

AI Risk Management Framework: Initial Draft

March 17, 2022

This initial draft of the Artificial Intelligence Risk Management Framework (AI RMF, or Framework) builds on the concept paper released in December 2021 and incorporates the feedback received. The AI RMF is intended for voluntary use in addressing risks in the design, development, use, and evaluation of AI products, services, and systems.

AI research and deployment is evolving rapidly. For that reason, the AI RMF and its companion documents will evolve over time. When AI RMF 1.0 is issued in January 2023, NIST, working with stakeholders, intends to have built out the remaining sections to reflect new knowledge, awareness, and practices.

Part I of the AI RMF sets the stage for why the AI RMF is important and explains its intended use and audience. Part II includes the AI RMF Core and Profiles. Part III includes a companion Practice Guide to assist in adopting the AI RMF.

That Practice Guide which will be released for comment includes additional examples and practices that can assist in using the AI RMF. The Guide will be part of a NIST AI Resource Center that is being established.

NIST welcomes feedback on this initial draft and the related Practice Guide to inform further development of the AI RMF. Comments may be provided at a [workshop on March 29-31, 2022](#), and also are strongly encouraged to be shared via email. NIST will produce a second draft for comment, as well as host a third workshop, before publishing AI RMF 1.0 in January 2023.

Please send comments on this initial draft to AIframework@nist.gov by April 29, 2022.



Comments are especially requested on:

1. Whether the AI RMF appropriately covers and addresses AI risks, including with the right level of specificity for various use cases.
2. Whether the AI RMF is flexible enough to serve as a continuing resource considering evolving technology and standards landscape.
3. Whether the AI RMF enables decisions about how an organization can increase understanding of, communication about, and efforts to manage AI risks.
4. Whether the functions, categories, and subcategories are complete, appropriate, and clearly stated.
5. Whether the AI RMF is in alignment with or leverages other frameworks and standards such as those developed or being developed by IEEE or ISO/IEC SC42.
6. Whether the AI RMF is in alignment with existing practices, and broader risk management practices.
7. What might be missing from the AI RMF.
8. Whether the soon to be published draft companion document citing AI risk management practices is useful as a complementary resource and what practices or standards should be added.
9. Others?

Note: This first draft does not include Implementation Tiers as considered in the concept paper. Implementation Tiers may be added later if stakeholders consider them to be a helpful feature in the AI RMF. Comments are welcome.

Table of Contents

Part 1: Motivation

1	OVERVIEW	1
2	SCOPE	2
3	AUDIENCE	3
4	FRAMING RISK	5
4.1	Understanding Risk and Adverse Impacts	5
4.2	Challenges for AI Risk Management	6
5	AI RISKS AND TRUSTWORTHINESS	7
5.1	Technical Characteristics	8
5.1.1	<i>Accuracy</i>	9
5.1.2	<i>Reliability</i>	9
5.1.3	<i>Robustness</i>	10
5.1.4	<i>Resilience or ML Security</i>	10
5.2	Socio-Technical Characteristics	10
5.2.1	<i>Explainability</i>	11
5.2.2	<i>Interpretability</i>	11
5.2.3	<i>Privacy</i>	11
5.2.4	<i>Safety</i>	12
5.2.5	<i>Managing Bias</i>	12
5.3	Guiding Principles	12
5.3.1	<i>Fairness</i>	13
5.3.2	<i>Accountability</i>	13
5.3.3	<i>Transparency</i>	13

Part 2: Core and Profiles

6	AI RMF CORE	14
6.1	Map	15
6.2	Measure	16
6.3	<i>Manage</i>	17
6.4	Govern	18
7	AI RMF PROFILES	20
8	EFFECTIVENESS OF THE AI RMF	20

Part 3: Practical Guide

9	PRACTICE GUIDE	20
----------	-----------------------	-----------

1 AI Risk Management Framework: Initial Draft -

2 Part 1: Motivation

3 1 Overview

4 Remarkable surges in artificial intelligence (AI) capabilities have led to a wide range of
5 innovations with the potential to benefit nearly all aspects of our society and economy –
6 everything from commerce and healthcare to transportation and cybersecurity. AI systems are
7 used for tasks such as informing and advising people and taking actions where they can have
8 beneficial impact, such as safety and housing.

9 AI systems sometimes do not operate as intended because they are making inferences from
10 patterns observed in data rather than a true understanding of what causes those patterns. Ensuring
11 that these inferences are helpful and not harmful in particular use cases – especially when
12 inferences are rapidly scaled and amplified – is fundamental to trustworthy AI. While answers to
13 the question of what makes an AI technology trustworthy differ, there are certain key
14 characteristics which support trustworthiness, including accuracy, explainability and
15 interpretability, privacy, reliability, robustness, safety, security (resilience) and mitigation of
16 harmful bias. There also are key guiding principles to take into account such as accountability,
17 fairness, and equity.

18 Cultivating trust and communication about how to understand and manage the risks of AI
19 systems will help create opportunities for innovation and realize the full potential of this
20 technology.

21 Many activities related to managing risk for AI
22 are common to managing risk for other types of
23 technology. An AI Risk Management Framework
24 (AI RMF, or Framework) can address challenges
25 unique to AI systems. This AI RMF is an initial
26 attempt to describe how the risks from AI-based
27 systems differ from other domains and to
28 encourage and equip many different stakeholders
29 in AI to address those risks purposefully.

It is important to note that the AI RMF is neither a checklist nor should be used in any way to certify an AI system. Likewise, using the AI RMF does not substitute for due diligence and judgment by organizations and individuals in deciding whether to design, develop, and deploy AI technologies – and if so, under what conditions.

30 This voluntary framework provides a flexible,
31 structured, and measurable process to address AI risks throughout the AI lifecycle, offering
32 guidance for the development and use of trustworthy and responsible AI. It is intended to
33 improve understanding of and to help organizations manage both enterprise and societal risks
34 related to the development, deployment, and use of AI systems. Adopting the AI RMF can assist
35 organizations, industries, and society to understand and determine their acceptable levels of risk.

1 In addition, it can be used to map compliance considerations beyond those addressed by this
2 framework, including existing regulations, laws, or other mandatory guidance.

3 Risks to any software or information-based system apply to AI; that includes important concerns
4 related to cybersecurity, privacy, safety, and infrastructure. This framework aims to fill the gaps
5 related specifically to AI. Rather than repeat information in other guidance, users of the AI RMF
6 are encouraged to address those non-AI specific issues via guidance already available.

7 Part 1 of this framework establishes the context for the
8 AI risk management process. Part 2 provides guidance
9 on outcomes and activities to carry out that process to
10 maximize the benefits and minimize the risks of AI.

11 Part 3 [yet to be developed] assists in using the AI
12 RMF and offers sample practices to be considered in
13 carrying out this guidance, before, during, and after AI products, services, and systems are
14 developed and deployed.

For the purposes of the NIST AI RMF the term *artificial intelligence* refers to algorithmic processes that learn from data in an automated or semi-automated manner.

15 The Framework, and supporting resources, will be updated and improved based on evolving
16 technology and the standards landscape around the globe. In addition, as the AI RMF is put into
17 use, additional lessons will be learned that can inform future updates and additional resources.

18 NIST’s development of the AI RMF in collaboration with the private and public sectors is
19 consistent with its broader AI efforts called for by the National AI Initiative Act of 2020 (P.L.
20 116-283), the National Security Commission on Artificial Intelligence recommendations, and the
21 Plan for Federal Engagement in AI Standards and Related Tools. Engagement with the broad AI
22 community during this Framework’s development also informs AI research and development
23 and evaluation by NIST and others.

24 2 Scope

25 The NIST AI RMF offers a process for managing risks related to AI systems across a wide
26 spectrum of types, applications, and maturity. This framework is organized and intended to be
27 understood and used by individuals and organizations, regardless of sector, size, or level of
28 familiarity with a specific type of technology. Ultimately, it will be offered in multiple formats,
29 including online versions, to provide maximum flexibility.

30 The AI RMF serves as a part of a broader NIST resource center containing documents,
31 taxonomy, suggested toolkits, datasets, code, and other forms of technical guidance related to the
32 development and implementation of trustworthy AI. Resources will include a knowledge base of
33 terminology related to trustworthy and responsible AI and how those terms are used by different
34 stakeholders.

35 The AI RMF is not a checklist nor a compliance mechanism to be used in isolation. It should be
36 integrated within the organization developing and using AI and be incorporated into enterprise

1 risk management; doing so ensures that AI will be treated along with other critical risks, yielding
2 a more integrated outcome and resulting in organizational efficiencies.

3 Attributes of the AI RMF

4 The AI RMF strives to:

- 5 1. Be risk-based, resource efficient, and voluntary.
- 6 2. Be consensus-driven and developed and regularly updated through an open, transparent process.
7 All stakeholders should have the opportunity to contribute to the AI RMF’s development.
- 8 3. Use clear and plain language that is understandable by a broad audience, including senior
9 executives, government officials, non-governmental organization leadership, and those who are
10 not AI professionals – while still of sufficient technical depth to be useful to practitioners. The AI
11 RMF should allow for communication of AI risks across an organization, between organizations,
12 with customers, and to the public at large.
- 13 4. Provide common language and understanding to manage AI risks. The AI RMF should offer
14 taxonomy, terminology, definitions, metrics, and characterizations for AI risk.
- 15 5. Be easily usable and mesh with other aspects of risk management. Use of the Framework should
16 be intuitive and readily adaptable as part of an organization’s broader risk management strategy
17 and processes. It should be consistent or aligned with other approaches to managing AI risks.
- 18 6. Be useful to a wide range of perspectives, sectors, and technology domains. The AI RMF should
19 be both technology agnostic and applicable to context-specific use cases.
- 20 7. Be outcome-focused and non-prescriptive. The Framework should provide a catalog of outcomes
21 and approaches rather than prescribe one-size-fits-all requirements.
- 22 8. Take advantage of and foster greater awareness of existing standards, guidelines, best practices,
23 methodologies, and tools for managing AI risks – as well as illustrate the need for additional,
24 improved resources.
- 25 9. Be law- and regulation-agnostic. The Framework should support organizations’ abilities to
26 operate under applicable domestic and international legal or regulatory regimes.
- 27 10. Be a living document. The AI RMF should be readily updated as technology, understanding, and
28 approaches to AI trustworthiness and uses of AI change and as stakeholders learn from
29 implementing AI risk management generally and this framework in particular.

30 3 Audience

31 AI risk management is a complex and relatively new area, and the list of individuals, groups,
32 communities, and organizations that can be affected by AI technologies is extensive. Identifying
33 and managing AI risks and impacts – both positive and adverse – requires a broad set of
34 perspectives and stakeholders.

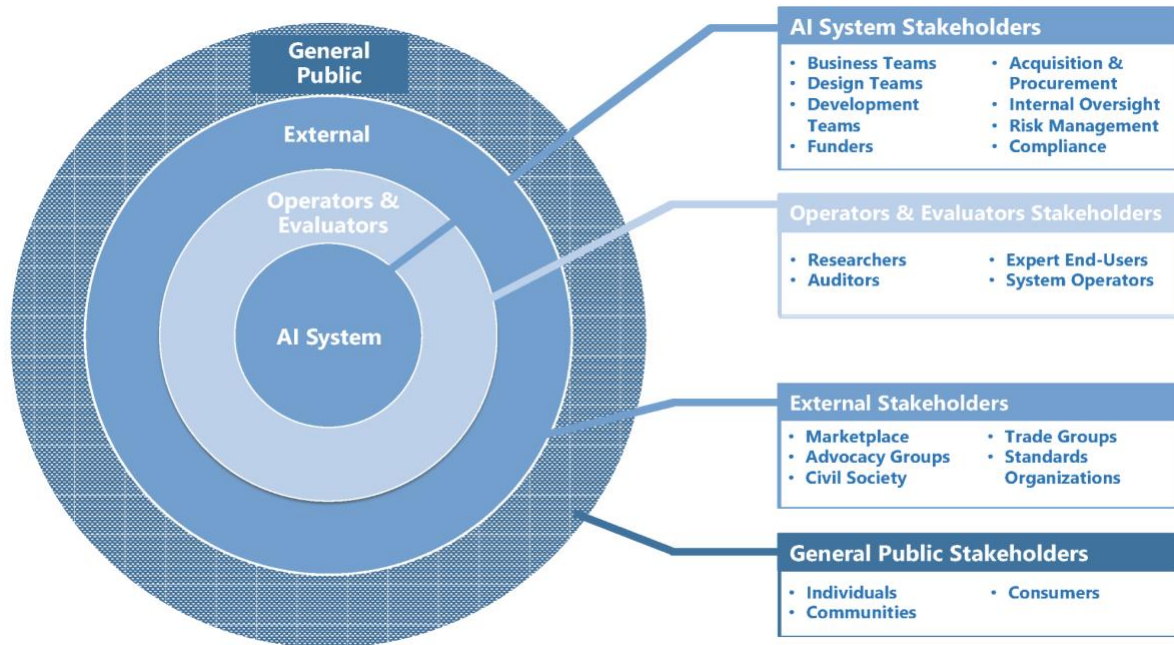


Figure 1: Key stakeholder groups associated with the AI RMF.

As Figure 1 illustrates, NIST has identified four stakeholder groups as intended audiences of this Framework: AI system stakeholders, operators and evaluators, external stakeholders, and the general public. Ideally, members of all stakeholder groups would be involved or represented in the risk management process, including those individuals and community representatives that may be affected by the use of AI technologies.

AI system stakeholders are those who have the most control and responsibility over the design, development, deployment, and acquisition of AI systems, and the implementation of AI risk management practices. This group comprises the primary adopters of the AI RMF. They may include individuals or teams within or among organizations with responsibilities to commission, fund, procure, develop, or deploy an AI system: **business teams**, design and development teams, internal risk management teams, and compliance teams. Small to medium-sized organizations face different challenges in implementing the AI RMF than large organizations.

Operators and evaluators provide monitoring and formal/informal test, evaluation, validation, and verification (TEVV) of system performance, relative to both technical and socio-technical requirements. These stakeholders, which include organizations which operate or **employ AI systems**, **use the output for decisions** or to evaluate their performance. This group can include users who interpret or **incorporate the output of AI systems in settings with a high potential for adverse impacts**. They might include academic, public, and private sector researchers; professional evaluators and auditors; system operators; and expert end users.

External stakeholders provide formal and/or **quasi-formal norms** or guidance for specifying and addressing AI risks. External to the primary adopters of the AI RMF, they can include trade

1 groups, standards developing organizations, advocacy groups, and civil society organizations.
2 Their actions can designate boundaries for operation (technical or legal) and **balance societal**
3 **values and priorities** related to civil liberties and rights, the economy, and security.

4 The **general public** is most likely to directly experience positive and adverse impacts of AI
5 technologies. They **may provide the motivation for actions taken by the other stakeholders** and
6 can include individuals, communities, and consumers in the context where an AI system is
7 developed or deployed.

8 **4 Framing Risk**

9 AI systems hold the potential to advance our quality of life and lead to new services, support,
10 and efficiencies for people, organizations, markets, and society. Identifying, mitigating, and
11 **minimizing risks and potential harms** associated with AI technologies are essential steps towards
12 the acceptance and widespread use of AI technologies. A risk management framework should
13 provide a structured, yet flexible, approach for managing enterprise and societal risk resulting
14 from the incorporation of AI systems into products, processes, organizations, systems, and
15 societies. Organizations managing an enterprise's AI risk also should be mindful of larger
16 societal AI considerations and risks. If a risk management framework can help to effectively
17 address and manage AI risk and adverse impacts, it can lead to more trustworthy AI systems.

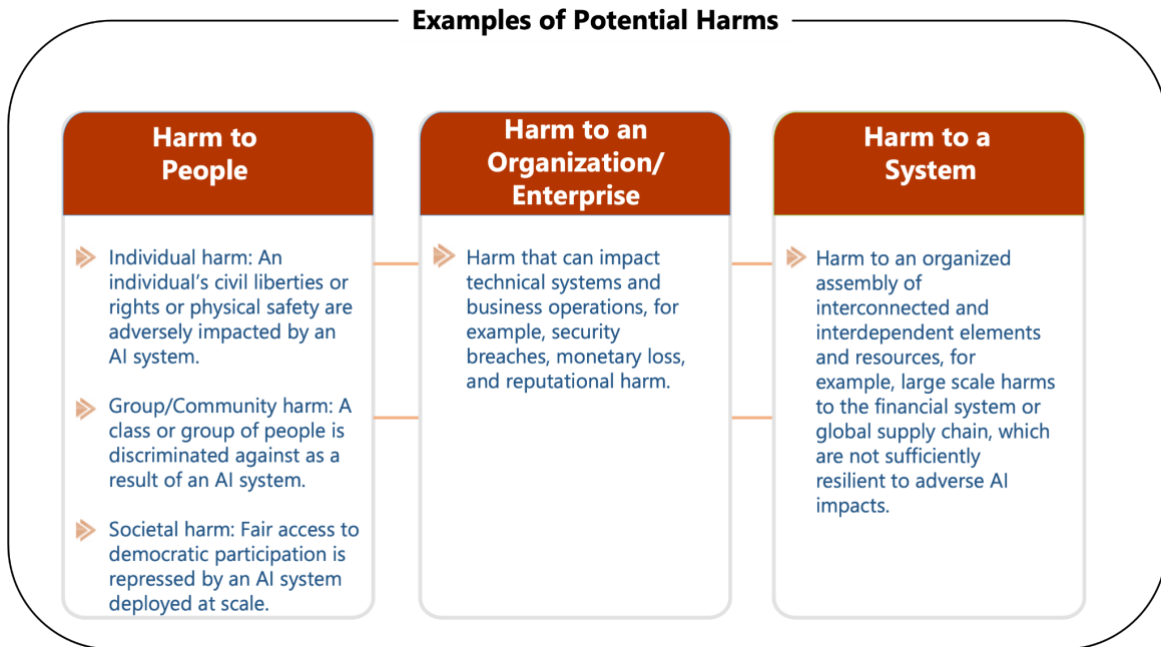
18 **4.1 Understanding Risk and Adverse Impacts**

19 Risk is a measure of the extent to which an entity is negatively influenced by a potential
20 circumstance or event. Typically, risk is a function of 1) the adverse impacts that could arise if
21 the circumstance or event occurs; and 2) the likelihood of occurrence. Entities can be individuals,
22 groups, or communities as well as systems, processes, or organizations.

23 The impact of AI systems can be positive, negative, or both and can address, create, or result in
24 opportunities **or threats**. According to the International Organization for Standardization (Guide
25 73:2009; IEC/ISO 31010), **certain risks can be positive**. While risk management processes
26 address adverse impacts, this framework intends to offer approaches to minimize anticipated
27 negative impacts of AI systems *and* identify opportunities to maximize positive impacts.
28 Additionally, this framework is designed to be responsive to new risks as they emerge rather than
29 enumerating all known risks in advance. This flexibility is particularly important where **impacts**
30 **are not easily foreseeable**, and applications are evolving rapidly. While AI benefits and some AI
31 risks are well-known, the AI community is only beginning to understand and classify incidents
32 and scenarios that result in harm. Figure 2 provides examples of potential harms from AI
33 systems.

34 Risk management can also drive AI developers and users to understand and account for the
35 inherent uncertainties and inaccuracy of their models and systems, which in turn can increase the

1 overall performance and trustworthiness of those models. Managing risk and adverse impacts
2 contributes to building trustworthy AI technologies and applications



3

4

Figure 2: Examples of potential harms from AI systems.

5 4.2 Challenges for AI Risk Management

6 4.2.1 Risk Measurement

7 AI risks and impacts that are not well-defined or adequately understood are difficult to measure
8 quantitatively or qualitatively. The presence of third-party data or systems may also complicate
9 risk measurement. Those attempting to measure the adverse impact on a population may not be
10 aware that certain demographics may experience harm differently than others.

11 AI risks can have a temporal dimension. Measuring risk at an earlier stage in the AI lifecycle
12 may yield different results than measuring risk at a later stage. Some AI risks may have a low
13 probability in the short term but have a high likelihood for adverse impacts. Other risks may be
14 latent at present but may increase in the long term as AI systems evolve.

15 Furthermore, inscrutable AI systems can complicate the measurement of risk. Inscrutability can
16 be a result of the opaque nature of AI technologies (lack of explainability or interpretability),
17 lack of transparency or documentation in AI system development or deployment, or inherent
18 uncertainties in AI systems.

19 4.2.2 Risk Thresholds

20 Thresholds refer to the values used to establish concrete decision points and operational limits
21 that trigger a response, action, or escalation. AI risk thresholds (sometimes referred to as Key
22 Risk Indicators) can involve both technical factors (such as error rates for determining bias) and
23 human values (such as social or legal norms for appropriate levels of transparency). These

1 factors and values can establish levels of risk (e.g., low, medium, or high) based on broad
2 categories of adverse impacts or harms.

3 Thresholds and values can also determine where AI systems present unacceptable risks to certain
4 organizations, systems, social domains, or demographics. In these cases, the question is not how
5 to better manage risk of AI, but whether an AI system should be designed, developed, or
6 deployed at all.

7 The AI RMF does not prescribe risk thresholds or values. Risk tolerance – the level of risk or
8 degree of uncertainty that is acceptable to organizations or society – is context and use case-
9 specific. Therefore, risk thresholds should be set through policies and norms that can be
10 established by AI system owners, organizations, industries, communities, or regulators (who
11 often are acting on behalf of individuals or societies). Risk thresholds and values are likely to
12 change and adapt over time as policies and norms change or evolve. In addition, different
13 organizations may have different risk thresholds (or tolerance) due to varying organizational
14 priorities and resource considerations. Even within a single organization there can be a balancing
15 of priorities and tradeoffs between technical factors and human values. Emerging knowledge and
16 methods for better informing these decisions are being developed and debated by business,
17 governments, academia, and civil society. To the extent that challenges for specifying risk
18 thresholds or determining values remain unresolved, there may be contexts where a risk
19 management framework is not yet readily applicable for mitigating AI risks and adverse impacts.
20 The AI RMF provides the opportunity for organizations to specifically define their risk
21 thresholds and then to manage those risks within their tolerances.

22 4.2.3 Organizational Integration

23 The AI RMF is not a checklist nor a compliance mechanism to be used in isolation. It should be
24 integrated within the organization developing and using AI technologies and be incorporated into
25 enterprise risk management; doing so ensures that AI will be treated along with other critical
26 risks, yielding a more integrated outcome and resulting in organizational efficiencies.

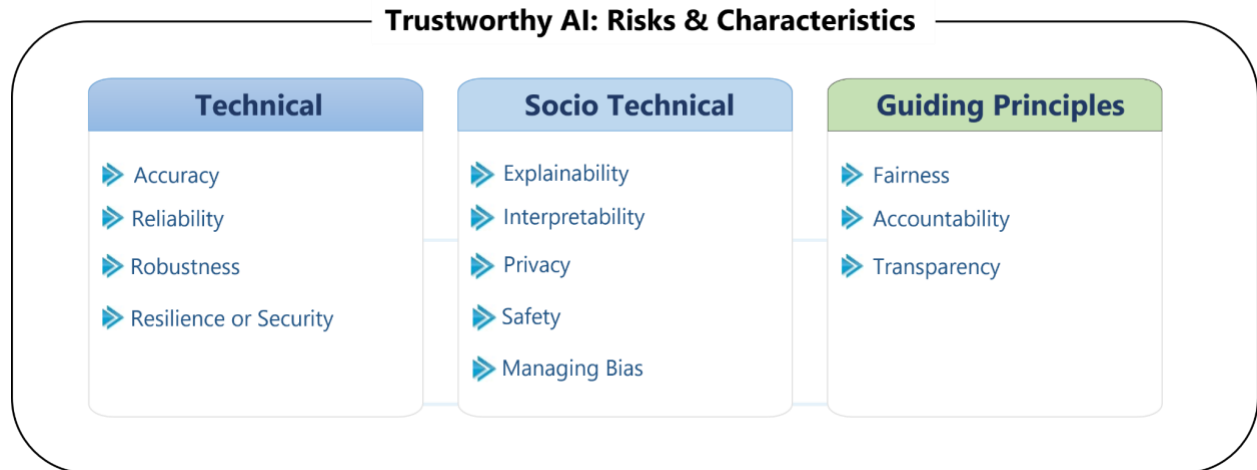
27 Organizations need to establish and maintain the appropriate accountability mechanisms, roles
28 and responsibilities, culture, and incentive structures for risk management to be effective. Use of
29 the AI RMF alone will not lead to these changes or provide the appropriate incentives. Effective
30 risk management needs organizational commitment at senior levels and may require significant
31 cultural change for an organization or industry.

32 Small to medium-sized organizations face different challenges in implementing the AI RMF than
33 large organizations.

34 5 AI Risks and Trustworthiness

35 The AI RMF uses a three-class taxonomy, depicted in Figure 3, to classify characteristics that
36 should be considered in comprehensive approaches for identifying and managing risk related to
37 AI systems: *technical characteristics, socio-technical characteristics, and guiding principles.*

1 This AI RMF taxonomy frames AI risk using characteristics that are aligned with trustworthy AI
 2 systems, in conjunction with contextual norms and values. Since AI trustworthiness and risk are
 3 inversely related, approaches which enhance trustworthiness can contribute to a reduction or
 4 attenuation of related risks. The AI RMF taxonomy articulates several key building blocks of
 5 trustworthy AI within each category, which are particularly suited to the examination of potential
 6 risk.



7
 8 **Figure 3:** AI Risks and Trustworthiness. The three-class taxonomy to classify characteristics that should
 9 be considered in comprehensive approaches for identifying and managing risk related to AI systems. The
 10 taxonomy articulates several key building blocks of trustworthy AI within each category, which are
 11 particularly suited to the examination of potential risk.

12 Figure 4 provides a mapping of the AI RMF taxonomy to the terminology used by the
 13 Organisation for Economic Co-operation and Development (OECD) in their Recommendation
 14 on AI, the European Union (EU) Artificial Intelligence Act, and United States Executive Order
 15 (EO) 13960.

16 5.1 Technical Characteristics

17 Technical characteristics in the AI RMF taxonomy refer to factors that are under the direct
 18 control of AI system designers and developers, and which may be measured using standard
 19 evaluation criteria. Technical characteristics include the tradeoff between convergent-
 20 discriminant validity (whether the data reflects what the user intends to measure and not other
 21 things) and statistical reliability (whether the data may be subject to high levels of statistical
 22 noise and measurement bias). Validity of AI, especially machine learning (ML) models, can be
 23 assessed using technical characteristics. Validity for deployed AI systems is often assessed with
 24 ongoing audits or monitoring that confirm that a system behaves as intended. It may be possible
 25 to utilize and automate explicit measures based on variations of standard statistical or ML
 26 techniques and specify thresholds in requirements. Data generated from experiments that are
 27 designed to evaluate system performance also fall into this category and might include tests of
 28 causal hypotheses and assessments of robustness to adversarial attack.

	Technical Design Characteristics	Socio-Technical Characteristics	Guiding Principles Contributing to Trustworthiness
AI RMF Taxonomy	<ul style="list-style-type: none"> • Accuracy • Reliability • Robustness • Resilience or ML Security 	<ul style="list-style-type: none"> • Explainability • Interpretability • Privacy • Safety • Managing Bias 	<ul style="list-style-type: none"> • Fairness • Accountability • Transparency
OECD AI Recommendation	<ul style="list-style-type: none"> • Robustness • Security 	<ul style="list-style-type: none"> • Safety • Explainability 	<ul style="list-style-type: none"> • Traceability to human values • Transparency and responsible disclosure • Accountability
EU AI Act	<ul style="list-style-type: none"> • Technical robustness 	<ul style="list-style-type: none"> • Safety • Privacy • Non-discrimination 	<ul style="list-style-type: none"> • Human agency and oversight • Data governance • Transparency • Diversity and fairness • Environmental and societal well-being • Accountability
EO 13960	<ul style="list-style-type: none"> • Purposeful and performance-driven • Accurate, reliable, and effective • Secure and resilient 	<ul style="list-style-type: none"> • Safe • Understandable by subject matter experts, users, and others, as appropriate 	<ul style="list-style-type: none"> • Lawful and respectful of our Nation's values • Responsible and traceable • Regularly monitored • Transparent • Accountable

Figure 4: Mapping of AI RMF taxonomy to AI policy documents.

The following technical characteristics lend themselves well to addressing AI risk: accuracy, reliability, robustness, and resilience (or ML security).

5.1.1 Accuracy

Accuracy indicates the degree to which the ML model is correctly capturing a relationship that exists within training data. Analogous to statistical conclusion validity, accuracy is examined via standard ML metrics (e.g., false positive and false negative rates, F1-score, precision, and recall), as well as assessment of model underfit or overfit (high testing errors irrespective of error rates in training). It is widely acknowledged that current ML methods cannot guarantee that the underlying model is capturing a causal relationship. Establishing internal (causal) validity in ML models is an active area of research. AI risk management processes should take into account the potential risks to the enterprise and society if the underlying causal relationship inferred by a model is not valid, calling into question decisions made on the basis of the model. Determining a threshold for accuracy that corresponds with acceptable risk is fundamental to AI risk management and highly context-dependent.

5.1.2 Reliability

Reliability indicates whether a model consistently generates the same results, within the bounds of acceptable statistical error. Techniques designed to mitigate overfitting (e.g., regularization) and to adequately conduct model selection in the face of the bias/variance tradeoff can increase model reliability. The definition of reliability is analogous to construct reliability in the social sciences, albeit without explicit reference to a theoretical construct. Reliability measures may give insight into the risks related to decontextualization, due to the common practice of reusing

1 ML datasets or models in ways that cause them to become **disconnected from the social contexts**
2 and time periods of their creation. As with accuracy, reliability provides an evaluation of the
3 validity of models, and thus can be a factor in determining thresholds for acceptable risk.

4 *5.1.3 Robustness*

5 Robustness is a measure of model sensitivity, indicating whether the model has minimum
6 sensitivity to variations in uncontrollable factors. A robust model will continue to function
7 despite the existence of faults in its components. The performance of the model may be
8 diminished or otherwise altered until the faults are corrected. Measures of robustness might
9 range from sensitivity of a model's outputs to small changes in its inputs, but might also include
10 error measurements on novel datasets. **Robustness contributes to sensitivity analysis** in the AI
11 risk management process.

12 *5.1.4 Resilience or ML Security*

13 A model that can withstand **adversarial attacks**, or more generally, unexpected changes in its
14 environment or use, may be said to be resilient or secure. This attribute has some relationship to
15 robustness except that it goes beyond the provenance of the data to encompass unexpected or
16 **adversarial use of the model or data. Other common ML security concerns relate to the**
17 **exfiltration of models, training data, or other intellectual property** through AI system endpoints.

18 **5.2 Socio-Technical Characteristics**

19 Socio-technical characteristics in the AI RMF taxonomy refer to how AI systems are used and
20 perceived in individual, group, and societal contexts. This includes mental representations of
21 models, whether the output provided is sufficient to evaluate compliance (transparency), whether
22 model operations can be **easily understood (explainability)**, whether they provide output that can
23 **be used to make a meaningful decision (interpretability)**, and whether the **outputs are aligned**
24 **with societal values. Socio-technical factors** are inextricably tied to human social and
25 organizational behavior, from the datasets used by ML processes and the decisions made by
26 those who build them, to the interactions with the humans who provide the insight and oversight
27 to make such systems actionable.

28 Unlike technical characteristics, **socio-technical characteristics** require significant human input
29 and cannot yet be measured through an automated process. Human judgment must be employed
30 when deciding on the specific metrics and the precise threshold values for these metrics. The
31 connection between human perceptions and interpretations, societal values, and enterprise and
32 societal risk is a key component of the kinds of cultural and organizational factors that will be
33 **necessary to properly manage AI risks. Indeed, input from a broad and diverse set of**
34 **stakeholders** is required throughout the AI lifecycle to ensure that risks arising in social contexts
35 are managed appropriately.

36 The following socio-technical characteristics have implications for addressing AI risk:
37 explainability, interpretability, privacy, safety, and managing bias.

1 5.2.1 Explainability

2 Explainability seeks to provide a programmatic, sometimes causal, description of how model
3 predictions are generated. Even given all the information required to make a model fully
4 transparent, a human must apply technical expertise if they want to understand how the model
5 works. Explainability refers to the user's perception of how the model works – such as what
6 output may be expected for a given input. Explanation techniques tend to summarize or visualize
7 model behavior or predictions for technical audiences. Explanations can be useful in promoting
8 human learning from machine learning, for addressing transparency requirements, or for
9 debugging issues with AI systems and training data. However, risks due to explainability may
10 arise for many reasons, including, for example, a lack of fidelity or consistency in explanation
11 methodologies, or if humans incorrectly infer a model's operation, or the model is not operating
12 as expected. Risk from lack of explainability may be managed by descriptions of how models
13 work to users' skill levels. Explainable systems can be more easily debugged and monitored, and
14 lend themselves to more thorough documentation, audit, and governance.

15 Explainability is related to transparency. Typically the more opaque a model is, the less it is
16 considered explainable. However, transparency does not guarantee explainability, especially if
17 the user lacks an understanding of ML technical principles.

18 5.2.2 Interpretability

19 Interpretability seeks to fill a meaning deficit. Although explainability and interpretability are
20 often used interchangeably, explainability refers to a representation of the mechanisms
21 underlying an algorithm's operation, whereas interpretability refers to the meaning of its output
22 in the context of its designed functional purpose. The underlying assumption is that perceptions
23 of risk stem from a lack of ability to make sense of, or contextualize, model output appropriately.
24 Model interpretability refers to the extent to which a user can determine adherence to this
25 function and the consequent implications of this output upon other consequential decisions for
26 that user. Interpretations are typically contextualized in terms of values and reflect simple,
27 categorical distinctions. For example, a society may value privacy and safety, but individuals
28 may have different determinations of safety thresholds. Risks to interpretability can often be
29 addressed by communicating the interpretation intended by model designers, although this
30 remains an open area of research. The prevalence of different interpretations can be readily
31 measured with psychometric instruments.

32 5.2.3 Privacy

33 Privacy refers generally to the norms and practices that help to safeguard values such as human
34 autonomy and dignity. These norms and practices typically address freedom from intrusion,
35 limiting observation, or individuals' control of facets of their identities (e.g., body, data,
36 reputation). Like safety and security, specific technical features of an AI system may promote
37 privacy, and assessors can identify how the processing of data could create privacy-related
38 problems. However, determinations of likelihood and severity of impact of these problems are
39 contextual and vary among cultures and individuals.

1 5.2.4 Safety

2 Safety as a concept is highly correlated with risk and generally denotes an absence (or
3 minimization) of failures or conditions that render a system dangerous. As AI systems interact
4 with humans more directly in factories and on the roads, for example, the safety of these systems
5 is a serious consideration for AI risk management. Safety is often – though not always –
6 considered through a legal lens. Practical approaches for AI safety often relate to rigorous
7 simulation and in-domain testing, real-time monitoring, and the ability to quickly shut down or
8 modify misbehaving systems.

9 5.2.5 Managing Bias

10 NIST has identified three major categories of bias in AI: systemic, computational, and human.
11 Managing bias in AI systems requires an approach that considers all three categories.
12 Bias exists in many forms, is omnipresent in society, and can become ingrained in the automated
13 systems that help make decisions about our lives. While bias is not always a negative
14 phenomenon, certain biases exhibited in AI models and systems can perpetuate and amplify
15 negative impacts on individuals, organizations, and society, and at a speed and scale far beyond
16 the traditional discriminatory practices that can result from implicit human or systemic biases.
17 Bias is tightly associated with the concepts of transparency and fairness in society. See NIST
18 publication “Towards a Standard for Identifying and Managing Bias in Artificial Intelligence.”

19 When managing risks in AI systems it is important to understand that the attributes
20 of the AI RMF risk taxonomy are interrelated. Highly secure but unfair systems,
21 accurate but opaque and uninterpretable systems, and inaccurate, but fair, secure,
22 privacy-protected, and transparent systems are all undesirable. It is possible for
23 trustworthy AI systems to achieve a high degree of risk control while retaining a
24 high level of performance quality. Achieving this difficult goal requires a
25 comprehensive approach to risk management, with tradeoffs among the technical
26 and socio-technical characteristics.

27 5.3 Guiding Principles

28 Guiding principles in the AI RMF taxonomy refer to broader societal norms and values that
29 indicate societal priorities. While there is no objective standard for ethical values, as they are
30 grounded in the norms and legal expectations of specific societies or cultures, it is widely agreed
31 that AI technologies should be developed and deployed in ways that meet contextual norms and
32 ethical values. When specified as policy, guiding principles can enable AI stakeholders to form
33 actionable, low-level requirements. Some requirements will be translated into quantitative
34 measures of performance and effectiveness, while some may remain qualitative in nature.

35 Guiding principles that are relevant for AI risk include fairness, accountability, and transparency.
36 Fairness in AI systems includes concerns for equality and equity by addressing socio-technical
37 issues such as bias and discrimination. Individual human operators and their organizations
38 should be answerable and held accountable for the outcomes of AI systems, particularly adverse

1 impacts stemming from risks. Absent transparency, users are left to guess about these factors and
2 may make unwarranted and unreliable assumptions regarding model provenance. Transparency
3 is often necessary for actionable redress related to incorrect and adverse AI system outputs.

4 5.3.1 Fairness

5 Standards of fairness can be complex and difficult to define because perceptions of fairness
6 differ among cultures. For one type of fairness, process fairness, AI developers assume that ML
7 algorithms are inherently fair because the same procedure applies regardless of user. However,
8 this perception has eroded recently as awareness of biased algorithms and biased datasets has
9 increased. Fairness is increasingly related to the existence of a harmful system, i.e., even if
10 demographic parity and other fairness measures are satisfied, sometimes the harm of a system is
11 in its existence. While there are many technical definitions for fairness, determinations of
12 fairness are not generally just a technical exercise. Absence of harmful bias is a necessary
13 condition for fairness.

14 5.3.2 Accountability

15 Determinations of accountability in the AI context are related to expectations for the responsible
16 party in the event that a risky outcome is realized. Individual human operators and their
17 organizations should be answerable and held accountable for the outcomes of AI systems,
18 particularly adverse impacts stemming from risks. The relationship between risk and
19 accountability associated with AI and technological systems more broadly differs across cultural,
20 legal, sectoral, and societal contexts. Grounding organizational practices and governing
21 structures for harm reduction, like risk management, can help lead to more accountable systems.

22 5.3.3 Transparency

23 Transparency seeks to remedy a common information imbalance between AI system operators
24 and AI system consumers. Transparency reflects the extent to which information is available to a
25 user when interacting with an AI system. Its scope spans from design decisions and training data
26 to model training, the structure of the model, its intended use case, how and when deployment
27 decisions were made and by whom, etc. Absent transparency, users are left to guess about these
28 factors and may make unwarranted and unreliable assumptions regarding model provenance.
29 Transparency is often necessary for actionable redress related to incorrect and adverse AI system
30 outputs. A transparent system is not necessarily a fair, privacy-protective, secure, or robust
31 system. However, it is difficult to determine whether an opaque system possesses such
32 desiderata, and to do so over time as complex systems evolve.

Part 2: Core and Profiles

6 AI RMF Core

The AI RMF Core provides outcomes and actions that enable dialogue, understanding, and activities to manage AI risks. The Core is composed of three elements: functions, categories, and subcategories. As illustrated in Figure 5, functions organize AI risk management activities at their highest level to map, measure, manage, and govern AI risks. Within each function, categories and subcategories subdivide the function into specific outcomes and actions.

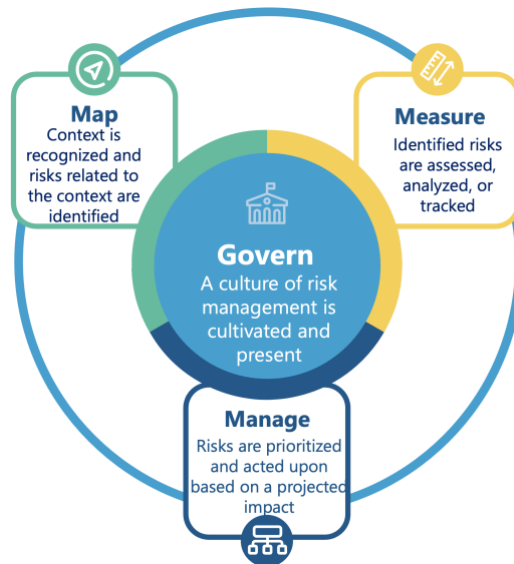


Figure 5: Functions organize AI risk management activities at their highest level to map, measure, manage, and govern AI risks. Governance is a cross-cutting function that is infused throughout and informs the other functions of the process.

Govern is a cross-cutting function that is infused throughout and informs the other functions of the process. Aspects of Govern, especially those related to compliance or evaluation, should be integrated into each of the other functions. Assuming a governance structure is in place, functions may be performed in any order across the AI lifecycle as deemed to add value by a user of the framework. In most cases, it will be more useful and effective to begin with Map before Measure and Manage. Regardless, the process should be iterative, with cross-referencing between functions as necessary. Similarly, there are categories and subcategories with elements that apply to multiple functions.

Technical and socio-technical characteristics and guiding principles of AI trustworthiness are essential considerations for each function. AI RMF core functions should be carried out in a way that reflects diverse and multidisciplinary perspectives, potentially including the views of stakeholders from outside the organization. Risk management should be performed throughout the AI system life cycle (Figure 6) to ensure it is continuous and timely.

1 On the following pages, Tables 1 through 4 provide the Framework Core listing.



2

3 **Figure 6:** Risk management should be performed throughout the AI system life cycle to ensure it is
 4 continuous and timely. Example activities for each stage of the AI lifecycle follow. *Pre-Design:* data
 5 collection, curation or selection, problem formulation, and identification of stakeholders. *Design &*
 6 *Development:* data analysis, data cleaning, model training, and requirement analysis. *Test & Evaluation:*
 7 technical validation and verification. *Deployment:* user feedback and override, post deployment
 8 monitoring, and decommissioning.

9 **6.1 Map**

10 The Map function establishes the context and applies the attributes of the AI RMF taxonomy
 11 (Figure 3) to frame risks related to an AI system. The information gathered while carrying out
 12 this function informs decisions about model management, including an initial decision about
 13 appropriateness or the need for an AI solution. Determination of whether AI use is appropriate or
 14 warranted can be considered in comparison to the status quo per a qualitative or more formal
 15 quantitative analysis of benefits, costs, and risks.

16 A companion document describes practices related to mapping AI risks. Table 1 lists the Map
 17 function’s categories and subcategories.

18 *Table 1: Example of categories and subcategories for Map function*

ID	Category	Subcategory
Map: Context is recognized and risks related to the context are identified		
1	Context is established and understood.	Intended purpose, setting in which the AI system will be deployed, the specific set of users along with their expectations, and impacts of system use are understood and documented as appropriate. The business purpose or context of use has been clearly defined or – in the case of assessing existing AI systems – re-evaluated. The organization’s mission and relevant goals for the AI technology are understood. Stakeholders are defined, a plan for continuous engagement/communication is developed, and outreach is conducted.

		System requirements are elicited and understood from relevant stakeholders (e.g., “the system shall respect the privacy of its users”). Design decisions take socio-technical implications into account for addressing AI risks.
		Risk tolerances are determined.
2	Classification of AI system is performed.	The specific task that the AI system will support is defined (e.g., recommendation, classification, etc.).
		Considerations related to data collection and selection are identified. (e.g., availability, representativeness, suitability).
		Detailed information is provided about the operational context in which the AI system will be deployed (e.g., human-machine teaming, etc.) and how output will be utilized.
3	AI capabilities, targeted usage, goals, and expected benefits and costs over status quo are understood.	Benefits of intended system behavior are examined.
		Cost (monetary or otherwise) of errors or unintended system behavior is examined.
		Targeted application scope is specified and narrowed to the extent possible based on established context and AI system classification.
4	Risks and harms to individual, organizational, and societal perspectives are identified.	Potential business and societal (positive or adverse) impacts of technical and socio-technical characteristics for potential users, the organizations, or society as a whole are understood.
		Potential harms of the AI system are elucidated along technical and socio-technical characteristics and aligned with guiding principles.
		Likelihood of each harm is understood based on expected use, past uses of AI systems in similar contexts, public incident reports or other data.
		Benefits of the AI system outweigh the risks, and risks can be assessed and managed. Ideally, this evaluation should be conducted by an independent third party or by experts who did not serve as front-line developers for the system, and who consults experts, stakeholders, and impacted communities.

1 6.2 Measure

2 The Measure function provides knowledge relevant to the risks associated with attributes of the
3 AI RMF taxonomy in Section 5. This includes analysis, quantitative or qualitative assessment,
4 and tracking the risk and its impact. Risk analysis and measurement may involve a detailed
5 consideration of uncertainties, tradeoffs, consequences, likelihood, events, controls, and their
6 effectiveness. An event can have multiple causes and consequences and can affect multiple
7 objectives.

8 Methods and metrics for quantitative or qualitative measurement rapidly evolve. Both qualitative
9 and quantitative methods should be used to track risks.

10 A companion document describes practices related to measuring AI risks. Table 2 lists the
11 Measure Function’s categories and subcategories.

1

Table 2: Example of categories and subcategories for Measure function

ID	Category	Subcategory
Measure: Identified risks are assessed, analyzed, or tracked		
1	Appropriate methods and metrics are identified and applied.	Elicited system requirements are analyzed.
		Approaches and metrics for quantitative or qualitative measurement of the enumerated risks, including technical measures of performance for specific inferences, are identified and selected for implementation.
		The appropriateness of metrics and effectiveness of existing controls is regularly assessed and updated.
2	Systems are evaluated.	Accuracy, reliability, robustness, resilience (or ML security), explainability and interpretability, privacy, safety, bias, and other system performance or assurance criteria are measured, qualitatively or quantitatively.
		Mechanisms for tracking identified risks over time are in place, particularly if potential risks are difficult to assess using currently available measurement techniques, or are not yet available.
3	Feedback from appropriate experts and stakeholders is gathered and assessed.	Subject matter experts assist in measuring and validating whether the system is performing consistently with their intended use and as expected in the specific deployment setting.
		Measurable performance improvements (e.g., participatory methods) based on consultations are identified.

2 6.3 Manage

3 This function addresses risks which have been mapped and measured and are managed in order
 4 to maximize benefits and minimize adverse impacts. These are risks associated with the
 5 attributes of the AI RMF taxonomy (Section 5). Decisions about this function take into account
 6 the context and the actual and perceived consequences to external and internal stakeholders. That
 7 includes interactions of the AI system with the status quo world and potential benefits or costs.

8 Management can take the form of deploying the system as is if the risks are deemed tolerable;
 9 deploying the system in production environments subject to increased testing or other controls;
 10 or decommissioning the system entirely if the risks are deemed too significant and cannot be
 11 sufficiently addressed. Like other risk management efforts, AI risk management must be
 12 ongoing.

13 Practices related to AI risk management are discussed in the companion document. Table 3 lists
 14 the Manage function's categories and subcategories.

15

Table 3: Example categories and subcategories for Manage function

ID	Category	Subcategory
Manage: Risks are prioritized and acted upon based on a projected impact		
1	Assessments of potential harms and results of analyses conducted via the map and measure functions are used to respond to and manage AI risks.	Assessment of whether the AI is the right tool to solve the given problem (e.g., if the system should be further developed or deployed).
		Identified risks are prioritized based on their impact, likelihood, resources required to address them, and available methods to address them.

		Responses to enumerated risks are identified and planned. Responses can include mitigating, transferring or sharing, avoiding, or accepting AI risks.
2	Priority actions to maximize benefits and minimize harm are planned, prepared, implemented, and communicated to internal and external stakeholders as appropriate (or required) and to the extent practicable.	Resources required to manage risks are taken into account, along with viable alternative systems, approaches, or methods, and related reduction in severity of impact or likelihood of each potential action.
		Plans are in place, both performance and control-related, to sustain the value of the AI system once deployed.
		Mechanisms are in place and maintained to supersede, disengage, or deactivate existing applications of AI that demonstrate performance or outcomes that are inconsistent with their intended use.
3	Responses to enumerated and measured risks are documented and monitored over time.	Plans related to post deployment monitoring of the systems are implemented, including mechanisms for user feedback, appeal and override, decommissioning, incident response, and change management.
		Measurable performance improvements (e.g., participatory methods) based on consultations are integrated into system updates.

1 6.4 Govern

2 The Govern function cultivates and implements a culture of risk management within
 3 organizations developing, deploying, or acquiring AI systems. Governance is designed to ensure
 4 risks and potential impacts are identified, measured, and managed effectively and consistently.

5 Governance processes focused on potential impacts of AI technologies are the backbone of risk
 6 management. Governance focuses on technical aspects of AI system design and development as
 7 well as on organizational practices and competencies that directly impact the individuals
 8 involved in training, deploying, and monitoring such systems. Governance should address supply
 9 chains, including third-party software or hardware systems and data as well internally developed
 10 AI systems.

11 Governance is a function that has relevance across all other functions, reflecting the importance
 12 of infusing governance considerations throughout risk management processes and procedures.
 13 Attention to governance is a continual and intrinsic requirement for effective AI risk
 14 management over an AI system’s entire lifespan. For example, compliance with internal and
 15 external policies or regulations is a universal aspect of the governance function in risk
 16 management. Similarly, governance provides a structure through which AI risk management
 17 functions can align with organizational policies and strategic priorities, including those not
 18 directly related to AI systems.

19 A companion document describes practices related to governance of AI risk management. Table
 20 4 lists Govern function’s categories and subcategories.

21

1

Table 4: Example categories and subcategories for Govern function

ID	Category	Subcategory
Govern: A culture of risk management is cultivated and present		
1	Policies, processes, procedures and practices across the organization related to the development, testing, deployment, use and auditing of AI systems are in place, transparent, and implemented effectively.	The risk management process and its outcomes are documented and traceable through transparent mechanisms, as appropriate and to the extent practicable.
		Ongoing monitoring and periodic review of the risk management process and its outcomes are planned, with responsibilities clearly defined.
		Methods for ensuring all dimensions of trustworthy AI are embedded into policies, processes, and procedures.
2	Accountability structures are in place to ensure that the appropriate teams and individuals are empowered, responsible, and trained for managing the risks of AI systems.	Roles and responsibilities and lines of communication related to identifying and addressing AI risks are clear to individuals and teams throughout the organization.
		The organization's personnel and partners are provided AI risk management awareness education and training to enable them to perform their duties and responsibilities consistent with related policies, procedures, and agreements.
		Executive leadership of the organization considers decisions about AI system development and deployment ultimately to be their responsibility.
3	Workforce diversity, equity and inclusion processes are prioritized.	Decision making throughout the AI lifecycle is informed by a demographically and disciplinarily diverse team, including internal and external personnel. Specifically, teams that are directly engaged with identifying design considerations and risks include a diversity of experience, expertise and backgrounds to ensure AI systems meet requirements beyond a narrow subset of users.
4	Teams are committed to a culture that considers and communicates risk.	Teams are encouraged to consider and document the impacts of the technology they design and to develop and communicate about these impacts more broadly.
		Organizational practices are in place to ensure that teams actively challenge and question steps in the design and development of AI systems to minimize harmful impacts.
5	Processes are in place to ensure that diversity, equity, inclusion, accessibility, and cultural considerations from potentially impacted individuals and communities are fully taken into account.	Organizational policies and practices are in place that prioritize the consideration and adjudication of external stakeholder feedback regarding the potential individual and societal harms posed by AI system deployment.
		Processes are in place to empower teams to make decisions about if and how to develop and deploy AI systems based on these considerations, and define periodic reviews of impacts, including potential harm.
6	Clear policies and procedures are in place to address AI risks arising from supply chain issues, including third-party software and data.	Policies and procedures include guidelines for ensuring supply chain and partner involvement and expectations regarding the value and trustworthiness of third-party data or AI systems.
		Contingency processes are in place to address potential issues with third-party data or AI systems.

2

1 7 AI RMF Profiles

2 Profiles are instantiations of the AI RMF Core for managing AI risks for context-specific use
3 cases. Using the AI RMF, profiles illustrate how risk can be managed at various stages of the AI
4 lifecycle or in sector, technology, or end-use applications. Profiles may state an “as is” and
5 “target” state of how an organization addresses AI risk management.

6 NOTE: Development of profiles is deferred until later drafts of the AI RMF are developed with
7 the community. NIST welcomes contributions of AI RMF profiles. These profiles will inform
8 NIST and the broader community about the usefulness of the AI RMF and likely lead to
9 improvements which can be incorporated into future versions of the framework.

10 8 Effectiveness of the AI RMF

11 The goal of the AI RMF is to offer a resource for improving the ability of organizations to
12 manage AI risks in order to maximize benefits and to minimize AI-related harms. Organizations
13 are encouraged to periodically evaluate whether the AI RMF has improved their ability to
14 manage AI risks, including but not limited to their policies, processes, practices, implementation
15 plans, indicators, and expected outcomes.

16 NOTE: NIST is deferring development of this section until later drafts of the AI RMF are
17 developed with the community.

18 Part 3: Practice Guide

19 9 Practice Guide

20 NOTE: NIST is developing a companion Practice Guide which will include additional
21 examples and practices that can assist in using the AI RMF. That Guide, which will reside
22 online only and will be updated regularly with contributions expected to come from many
23 stakeholders, will be part of the NIST AI Resource Center that is being established.