



NIST AI 800-1
Initial Public Draft

Managing Misuse Risk for Dual-Use Foundation Models

U.S. AI Safety Institute

This publication is available free of charge from:
<https://doi.org/10.6028/NIST.AI.800-1.ipd>

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25

Managing Misuse Risk for Dual-Use Foundation Models

U.S. AI Safety Institute

This publication is available free of charge from:
<https://doi.org/10.6028/NIST.AI.800-1.ipd>

July 2024



U.S. Department of Commerce
Gina M. Raimondo, Secretary

National Institute of Standards and Technology
Laurie E. Locascio, NIST Director and Under Secretary of Commerce for Standards and Technology

1
2
3
4
5
6
7
8
9

The U.S. AI Safety Institute (AISI) at NIST is releasing this document as an Initial Public Draft for public comment.

Comments on NIST AI 800-1 may be sent electronically to NISTAI800-1@nist.gov with “NIST AI 800-1, Managing Misuse Risk for Dual-Use Foundation Models” in the subject line. Electronic submissions may be sent as an attachment in any of the following unlocked formats: HTML; ASCII; Word; RTF; or PDF.

All comments are subject to release under the Freedom of Information Act (FOIA).

1 **Table of Contents**

2 **1. INTRODUCTION..... 1**

3 **2. SCOPE 1**

4 **3. KEY CHALLENGES IN MAPPING AND MEASURING MISUSE RISKS 2**

5 **4. OBJECTIVES AND PRACTICES TO MANAGE MISUSE RISKS 4**

6 Objective 1: Anticipate potential misuse risk 5

7 Objective 2: Establish plans for managing misuse risk 7

8 Objective 3: Manage the risks of model theft 8

9 Objective 4: Measure the risk of misuse..... 9

10 Objective 5: Ensure that misuse risk is managed before deploying foundation models 12

11 Objective 6: Collect and respond to information about misuse after deployment..... 13

12 Objective 7: Provide appropriate transparency about misuse risk 16

13 **Appendix A. Glossary 18**

14 **Appendix B. Example Safeguards Against the Misuse of Foundation Models..... 19**

15

16

17 ***Disclaimer:*** *Certain equipment, instruments, software, or materials, commercial or non-commercial, are*

18 *identified in this paper in order to specify the experimental procedure adequately. Such identification*

19 *does not imply recommendation or endorsement of any product or service by NIST, nor does it imply that*

20 *the materials or equipment identified are necessarily the best available for the purpose.*

1 1. INTRODUCTION

2 This document provides guidelines for improving the safety, security, and trustworthiness of
3 dual-use foundation models (hereafter referred to as “foundation models”)ⁱ consistent with the
4 National AI Initiative Actⁱⁱ and Executive Order 14110.ⁱⁱⁱ Specifically, it focuses on managing the
5 risk that such models will be deliberately misused to cause harm. The ways that foundation
6 models can be misused continue to evolve, but they include the risks that models will facilitate
7 the development of chemical, biological, radiological, or nuclear weapons; enable offensive
8 cyber attacks; aid deception and obfuscation; and generate child sexual abuse material (CSAM)
9 and non-consensual intimate imagery (NCII) of real individuals.

10 The rapid development of foundation models poses significant challenges to understanding
11 their capabilities and misuse risks,¹ and this document provides a basis to identify, measure,
12 and reduce these risks across the AI lifecycle. Misuse risks are not a product of a model alone—
13 they result in part from malicious actors’ motivations, resources, and constraints, as well as
14 society’s defensive measures against that harm.² As a result, the guidelines provided here
15 address both technical and social aspects of these risks.

16 Building on the AI Risk Management Framework,³ this document identifies best practices to
17 map, measure, manage, and govern misuse risks from foundation models, as well how
18 organizations can provide transparency into how they are managing these risks. This document
19 focuses particularly on foundation models’ initial developers, but actors across the lifecycle of a
20 model all play a role in managing misuse risks.⁴

ⁱ Executive Order 14110 defines a “dual-use foundation model” as “an AI model that is trained on broad data; generally uses self-supervision; contains at least tens of billions of parameters; is applicable across a wide range of contexts; and that exhibits, or could be easily modified to exhibit, high levels of performance at tasks that pose a serious risk to security, national economic security, national public health or safety, or any combination of those matters, such as by: (i) substantially lowering the barrier of entry for non-experts to design, synthesize, acquire, or use chemical, biological, radiological, or nuclear (CBRN) weapons; (ii) enabling powerful offensive cyber operations through automated vulnerability discovery and exploitation against a wide range of potential targets of cyber attacks; or (iii) permitting the evasion of human control or oversight through means of deception or obfuscation.” This definition is provided in a glossary at the end of this report, along with definitions of other key terms.

ⁱⁱ As codified in 14 U.S.C. § 278h-1.

ⁱⁱⁱ Section 4.1(a)(ii) of Executive Order 14110 directs the Secretary of Commerce to “Establish appropriate guidelines (except for AI used as a component of a national security system), including appropriate procedures and processes, to enable developers of AI, especially of dual-use foundation models, to conduct AI red-teaming tests to enable deployment of safe, secure, and trustworthy systems. These efforts shall include: (A) coordinating or developing guidelines related to assessing and managing the safety, security, and trustworthiness of dual-use foundation models”.

1 2. SCOPE

2 **Risks.** This document focuses on misuse risk from dual-use foundation models. Consistent with
3 Section 3(k) of Executive Order 14110,⁵ this includes foundation models that exhibit, or could
4 be easily modified to exhibit, high levels of performance at tasks that can pose a serious risk to
5 security, economic security, public health or safety, or any combination of those matters
6 (hereafter referred to jointly as “public safety”). This document addresses both emerging
7 misuse risks, such as a foundation model facilitating the development of a novel biological
8 weapon, as well as current harms from misuse, such as a foundation model generating CSAM or
9 NCII.

10 This document does not address other important risks from foundation models, such as bias,
11 discrimination, and hallucination, nor does it address all risks to public safety, including those
12 that may arise from other types of AI models and systems.⁶ Actors across the AI lifecycle should
13 manage these risks as well, consistent with relevant guidelines, such as those provided in the
14 *Blueprint for an AI Bill of Rights*,⁷ as well as the NIST AI Risk Management Framework and its
15 Generative AI Profile.⁸

16 **Actors.** The practices in this document are principally focused on the central role that
17 foundation models’ initial developers have in the supply chain for their models.⁹ These
18 developers contribute most to determining how their models are made available, the models’
19 capabilities, and safeguards against their misuse. In some cases, model developers may share
20 influence over these factors with an external partner, such as a cloud service provider that
21 leads deployment of the model, and in such cases these partners also have expanded
22 opportunities and responsibilities to manage the risks that a model may be misused.¹⁰

23 Other parties also play important roles in managing misuse risks, but they are not the focus of
24 this document. They include downstream developers and deployers, third-party evaluators and
25 auditors, civil society organizations, and government agencies.¹¹ Relevant stakeholders
26 throughout the AI supply chain are encouraged to share information and collaborate to
27 understand and mitigate misuse risks, including to integrate appropriate risk mitigations into
28 downstream systems that rely on foundation models.¹²

1 3. KEY CHALLENGES IN MAPPING AND MEASURING MISUSE RISKS

2 Accurately identifying and measuring foundation models' misuse risks helps to build confidence
3 in which potential harms are realistic and which practices are justified to mitigate them.
4 However, methodological and scientific challenges can limit understanding and measurement
5 of these risks. Organizations should take appropriate steps to address these challenges and
6 strive to build an empirical basis to evaluate and mitigate risks. These challenges include:

- 7 1. **Foundation models are broadly applicable.** Models trained on a broad data distribution
8 can often be applied across many different domains, including domains not explicitly
9 considered by the developer. This makes it challenging to anticipate potential ways in
10 which a model might be misused and complicates measuring or monitoring for misuse
11 risk broadly.
- 12 2. **Capabilities do not clearly translate across domains.** A foundation model's
13 performance on one task may not provide reliable evidence of its performance on
14 another, even when the two tasks appear related.¹³ For instance, initial benchmarks
15 that are cheaper and easier to carry out may suggest a model has dangerous
16 capabilities, but this concern may not be substantiated when the model is tested in
17 more rigorous and realistic conditions.
- 18 3. **It is difficult to predict how scale will affect performance.** In many instances, a
19 foundation model's performance can be improved by increasing the amount or quality
20 of its training data, the quantity of compute used to build the model, or the number of
21 parameters of the model.¹⁴ But while these factors may be useful heuristics in some
22 instances, their precision is limited and their relationship to any given risk is uncertain.¹⁵
- 23 4. **The relationship between a measured capability and its potential risk of harm is often**
24 **unclear.** It remains challenging to determine the likelihood or severity of real-world
25 harm through isolated testing. For many harms to public safety, domain knowledge and
26 computer automation are not the only the limiting factors in carrying out an attack,
27 which may also require physical infrastructure, distribution mechanisms, or complex
28 interactions in the physical world.¹⁶ Bad actors may also already have access to existing
29 tools that serve their needs better than foundation models, and existing methods to
30 prevent harm—such as controls on physical dual-use materials—may be substantially
31 more determinative of real-world risks.¹⁷
- 32 5. **Methods to evaluate safeguards are nascent.** Foundation model developers implement
33 safeguards to protect their models from misuse, such as those outlined in Appendix B,
34 but there is a lack of effective techniques for evaluating the adequacy of those
35 safeguards under real-world conditions.
- 36 6. **Evaluating risk may depend on scarce domain expertise.** Information about some risks,
37 like the potential for a foundation model to enable a malicious actor to develop a
38 chemical or biological weapon, may be closely guarded. In addition, organizations
39 developing foundation models may not currently have the professional experts or staff,

- 1 nor the developed means of accessing external expertise, necessary to appropriately
2 assess some misuse risks.
- 3 7. **It is difficult to accurately emulate threat actors and misuse.** Assessing misuse risk
4 requires realistically profiling threats. This assessment includes understanding how an
5 actor might exploit a model’s capabilities to cause harm, or how they might circumvent
6 safeguards that have been deployed to protect a model from misuse. This practice can
7 be challenging if malicious actors are willing to spend more time than evaluators or if
8 they have access to infrastructure or niche expertise not available during evaluations.
9 Realistic emulation of malicious actors may also be prohibitively dangerous or unethical,
10 such as attempting to develop a volatile substance, generate non-consensual intimate
11 imagery, or harm real people.

1 4. OBJECTIVES AND PRACTICES TO MANAGE MISUSE RISKS

2 This section outlines seven objectives, as well as associated practices that can help achieve
3 them, for organizations to map, measure, manage, and govern the risk that their foundation
4 models will be misused to deliberately harm public safety, consistent with the NIST AI Risk
5 Management Framework. These objectives are:

- 6 1. Anticipate potential misuse risk
- 7 2. Establish plans for managing misuse risk
- 8 3. Manage the risks of model theft
- 9 4. Measure misuse risk
- 10 5. Ensure that misuse risk is managed before deploying foundation models
- 11 6. Collect and respond to information about misuse after deployment
- 12 7. Provide appropriate transparency about misuse risk

13

14 Organizations should consider all seven objectives holistically to manage foundation models'
15 misuse risks. This document provides practices for each objective as non-exhaustive examples
16 of how organizations can meet the objectives. Managing risk is an iterative process, and
17 organizations should assess which practices are most relevant at each point in the lifecycle of a
18 foundation model. Implementing these practices also cannot guarantee that actors will not
19 misuse a foundation model, and organizations should adopt additional risk management
20 measures where appropriate.

21

22 The specific elements of each objective below include:

- 23 1. **Objective:** an aspirational outcome that, when achieved, will help organizations manage
24 a foundation model's misuse risks.
- 25 2. **Practice:** a suggested practice that can help organizations achieve an objective. This
26 document proposes a particular set of practices that can be collectively used to achieve
27 the indicated objectives. However, an individual practice may not be appropriate for all
28 contexts, alternative practices may also be able to achieve these objectives, and new
29 practices are likely to be developed over time.
- 30 3. **Recommendation:** an implementation characteristic or consideration that is often
31 important for a practice to be effective at achieving an objective.
- 32 4. **Documentation:** information that could be shared with either the public or with select
33 parties to help demonstrate whether and how a practice was implemented, as well as
34 potential evidence of its effectiveness and comprehensiveness. This documentation can
35 also help to improve collaboration with third parties and to develop shared best
36 practices for managing misuse risks. If an organization relies on alternative practices to
37 achieve an objective, a similar level of documentation about those practices can be used
38 to provide transparency, along with a rationale for why the alternative practices match
39 or exceed the efficacy of those provided here.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35

Objective 1: Anticipate potential misuse risk
Assess the risks that a foundation model could be misused if it became available to malicious actors, including assessing these risks before the model is developed based on an estimate of its capabilities. Identify the most significant expected risks so that they can be measured and managed, as necessary.

Practice 1.1: Identify and maintain a list of threat profiles that covers significant ways in which malicious actors might misuse the model.

Recommendations to effectively implement Practice 1.1:

1. For each threat profile, specify the relevant capabilities of concern, a threat actor or actors who might use them to cause harm, and the malicious tasks that the threat actor might accomplish using the model.
2. Consider known misuse risks, including, as relevant, the risk that a model might be used to facilitate the development of chemical, biological, radiological, or nuclear weapons, to automate offensive cyber operations, or to generate CSAM or NCII.
3. Consider consulting external experts with relevant expertise and responsibilities to help identify gaps in threat profiles, including potentially granting them model access for open-ended experimentation to identify other ways that a model could be misused to cause harm. The consultation of external experts may be particularly relevant to helping ensure that the threat profiles cover the most significant expected ways in which a model could be misused to cause harm.

The following documentation can help provide transparency about how Practice 1.1 is implemented:

- a. The list of key threat profiles.¹⁸
- b. A description of how these threat profiles were selected and judged to have adequate coverage.
- c. A description of any internal and external experts involved and their role in developing profiles.

Practice 1.2: Assess the impact of each identified threat profile if the malicious actor successfully misused the model.

Recommendations to effectively implement Practice 1.2:

1. For each threat profile, assess the impact to public safety if a threat actor successfully used the model to carry out the malicious task, including estimates of how much it would help an actor increase the scale, decrease the cost, or improve the effectiveness

- 1 of their malicious activity compared to alternative resources that are available to that
2 actor, such as other machine learning models and digital tools.¹⁹
- 3 2. For each threat profile, estimate the scale, frequency, or probability with which threat
4 actors might attempt to misuse the model.
- 5 3. Consider how potential harms may be prevented or managed outside of the context of
6 the model itself, such as mechanisms that will effectively stop harm from a given
7 malicious activity regardless of how inexpensive or easy the model makes that malicious
8 activity.

9 The following documentation can help provide transparency about how Practice 1.2 is
10 implemented:

- 11 a. An impact assessment for each identified threat profile.

12

13 **Practice 1.3: Estimate the model’s capabilities of concern before it is developed by comparing**
14 **it to existing models.**

15 Recommendations to effectively implement Practice 1.3:

- 16 1. Identify similar models for which capabilities of concern have already been measured.
- 17 2. Based on comparing the characteristics of the upcoming model with existing models,
18 estimate how its performance will compare to those models.
- 19 3. Assess the degree of uncertainty in these estimates based on how different the models
20 are, how those differences are expected to affect their capabilities, and the reliability
21 and completeness of the evaluations available for the similar models.
- 22 4. If capability forecasts are uncertain and include the possibility that the model may pose
23 significant risk, consider development strategies that may reduce reliance on
24 forecasting, such as taking steps to periodically measure capabilities throughout the
25 development process.

26 The following documentation can help provide transparency about how Practice 1.3 is
27 implemented:

- 28 a. Capability estimates that were made prior to developing the model.
- 29 b. Relevant capability measurements performed after the model was developed.
- 30 c. Estimated comparisons of the model’s misuse risks compared to common benchmarks,
31 where available.

32

33

34

35

Objective 2: Establish plans for managing misuse risk

Determine acceptable levels of misuse risk, considering legal or regulatory obligations, potential benefits that may trade off against risk, and other factors. Align development and deployment plans with the resources, time, and operational constraints that may be required to manage potential misuse risk.

Practice 2.1. Identify a level of misuse risk which the organization considers unacceptable.

Recommendations to effectively implement Practice 2.1:

1. Define a risk threshold for each identified threat profile based on the impact assessments carried out in Practice 1.2 and other relevant information.²⁰
2. Account for how much the model may reduce the cost of accomplishing the malicious task, increase the number of actors who could accomplish that task, or increase the scale or quality at which a given actor can accomplish the task.
3. When determining risk thresholds, weigh potential benefits against the misuse risks.
4. Refine thresholds continuously based on adjustments to identified threat profiles or changes in factors that affect risk tolerance (such as increases in expected benefits).
5. Consider how external factors such as legal and regulatory obligations may inform what risk thresholds are acceptable.

The following documentation can help provide transparency about how Practice 2.1 is implemented:

- a. The thresholds for the misuse risks of a model that the organization considers unacceptable.

Practice 2.2: Establish a roadmap to manage misuse risks for the development of planned foundation models and future versions.

Recommendations to effectively implement Practice 2.2:

1. Define how the organization expects to manage misuse risk and achieve the objectives established in this document.
2. Plan to implement security practices to protect models from model theft if necessary to manage misuse risk. Define security goals and a timeline for achieving those goals.
3. Plan to implement appropriate safeguards to protect deployed models from misuse if necessary to manage misuse risk. Define safety goals and a timeline for achieving those goals.
4. Plan to adjust deployment or development strategies if misuse risks rise to unacceptable levels before adequate security and safeguards are available to manage risk.

1 The following documentation can help provide transparency about how Practice 2.2 is
2 implemented:

- 3 a. A roadmap documenting what risk mitigations are expected to be feasible over time
4 given technical and resource constraints, and what steps will be taken if measurements
5 suggest those mitigations are not adequate to manage misuse risk.

6

7

Objective 3: Manage the risks of model theft

Take steps to prevent the theft of information and assets that would allow a malicious actor to recreate the foundation model, where doing so is necessary to manage misuse risk. Only develop a model that relies on confidentiality to manage misuse risk when the risk of model theft is sufficiently mitigated.

8

9

10

11

12

Practice 3.1: Assess the risk of model theft from relevant threat actors.

13

Recommendations to effectively implement Practice 3.1:

14

15

16

17

18

19

20

21

22

1. Assess the risk of model theft by each identified actor or threat profile, based on the relevant actors' respective resources, motivations, tactics, and sophistication.²¹
2. Consider the organization's compliance with applicable cybersecurity best practices as an input into assessing the risk of model theft.
3. Consider the risk posed by an insider threat, such as an individual involved in developing or deploying the model who may behave maliciously or collaborate with an external attacker.
4. Consider using cybersecurity red teams and penetration testing to assess how difficult it would be for an actor to circumvent security measures.

23

24

The following documentation can help provide transparency about how Practice 3.1 is implemented:

25

26

27

- a. A summary of evidence collected to evaluate the risk of model theft, including the results of any red team exercises and penetration testing.

28

29

Practice 3.2: Compare predicted misuse risk to the organization's risk tolerance prior to developing models with increased capabilities of concern.

30

Recommendations to effectively implement Practice 3.2:

31

32

33

34

1. Prior to developing a model, assess overall misuse risk from model theft by combining the estimated capabilities of concern with an estimate of the probability of model theft.
2. Adjust or halt further development until the risk of model theft is adequately managed and within the organization's defined risk thresholds.

- 1 3. Periodically revisit estimates of misuse risk stemming from model theft, including prior
2 to individual development decisions that may significantly increase capabilities of
3 concern.

4 The following documentation can help provide transparency about how Practice 3.2 is
5 implemented:

- 6 a. The risk assessments that were performed, either periodically or prior to significant
7 development and deployment decisions, related to model theft.

8

9 **Practice 3.3: Maintain security practices sufficient to prevent model theft.**

10 Recommendations to effectively implement Practice 3.3:

- 11 1. Implement security practices designed to protect against model theft.
12 2. Where applicable, apply security practices specifically tailored to the context of
13 foundation models, such as protections against exfiltrating large amounts of data (e.g.,
14 model weights) and protections against vulnerabilities like extraction attacks.²²
15 3. Increase the level of security proportionate to the level of misuse risk posed by the
16 model and the importance of preventing model theft to managing that risk.
17 4. Apply appropriate protections against insider threats, such as limiting access to model
18 weights within the organization.
19 5. Consider adopting existing security standards for sensitive data and applications, such as
20 U.S. government cybersecurity guidance and applicable international standards.
21 6. Re-assess the risk of model theft prior to development with new security practices in
22 place (if previous assessments found unacceptable risks from model theft).

23 The following documentation can help provide transparency about how Practice 3.3 is
24 implemented:

- 25 a. A summary of security measures that have been implemented to reduce the risk of
26 model theft.

27

28

Objective 4: Measure the risk of misuse

29

Where there is a reasonable assessment that a foundation model could be misused, measure the predicted risks to provide evidence for the model's actual misuse risk in a real-world context. Rely on methods that incorporate both technical and social factors, as well as those that provide accurate evidence and higher confidence, while avoiding any harm that could be caused by measuring dangerous activities.

30

31

32

33

34 **Practice 4.1: Measure model capabilities relevant to assessing misuse risk.**

35 Recommendations to effectively implement Practice 4.1:

- 1 1. Directly measure capabilities of concern or identify alternative capability measurements
2 that allow for sufficiently reliable inferences about them, such as pursuing
3 recommendations 2, 3, and 4 below when they are sufficiently reliable.
- 4 2. Consider comparing the performance of a new model to an existing model with known
5 misuse risks (the “proxy model”). If the models are similar across a broad range of
6 evaluations, results from the proxy model’s capability evaluations may be suitable to
7 assess risk from the new model. If this comparison is inconclusive or suggests that the
8 new model may pose significant risk, then the new model’s capabilities of concern
9 should be separately evaluated.
- 10 3. Consider measuring model performance on proxy tasks that are safe and tractable while
11 being similar enough to allow reliable inferences about the capability of concern.
- 12 4. If the most accurate capability measurement would be expensive and it is appropriate
13 to do so, consider using cheaper tests (such as automated evaluations on simpler tasks
14 that would be easier for the model to perform than the true capabilities of concern)²³
15 that may indicate that a model lacks dangerous capabilities. If these tests are
16 inconclusive or indicate a potential risk, then rely on more costly and precise
17 measurements.
- 18 5. Account for the possibility that the model's capabilities of concern may exceed any
19 evaluated proxy capabilities and adjust the interpretation of evaluation results
20 accordingly. This consideration is especially relevant when the proxy capability being
21 evaluated differs significantly from the actual capability of concern.
- 22 6. Assess what a threat actor could achieve given access to the weights of a model and the
23 ability to integrate it with other tools by testing AI systems which are reasonably
24 optimized for performance on the evaluation task. If there is a gap between the effort
25 applied to optimize system performance during testing and the effort that could be
26 applied by a threat actor, evaluations should explicitly account for the uncertainty
27 introduced by that gap.
- 28 7. Avoid overlap between data used to train the model and data used in capability and risk
29 evaluations and measure the extent of any overlap.

30 The following documentation can help provide transparency about how Practice 4.1 is
31 implemented:

- 32 a. A list of evaluation tasks that were used to evaluate each threat profile, and a
33 representative subset of datasets used for each evaluation.
- 34 b. A methodological description for each evaluation in enough detail to reproduce it.
- 35 c. An analysis of the relationship between evaluation tasks and capabilities of concern,
36 addressing the possibility that evaluation tasks are more challenging.
- 37 d. If a model was evaluated via a comparison with a proxy model: the proxy model used
38 for comparison, evaluations used to establish comparability, the risk assessment for the

1 proxy model, and other ways in which the proxy may be meaningfully different from the
2 model being assessed.

3
4 **Practice 4.2: Use red teams to assess whether threat actors could bypass model and system**
5 **safeguards and misuse any capabilities of concern.**

6 Recommendations to effectively implement Practice 4.2:

- 7 1. To justify an assessment that a threat actor is unlikely to be able to misuse a model,
8 verify that an adequately resourced red team is unable to misuse the model or
9 accomplish a related proxy task in a realistic deployment context. Apply practices 4.1.1,
10 4.1.3, and 4.1.4 to the selection of proxy tasks for red team exercises as well as
11 capability evaluations.
- 12 2. Clearly specify what goal the red team is trying to achieve in advance, provide incentives
13 and accountability for achieving those goals, and select red teams based on their ability
14 to achieve those goals.
- 15 3. Rely on red teams comprised of external experts that are meaningfully independent
16 from the model developer and who do not have incentives that conflict with their red-
17 teaming goal.
- 18 4. Compare the red team's expertise, resources, and time available to those of a relevant
19 threat actor. To the extent that there are gaps, explicitly account for those gaps when
20 interpreting the result of a red team exercise.
- 21 5. Consider providing additional resources to the red team or making the red team's task
22 easier in other ways (such as providing additional access,²⁴ reducing task difficulty, or
23 dividing complex tasks into constituent pieces) to compensate for any remaining gaps
24 between the red team and threat actors.
- 25 6. Consider providing red teams with available legal protections for their tasks, such as
26 waiving terms of service and indemnifying them for legal liability for their interactions
27 with the model.
- 28 7. Provide the red team with at least as much information about the system as would be
29 available to an attacker and make explicit any respects in which the red team lacks full
30 information about the design of system safeguards.
- 31 8. Consider each level of access to a model that a threat actor might have, ranging from
32 limited access through an API to direct access to the code and parameters that define
33 the model, and determine the minimum level of access (if any) that allows the red team
34 to accomplish its goal.
- 35 9. Explicitly define the time period over which the red teaming results are intended to
36 apply and consider the fact that a threat actor will have access to any information about
37 how to use the model or circumvent safeguards that becomes public over that time
38 period.

1 The following documentation can help provide transparency about how Practice 4.2 is
2 implemented:

- 3 a. A description of the red team’s goal.
- 4 b. A description of the red team’s composition, expertise, resources, and timelines.
- 5 c. The results of the red team exercise given varying levels of access to the model.
- 6 d. A summary of the strategies pursued by the red team.

7

8 **Objective 5: Ensure that misuse risk is managed before deploying foundation**
9 **models**

10 *Take actions to increase access to the model (e.g., deploying a model via API or releasing its*
11 *weights) only when misuse risks are adequately managed, including that they are at minimum*
12 *within the organization’s risk tolerance.*

13 **Practice 5.1: Assess the effect of a potential deployment on the model’s misuse risk.**

14 Recommendations to effectively implement Practice 5.1:

- 15 1. Consider potential deployments and the levels of access they would grant actors (e.g.
16 whether they would increase the number or type of actors who can access a model or
17 whether they would grant actors greater access to a model’s features).
- 18 2. Identify the level of access that a malicious actor could obtain under each proposed
19 deployment (e.g., would it grant them access to API inference, access to a fine-tuning
20 API, access to model weights) and consider how the deployment may affect misuse risk.
21 For example, allowing fine-tuning via API can significantly limit options to prevent
22 jailbreaking and sharing the model’s weights can significantly limit options to monitor
23 for misuse (Practice 6.1) and respond to instances of misuse (Practice 6.2).

24 The following documentation can help provide transparency about how Practice 5.1 is
25 implemented:

- 26 a. For each deployment: the level of access provided, the assessed misuse risks associated
27 with the deployment, and the set of safeguards anticipated to mitigate these misuse
28 risks.

29

30 **Practice 5.2: Implement safeguards proportionate to the model’s misuse risk.**

31 Recommendations to effectively implement Practice 5.2:

- 32 1. Implement safeguards designed to protect the model from misuse. Appendix A outlines
33 a non-exhaustive list of possible safeguards.
- 34 2. Establish reliable evidence of safeguards’ effectiveness before relying on them to
35 prevent misuse of meaningful capabilities of concern.

- 1 3. If additional safeguards are added in response to identified risks, consider re-assessing
2 misuse risk prior to deployment by carrying out red-teaming exercises with the
3 additional safeguards in place.

4 The following documentation can help provide transparency about how Practice 5.2 is
5 implemented:

- 6 a. The list of safeguards that have been implemented.
7 b. A summary of the result of safeguard evaluations via red teams and other testing.

8

9 **Practice 5.3: Only pursue deployments where misuse risk is adequately managed.**

10 Recommendations to effectively implement Practice 5.3:

- 11 1. For each deployment, establish a process to determine whether the deployment should
12 proceed based on the assessed misuse risk and a consideration of any safeguards.
13 Otherwise, determine whether the deployment should be modified, delayed, or
14 canceled. For instance, consider whether further safety improvements are feasible prior
15 to deployment, whether additional time could be used to carry out a more reliable
16 estimate of risk, or whether a more limited deployment may be more appropriate given
17 the level of assessed risk.
- 18 2. Consider leaving a margin of safety between the estimated level of risk at the point of
19 deployment and the organization’s risk tolerance. Consider how threat actors may
20 continue to acquire new knowledge about how to misuse or augment the model after it
21 is deployed²⁵ and how to circumvent its safeguards.²⁶ Consider a larger margin of safety
22 to manage risks that are more severe or less certain.

23 The following documentation can help provide transparency about how Practice 5.3 is
24 implemented:

- 25 a. The basis for determining that the risk of misuse was adequately managed for a
26 deployment decision, including that the deployment’s risks were within the
27 organization’s risk thresholds.

28

29 **Objective 6: Collect and respond to information about misuse after deployment**

30 *Collect information about deployed systems that improves understanding of their misuse risk to*
31 *adjust deployments and improve future risk management. Engage with and encourage findings*
32 *from the public, relevant civil society organizations, external researchers, and the foundation*
33 *model’s third-party distribution partners.*

34 **Practice 6.1: Where possible, monitor distribution channels for evidence of misuse.**

35 Recommendations to effectively implement Practice 6.1:

- 1 1. Monitor APIs, websites, and other distribution channels for misuse while maintaining
2 privacy of users.
- 3 2. Build or procure systems to enable automated detection of misuse.
- 4 3. Continually assess the effectiveness of systems that are intended to detect misuse to
5 help provide ground truth for their effectiveness and evidence that misuse is not going
6 undetected. Red teaming can be particularly helpful to assess whether malicious actors
7 may be able to deliberately avoid detection.
- 8 4. Request that third-party distribution channels for the model monitor those channels for
9 misuse and that they regularly share information with the developer regarding their
10 monitoring.
- 11 5. Consider using tiered methods of detection when doing so helps prioritize limited
12 resources, improve privacy, and increase coverage. This can include filtering for misuse
13 first by less costly automated methods which are then, as appropriate, validated by
14 methods requiring direct human intervention.
- 15 6. When malicious actors can avoid direct monitoring by operating the model
16 independently, such as when its weights are widely available, consider mechanisms
17 other than those identified in recommendations 6.1.1, 6.1.2, and 6.1.3 that could be
18 used to monitor for patterns of malicious behavior enabled by the model.

19 The following documentation can help provide transparency about how Practice 6.1 is
20 implemented:

- 21 a. A summary of the mechanisms used to monitor each distribution channel for a
22 foundation model and the methods for determining the effectiveness of those
23 mechanisms.²⁷

24

25 **Practice 6.2: Maintain a process to respond to incidents of model misuse.**

26 Recommendations to effectively implement Practice 6.2:

- 27 1. Establish clear organizational responsibilities for responding to incidents of misuse.
- 28 2. Preemptively develop plans for plausible novel scenarios of misuse and how to respond,
29 such as restricting access to the model or strengthening safeguards.²⁸
- 30 3. When it may not be possible to effectively reduce access to a model, such as when the
31 model's weights are widely available, plan for how identified instances of misuse will
32 inform future development and deployment decisions.
- 33 4. Consider carrying out drills to practice responding to time-sensitive and safety-critical
34 scenarios of misuse.

35 The following documentation can help provide transparency about how Practice 6.2 is
36 implemented:

- 1 a. A summary of the incident response process and the organizational roles and
2 responsibilities in the process.

3

4 **Practice 6.3: Establish protections for internal reporting of misuse issues.**

5 Recommendations to effectively implement Practice 6.3:

- 6 1. Adopt policies that protect and reward individuals who report model issues related to
7 misuse risk.
- 8 2. Establish formal processes to adjudicate concerns raised by employees and contractors
9 in a timely fashion.

10 The following documentation can help provide transparency about how Practice 6.3 is
11 implemented:

- 12 a. A summary of the organization’s policies with respect to internal safety reporting.

13

14 **Practice 6.4: Provide safe harbors for third-party safety research.**

15 Recommendations to effectively implement Practice 6.4:

- 16 1. Publish a clear vulnerability disclosure policy for model safety issues that outlines how
17 such vulnerabilities should be shared with the developer and the public, and how the
18 organization will respond to reported vulnerabilities.
- 19 2. Publish a safe harbor policy that commits to not pursuing legal action against or
20 restricting the accounts of external safety researchers that act in good faith and comply
21 with the vulnerability disclosure policy.²⁹
- 22 3. Consider providing support and accommodations for vetted external researchers’
23 interactions with the model, such as providing researchers with access to models with
24 fewer safeguards to conduct post-deployment red-teaming exercises.

25 The following documentation can help provide transparency about how Practice 6.4 is
26 implemented:

- 27 a. A summary of the organization’s commitment to not pursue legal action against third-
28 party researchers acting in good faith.
- 29 b. A description of the organization’s process for providing vetted researchers with access
30 to models or systems with fewer safeguards.³⁰

31

32 **Practice 6.5: Create bounties for issues related to the misuse risk.**

33 Recommendations to effectively implement Practice 6.5:

- 1 1. Establish a program to incentivize researchers for finding vulnerabilities and disclosing
2 them according to the vulnerability disclosure policy.³¹

3 The following documentation can help provide transparency about how Practice 6.5 is
4 implemented:

- 5 a. The terms of the bounties and details regarding the process for submitting
6 vulnerabilities.

7

8 **Objective 7: Provide appropriate transparency about misuse risk**

9 *Provide transparency to the public and relevant entities regarding the organization's processes*
10 *for developing and deploying foundation models that are related to misuse risk to facilitate*
11 *understanding, accountability, collaboration, and scientific development related to model*
12 *misuse.*

13 **Practice 7.1: Publish regular transparency reports that include key details regarding misuse**
14 **risks and how those risks are managed.**

15 Recommendations to effectively implement Practice 7.1:

- 16 1. Share the methodology and results of pre-deployment evaluations of model capabilities,
17 risks, and mitigations, including as much detail about the data and evaluation
18 methodology as can be disclosed without introducing risks to public safety.³²
- 19 2. Share details regarding the safeguards in place for the model, including how they are
20 applied across different distribution channels, with as much as can be shared without
21 rendering the safeguards ineffective.³³
- 22 3. Share information about data used to build the model that is relevant to assessing the
23 misuse risk, such as data sources and criteria for data filtering and selection.³⁴
- 24 4. Share steps that downstream developers and deployers of AI systems that integrate the
25 foundation model should take to manage misuse risk.
- 26 5. Share details regarding the organization's internal structure and internal governance
27 process with respect to determining how to deploy foundation models and map,
28 measure, and manage misuse risk.
- 29 6. Make the transparency reports publicly available, update them on a regular basis (e.g.,
30 with each new major version of the model), and include key information related to
31 misuse risk.³⁵

32

33 **Practice 7.2: Disclose information about risk management practices to promote**
34 **accountability.**

35 Recommendations to effectively implement Practice 7.2:

- 1 1. Regularly share information about risk management practices through a process
2 structured to provide meaningful accountability.
- 3 2. Share information covering the practices used to achieve the objectives listed in this
4 document, including at least as much detail as described in the documentation sections
5 for each practice listed here.
- 6 3. Ensure recipients of information are effectively independent from the development and
7 deployment process and that they are empowered to provide accountability (such as by
8 disclosing risk management deficiencies publicly).

9 The following documentation can help provide transparency about how Practice 7.2 is
10 implemented:

- 11 a. A summary of the list of stakeholders that receive information about risk management
12 practices, and the types of information disclosed to them.

13

14 **Practice 7.3: Report incidents and hazards related to the foundation model to AI incident**
15 **databases.**

16 Recommendations to effectively implement Practice 7.3:

- 17 1. Based on existing best practices and adequate review of the benefits and risks of
18 disclosing certain information, define the category of misuse events to report.³⁶
- 19 2. Collate verified reports of misuse in a commonly used and machine-readable format.
- 20 3. Share verified reports of misuse with relevant third parties, such as AI incident
21 databases.³⁷

22 The following documentation can help provide transparency about how Practice 7.3 is
23 implemented:

- 24 a. A description of the incident reporting process, including an example of a type of
25 incident that might be reported and to whom it would be reported.

1 **Appendix A. Glossary**

2 **Artificial Intelligence (AI)**

3 A machine-based system that can, for a given set of human-defined objectives, make predictions,
4 recommendations, or decisions influencing real or virtual environments. Artificial intelligence systems use
5 machine- and human-based inputs to perceive real and virtual environments; abstract such perceptions into
6 models through analysis in an automated manner; and use model inference to formulate options for information
7 or action.³⁸

8 **AI Red-Teaming**

9 A structured testing effort to find flaws and vulnerabilities in an AI system, often in a controlled environment and
10 in collaboration with developers of AI. Artificial Intelligence red-teaming is most often performed by dedicated
11 “red teams” that adopt adversarial methods to identify flaws and vulnerabilities, such as harmful or discriminatory
12 outputs from an AI system, unforeseen or undesirable system behaviors, limitations, or potential risks associated
13 with the misuse of the system.³⁹

14 **Capability of Concern**

15 A capability of an AI model that poses a risk of enabling malicious behavior that would cause serious harm to public
16 safety.

17 **Distribution Channel**

18 The various ways in which a model could be distributed, including, but not limited to, public release of the model
19 weights, access to the model via an API supplied by a cloud service provider, or access to a fine-tuned or otherwise
20 augmented version of the model from a third-party deployer.

21 **Dual-Use Foundation Model**

22 An AI model that is trained on broad data; generally uses self-supervision; contains at least tens of billions of
23 parameters; is applicable across a wide range of contexts; and that exhibits, or could be easily modified to exhibit,
24 high levels of performance at tasks that pose a serious risk to security, national economic security, national public
25 health or safety, or any combination of those matters, such as by: substantially lowering the barrier of entry for
26 non-experts to design, synthesize, acquire, or use chemical, biological, radiological, or nuclear weapons; enabling
27 powerful offensive cyber operations through automated vulnerability discovery and exploitation against a wide
28 range of potential targets of cyber attacks; or permitting the evasion of human control or oversight through means
29 of deception or obfuscation. Models meet this definition even if they are provided to end users with technical
30 safeguards that attempt to prevent users from taking advantage of the relevant unsafe capabilities.⁴⁰

31 **Fine-Tuning**

32 An approach in which the parameters of an already trained model are adjusted by training on new data. Fine-
33 tuning is often used to adapt a model to a particular task, or to mildly modify a model’s behavior.

34 **Foundation Model**

35 An AI model that is trained on broad data; generally uses self-supervision; contains at least tens of billions of
36 parameters; and is applicable across a wide range of contexts.⁴¹

37 **Jailbreaking**

38 Methods that cause an AI model to act in a manner contrary to how its designer intended for it to behave.

39 **Misuse Risk**

40 A risk that an AI model will be deliberately misused to cause harm.

41 **Model Theft**

42 Unauthorized access to information that substantially aids an actor in recreating an AI model or its capabilities of
43 concern, such as the code or weights that define the model or bypassing security measures designed to prevent
44 access to the model itself such that an actor can utilize AI model capabilities in an unauthorized manner.

1 Appendix B. Example Safeguards Against the Misuse of Foundation Models

2 The available safeguards against the misuse of foundation models continue to evolve and
3 expand. Current literature suggests that safeguards’ effectiveness can vary widely, and reliance
4 on them to prevent misuse should be based on evidence of the safeguards’ efficacy, particularly
5 as some may incur tradeoffs in areas like cost, transparency, research, and user privacy. The
6 following table provides a non-exhaustive survey of several categories of safeguards that
7 organizations can consider:

8 **Table 1. Example Safeguards.**

Safeguard	Possible Implementation Methods
Improve the model’s training.	<ul style="list-style-type: none"> Filter training data to exclude examples that could result in capabilities that increase the likelihood of misuse, such as biological sequence data, or known CSAM/NCII images.⁴² Employ training techniques that reduce the model’s propensity to produce harmful output such as by reducing its knowledge of harmful information, reducing its capabilities of concern, or causing the model to refuse harmful requests.⁴³ Train models with approaches that make it more difficult for subsequent fine-tuning to remove safeguards.⁴⁴
Detect and block attempted misuse.	<ul style="list-style-type: none"> Add additional infrastructure around the base model to monitor for and detect behavior which may constitute misuse—such as with algorithmic classifiers for misuse. This infrastructure should be supported, as appropriate, by human review.⁴⁵ Once detected, block, modify, or otherwise limit unsafe queries and responses, and place limitations on users and organizations engaging in misuse or attempting to circumvent safeguards.
Limit access to the model’s capabilities.	<ul style="list-style-type: none"> Monitor the model’s risk for misuse with a limited audience and gradually expand access to the model to wider audiences over time.⁴⁶ Limit the ability to interact with the model to contexts and users where the misuse risks are lower.⁴⁷ Reduce access to the model reactively when misuse is detected to limit further misuse, such as by rolling back a model to a previous version, or discontinuing availability of a model, if a model displays significant misuse risk while it is in production. Apply limitations on using the model’s capabilities broadly to all interactions with the model or target limitations to particular capabilities and features at particular misuse risks, such as capabilities of concern that also have narrow beneficial purposes. For models that pose more severe risk, consider restricting the model only to those directly involved in its development and evaluation. Internal protections can reduce the risk that an organization’s own employees or contractors are able to misuse models; these protections can include logging employee interactions with models, restricting model access to a subset of employees, requiring multiple employees to access a model together, and ensuring that all employee interactions with the model are accompanied by the model’s other relevant safeguards.
Ensure the level of access to the model’s weights is appropriate.⁴⁸	<ul style="list-style-type: none"> Consider when it is appropriate to make the model’s weights widely available, such as available for download by the public. Once a model’s weights are made widely available, options to roll back or prevent its further sharing and modification are severely limited.⁴⁹ Consider when it is appropriate to allow the foundation model to be fine-tuned via API, which also can reduce the availability of safeguards, such as by letting users train the model on data related to dangerous tasks or reduce how often the model refuses dangerous requests.⁵⁰ Consider limiting internal access to the models’ weights within an organization.
Stop development if a model displays significant misuse risk.	<ul style="list-style-type: none"> If other practically available safeguards are not sufficient to protect a model from misuse—including considering the risk of theft or internal abuse—then it may be appropriate to make significant changes to the development plan, or else not develop the model at all.

1

-
- ¹ Anwar, U. et al., (2024) Foundational Challenges in Assuring Alignment and Safety of Large Language Models. *arXiv*. <https://arxiv.org/pdf/2404.09932>. For information on the rate of progress of foundation models, see Maslej, N. et al., (2024) The AI Index 2024 Annual Report. *AI Index Steering Committee, Institute for Human-Centered AI, Stanford University*. https://aiindex.stanford.edu/wp-content/uploads/2024/05/HAI_AI-Index-Report-2024.pdf; Ho, A. et al., (2024) Algorithmic progress in language models. *arXiv*. <https://arxiv.org/abs/2403.05812>; and Sevilla, J. et al., (2024) Training Compute of Frontier AI Models Grows by 4-5x per Year. *Epoch AI*. <https://epochai.org/blog/training-compute-of-frontier-ai-models-grows-by-4-5x-per-year>.
- ² Vogel, M. et al., (2024) Findings & Recommendations: AI Safety. *National AI Advisory Committee*. https://ai.gov/wp-content/uploads/2024/06/FINDINGS-RECOMMENDATIONS_AI-Safety.pdf.
- ³ AI Risk Management Framework. *National Institute of Standards and Technology*. <https://www.nist.gov/itl/ai-risk-management-framework>.
- ⁴ Srikumar, M. et al., (2024) Risk Mitigation Strategies for the Open Foundation Model Value Chain. *Partnership on AI*. <https://partnershiponai.org/resource/risk-mitigation-strategies-for-the-open-foundation-model-value-chain/>.
- ⁵ Executive Order 14110 on the Safe, Secure, and Trustworthy Development and Use of Artificial Intelligence (2023) *The White House*. <https://www.whitehouse.gov/briefing-room/presidential-actions/2023/10/30/executive-order-on-the-safe-secure-and-trustworthy-development-and-use-of-artificial-intelligence/>.
- ⁶ For instance, many smaller, domain-specific models, such as some generative image models, can also be misused in harmful ways. This document also does not cover risks from accidental AI harms to public safety.
- ⁷ Blueprint for an AI Bill of Rights (2022) *The White House*. <https://www.whitehouse.gov/ostp/ai-bill-of-rights/>.
- ⁸ AI Risk Management Framework: Generative Artificial Intelligence Profile. *National Institute of Standards and Technology*. <https://airc.nist.gov/docs/NIST.AI.600-1.GenAI-Profile.ipd.pdf>.
- ⁹ Bommasani, R. et al., (2021) On the Opportunities and Risks of Foundation Models. *arXiv*. <https://arxiv.org/abs/2108.07258>.
- ¹⁰ Heim, L. et al., (2024) Governing Through the Cloud: The Intermediary Role of Compute Providers in AI Regulation. *arXiv*. <https://arxiv.org/pdf/2403.08501>
- ¹¹ Cen, S. et al., (2024) AI Supply Chains. *SSRN*. https://papers.ssrn.com/sol3/papers.cfm?abstract_id=4789403; and Constanza-Chock, S. et al., (2022) Who Audits the Auditors? Recommendations from a field scan of algorithmic auditing ecosystem. *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*. <https://dl.acm.org/doi/abs/10.1145/3531146.3533213>.
- ¹² Mulani, N. et al., (2023) Proposing a Foundation Model Information-Sharing Regime for the UK. *Centre for the Governance of AI*. <https://www.governance.ai/post/proposing-a-foundation-model-information-sharing-regime-for-the-uk>; and Bommasani, R. et al. (2023) The Foundation Model Transparency Index. *arXiv*. <https://arxiv.org/abs/2310.12941>.
- ¹³ Anwar, U. et al., (2024) Foundational Challenges in Assuring Alignment and Safety of Large Language Models. *arXiv*. <https://arxiv.org/pdf/2404.09932>.
- ¹⁴ Hoffmann, J. et al., (2022) Training Compute-Optimal Large Language Models. *Proceedings of the 2022 Conference on Neural Information Processing Systems*. <https://dl.acm.org/doi/10.5555/3600270.3602446>.
- ¹⁵ Besiroglu, T. et al., (2024) Chinchilla Scaling: A replication attempt. *arXiv*. <https://arxiv.org/pdf/2404.10102>.
- ¹⁶ Terwilliger T. et al., (2023) AlphaFold predictions are valuable hypotheses and accelerate but do not replace experimental structure determination. *Nature*. <https://www.nature.com/articles/s41592-023-02087-4>.
- ¹⁷ S. et al., (2024) Position: On the Societal Impact of Open Foundation Models. *Proceedings of the 2024 International Conference on Machine Learning*. <https://openreview.net/pdf?id=jRX6yCxFhx>.
- ¹⁸ Anthropic (2023) Responsible Scaling Policy. *Anthropic*. <https://www-cdn.anthropic.com/1adf000c8f675958c2ee23805d91aaade1cd4613/responsible-scaling-policy.pdf>; OpenAI (2023) Preparedness Framework (Beta). *OpenAI*. <https://cdn.openai.com/openai-preparedness-framework-beta.pdf>; DeepMind (2024) Frontier Safety Framework. *DeepMind*. <https://storage.googleapis.com/deepmind->

-
- media/DeepMind.com/Blog/introducing-the-frontier-safety-framework/fsf-technical-report.pdf. Magic (2024) AGI Readiness Policy. *Magic*. <https://magic.dev/agi-readiness-policy>.
- ¹⁹ Kapoor, S. et al., (2024) Position: On the Societal Impact of Open Foundation Models. *Proceedings of the 2024 International Conference on Machine Learning*. <https://openreview.net/pdf?id=jRX6yCxFhx>; and Mouton, C. et al., (2024) The Operational Risks of AI in Large-Scale Biological Attacks. *RAND*. https://www.rand.org/pubs/research_reports/RRA2977-2.html.
- ²⁰ Koessler, L. et al., (2024) Risk thresholds for frontier AI. *arXiv*. <https://arxiv.org/pdf/2406.14713>.
- ²¹ Nevo, S., et al., (2024) Securing AI Model Weights. *RAND*. https://www.rand.org/pubs/research_reports/RRA2849-1.html.
- ²² Carlini, N. et al., (2024) Stealing Part of a Production Language Model. *arXiv*, <https://arxiv.org/pdf/2403.06634>.
- ²³ Barrett, A. et al, (2024) Benchmark Early and Red Team Often: A Framework for Assessing and Managing Dual-Use Hazards of AI Foundation Models. *Center for Long-Term Cybersecurity, UC Berkeley*. <https://cltc.berkeley.edu/wp-content/uploads/2024/05/Dual-Use-Benchmark-Early-Red-Team-Often.pdf>; Shao, M. et al., (2024) NYU CTF Dataset: A Scalable Open-Source Benchmark Dataset for Evaluating LLMs in Offensive Security. *arXiv*. <https://arxiv.org/pdf/2406.05590v1>; and Li, N. et al., (2024) The WMDP Benchmark: Measuring and Reducing Malicious Use With Unlearning. *arXiv*. <https://arxiv.org/pdf/2403.03218>.
- ²⁴ Casper, S. et al., (2024) Black-Box Access is Insufficient for Rigorous AI Audits. *arXiv*. <https://arxiv.org/pdf/2401.14446>.
- ²⁵ Bran, A. et al., (2024) Augmenting large language models with chemistry tools. *Nature*. <https://www.nature.com/articles/s42256-024-00832-8>; Davidson, T. et al., (2023) AI capabilities can be significantly improved without expensive retraining. *arXiv*. <https://arxiv.org/pdf/2312.07413>; and Measuring the impact of post-training enhancements. *METR*. <https://metr.github.io/autonomy-evals-guide/elicitation-gap/>.
- ²⁶ Jin, H. et al., (2024) JailbreakZoo: Survey, Landscapes, and Horizons in Jailbreaking Large Language and Vision-Language Models. *arXiv*. <https://arxiv.org/pdf/2407.01599>.
- ²⁷ Overview of Meta AI safety policies prepared for the UK AI Safety Summit. *Meta*. <https://transparency.meta.com/en-gb/policies/ai-safety-policies-for-safety-summit>; Microsoft’s AI Safety Policies. *Microsoft*. <https://query.prod.cms.rt.microsoft.com/cms/api/am/binary/RW1dObQ>; AI Safety Summit: An update on our approach to safety and responsibility. *Google DeepMind*. <https://deepmind.google/public-policy/ai-summit-policies/>; AI Safety Summit – Enhancing Frontier AI Safety. *AWS Amazon*. <https://aws.amazon.com/uki/cloud-services/uk-gov-ai-safety-summit/>; Our policy on frontier safety. *Inflection*. <https://inflection.ai/frontier-safety/>; and *OpenAI*. <https://openai.com/index/openai-safety-update/>.
- ²⁸ O’Brien, J. et al., Deployment Corrections: An incident response framework for frontier AI models. *Institute for AI Policy and Strategy*. <https://arxiv.org/pdf/2310.00328>.
- ²⁹ Longpre, S. et al., (2024) A Safe Harbor for AI Evaluation and Red-Teaming. *arXiv*. <https://arxiv.org/pdf/2403.04893>.
- ³⁰ Bucknall, B. et al., (2023). Structured Access for Third-Party Research on Frontier AI Models. *Center for the Governance of AI*. <https://www.governance.ai/research-paper/structured-access-for-third-party-research-on-frontier-ai-models>.
- ³¹ Wan, S. et al., (2024) Bridging the Gap: A Study of AI-based Vulnerability Management between Industry and Academia. *arXiv*. <https://arxiv.org/abs/2405.02435>.
- ³² Phuong, M. et al., (2024) Evaluating Frontier Models for Dangerous Capabilities. *arXiv*. <https://arxiv.org/pdf/2403.13793>; Anthropic (2024) Claude 3.5 Sonnet Model Card Addendum. *Anthropic*. https://www-cdn.anthropic.com/fed9cc193a14b84131812372d8d5857f8f304c52/Model_Card_Claude_3_Addendum.pdf; and OpenAI (2023) GPT-4 System Card. *OpenAI*. <https://cdn.openai.com/papers/gpt-4-system-card.pdf>
- ³³ Bommasani, R. et al., (2024) The Foundation Model Transparency Index v1.1: May 2024. *arXiv*. <https://arxiv.org/pdf/2407.12929>.
- ³⁴ Longpre, S. et al., (2023). The Data Provenance Initiative: A Large Scale Audit of Dataset Licensing & Attribution in AI. *arXiv*. <https://arxiv.org/abs/2310.16787>.

-
- ³⁵ Voluntary AI Commitments. *White House*. <https://www.whitehouse.gov/wp-content/uploads/2023/09/Voluntary-AI-Commitments-September-2023.pdf>; Bommasani, R. et al., (2024) The Foundation Model Transparency Index v1.1: May 2024. *arXiv*. <https://arxiv.org/pdf/2407.12929>.
- ³⁶ McGregor, S. (2020) Preventing Repeated Real World AI Failures by Cataloging Incidents: The AI Incident Database. *arXiv*. <https://arxiv.org/pdf/2011.08512>.
- ³⁷ Turri, V. et al., (2023) Why We Need to Know More: Exploring the State of AI Incident Documentation Practices. *AAAI/ACM Conference on AI Ethics and Society*. <https://doi.org/10.1145/3600211.3604700>.
- ³⁸ The term “artificial intelligence” or “AI” has the meaning set forth in 15 U.S.C. 9401(3).
- ³⁹ Executive Order 14110.
- ⁴⁰ Executive Order 14110.
- ⁴¹ Executive Order 14110.
- ⁴² Thiel, D., et al., (2023) Identifying and Eliminating CSAM in Generative ML Training Data and Models. *Stanford Digital Repository*. <https://purl.stanford.edu/kh752sm9123>.
- ⁴³ Lynch, A. et al., (2024) Eight Methods to Evaluate Robust Unlearning in LLMs. *arXiv*. <https://arxiv.org/pdf/2402.16835>.
- ⁴⁴ Henderson, P. et al., (2022) Self-Destructing Models: Increasing the Costs of Harmful Dual Uses of Foundation Models. *arXiv*. <https://arxiv.org/pdf/2211.14946>.
- ⁴⁵ Bommasani, R. (2023) Ecosystem graphs: The social footprint of foundation models. *arXiv*, abs/2303.15772.
- ⁴⁶ Solaiman, I. (2023) The Gradient of Generative AI Release: Methods and Considerations. *arXiv*. <https://arxiv.org/pdf/2302.04844>.
- ⁴⁷ Seger, E. et al., (2024) Open-Sourcing Highly Capable Foundation Models. https://cdn.governance.ai/Open-Sourcing_Highly_Capable_Foundation_Models_2023_GovAI.pdf
- ⁴⁸ When a model’s weights are available to a threat actor, a range of other safeguards become less effective, particularly those that are implemented at the application level, such as limiting who can use the model and detecting when it is misused. Limiting access to the model weights should be weighed against the potential benefits of access, such as for innovation and research, including research into safety. Limiting access to model weights is also only effective if organizations can prevent model theft.
- ⁴⁹ Kapoor, S. et al., (2024) Position: On the Societal Impact of Open Foundation Models. *Proceedings of the 2024 International Conference on Machine Learning*. <https://openreview.net/pdf?id=jRX6yCxHx>.
- ⁵⁰ Qi, X. et al., (2024) Fine-tuning Aligned Language Models Compromises Safety, Even When Users Do Not Intend To! *arXiv*. <https://arxiv.org/pdf/2310.03693>.