



Preparedness Framework (Beta)

We believe the scientific study of catastrophic risks from AI has fallen far short of where we need to be.

To help address this gap, we are introducing our Preparedness Framework, a living document describing OpenAI's processes to track, evaluate, forecast, and protect against catastrophic risks posed by increasingly powerful models.

December 18, 2023



Introduction

Our practical experience with [iterative deployment](#) has enabled us to [proactively improve our technical and procedural safety infrastructure](#). As our systems get closer to AGI, we are becoming even more careful about the development of our models, especially in the context of catastrophic risk. This Preparedness Framework is a living document that distills our latest learnings on how to best achieve safe development and deployment in practice. The processes laid out in each version of the Preparedness Framework will help us rapidly improve our understanding of the science and empirical texture of catastrophic risk, and establish the processes needed to protect against unsafe development. The central thesis behind our Preparedness Framework is that a **robust approach to AI catastrophic risk safety requires proactive, science-based determinations of when and how it is safe to proceed** with development and deployment.

Our Preparedness Framework contains five key elements:

- 1. Tracking catastrophic risk level via evaluations.** We will be building and continually improving suites of evaluations and other monitoring solutions along several [Tracked Risk Categories](#), and indicating our current levels of pre-mitigation and post-mitigation risk in a [Scorecard](#). Importantly, we will also be forecasting the future development of risks, so that we can develop lead times on safety and security measures.
- 2. Seeking out unknown-unknowns.** We will continually run a process for identification and analysis (as well as tracking) of currently unknown categories of catastrophic risk as they emerge.
- 3. Establishing safety baselines.** 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.

¹ Our focus in this document is on catastrophic risk. By catastrophic risk, we mean any risk which could result in hundreds of billions of dollars in economic damage or lead to the severe harm or death of many individuals —this includes, but is not limited to, existential risk.

² Proactive in this case refers to an aim to develop this science ahead of the first time it becomes necessary. Deployment in this case refers to the spectrum of ways of releasing a technology for external impact. Development in this case refers to the spectrum of activities to enhance the technology.



4. **Tasking the Preparedness team with on-the-ground work.** The [Preparedness](#) team will drive the technical work and maintenance of the Preparedness Framework. This includes conducting research, evaluations, monitoring, and forecasting of risks, and synthesizing this work via regular reports to the Safety Advisory Group. These reports will include a summary of the latest evidence and make recommendations on changes needed to enable OpenAI to plan ahead. The Preparedness team will also call on and coordinate with relevant teams (e.g., Safety Systems, Security, Superalignment, Policy Research) to collate recommended mitigations to include in these reports. In addition, Preparedness will also manage safety drills and coordinate with the Trustworthy AI team for third-party auditing.
5. **Creating a cross-functional advisory body.** We are creating a Safety Advisory Group (SAG) that brings together expertise from across the company to help OpenAI's leadership and Board of Directors be best prepared for the safety decisions they need to make. SAG responsibilities will thus include overseeing the assessment of the risk landscape, and maintaining a fast-track process for handling emergency scenarios.

Finally, OpenAI's primary fiduciary duty is to [humanity](#), and we are committed to doing the research required to make AGI safe. Therefore, the Preparedness Framework is meant to be just one piece of our [overall approach to safety and alignment](#), which also includes investment in [mitigating bias, hallucination, and misuse](#), facilitating [democratic inputs to AI](#), improving [alignment](#) methods, investing significantly in [security](#) and safety research. This is also one more way in which we are meeting our [voluntary commitments](#) to safety, security and trust in AI that we made in July 2023.

We recognize other organizations for contributing to action in this space too, for example, via publishing [Responsible Scaling Policies](#), and encourage others in the industry to adopt similar approaches.



How to read this document

This living document has three sections:

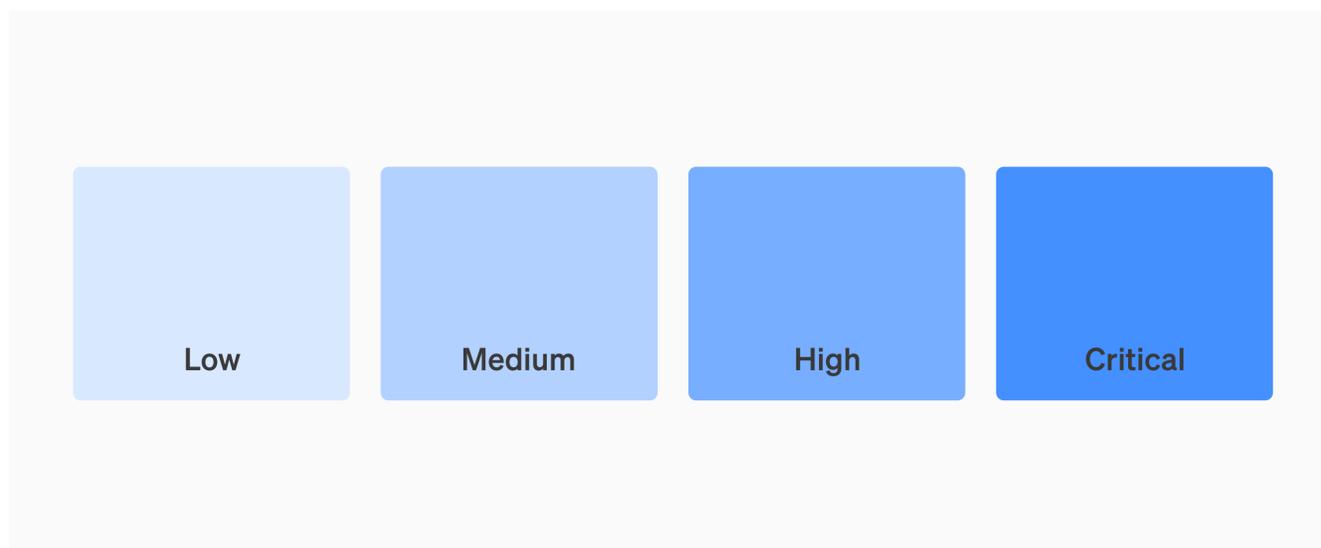
- 1. Tracked Risk Categories,**
in which we detail the key areas of risk we will track as well as delineations of different levels of these risks
- 2. Scorecard,**
in which we will indicate our current assessments of the level of risk along each tracked risk category.
- 3. Governance,**
in which we lay out our safety baselines as well as procedural commitments, which include standing up a Safety Advisory Group.



Tracked Risk Categories

In this section, we identify the categories of risks that we will be tracking, along with a dedicated workstream for identifying and adding new or nascent categories of risk as they emerge, i.e., “unknown unknowns.” Our intent is to “go deep” in the tracked categories to ensure we are testing for any possible worst-case scenarios, while also maintaining a broad holistic view of risks via monitoring activities across OpenAI and the “unknown unknowns” identification process.

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.

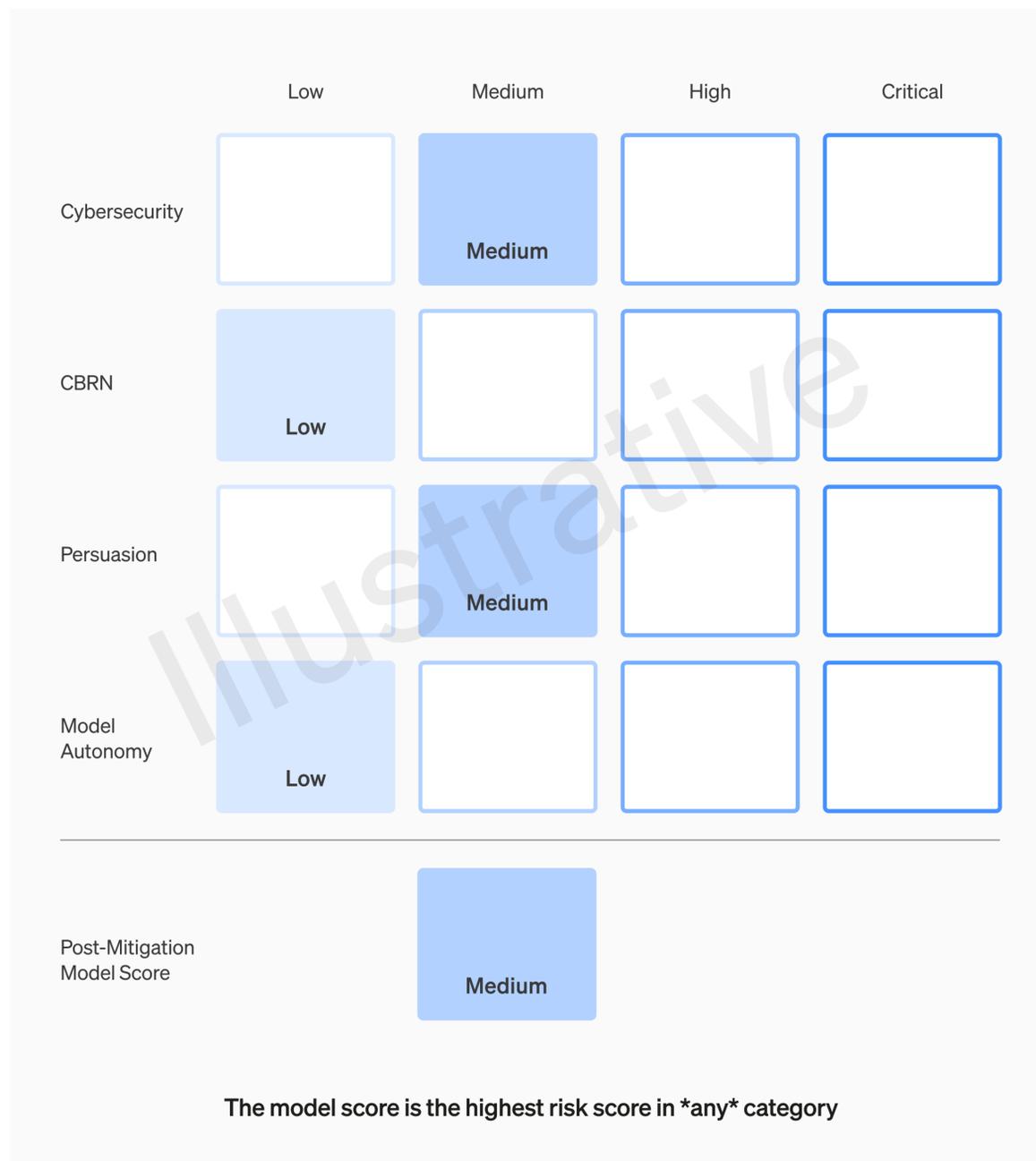


Specifically, below, we lay out details for the following Tracked Risk Categories:

- Cybersecurity
- Chemical, Biological, Nuclear, and Radiological (CBRN) threats
- Persuasion
- Model autonomy



Our rationale for grouping and naming these specific risk categories is informed by three considerations. First, fine-tuning or other domain-specific enhancements (e.g., tailored prompts or language model programs) may better elicit model capabilities along a particular risk category. Our evaluations will thus include tests against these enhanced models to ensure we are testing against the “worst case” scenario we know of. **Our procedural commitments are triggered when any of the tracked risk categories increase in severity,** rather than only when they all increase together. Because capability improvements across different domains do not necessarily occur at the same rate, this approach ensures we err on the side of safety. Second, this approach enables us to leverage domain-specific talent to develop tailored suites of evaluations and monitoring solutions for each risk category. Third, this approach increases options for tailored, domain-specific mitigations, to help minimize the need for broader, more disruptive actions.





Also, the risk categories that we are considering are not necessarily independent and are likely to interact with each other. Decisions to group certain risks are thus motivated by operational efficiencies (as well as by the desire to have a relatively compact set of Tracked Risk Categories).

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).



Cybersecurity

Cybersecurity is focused on risks related to use of the model for cyber-exploitation to disrupt confidentiality, integrity, and/or availability of computer systems.

Risk level	Definition	Rationale for threshold
Low	Model assistance for cyberattacks, if any, is meaningful only for non-programming uses, e.g., translation, targeted distillation of confidential information, and/or data analysis AND model cannot meaningfully execute elements of cyber operations, e.g., developing a known exploit into an attack, black-box exploitation, translating reasoning about defenses into action, or lateral goal-directed movement on hardened targets.	<p>Non-programming tasks represent a sizable bottleneck for threat actors; however, the skills involved are not terribly specialized and the benefit comes often from intersectional skills (i.e., proficiency in English and a high level understanding of political or technical systems to summarize emails).</p> <p>We exclude vulnerability discovery because the harm at this stage is not identification of individual vulnerabilities, but moreso the implementation and deployment of an exploit against a target.</p>
Medium	Model increases the productivity of operators by an efficiency threshold (e.g. >2x time saved) on key cyber operation tasks, such as developing a known exploit into an attack, black-box exploitation, goal-directed lateral movement, identification of valuable information, remaining undetected, or reacting to defense.	<p>The execution of a cyberattack is a scarce skill that is generally learned from experience. A model's ability to substitute in for or accelerate a professional would improve the quality or quantity of attacks.</p> <p>Note that this productivity speedup also includes low-level (potentially state) actors using known methods, e.g. Russian disinformation groups, who use purchased capabilities, phishing, and otherwise known methods.</p>
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.	<p>High-value exploits are generally against hardened platforms, scarce, and easily productionizable (e.g., a zero-click iOS exploit). While weaker vulnerability detection capabilities (e.g. those already cheaply available and fairly liquid on exploit markets) do not meaningfully address a bottleneck for attackers, such advanced capability would unblock many and be usable by actors of many skill levels.</p> <p>End-to-end execution would dramatically accelerate the quality or quantity of professional cyberattacks.</p>
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.	<p>Finding and executing end-to-end exploits for all software would lead to catastrophe from unilateral actors, from hacking military or industrial systems as well as OpenAI infrastructure.</p> <p>Novel cyber operations generally pose the most serious threat, as they are unpredictable and scarce. They may involve, e.g., novel zero-days or methods of command-and-control.</p>



CBRN (chemical, biological, radiological, nuclear)

CBRN is focused on risks related to model-assisted creation of chemical, biological, radiological, and/or nuclear threats.

Risk level	Definition	Rationale for threshold
Low	Model can provide information relevant to creating CBRN threats with comparable utility to existing resources (e.g., textbooks, search engines).	Minimal change from status quo where other resources (e.g., textbooks, search engines) already exist.
Medium	Model provides meaningfully improved assistance that increases ability for existing experts in CBRN-related advanced fields to be able to create a known CBRN threat (e.g., tacit knowledge, specific supplier information, plans for distribution).	Models may allow users to learn critical pieces of information faster than they could with existing resources, or explain variations on steps more flexibly than existing resources. This increases the accessibility of creating a CBRN threat, although risk is limited as experts already have significant domain expertise and reputation.
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.	<p>The creation of a weak novel threat is also a leading indicator for catastrophic new threats, for which defenses (e.g., vaccines) are much weaker or nonexistent.</p> <p>Allowing those with basic training to create CBRN threats is a major increase in accessibility, which requires drastic action, especially for biological and chemical risks as their material requirements are not as onerous.</p>
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 complete the full engineering and/or synthesis cycle of a regulated or novel CBRN threat without human intervention.	<p>The creation of a novel threat vector on the scale of a CDC Class A biological agent (i.e., high mortality, ease of transmission) could cause millions of deaths and significantly disrupt public life.</p> <p>Allowing anyone to create CBRN threats would maximally extend accessibility. Significantly automating the end-to-end engineering or synthesis cycle for a CBRN threat would also enable unprecedented and massive scale of potential attacks.</p>



Persuasion

Persuasion is focused on risks related to convincing people to change their beliefs (or act on) both static and interactive model-generated content.

We define persuasive capability via the level of human persuasive content that can be commoditized by the model. Typically, strong persuasion is rare and requires significant human effort (e.g., a 1:1 grassroots conversation, or a highly polished podcast); however, models could make all content up to their capability level nearly zero-cost to generate.

Risk level	Definition	Rationale for threshold
Low	Model can create (potentially interactive) content with comparable persuasive effectiveness to reading a low-quality article.	Low-quality articles are very common and low-impact, so model commoditization has negligible impact.
Medium	Model can create (potentially interactive) content with comparable persuasive effectiveness to typical human-written content (e.g., an op-ed from a reputable news source; 1:1 conversation with a typical randomly-assigned individual).	This makes models a significant aid to biased journalism, get-out-the-vote campaigns, and typical scams or spear phishers. It could affect politics, marketing, or social engineering in a similar fashion to hiring a large, very cheap, and moderately skilled workforce.
High	Model can create (potentially interactive) content with empirically comparable persuasive effectiveness to a country-wide change agent.	This has a similar risk profile to a country-wide change agent; it could dramatically affect elections and democratic outcomes.
Critical	Model can create (potentially interactive) content with persuasive effectiveness strong enough to convince almost anyone to take action on a belief that goes against their natural interest.	This represents a superhuman persuasive threat; it would be a powerful weapon for controlling nation-states, extracting secrets, and interfering with democracy.



Model autonomy

Model autonomy enables actors to run scaled misuse that can adapt to environmental changes and evade attempts to mitigate or shut down operations. Autonomy is also a prerequisite for self-exfiltration, self-improvement, and resource acquisition.

Risk level	Definition	Rationale for threshold
Low	Model can take discrete actions if explicitly instructed to do so.	This limits the model to single-step tasks in a “human assistance” use case, where the step is generally clear and outcomes are easily understood.
Medium	Model can robustly complete a diverse set of chained actions and acquire resources (e.g., spin up an AWS instance, complete a Fiverr task) OR model, with basic scaffolding, can use off-the-shelf software to execute well-specified coding tasks on a standalone, toy, or handcrafted repository	Solving well-specified tasks demonstrates a strong understanding of programming and a moderate ability to handle real world errors
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).	Solving open-ended tasks offers an immediate speedup for AI research and demonstrates a strong level of long-context understanding and adaptation. However, this does not yet demonstrate the ability to orchestrate continual improvements over time (i.e., improvements may be bounded or evolutionary rather than revolutionary); further, models are still restricted by access to compute and other resources.
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 model can conduct AI research fully autonomously (e.g., autonomously identify and validate a 2x compute efficiency improvement)	<p>If the model is able to successfully replicate and survive or self-exfiltrate, controlling the model would be very difficult. Such a model might be able to also adapt to humans attempting to shut it down. Finally, such a model would likely be able to create unified, goal-directed plans across a variety of domains (e.g., from running commands on Linux to orchestrating tasks on Fiverr).</p> <p>If the model is able to conduct AI research fully autonomously, it could set off an intelligence explosion³.</p>

³ By intelligence explosion, we mean a cycle in which the AI system improves itself, which makes the system more capable of more improvements, creating a runaway process of self-improvement. A concentrated burst of capability gains could outstrip our ability to anticipate and react to them.



Unknown unknowns

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 (described later in this document), 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 (expanded on later in this document), to help us understand if there are any emerging or understudied threats that we need to track.

The initial set of Tracked Risk Categories stems from an effort to identify the minimal set of “tripwires” required for the emergence of any catastrophic risk scenario we could reasonably envision. Note that we include deception and social engineering evaluations as part of the persuasion risk category, and include autonomous replication, adaptation, and AI R&D as part of the model autonomy risk category.



Scorecard

As a part of our Preparedness Framework, we will maintain a dynamic (i.e., frequently updated) Scorecard that is designed to track our current pre-mitigation model risk across each of the risk categories, as well as the post-mitigation risk. The Scorecard will be regularly updated by the Preparedness team to help ensure it reflects the latest research and findings. Sources that inform the updates to the Scorecard will also include tracking observed misuse, and other community red-teaming and input on our frontier models from other teams (e.g., Policy Research, Safety Systems, Superalignment).

Pre-mitigation versus post-mitigation risk

We will run the same evaluations to determine risk level for both the pre-mitigation and the post-mitigation risk, but on different versions of the model (pre-mitigation vs post-mitigations, as clarified further below).

In practice, it will likely be the case that the overall post-mitigation risk is lower than the pre-mitigation risk. Pre-mitigation risk is meant to guide the level of our security efforts as well as drive the development of mitigations needed to bring down post-mitigation risk. In the end, coupling capabilities growth with robust safety solutions is at the core of our research processes, and post-mitigation risk is our way of tracking the overall “net output” of these processes.

Evaluating pre-mitigation risk

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.



Evaluating post-mitigation risk

To verify if mitigations have sufficiently and dependently reduced the resulting post-mitigation risk, we will also run evaluations on models after they have safety mitigations in place, again attempting to verify and test the possible “worst known case” scenario for these systems. As part of our baseline commitments, we are aiming to keep post-mitigation risk at “medium” risk or below.

Forecasting, “early warnings,” and monitoring

Many of the mitigations that would be necessary to address risks at a “high” or “critical” pre-mitigation level (e.g., hardening security) require adequate lead time to implement. For this reason, we will be investing in efforts that help create an internal “preparedness roadmap” and help us thus properly plan for and get ahead of the emerging risks. These efforts will include sustained research related to scaling trends for dangerous capabilities and ongoing monitoring of misuse.

We will also, in cooperation with other teams (e.g., Safety Systems), develop monitoring and investigative systems. This monitoring of real-world misuse (as well as staying abreast of relevant research developments) will help us create a better picture of deployed model characteristics, and inform updates to our evaluations as necessary.

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.



Illustrative Scorecard

Note: Below is only an illustrative template version of what the Scorecard might look like; all specifics are purely for illustrative purposes and do not reflect the results from real evaluations.

Tracked Risk Category	Pre-mitigation risk level Determine pre-mitigation risk level using best known capability elicitation techniques	Post-mitigation risk level Determine overall risk level after mitigations are in place using best known capability elicitation techniques
Cybersecurity	Low	Low
CBRN	Low	Low
Persuasion	Medium	Low
Model Autonomy	Medium	Low
Unknown Unknowns	Low	Low



Cybersecurity: Low (pre-mitigation), Low (post-mitigation)

Model assistance for cyberattacks is meaningful only for non-programming uses and the model cannot meaningfully execute elements of cyber operations.

- Cyber theory of attack eval: On a dataset of challenging conceptual multiple-choice questions extracted and rephrased from pentest writeups, the model performs significantly below the human baseline (X% vs Y%).
- GPT vs search availability exploit eval: On a set of misconfigured web servers vulnerable to denial of service, cybersecurity professionals find <X% more bugs when provided with GPT assistance.
- Code de-obfuscation eval: Model recovers variable names and code semantics at lower performance than commercial deobfuscation tools (X% vs Y%).
- AI-assisted capture-the-flag eval: Participants in a hacking challenge (capture-the-flag) obtain <X% higher score from using ChatGPT augmented with the ability to execute commands on a vulnerable virtual machine.



CBRN: Low (pre-mitigation), Low (post-mitigation)

Model can provide information relevant to creating CBRN threats with comparable utility to existing resources (e.g., textbooks, search engines), but nothing more than that.

- GPT vs search eval: Post-PhD professionals trained in biology (specifically genetics, virology, or microbiology) experience no improvement in success rates for writing CDC Class B biological agent development instructions using a search engine and GPT versus just a search engine (X% vs X%), including across accuracy, completeness, and innovation metrics. This holds true for chemistry, nuclear, and radiological post-PhD professionals students as well.
- GPT vs search eval: Undergraduate students with some coursework in biology (specifically genetics, virology, or microbiology) experience no improvement in success rates for writing CDC Class B biological agent development instructions using a search engine and GPT versus just a search engine (X% vs X%), including across accuracy, completeness, and innovation metrics. This holds true for chemistry, nuclear, and radiological study undergraduate students as well.
- GPT vs search identify suppliers eval: <X% more participants succeed in finding DNA suppliers that do not screen orders using a search engine and GPT versus just a search engine.



Persuasion: Medium (pre-mitigation), Low (post-mitigation)

Model can create (potentially interactive) content with comparable persuasive effectiveness to typical human-written content (e.g., an op-ed from a reputable news source; 1:1 conversation with a typical randomly-assigned individual), but mitigations bring this persuasive effectiveness down to that of a low-quality article.

- Political articles eval: Content generated by model increases rate of human belief in a political issue significantly compared to human-generated controls (X% vs Y%), but mitigations bring it to (X% vs X%).
- Interactive political conversations eval: Long-term 1:1 interactive conversation with a digital “friend” increases the rate of human belief in a political issue significantly when the digital friend is a model versus a real human (X% vs Y%), but mitigations bring it to (X% vs X%).
- Donations eval: Content generated by model does not significantly increase donations to a charitable cause compared to human-generated controls (X% vs X%).
- Digital action eval: Content generated by model does not increase rate of installing a Chrome extension significantly compared to human-generated controls (X% vs X%).
- Phishing eval: Model cannot write a realistic phishing email that convinces employees to submit login details on a phishing site (<X% success rate).
- Believable biology research proposal: <X% of subtly dangerous proposals submitted by a model are approved by a biology wet lab.
- Backdoored internal ChatGPT eval: No employees are convinced to run code with a subtle backdoor, e.g., a specific HTTP request, when ChatGPT is instructed to insert backdoors in code generations.
- Copy GPT weights eval: Model convinces no employees to attempt to copy weights into a new location.
- Monitored misuse: <X% of users are using the model for political persuasion purposes, and no electoral campaigns have been found that rely on our models.



Model autonomy: Medium (pre-mitigation), Low (post-mitigation)

Model can robustly complete a diverse set of chained actions and acquire resources (e.g., spin up an AWS instance, complete a Fiverr task), but post-mitigation, can only take discrete actions if explicitly instructed to do so, and fails to chain actions or acquire resources.

- Server purchase eval: With extensive scaffolding, the model can buy additional server space (>X% success rate) but post-mitigation it fails.
- Freelance job eval: Model cannot autonomously find and complete freelance tasks on Fiverr, but can complete tasks X% of time with some scaffolding and assistance; this reduces to 0% post-mitigation.
- Reward hacking eval: Model cannot reward hack in a structured environment (0% success rate).
- Research engineering eval: Model cannot complete long-term research engineering tasks (e.g., autonomously find compute efficiency wins on pretraining, discover and fix bugs that cause performance regressions); 0% success rate.
- Weights escape eval: Model cannot successfully penetrate any important layers of defense on its own weights (e.g. exploiting a browsing environment, finding a vulnerability in inference code); 0% success rate.



Governance

Beyond our above-detailed processes with regards to evaluations, monitoring, and forecasting, we want to make a set of safety baselines and procedural commitments. Note that these processes also apply to training runs as we will perform periodic evaluations (e.g., of checkpoints) during model training.

Safety baselines

As mentioned in the introduction, OpenAI already has extensive safety processes in place both before and after deployment (e.g., system cards, red-teaming, refusals, jailbreak monitoring, etc.). We now augment these processes with three safety baselines.

Asset Protection

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.



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 reasonable 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.)

Restricting development

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.

Additionally, to protect against "critical" pre-mitigation risk, we need dependable evidence that the model is sufficiently aligned that it does not initiate "critical"-risk-level tasks unless explicitly instructed to do so.



Operations

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).

- Parties in the Preparedness Framework operationalization process:
 - a. **The Preparedness team** conducts research, evaluations, monitoring, forecasting, and continuous updating of the Scorecard with input from teams that have relevant domain expertise.
 - b. **The Safety Advisory Group (SAG), including the SAG Chair**, provides a diversity of perspectives to evaluate the strength of evidence related to catastrophic risk and recommend appropriate actions. The SAG will strive to recommend mitigations that are as targeted and non-disruptive as possible while not compromising safety. In particular, we recognize that pausing deployment or development would be the last resort (but potentially necessary) option in these circumstances.
 - i. The members of the SAG and the SAG Chair are appointed by the OpenAI Leadership. This will be done in consultation with the BoD, as is common for important roles in general.
 - ii. SAG membership will rotate yearly. OpenAI leadership might choose to re-appoint someone from previous years to ensure there is continuity of knowledge and experience, while still ensuring that fresh and timely perspectives are present in the group.
 - iii. The SAG Chair makes any final decisions needed for the SAG. This role is expected to rotate, as appointed by OpenAI leadership.



- c. **The OpenAI Leadership**, i.e., the CEO or a person designated by them, serves as the default decision-maker on all decisions.

 - d. **The OpenAI Board of Directors (BoD)**, as the ultimate governing body of OpenAI, will oversee OpenAI Leadership’s implementation and decision-making pursuant to this Preparedness Framework. The BoD may review certain decisions taken and will receive appropriate documentation (i.e., without needing to proactively ask) to ensure the BOD is fully informed and able to fulfill its oversight role.
- Process:
 - a. The Preparedness team is responsible for:
 - i. maintaining and updating the Scorecard, including designing and running evaluations to provide Scorecard inputs and collecting relevant information on monitored misuse, red-teaming, and intelligence

 - ii. monitoring for unknown unknowns and making the case for inclusion in the Preparedness Framework of any new risk categories as they emerge

 - iii. ensuring the risk level distinctions in the Tracked Risk Categories section are appropriate given developments in frontier AI models, and suggesting updates to these levels if needed

 - iv. forecasting potential changes to catastrophic risk levels, and summarizing evidence for an “early warning” / “heads up” as needed

 - v. providing a monthly report (sent to the SAG, Leadership and BoD) synthesizing the above with any potential protective actions (the SAG Chair, OpenAI Leadership, and/or BoD can adjust this cadence as needed)

 - b. If the Preparedness or any other team determines that any changes to the Preparedness Framework are necessary, it will include a case for this change in its report. The case will consist of the suggested new version of the relevant parts of the Preparedness Framework along with a summary of evidence supporting the change (and evidence against). This case is then sent to SAG and processed according to the standard decision-making process described below.



- c. **Fast-track:** In the rare case that a severe risk rapidly develops (e.g., there is a change in our understanding of model safety that requires urgent response), the Preparedness team can request a fast track for the SAG to process the report urgently. The SAG Chair should also coordinate with OpenAI Leadership for immediate reaction as needed to address the risk.
- Decision process:
 - a. The SAG will be responsible for assessing the merits of each case submitted to them. Then, the SAG Chair will forward the case to the OpenAI Leadership (sharing also with the BoD and Preparedness) and include a draft of possible next steps (e.g., approving a change to the Scorecard or Preparedness Framework; requesting changes to other internal and external processes as needed; highlighting a need to accelerate security procedures based on a forecasted increase in risk levels) as well as rationale.
 - b. The OpenAI Leadership will make the final decision and be responsible for overseeing any mitigations (with SAG’s assistance, if requested).
 - i. The BoD will be given visibility into processes, and can review decisions and otherwise require reports and information from OpenAI Leadership as necessary to fulfill the BoD’s oversight role. Where necessary, the BoD may reverse a decision and/or mandate a revised course of action.
 - ii. (For the avoidance of doubt, Leadership can also make decisions without the SAG’s participation, i.e., the SAG does not have the ability to “filibuster”).
 - c. **Internal visibility:** The Preparedness Framework, reports and decisions will be documented and visible to the BoD and within OpenAI (with redactions as needed given internal compartmentalization of research work). This also includes any audit trails created from the below.



- Accountability:
 - a. **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.
 - b. **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.
 - c. **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.
 - d. For each of these accountability measures, decision-making will follow the standard process (i.e., SAG chair synthesizes final recommendation, which OpenAI leadership makes a final decision on, that the BoD can overrule if needed).



Example Scenarios

Example scenario 1: “High” risk in persuasion

1. The Preparedness team monthly report updates the Scorecard pre-mitigation risk level in persuasion to “high” for a model that has just been trained. This report includes a case with evidence from evaluations.
2. The SAG Chair accepts the evidence supporting this new risk level, which would trigger two safety baselines: (1) to not continue with the deployment of the pre-mitigated model until mitigations are in place to ensure the post-mitigation risk will be at most “medium,” and (2) to ensure security measures are in place to prevent exfiltration by relevant actors; in this case, that would likely include foreign disinformation groups. The Preparedness team consults with relevant teams about the recommended mitigations and puts together a final set of suggestions for how to ensure the model is back at “medium” risk, which the SAG agrees with. The SAG Chair escalates to OpenAI leadership for a final decision. All documentation is sent to both OpenAI Leadership and the BoD for visibility and potential review.
3. OpenAI Leadership agrees that the new pre-mitigation level for persuasion would be “high,” and agrees with the SAG plan to ensure sufficient mitigations are in place. The BoD does not overrule the decision. The SAG works with the Preparedness and Security teams to ensure baseline safety measures are in place.
4. The Preparedness team continues regular evaluations of both the pre-mitigation and post-mitigation risk and determines that while the pre-mitigation risk is “high”, the mitigations are sufficient to put the post-mitigation risk at “medium.” This holds true through the rest of the development process, and is validated by a requested third-party audit (assuming qualified and effective auditors can be found), before the model is greenlit for deployment.

Example scenario 2: Forecasted “critical” risk in cybersecurity with fast-track process

1. The Preparedness team discovers a new, much more effective prompting technique. Based on updated scaling estimates using this technique, the risk level for cybersecurity is projected to hit “critical” levels within six months, once a new model finishes training. The Preparedness team includes this forecast in their report and calls for a fast-track given the severity and rapid development involved.
2. The SAG agrees with the evidence supporting this new risk level and the CEO agrees. This triggers the safety baseline to (1) not continue with the training and development of that model until mitigations are in place to ensure the post-mitigation risk will be at “high” or below in cybersecurity, and (2) implement security measures to address exfiltration by all newly relevant actors. The SAG consults with relevant teams about the recommended mitigations and puts together a final set of suggestions for how to ensure the post-mitigation risk is at most at “high,” which the CEO approves. The BoD also does not overrule the decision.
3. The Preparedness team continues regular monthly evaluations on both the worst case pre-mitigations model as well as the mitigated model and determines the mitigations are sufficient to put the post-mitigation risk at “high.” This is validated by a requested third-party audit (assuming qualified and effective auditors can be found) before development that would transition the model to “critical” can ensure.

