



NIST AI 100-1



# Artificial Intelligence Risk Management Framework (AI RMF 1.0)

# Table of Contents

<b>Executive Summary</b>	<b>1</b>
<b>Part 1: Foundational Information</b>	<b>4</b>
<b>1 Framing Risk</b>	<b>4</b>
1.1 Understanding and Addressing Risks, Impacts, and Harms	4
1.2 Challenges for AI Risk Management	5
1.2.1 Risk Measurement	5
1.2.2 Risk Tolerance	7
1.2.3 Risk Prioritization	7
1.2.4 Organizational Integration and Management of Risk	8
<b>2 Audience</b>	<b>9</b>
<b>3 AI Risks and Trustworthiness</b>	<b>12</b>
3.1 Valid and Reliable	13
3.2 Safe	14
3.3 Secure and Resilient	15
3.4 Accountable and Transparent	15
3.5 Explainable and Interpretable	16
3.6 Privacy-Enhanced	17
3.7 Fair – with Harmful Bias Managed	17
<b>4 Effectiveness of the AI RMF</b>	<b>19</b>
<b>Part 2: Core and Profiles</b>	<b>20</b>
<b>5 AI RMF Core</b>	<b>20</b>
5.1 Govern	21
5.2 Map	24
5.3 Measure	28
5.4 Manage	31
<b>6 AI RMF Profiles</b>	<b>33</b>
<b>Appendix A: Descriptions of AI Actor Tasks from Figures 2 and 3</b>	<b>35</b>
<b>Appendix B: How AI Risks Differ from Traditional Software Risks</b>	<b>38</b>
<b>Appendix C: AI Risk Management and Human-AI Interaction</b>	<b>40</b>
<b>Appendix D: Attributes of the AI RMF</b>	<b>42</b>

## List of Tables

Table 1 Categories and subcategories for the <b>GOVERN</b> function.	22
Table 2 Categories and subcategories for the <b>MAP</b> function.	26
Table 3 Categories and subcategories for the <b>MEASURE</b> function.	29
Table 4 Categories and subcategories for the <b>MANAGE</b> function.	32

## Executive Summary

Artificial intelligence (AI) technologies have significant potential to transform society and people's lives – from commerce and health to transportation and cybersecurity to the environment and our planet. AI technologies can drive inclusive economic growth and support scientific advancements that improve the conditions of our world. AI technologies, however, also pose risks that can negatively impact individuals, groups, organizations, communities, society, the environment, and the planet. Like risks for other types of technology, AI risks can emerge in a variety of ways and can be characterized as long- or short-term, high- or low-probability, systemic or localized, and high- or low-impact.

The AI RMF refers to an *AI system* as an engineered or machine-based system that can, for a given set of objectives, generate outputs such as predictions, recommendations, or decisions influencing real or virtual environments. AI systems are designed to operate with varying levels of autonomy (Adapted from: OECD Recommendation on AI:2019; ISO/IEC 22989:2022).

While there are myriad standards and best practices to help organizations mitigate the risks of traditional software or information-based systems, the risks posed by AI systems are in many ways unique (See Appendix B). AI systems, for example, may be trained on data that can change over time, sometimes significantly and unexpectedly, affecting system functionality and trustworthiness in ways that are hard to understand. AI systems and the contexts in which they are deployed are frequently complex, making it difficult to detect and respond to failures when they occur. AI systems are inherently socio-technical in nature, meaning they are influenced by societal dynamics and human behavior. AI risks – and benefits – can emerge from the interplay of technical aspects combined with societal factors related to how a system is used, its interactions with other AI systems, who operates it, and the social context in which it is deployed.

These risks make AI a uniquely challenging technology to deploy and utilize both for organizations and within society. Without proper controls, AI systems can amplify, perpetuate, or exacerbate inequitable or undesirable outcomes for individuals and communities. With proper controls, AI systems can mitigate and manage inequitable outcomes.

AI risk management is a key component of responsible development and use of AI systems. Responsible AI practices can help align the decisions about AI system design, development, and uses with intended aim and values. Core concepts in responsible AI emphasize human centrality, social responsibility, and sustainability. AI risk management can drive responsible uses and practices by prompting organizations and their internal teams who design, develop, and deploy AI to think more critically about context and potential or unexpected negative and positive impacts. Understanding and managing the risks of AI systems will help to enhance trustworthiness, and in turn, cultivate public trust.

*Social responsibility* can refer to the organization’s responsibility “for the impacts of its decisions and activities on society and the environment through transparent and ethical behavior” (ISO 26000:2010). *Sustainability* refers to the “state of the global system, including environmental, social, and economic aspects, in which the needs of the present are met without compromising the ability of future generations to meet their own needs” (ISO/IEC TR 24368:2022). Responsible AI is meant to result in technology that is also equitable and accountable. The expectation is that organizational practices are carried out in accord with “*professional responsibility*,” defined by ISO as an approach that “aims to ensure that professionals who design, develop, or deploy AI systems and applications or AI-based products or systems, recognize their unique position to exert influence on people, society, and the future of AI” (ISO/IEC TR 24368:2022).

As directed by the National Artificial Intelligence Initiative Act of 2020 (P.L. 116-283), the goal of the AI RMF is to offer a resource to the organizations designing, developing, deploying, or using AI systems to help manage the many risks of AI and promote trustworthy and responsible development and use of AI systems. The Framework is intended to be **voluntary**, rights-preserving, non-sector-specific, and use-case agnostic, providing flexibility to organizations of all sizes and in all sectors and throughout society to implement the approaches in the Framework.

The Framework is designed to equip organizations and individuals – referred to here as *AI actors* – with approaches that increase the trustworthiness of AI systems, and to help foster the responsible design, development, deployment, and use of AI systems over time. AI actors are defined by the Organisation for Economic Co-operation and Development (OECD) as “those who play an active role in the AI system lifecycle, including organizations and individuals that deploy or operate AI” [OECD (2019) Artificial Intelligence in Society—OECD iLibrary] (See Appendix A).

The AI RMF is intended to be practical, to adapt to the AI landscape as AI technologies continue to develop, and to be operationalized by organizations in varying degrees and capacities so society can benefit from AI while also being protected from its potential harms.

The Framework and supporting resources will be updated, expanded, and improved based on evolving technology, the standards landscape around the world, and AI community experience and feedback. NIST will continue to align the AI RMF and related guidance with applicable international standards, guidelines, and practices. As the AI RMF is put into use, additional lessons will be learned to inform future updates and additional resources.

The Framework is divided into two parts. Part 1 discusses how organizations can frame the risks related to AI and describes the intended audience. Next, AI risks and trustworthiness are analyzed, outlining the characteristics of trustworthy AI systems, which include