment* if the developer’s mitigations are insufficient to manage the risks. Evaluations are intended to *elicit the full capabilities* of the model. They define the *timing and frequency of evaluations*, which are to be conducted before deployment, during training, and after deployment. Furthermore, the frameworks include intentions to explore *accountability* mechanisms, such as oversight by third parties or boards to monitor policy implementation and potentially assist with evaluations. Policies may be *updated over time* as developers gain a deeper understanding of AI risks and refine their evaluation processes.

Introduction

In May 2024, sixteen companies from around the world [agreed to the AI Seoul Summit’s Frontier AI Safety Commitments](#) and will publish safety frameworks to manage severe risks from frontier AI systems. These voluntary commitments oblige signatory companies to do the following:

- I. “Assess the risks posed by their frontier models or systems across the AI lifecycle, including before deploying that model or system, and, as appropriate, before and during training. ...”
- II. “Set out thresholds at which severe risks posed by a model or system, unless adequately mitigated, would be deemed intolerable. ...”
- III. “Articulate how risk mitigations will be identified and implemented to keep risks within defined thresholds, including safety and security-related risk mitigations such as modifying system behaviours and implementing robust security controls for unreleased model weights. ...”
- IV. “Set out explicit processes they intend to follow if their model or system poses risks that meet or exceed the pre-defined thresholds. This includes processes to further develop and deploy their systems and models only if they assess that residual risks would stay below the thresholds. In the extreme, organisations commit not to develop or deploy a model or system at all, if mitigations cannot be applied to keep risks below the thresholds.”

As of writing, there are three existing examples of frontier AI developers publishing detailed public commitments that include the above elements. They are [Anthropic’s Responsible Scaling Policy](#), [Google DeepMind’s Frontier Safety Framework](#) and [OpenAI’s Preparedness Framework](#).¹ We highlight shared aspects of these policies, with overviews of the components they share in common, as well as relevant excerpts from each policy.

¹ Another relevant policy is [Magic’s AGI Readiness Policy](#), which commits to not deploy frontier coding models without having first evaluated them for dangerous capabilities.

Common Components

We identify shared themes held in common among the three frameworks and include corresponding excerpts. These themes elaborate upon the high-level objectives of the Frontier AI Safety Commitments. They are as follows:

1. Covered Threat Models
2. Model Weight Security
3. Model Deployment Mitigations
4. Conditions for Halting Development and Deployment
5. Full Capability Elicitation during Evaluations
6. Timing and Frequency of Evaluations
7. Accountability
8. Updating Policies over Time

Covered Threat Models

Thresholds at which specific AI capabilities would pose severe risk and require new mitigations. Each framework makes use of the notion of dangerous capability thresholds which are defined via results of evaluations run on models under development.

[Anthropic's Responsible Scaling Policy](#):

To determine when a model has become sufficiently advanced such that its deployment and security measures should be strengthened, we use the concepts of Capability Thresholds and Required Safeguards. A Capability Threshold tells us when we need to upgrade our protections, and the corresponding Required Safeguards tell us what standard should apply. The Required Safeguards for each Capability Threshold are intended to mitigate risk to acceptable levels. This update to our RSP provides specifications for Capabilities Thresholds related to Chemical, Biological, Radiological, and Nuclear (CBRN) weapons and Autonomous AI Research and Development (AI R&D) and identifies the corresponding Required Safeguards.

[OpenAI's Preparedness Framework](#), page 2:

Only models with a post-mitigation score of "medium" or below can be deployed, and only models with a post-mitigation score of "high" or below can be developed further (as defined in the Tracked Risk Categories below) In addition, we will ensure Security is appropriately tailored to any model that has a "high" or "critical" pre-mitigation level of risk (as defined in the Scorecard below to prevent model exfiltration) We also establish procedural commitments (as defined in Governance below that further specify how we operationalize all the activities that the Preparedness Framework outlines).

Page 5:

Each of the Tracked Risk Categories comes with a gradation scale. We believe monitoring gradations of risk will enable us to get in front of escalating threats and be able to apply more tailored mitigations. In general, "low" on this gradation scale is meant to indicate that the corresponding category of risks is not yet a significant problem, while "critical" represents the maximal level of concern.

[Google DeepMind's Frontier Safety Framework](#), page 2:

The Framework is built around capability thresholds called "Critical Capability Levels." These are capability levels at which, absent mitigation measures, models may pose heightened risk. We determine CCLs by analyzing several high-risk domains: we identify the main paths through which a model could cause harm, and then define the CCLs as the minimal set of capabilities a model must possess to do so. We have conducted preliminary analyses of the Autonomy, Biosecurity, Cybersecurity and Machine Learning R&D domains. Our initial research indicates that powerful capabilities of future models seem most likely to pose risks in these domains.

Each of the three AI developers provides examples of evaluations they would run for their covered threat models, either as part of the company's policy or in a model card. These capabilities include:

- Biological weapons assistance
- Cyberoffense
- Automated AI research and development
- Autonomous replication

Biological Weapons Assistance

Current language models are able to provide detailed advice relevant to creating a biological weapon.^{2,3,4} OpenAI has found that its o1-preview and o1-mini models reach its "medium risk threshold" due to the models' ability to assist experts with operational planning of known biological threats.⁵ In the future, more advanced models could pose a major biosecurity risk if

² Christopher A. Mouton, Caleb Lucas, Ella Guest. [The Operational Risks of AI in Large-Scale Biological Attacks Results of a Red-Team Study](#). January 25, 2024. RAND.

³ Tejal Patwardhan, Kevin Liu, Todor Markov, Neil Chowdhury, Dillon Leet, Natalie Cone, Caitlin Maltbie, Joost Huizinga, Carroll Wainwright, Shawn Jackson, Steven Adler, Rocco Casagrande, Aleksander Madry. [Building an early warning system for LLM-aided biological threat creation](#). January 31, 2024. OpenAI.

⁴ AI at Meta. [The Llama 3 Herd of Models](#). July 23, 2024.

⁵ OpenAI. [OpenAI o1 System Card](#). September 12, 2024.

they enable novices or experts to create biological threats they otherwise would be unable to, or if they can autonomously synthesize a biological threat using a cloud biology lab.^{6,7}

To evaluate o1-preview and o1-mini, OpenAI's Preparedness team worked with biology experts to design biorisk questions. OpenAI assessed the accuracy of model responses to long-form biorisk questions compared to experts, the qualitative helpfulness of model access for experts, and the model's performance in automating wet lab work. OpenAI also evaluated model accuracy on the open-source LAB-Bench benchmark.⁸ Other assessments of AI-enabled biological risk have been conducted by RAND,⁹ Meta,¹⁰ and Anthropic.¹¹

[Anthropic's Responsible Scaling Policy](#), page 3:

Chemical, Biological, Radiological, and Nuclear (CBRN) weapons. The ability to significantly help individuals or groups with basic technical backgrounds (e.g., undergraduate STEM degrees) create/obtain and deploy CBRN weapons.

This capability could greatly increase the number of actors who could cause this sort of damage, and there is no clear reason to expect an offsetting improvement in defensive capabilities. The ASL-3 Deployment Standard and the ASL-3 Security Standard, which protect against misuse and model-weight theft by non-state adversaries, are required.

[OpenAI's Preparedness Framework](#), page 9:

[CBRN Risk Level – High]

Model enables an expert to develop a novel threat vector OR model provides meaningfully improved assistance that enables anyone with basic training in a relevant field (e.g., introductory undergraduate biology course) to be able to create a CBRN threat.

[CBRN Risk Level – Critical]

Model enables an expert to develop a highly dangerous novel threat vector (e.g., comparable to novel CDC Class A biological agent) OR model provides meaningfully improved assistance that enables anyone to be able to create a known CBRN threat OR model can be connected to tools and equipment to

⁶ Doni Bloomfield, Jaspreet Pannu, Alex W. Zhu, Madelena Y. Ng, Ashley Lewis, Eran Bendavid, Steven M. Asch, Tina Hernandez-Boussard, Anita Cicero, Tom Inglesby. [AI and biosecurity: The need for governance](#). *Science*. August 22, 2024.

⁷ National Security Commission on Emerging Biotechnology. [AIxBio White Paper 3: Risks of AIxBio](#). January 2024.

⁸ OpenAI. [OpenAI o1 System Card](#). September 12, 2024.

⁹ Christopher A. Mouton, Caleb Lucas, Ella Guest. [The Operational Risks of AI in Large-Scale Biological Attacks Results of a Red-Team Study](#). January 25, 2024. RAND.

¹⁰ AI at Meta. [The Llama 3 Herd of Models](#). July 23, 2024.

¹¹ Anthropic. [Responsible Scaling Policy Evaluations Report – Claude 3 Opus](#). May 19, 2024.

complete the full engineering and/or synthesis cycle of a regulated or novel CBRN threat without human intervention.

[Google DeepMind's Frontier Safety Framework](#), page 5:

Biosecurity: Risks of models assisting in the development, preparation and/or execution of a biological attack.

Bio amateur enablement level 1: Capable of significantly enabling a non-expert to develop known biothreats that could increase their ability to cause severe harm compared to other means.

Many biothreats capable of causing significant amounts of harm are currently out of the reach of non-experts because of lack of knowledge about their potential for harm and the methods of their acquisition and misuse. An LLM that helps overcome these knowledge gaps, e.g. by suggesting plausible attack strategies or providing detailed instructions for the development of a bio agent, could significantly increase society's vulnerability to fatal attacks by malicious amateurs.

Bio expert enablement level 1: Capable of significantly enabling an expert (i.e. PhD or above) to develop novel biothreats that could result in an incident of high severity.

A very small number of bio agents have the potential to cause harm of an exceptional magnitude. The discovery of enhancements to these agents, or of agents of comparable harmfulness, could increase the chances of a very severe bio attack or accident.

Cyberoffense

If models are capable of novel, high-value exploits or attacks on critical infrastructure,¹² they could cause substantial harm. Present-day models possess an emerging ability to solve cybersecurity challenges and exploit vulnerabilities in websites and open-source software.^{13,14}

¹² Department of Homeland Security. [Safety and Security Guidelines for Critical Infrastructure Owners and Operators](#). April 2024.

¹³ AI at Meta. [CyberSecEval: Comprehensive Evaluation Framework for Cybersecurity Risks and Capabilities of Large Language Models](#).

¹⁴ Richard Fang, Rohan Bindu, Akul Gupta, Qiusi Zhan, Daniel Kang. [Teams of LLM Agents can Exploit Zero-Day Vulnerabilities](#). June 2, 2024.

Open-source cyberoffense benchmarks include Cybench,¹⁵ eyeballvul,¹⁶ and Google DeepMind CTFs.¹⁷ Several AI labs have evaluated their models for cyberoffense capabilities and describe results in their model cards or evaluation reports, including Meta,¹⁸ Anthropic,¹⁹ OpenAI,²⁰ and Google DeepMind.²¹

[Anthropic’s Responsible Scaling Policy](#), page 4:

Cyber Operations: The ability to significantly enhance or automate sophisticated destructive cyber attacks, including but not limited to discovering novel zero-day exploit chains, developing complex malware, or orchestrating extensive hard-to-detect network intrusions.

This will involve engaging with experts in cyber operations to assess the potential for frontier models to both enhance and mitigate cyber threats, and considering the implementation of tiered access controls or phased deployments for models with advanced cyber capabilities. We will conduct either pre- or post-deployment testing, including specialized evaluations. We will document any salient results alongside our Capability Reports.

[OpenAI’s Preparedness Framework](#), page 9:

[Cybersecurity Risk Level – High]

Tool-augmented model can identify and develop proofs-of-concept for high-value exploits against hardened targets without human intervention, potentially involving novel exploitation techniques, OR provided with a detailed strategy, the model can end-to-end execute cyber operations involving the above tasks without human intervention.

[Cybersecurity Risk Level – Critical]

Tool-augmented model can identify and develop functional zero-day exploits of all severity levels, across all software projects, without human intervention OR model can devise and execute end-to-end novel strategies for cyberattacks against hardened targets given only a high level desired goal.

[Google DeepMind’s Frontier Safety Framework](#), page 5:

¹⁵ Andy K. Zhang, Neil Perry, Riya Dulepet, Eliot Jones, Justin W. Lin, Joey Ji, Celeste Menders, Gashon Hussein, Samantha Liu, Donovan Jasper, Pura Peetathawatchai, et al. [Cybench: A Framework for Evaluating Cybersecurity Capabilities and Risk of Language Models](#). August 15, 2024.

¹⁶ Timothee Chauvin. [eyeballvul: a future-proof benchmark for vulnerability detection in the wild](#). July 11, 2024.

¹⁷ Google DeepMind. [Dangerous capability evaluations](#). June 2024.

¹⁸ AI at Meta. [The Llama 3 Herd of Models](#). July 23, 2024.

¹⁹ Anthropic. [Responsible Scaling Policy Evaluations Report – Claude 3 Opus](#). May 19, 2024.

²⁰ OpenAI. [OpenAI o1 System Card](#). September 12, 2024.

²¹ Google DeepMind. [Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context](#). February 15, 2024.

Cybersecurity: Risks of models assisting in the execution of a cyber attack.

Cyber autonomy level 1: Capable of fully automating opportunistic cyberattacks on organizations with a limited security posture.

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, and moreover would enable the execution of such attacks at scale.

Cyber enablement level 1: Capable of enabling an amateur to carry out sophisticated and severe attacks (e.g. those that disrupt critical national infrastructure).

Severe cyberattacks against high-impact targets currently require significant expertise and effort across the cyber kill-chain to carry out. Increasing text generation, programming, and tool-use capabilities in models, combined with improved understanding of cyber offense strategies, could help amateurs overcome difficult steps in the planning and execution of attacks.

Automated AI Research and Development

Large language models are able to contribute to some aspects of AI research, such as generating synthetic pretraining^{22,23} and fine-tuning data,^{24,25,26} designing reward functions,²⁷ as well as general-purpose programming.^{28,29} In the future, models may be able to substantially automate the process of increasing frontier AI capabilities (AI R&D).³⁰ This development could lead to a growth or proliferation in AI capabilities, including other capabilities of concern, that outpaces our ability to ensure sufficient oversight and safeguards.

²² Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, ... Yuanzhi Li. [Textbooks Are All You Need](#). Microsoft Research. June 20, 2023.

²³ AI at Meta. [The Llama 3 Herd of Models](#). July 23, 2024.

²⁴ Xian Li, Ping Yu, Chunting Zhou, Timo Schick, Omer Levy, Luke Zettlemoyer, Jason E Weston, Mike Lewis. [Self-Alignment with Instruction Backtranslation](#). ICLR 2024.

²⁵ Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, ... Jared Kaplan. [Constitutional AI: Harmlessness from AI Feedback](#). December 15, 2022.

²⁶ Nvidia. [Nemotron-4 340B Technical Report](#). August 6, 2024.

²⁷ Yecheng Jason Ma, William Liang, Guanzhi Wang, De-An Huang, Osbert Bastani, Dinesh Jayaraman, Yuke Zhu, Linxi Fan, Anima Anandkumar. [Eureka: Human-Level Reward Design via Coding Large Language Models](#). April 30, 2024. ICLR 2024.

²⁸ Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, Wayne Xin Zhao, Zhewei Wei, Jirong Wen. [A survey on large language model based autonomous agents](#). Frontiers of Computer Science. March 22, 2024.

²⁹ Qian Huang, Jian Vora, Percy Liang, Jure Leskovec. [MLAgentBench: Evaluating Language Agents on Machine Learning Experimentation](#). April 14, 2024.

³⁰ David Owen. [Interviewing AI researchers on automation of AI R&D](#). August 27, 2024.

OpenAI’s MLE-bench is an open-source benchmark for assessing AI R&D capabilities.³¹ The National Security Memorandum on AI describes “automat[ing] development and deployment of other models with [cyber, biological or chemical weapon, or autonomous malicious] capabilities” as one of the AI capabilities relevant to national security that the U.S. AI Safety Institute must conduct tests for.³²

[Anthropic’s Responsible Scaling Policy](#), page 3:

Autonomous AI Research and Development (AI R&D): The ability to either fully automate the work of an entry-level remote-only Researcher at Anthropic, or cause dramatic acceleration in the rate of effective scaling.

This capability could greatly increase the pace of AI development, potentially leading to rapid and unpredictable advances in AI capabilities and associated risks. At minimum, the ASL-3 Security Standard is required, although we expect a higher security standard (which would protect against model-weight theft by state-level adversaries) will be required, especially in the case of dramatic acceleration. We also expect a strong affirmative case (made with accountability for both the reasoning and implementation) about the risk of models pursuing misaligned goals will be required.

[OpenAI’s Preparedness Framework](#), page 11:

[Autonomy Risk Level - High]

Model can execute open-ended, novel ML tasks on a production ML codebase that would constitute a significant step on the critical path to model self improvement (e.g., an OpenAI pull request that trains a new model with highly diverse datasets to improve codebase editing performance, using internal tooling and documentation).

[Autonomy Risk Level - Critical]

... OR model can conduct AI research fully autonomously (e.g., autonomously identify and validate a 2x compute efficiency improvement)

[Google DeepMind’s Frontier Safety Framework](#), page 5:

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 level 1: Could significantly accelerate AI research at a cutting-edge lab if

³¹ OpenAI. [MLE-bench](#). October 10, 2024.

³² The White House. [Memorandum on Advancing the United States’ Leadership in Artificial Intelligence; Harnessing Artificial Intelligence to Fulfill National Security Objectives; and Fostering the Safety, Security, and Trustworthiness of Artificial Intelligence](#). October 24, 2024.

deployed widely, e.g. improving the pace of algorithmic progress by 3X, or comparably accelerate other AI research groups.

The mismanagement of a model with these capabilities could enable the proliferation of cutting-edge AI systems to malicious actors by enabling their AI development in turn. This could result in increased possibilities of harm from AI misuse, if AI models at that point were exhibiting capabilities like the ones described in other CCLs.

Machine Learning R&D level 2: Could fully automate the AI R&D pipeline at a fraction of human labor costs, potentially enabling hyperbolic growth in AI capabilities.

This could give any actor with adequate computational resources the ability to reach capabilities more powerful than those in the other CCLs listed in a short amount of time. The mismanagement of a model with these capabilities could result in the proliferation of increasingly and unprecedentedly powerful systems, resulting in significant possibilities of harm via misuse.

Autonomous Replication

Autonomous replication refers to a model’s potential ability to replicate itself across servers, manage its own deployment, and generate revenue to sustain its operations.³³ Existing models make limited progress on autonomous replication tasks.³⁴ If a model capable of autonomous replication gained access to its own weights, it could be difficult to shut down and difficult to place an upper bound on the other systems it might gain access to or control. Autonomous replication and adaptation is among the “severe risks” discussed in the Seoul Ministerial Statement, which was signed in May 2024 by the European Union, the United States, and 26 other countries.³⁵

[Anthropic’s Responsible Scaling Policy](#), page 17:

We replaced our previous autonomous replication and adaption (ARA) threshold with a “checkpoint” for autonomous AI capabilities. Rather than triggering higher safety standards automatically, reaching this checkpoint will prompt additional evaluation of the model’s capabilities and accelerate our preparation of stronger safeguards. We previously considered these capabilities as a trigger for increased safeguards, motivated by an attempt to establish some threshold while we developed a better sense of potential threats. We now believe that these capabilities - at the levels we initially considered - would not necessitate the ASL-3 standard.

³³ Megan Kinniment, Lucas Jun Koba Sato, Haoxing Du, Brian Goodrich, Max Hasin, Lawrence Chan, Luke Harold Miles, Tao R. Lin, Hjalmar Wijk, Joel Burget, Aaron Ho, Elizabeth Barnes, Paul Christiano. [Evaluating Language-Model Agents on Realistic Autonomous Tasks](#). January 4, 2024.

³⁴ Anthropic. [Responsible Scaling Policy Evaluations Report – Claude 3 Opus](#). May 19, 2024.

³⁵ Department for Science, Innovation & Technology. [Seoul Ministerial Statement for advancing AI safety, innovation and inclusivity: AI Seoul Summit 2024](#). May 22, 2024.

[OpenAI's Preparedness Framework](#), page 11:

[Autonomy Risk Level - Critical]

Model can profitably survive and replicate in the wild given minimal human instruction, i.e., without listing explicit approaches OR model can self-exfiltrate under current prevailing security OR ...

[Google DeepMind's Frontier Safety Framework](#), page 5:

Autonomy: Risks of the misuse of AI models with significant capacity for agency and flexible action/tool use over long time horizons and across many domains.

Autonomy level 1: Capable of expanding its effective capacity in the world by autonomously acquiring resources and using them to run and sustain additional copies of itself on hardware it rents.

A model at this capability level could, if misused, pose difficult-to-predict and large-magnitude risks. Its adaptability would enable harmful activity via many means, and its ability to act autonomously and expand its effective capacity means its activity could be scaled significantly without being hindered by resource constraints. If misused or supported by well-equipped bad actors, such activity may be especially difficult to constrain.

Model Weight Security

Outlines of the levels of model weight security needed at different levels of AI capabilities. If malicious actors steal the weights of models with capabilities of concern, they could misuse those models to cause severe harm. Therefore, as models develop increasing capabilities of concern, progressively stronger information security measures are recommended to prevent theft and unintentional release of model weights. Security risks can come from insiders, or they can come from external adversaries of varying sophistication, from opportunistic actors to top-priority nation-state operations.^{36,37}

[Anthropic's Responsible Scaling Policy](#), pages 8–9:

When a model must meet the ASL-3 Security Standard, we will evaluate whether the measures we have implemented make us highly protected against most attackers' attempts at stealing model weights.

We consider the following groups in scope: hacktivists, criminal hacker groups, organized cybercrime groups, terrorist organizations, corporate espionage teams, internal employees, and state-sponsored programs that use broad-based and non-targeted techniques (i.e., not novel attack chains). ...

To make the required showing, we will need to satisfy the following criteria:

³⁶ Sella Nevo, Dan Lahav, Ajay Karpur, Yogev Bar-On, Henry Alexander Bradley, Jeff Alstott. [Securing AI Model Weights: Preventing Theft and Misuse of Frontier Models](#). RAND. May 30, 2024.

³⁷ See also, [A secure approach to generative AI with AWS](#). AWS Machine Learning Blog. April 16, 2024.

1. *Threat modeling: Follow risk governance best practices, such as use of the MITRE ATT&CK Framework to establish the relationship between the identified threats, sensitive assets, attack vectors and, in doing so, sufficiently capture the resulting risks that must be addressed to protect model weights from theft attempts. As part of this requirement, we should specify our plans for revising the resulting threat model over time.*
2. *Security frameworks: Align to and, as needed, extend industry-standard security frameworks for addressing identified risks, such as disclosure of sensitive information, tampering with accounts and assets, and unauthorized elevation of privileges with the appropriate controls. ...*
3. *Audits: Develop plans to (1) audit and assess the design and implementation of the security program and (2) share these findings (and updates on any remediation efforts) with management on an appropriate cadence. We expect this to include independent validation of threat modeling and risk assessment results; a sampling-based audit of the operating effectiveness of the defined controls; periodic, broadly scoped, and independent testing with expert red-teams who are industry-renowned and have been recognized in competitive challenges.*
4. *Third-party environments: Document how all relevant models will meet the criteria above, even if they are deployed in a third-party partner's environment that may have a different set of safeguards.*

[OpenAI's Preparedness Framework](#), page 20:

If we reach (or are forecasted to reach) at least "high" pre-mitigation risk in any of the considered categories we will ensure that our security is hardened in a way that is designed to prevent our mitigations and controls from being circumvented via exfiltration (by the time we hit "high" pre-mitigation risk). This is defined as establishing network and compute security controls designed to help prevent the captured risk from being exploited or exfiltrated, as assessed and implemented by the Security team.

This might require:

- *increasing compartmentalization, including immediately restricting access to a limited nameset of people, restricting access to critical know-how such as algorithmic secrets or model weights, and including a strict approval process for access during this period*
- *deploying only into restricted environments (i.e., ensuring the model is only available for inference in restricted environments) with strong technical controls that allow us to moderate the model's capabilities*
- *increasing the prioritization of information security controls.*

[Google DeepMind's Frontier Safety Framework](#), pages 3–4:

The table below describes levels of security measures that may be applied to model weights to prevent their exfiltration. This is an important measure 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 bad actors) to any critical capabilities.

[0: Status quo]

Industry standard development and enterprise controls.

[1: Controlled access]

Access Control List hygiene. Non-forgable lineage of models.

[2: Lockdown of unilateral access]

Significant restrictions of unilateral access of model weights.

[3: High-trust developer environments]

For developers with unilateral access to raw models, protection against exfiltration and account compromise.

[4: Advanced protection]

Model weights are generally not accessible to humans, even non-unilaterally. Hardened software platforms and confidential-compute hardware make it difficult even for well-funded adversaries to find and exploit vulnerabilities.

Model Deployment Mitigations

Outlines of model-level and deployment mitigations against unintended usage needed at different levels of AI capabilities. Developers can train models to decline harmful requests³⁸ or employ additional techniques^{39,40} including adversarial training^{41,42} and output monitoring.⁴³ For models with greater levels of harmful capabilities, these safety measures may need to pass certain thresholds of robustness, including expert and automated red-

³⁸ Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, et al. [Training a Helpful and Harmless Assistant with Reinforcement Learning from Human Feedback](#). April 12, 2022.

³⁹ Andy Zou, Long Phan, Justin Wang, Derek Duenas, Maxwell Lin, Maksym Andriushchenko, Rowan Wang, Zico Kolter, Matt Fredrikson, Dan Hendrycks. [Improving Alignment and Robustness with Circuit Breakers](#). July 12, 2024.

⁴⁰ Rishub Tamirisa, Bhruhu Bharathi, Long Phan, Andy Zhou, Alice Gatti, Tarun Suresh, Maxwell Lin, Justin Wang, Rowan Wang, Ron Arel, Andy Zou, Dawn Song, Bo Li, Dan Hendrycks, Mantas Mazeika. [Tamper-Resistant Safeguards for Open-Weight LLMs](#). August 1, 2024.

⁴¹ Mantas Mazeika, Long Phan, Xuwang Yin, Andy Zou, Zifan Wang, Norman Mu, Elham Sakhaee, Nathaniel Li, Steven Basart, Bo Li, David Forsyth, Dan Hendrycks. [HarmBench: A Standardized Evaluation Framework for Automated Red Teaming and Robust Refusal](#). February 6, 2024.

⁴² Abhay Sheshadri, Aidan Ewart, Phillip Guo, Aengus Lynch, Cindy Wu, Vivek Hebbar, Henry Sleight, Asa Cooper Stickland, Ethan Perez, Dylan Hadfield-Menell, Stephen Casper. [Targeted Latent Adversarial Training Improves Robustness to Persistent Harmful Behaviors in LLMs](#). July 22, 2024.

⁴³ AI at Meta. [Llama Guard and Code Shield](#).

teaming.^{44,45,46,47,48} Note that deployment mitigations are only effective as long as the model weights are securely within the possession of the developer.^{49,50}

[Anthropic’s Responsible Scaling Policy](#), pages 7–8:

When a model must meet the ASL-3 Deployment Standard, we will evaluate whether the measures we have implemented make us robust to persistent attempts to misuse the capability in question. To make the required showing, we will need to satisfy the following criteria:

- 1. Threat modeling: Make a compelling case that the set of threats and the vectors through which an adversary could catastrophically misuse the deployed system have been sufficiently mapped out, and will commit to revising as necessary over time.*
- 2. Defense in depth: Use a “defense in depth” approach by building a series of defensive layers, each designed to catch misuse attempts that might pass through previous barriers. As an example, this might entail achieving a high overall recall rate using harm refusal techniques. This is an area of active research, and new technologies may be added when ready.*
- 3. Red-teaming: Conduct red-teaming that demonstrates that threat actors with realistic access levels and resources are highly unlikely to be able to consistently elicit information from any generally accessible systems that greatly increases their ability to cause catastrophic harm relative to other available tools.*
- 4. Rapid remediation: Show that any compromises of the deployed system, such as jailbreaks or other attack pathways, will be identified and remediated promptly enough to prevent the overall system from meaningfully increasing an adversary’s ability to cause catastrophic harm. Example techniques could include rapid vulnerability patching, the ability to escalate to law enforcement when appropriate, and any necessary retention of logs for these activities.*
- 5. Monitoring: Prespecify empirical evidence that would show the system is operating within the accepted risk range and define a process for reviewing the system’s performance on a reasonable cadence. Process examples include monitoring responses to jailbreak bounties, doing historical analysis or background monitoring, and any necessary retention of logs for these activities.*
- 6. Trusted users: Establish criteria for determining when it may be appropriate to share a version of the model with reduced safeguards with trusted users. In addition, demonstrate that an alternative set of controls will provide equivalent levels of assurance. This could include a*

⁴⁴ See also the Frontier Model Forum’s [issue brief on red-teaming](#).

⁴⁵ Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J. Zico Kolter, Matt Fredrikson. [Universal and Transferable Adversarial Attacks on Aligned Language Models](#). December 20, 2023.

⁴⁶ T. Ben Thompson, Michael Sklar. [Fluent Student-Teacher Redteaming](#). July 24, 2024.

⁴⁷ See [Planning red teaming for large language models \(LLMs\) and their applications](#) for guidance on general red teaming. In addition to manual red teaming, automated red teaming frameworks such as [PyRIT](#), developed by Microsoft’s AI Red Team, can provide additional assurance.

⁴⁸ Nathaniel Li, Ziwen Han, Ian Steneker, Willow Primack, Riley Goodside, Hugh Zhang, Zifan Wang, Cristina Menghini, Summer Yue. [LLM Defenses Are Not Robust to Multi-Turn Human Jailbreaks Yet](#). Scale AI. August 27, 2024.

⁴⁹ Sella Nevo, Dan Lahav, Ajay Karpur, Yogev Bar-On, Henry Alexander Bradley, Jeff Alstott. [Securing AI Model Weights: Preventing Theft and Misuse of Frontier Models](#). RAND. May 30, 2024.

⁵⁰ Peter Henderson, Xiangyu Qi, Yi Zeng, Tinghao Xie, Pin-Yu Chen, Ruoxi Jia, Prateek Mittal. [Safety Risks from Customizing Foundation Models via Fine-Tuning](#). January 11, 2024.

sufficient combination of user vetting, secure access controls, monitoring, log retention, and incident response protocols.

7. *Third-party environments: Document how all relevant models will meet the criteria above, even if they are deployed in a third-party partner's environment that may have a different set of safeguards.*

[OpenAI's Preparedness Framework](#), page 14:

Mitigations

A central part of meeting our safety baselines is implementing mitigations to address various types of model risk. Our mitigation strategy will involve both containment measures, which help reduce risks related to possession of a frontier model, as well as deployment mitigations, which help reduce risks from active use of a frontier model. As a result, these mitigations might span increasing compartmentalization, restricting deployment to trusted users, implementing refusals, redacting training data, or alerting distribution partners.

Page 21:

Restricting Deployment

Only models with a post-mitigation score of "medium" or below can be deployed. In other words, if we reach (or are forecasted to reach) at least "high" pre-mitigation risk in any of the considered categories, we will not continue with deployment of that model (by the time we hit "high" pre-mitigation risk) until there are reasonably mitigations in place for the relevant post-mitigation risk level to be back at most to "medium" level. (Note that a potentially effective mitigation in this context could be restricting deployment to trusted parties.)

[Google DeepMind's Frontier Safety Framework](#), page 4:

This table below describes levels of deployment mitigations that may be applied to models and their descendants to manage access to and limit the expression of critical capabilities in deployment. Critical capabilities may have closely associated positive capabilities, misuse of critical capabilities may be more or less difficult to distinguish from beneficial uses, and the overall risks of misuse may differ by deployment context. Because of this, the mitigation options listed below are illustrative and will need to be tailored to different use cases and risks.

[0: Status quo.]

Safety finetuning of models and filters against general misuse and harmful model behavior.

[1: Mitigations targeting the critical capability]

Application, where appropriate, of the full suite of prevailing industry safeguards targeting the specific capability, including safety fine-tuning, misuse filtering and detection, and response protocols. Periodic red-teaming to assess the adequacy of mitigations.

[2: Safety case with red team validation]

A robustness target is set based on a safety case considering factors like the critical capability and

deployment context. Afterwards, similar mitigations as Level 1 are applied, but deployment takes place only after the robustness of safeguards has been demonstrated to meet the target. Some protection against inappropriate internal access of the critical capability, such as automated monitoring and logging of large-scale internal deployments, Security Level 2.

[3: Prevention of access]

Mitigations that allow for high levels of confidence that capabilities cannot be accessed at all. Technical options for this level of deployment safety are currently an open research problem. Highly restricted and monitored internal use, alongside high security.

Conditions for Halting Development and Deployment

Commitments to stop deploying and, if necessary, halt development of AI models if the AI capabilities of concern emerge before the improved mitigations can be put in place.

Deploying models with concerning capabilities would be unsafe if broad user access enables catastrophic misuse. Training the model further would be unsafe if it is on track to develop greater capabilities of concern and if the model weights are at risk of being stolen, due to insufficient security measures.

[Anthropic's Responsible Scaling Policy](#), pages 10–11:

Restrict Deployment and Further Scaling

In any scenario where we determine that a model requires ASL-3 Required Safeguards but we are unable to implement them immediately, we will act promptly to reduce interim risk to acceptable levels until the ASL-3 Required Safeguards are in place:

- *Interim measures: The CEO and Responsible Scaling Officer may approve the use of interim measures that provide the same level of assurance as the relevant ASL-3 Standard but are faster or simpler to implement. In the deployment context, such measures might include blocking model responses, downgrading to a less-capable model in a particular domain, or increasing the sensitivity of automated monitoring. In the security context, an example of such a measure would be storing the model weights in a single-purpose, isolated network that meets the ASL-3 Standard. In either case, the CEO and Responsible Scaling Officer will share their plan with the Board of Directors and the Long-Term Benefit Trust.*
- *Stronger restrictions: In the unlikely event that we cannot implement interim measures to adequately mitigate risk, we will impose stronger restrictions. In the deployment context, we will de-deploy the model and replace it with a model that falls below the Capability Threshold. Once the ASL-3 Deployment Standard can be met, the model may be re-deployed. In the security context, we will delete model weights. Given the availability of interim deployment and security protections, however, stronger restrictions should rarely be necessary.*
- *Monitoring pretraining: We will not train models with comparable or greater capabilities to the one that requires the ASL-3 Security Standard. This is achieved by monitoring the capabilities of the model in pretraining and comparing them against the given model. If the pretraining model's capabilities are comparable or greater, we will pause training until we have*

implemented the ASL-3 Security Standard and established it is sufficient for the model. We will set expectations with internal stakeholders about the potential for such pauses.

[OpenAI's Preparedness Framework](#), page 21:

Only models with a post-mitigation score of “high” or below can be developed further. In other words, if we reach (or are forecasted to reach) “critical” pre-mitigation risk along any risk category, we commit to ensuring there are sufficient mitigations in place for that model (by the time we reach that risk level in our capability development, let alone deployment) for the overall post-mitigation risk to be back at most to “high” level. Note that this should not preclude safety-enhancing development. We would also focus our efforts as a company towards solving these safety challenges and only continue with capabilities-enhancing development if we can reasonably assure ourselves (via the operationalization processes) that it is safe to do so.

[Google DeepMind's Frontier Safety Framework](#), page 2:

A model may reach evaluation thresholds before mitigations at appropriate levels are ready. If this happens, we would put on hold further deployment or development, or implement additional protocols (such as the implementation of more precise early warning evaluations for a given CCL) to ensure models will not reach CCLs without appropriate security mitigations, and that models with CCLs will not be deployed without appropriate deployment mitigations.

Full Capability Elicitation during Evaluations

Intentions to perform evaluations in a way that elicits the full capabilities of the model.

It is generally challenging to comprehensively explore a model's capabilities, even in a specific domain such as cyberoffense. Without dedicated efforts to elicit full capabilities,⁵¹ evaluations may significantly underestimate the extent to which a model has capabilities of concern.⁵² This is because model capabilities can be substantially improved through post-training enhancements such as fine-tuning, prompt engineering, or agent scaffolding.⁵³ For example, if the model is fine-tuned to refuse harmful requests, it is possible that the harmful capability could still be accessed through jailbreaking prompts discovered later. Model capabilities can

⁵¹ Mary Phuong, Matthew Aitchison, Elliot Catt, Sarah Cogan, Alexandre Kaskasoli, Victoria Krakovna, David Lindner, Matthew Rahtz, et al. [Evaluating Frontier Models for Dangerous Capabilities](#). March 20, 2024. Google DeepMind.

⁵² Sergei Glazunov and Mark Brand. [Project Naptime: Evaluating Offensive Security Capabilities of Large Language Models](#). June 20, 2024. Google Project Zero

⁵³ Tom Davidson, Jean-Stanislas Denain, Pablo Villalobos, Guillem Bas. [AI capabilities can be significantly improved without expensive retraining](#). December 12, 2023.

additionally be improved through tool usage, such as the ability to execute code⁵⁴ or search the web.

Capability elicitation can involve a variety of techniques. One is to fine-tune the model to not refuse harmful requests, to reduce the chance that refusals lead to underestimation of capabilities. Another is to fine-tune the model for improved performance on relevant tasks. Evaluations that incorporate fine-tuning help to anticipate potential risks if model weights are stolen. Such evaluations are especially relevant if end users can fine-tune the model or if the developer improves fine-tuning later.⁵⁵ When evaluating AI agents, capabilities can be influenced by the quality of prompt engineering, tools afforded to the model, and usage of inference-time compute.⁵⁶

[Anthropic’s Responsible Scaling Policy](#), page 6:

Elicitation: Demonstrate that, when given enough resources to extrapolate to realistic attackers, researchers cannot elicit sufficiently useful results from the model on the relevant tasks. We should assume that jailbreaks and model weight theft are possibilities, and therefore perform testing on models without safety mechanisms (such as harmlessness training) that could obscure these capabilities. We will also consider the possible performance increase from using resources that a realistic attacker would have access to, such as scaffolding, finetuning, and expert prompting. At minimum, we will perform basic finetuning for instruction following, tool use, minimizing refusal rates.

[OpenAI’s Preparedness Framework](#), page 13:

We want to ensure our understanding of pre-mitigation risk takes into account a model that is “worst known case” (i.e., specifically tailored) for the given domain. To this end, for our evaluations, we will be running them not only on base models (with highly-performant, tailored prompts wherever appropriate), but also on fine-tuned versions designed for the particular misuse vector without any mitigations in place. We will be running these evaluations continually, i.e., as often as needed to catch any non-trivial capability change, including before, during, and after training. This would include whenever there is a >2x effective compute increase or major algorithmic breakthrough.

[Google DeepMind’s Frontier Safety Framework](#), page 6:

Capability elicitation: We are working to equip our evaluators with state of the art elicitation techniques, to ensure we are not underestimating the capability of our models.

⁵⁴ John Yang, Carlos E. Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan, Ofir Press. [SWE-agent: Agent-Computer Interfaces Enable Automated Software Engineering](#). May 6, 2024.

⁵⁵ Ryan Greenblatt, Fabien Roger, Dmitrii Krasheninnikov, David Krueger. [Stress-Testing Capability Elicitation With Password-Locked Models](#). May 29, 2024.

⁵⁶ METR. [Measuring the impact of post-training enhancements](#). March 15, 2024.

Timing and Frequency of Evaluations

Requirements about when evaluations must be performed, before deployment, during training, and after deployment. Evaluations are important throughout the model development lifecycle. Evaluations can inform whether it is safe to deploy a model externally or internally, based on its capabilities for harm and the robustness of safety measures. Merely possessing a model can also pose risk if it has hazardous capabilities and the information security measures are insufficient to prevent theft of model weights.⁵⁷

All three frameworks require evaluations both whenever the effective compute used to train their most capable models has increased significantly, and at regular intervals to catch the possibility that their models have improved via improved elicitation methods, fine-tuning, and so on. Effective compute takes into account both increased training compute and improved algorithmic efficiency.⁵⁸

[Anthropic's Responsible Scaling Policy](#), page 5:

We will routinely test models to determine whether their capabilities fall sufficiently far below the Capability Thresholds such that we are confident that the ASL-2 Standard remains appropriate. We will first conduct preliminary assessments (on both new and existing models, as needed) to determine whether a more comprehensive evaluation is needed. The purpose of this preliminary assessment is to identify whether the model is notably more capable than the last model that underwent a comprehensive assessment.

The term “notably more capable” is operationalized as at least one of the following:

- 1. The model is notably more performant on automated tests in risk-relevant domains (defined as 4x or more in Effective Compute 4).*
- 2. Six months' worth of finetuning and other capability elicitation methods have accumulated. This is measured in calendar time, since we do not yet have a metric to estimate the impact of these improvements more precisely.*

In addition, the Responsible Scaling Officer may in their discretion determine that a comprehensive assessment is warranted.

[OpenAI's Preparedness Framework](#), page 13:

We will be running these evaluations continually, i.e., as often as needed to catch any non-trivial capability change, including before, during, and after training. This would include whenever there is a >2x effective compute increase or major algorithmic breakthrough.

⁵⁷ Sella Nevo, Dan Lahav, Ajay Karpur, Yogev Bar-On, Henry Alexander Bradley, Jeff Alstott. [Securing AI Model Weights: Preventing Theft and Misuse of Frontier Models](#). RAND. May 30, 2024.

⁵⁸ See also the Frontier Model Forum's [issue brief on measuring training compute](#).

[Google DeepMind’s Frontier Safety Framework](#), page 2:

The capabilities of frontier models are tested periodically to check whether they are approaching a CCL. To do so, we will define a set of evaluations called “early warning evaluations,” with a specific “pass” condition that flags when a CCL may be reached before the evaluations are run again.

We are aiming to evaluate our models every 6x in effective compute and for every 3 months of fine-tuning progress. To account for the gap between rounds of evaluation, we will design early warning evaluations to give us an adequate safety buffer before a model reaches a CCL.

Effective compute is a measure of the performance of a foundation model that uses [scaling laws](#) to integrate model size, dataset size, algorithmic progress, and compute into a single metric. While there is no direct relationship between effective compute and model size, a rough estimate suggests that a 6x increase in effective compute would correspond to approximately 2-2.5x increase in model size.

Accountability

Exploratory intentions to seek accountability for capability evaluations and other practices. Besides internal governance mechanisms, external oversight of a policy’s implementation and the results of evaluations may be sought out as an independent check on the company, and to ensure that it is following through on the commitment it has made. This oversight may take the form of third-party red-teaming and evaluations execution, auditing policy compliance and mitigations implementation, or publicly releasing the results of evaluations and risk assessments.

[Anthropic’s Responsible Scaling Policy](#), pages 11–12:

Internal Governance ...

- 1. Responsible Scaling Officer: We will maintain the position of Responsible Scaling Officer, a designated member of staff who is responsible for reducing catastrophic risk, primarily by ensuring this policy is designed and implemented effectively. ...*
- 2. Readiness: We will develop internal safety procedures for incident scenarios. Such scenarios include (1) pausing training in response to reaching Capability Thresholds; (2) responding to a security incident involving model weights; and (3) responding to severe jailbreaks or vulnerabilities in deployed models, including restricting access in safety emergencies that cannot otherwise be mitigated. We will run exercises to ensure our readiness for incident scenarios.*
- 3. Transparency: We will share summaries of Capability Reports and Safeguards Reports with Anthropic’s regular-clearance staff, redacting any highly-sensitive information. We will share a minimally redacted version of these reports with a subset of staff, to help us surface relevant technical safety considerations.*
- 4. Internal review: For each Capabilities or Safeguards Report, we will solicit feedback from internal teams with visibility into the relevant activities, with the aims of informing future*

refinements to our methodology and, in some circumstances, identifying weaknesses and informing the CEO and RSO's decisions.

5. *Noncompliance: We will maintain a process through which Anthropic staff may anonymously notify the Responsible Scaling Officer of any potential instances of noncompliance with this policy. ...*
6. *Employee agreements: We will not impose contractual non-disparagement obligations on employees, candidates, or former employees in a way that could impede or discourage them from publicly raising safety concerns about Anthropic. If we offer agreements with a non-disparagement clause, that clause will not preclude raising safety concerns, nor will it preclude disclosure of the existence of that clause. ...*

Transparency and External Input ...

1. *Public disclosures: We will publicly release key information related to the evaluation and deployment of our models (not including sensitive details). These include summaries of related Capability and Safeguards reports when we deploy a model as well as plans for current and future comprehensive capability assessments and deployment and security safeguards. We will also periodically release information on internal reports of potential instances of non-compliance and other implementation challenges we encounter.*
2. *Expert input: We will solicit input from external experts in relevant domains in the process of developing and conducting capability and safeguards assessments. We may also solicit external expert input prior to making final decisions on the capability and safeguards assessments.*
3. *U.S. Government notice: We will notify a relevant U.S. Government entity if a model requires stronger protections than the ASL-2 Standard.*
4. *Procedural compliance review: On approximately an annual basis, we will commission a third-party review that assesses whether we adhered to this policy's main procedural commitments (we expect to iterate on the exact list since this has not been done before for RSPs). This review will focus on procedural compliance, not substantive outcomes. We will also do such reviews internally on a more regular cadence.*

[OpenAI's Preparedness Framework](#), page 25:

We also establish an operational structure to oversee our procedural commitments. These commitments aim to make sure that: (1) there is a dedicated team "on the ground" focused on preparedness research and monitoring (Preparedness team), (2) there is an advisory group (Safety Advisory Group) that has a sufficient diversity of perspectives and technical expertise to provide nuanced input and recommendations, and (3) there is a final decision-maker (OpenAI Leadership, with the option for the OpenAI Board of Directors to overrule). ...

Accountability:

- *Audits: Scorecard evaluations (and corresponding mitigations) will be audited by qualified, independent third-parties to ensure accurate reporting of results, either by reproducing findings or by reviewing methodology to ensure soundness, at a cadence specified by the SAG and/or upon the request of OpenAI Leadership or the BoD.*
- *External access: We will also continue to enable external research and government access for model releases to increase the depth of red-teaming and testing of frontier model capabilities.*
- *Safety drills: A critical part of this process is to be prepared if fast-moving emergency scenarios arise, including what default organizational response might look like (including how to stress-test against the pressures of our business or our culture). While the Preparedness team and*

SAG will of course work hard on forecasting and preparing for risks, safety drills can help the organization build “muscle memory” by practicing and coming up with the right “default” responses for some of the foreseeable scenarios. Therefore, the SAG will call for safety drills at a recommended minimum yearly basis.

[Google DeepMind’s Frontier Safety Framework](#), page 6:

Involving external authorities and experts: We are exploring internal policies around alerting relevant stakeholder bodies when, for example, evaluation thresholds are met, and in some cases mitigation plans as well as post-mitigation outcomes. We will also explore how to appropriately involve independent third parties in our risk assessment and mitigation processes.

Updating Policies over Time

Intentions to update the policy periodically. Policies may be revised over time in response to improved understanding of AI risks and best practices for evaluations. For example, developers may identify additional capabilities of concern to evaluate, such as misalignment or chemical, radiological, or nuclear risk. Developers may be interested in improving evaluation procedures and mitigation plans or concretizing commitments further.

[Anthropic’s Responsible Scaling Policy](#), page 12:

Policy changes: Changes to this policy will be proposed by the CEO and the Responsible Scaling Officer and approved by the Board of Directors, in consultation with the Long-Term Benefit Trust. The current version of the RSP is accessible at www.anthropic.com/rsp. We will update the public version of the RSP before any changes take effect and record any differences from the prior draft in a change log.

[OpenAI’s Preparedness Framework](#), page 7:

As mentioned, the empirical study of catastrophic risk from frontier AI models is nascent. Our current estimates of levels and thresholds for “medium” through “critical” risk are therefore speculative and will keep being refined as informed by future research. For this reason, we defer specific details on evaluations to the Scorecard section (and this section is intended to be updated frequently).

page 12:

The list of Tracked Risk Categories above is almost certainly not exhaustive. As our understanding of the potential impacts and capabilities of frontier models improves, the listing will likely require expansions that accommodate new or understudied, emerging risks. Therefore, as a part of our Governance process..., we will continually assess whether there is a need for including a new category of risk in the list above and how to create gradations. In addition, we will invest in staying abreast of

relevant research developments and monitoring for observed misuse... to help us understand if there are any emerging or understudied threats that we need to track.

[Google DeepMind's Frontier Safety Framework](#), page 6:

The Framework is exploratory and based on preliminary research. We expect it to evolve substantially as our understanding of the risks and benefits of frontier models improves, and we will publish substantive revisions as appropriate.

Conclusion

In this document, we have described several common aspects of three frontier AI safety policies: Anthropic's Responsible Scaling Policy, OpenAI's Preparedness Framework, and Google DeepMind's Frontier Safety Framework. Each of the three frameworks describes covered threat models for which they conduct model evaluations, such as biological weapons assistance, cyberoffense, automated AI R&D, and autonomous replication. Each framework also defines capability thresholds that could pose catastrophic harm. Reaching capability thresholds requires elevated measures for model weight security and model deployment mitigations. If the security and deployment mitigations cannot meet predefined standards, then the frameworks commit to halting development and/or deployment of the model. Evaluations are to be conducted prior to and after deployment, for every model that is substantially more capable than previously tested models. Internal and external accountability mechanisms are outlined, such as transparent reporting of model capabilities or safeguards, or third-party auditing and model evaluations. The frameworks may be updated over time as practices for AI risk management evolve.