#### **JUNE 2024**

# **AI POLICY TEMPLATE**

**Build Your Foundational Organizational AI Policy** 

**Responsible AI Institute** 

Version 1.0

Last updated: June 6, 2024

#### Al Policy Template Disclaimer

This Template [AI Policy] includes various provisions throughout that must be reviewed and, potentially, revised based on the specifics of your business and your use of artificial intelligence technologies. You are advised to confirm that all prepopulated information is accurate and appropriate for your business.

This Policy Template was developed to reflect sound practices for responsible Al management at the time it was created. The Policy is not intended to, and does not purport to, satisfy particular legal requirements in every jurisdiction. In all cases, you are advised to consult an attorney for guidance on the laws and regulations applicable to your business and to determine whether this Policy is adequate to address such legal requirements.

Artificial intelligence is a rapidly evolving field, and the best practices for responsible Al management are likewise evolving. We recommend revisiting your policy regularly (at least annually) and updating it as needed to reflect current standards or legal requirements.

You use this Policy Template at its own risk. This Policy Template does not constitute legal advice, and by using all or any part of this Policy Template, you agree to this disclaimer. You are advised to (i) consult independent legal advice before adopting or publishing your policy; (ii) read this Policy with care and modify, delete or add all and any areas as necessary; and (iii) not rely on this Policy for any purpose without seeking legal advice from a licensed attorney in your jurisdiction. This Policy Template is provided only for informational purposes and may or may not reflect the most current legal developments; accordingly, it is not promised or guaranteed to be correct or complete.

# **Table of Contents**

Introduction to the AI Policy Template	3
I. Purpose and Scope	5
II. Al Principles	5
III. AI Objectives and Strategy	7
IV. Governance	9
V. Data Management	13
VI. Risk Management	18
VII. Project Management	28
VIII. Stakeholder Management and Engagement	35
IX. Workforce Management	38
X. Regulatory Compliance	39
XI. Al Procurement	40
XII. Documentation Management	44
XIII. Review and Enforcement of the Al Policy	45
XIV. Conclusion/Acknowledgement	46
Appendix A	47

## **Introduction to the AI Policy Template**

As artificial intelligence (AI) unlocks opportunities for businesses to increase efficiency and generate novel insights and offerings, organizations seek to build responsible AI practices that align with ethical principles, mitigate potential risks, and foster trust with stakeholders. For many, an organizational-level policy can be a foundational document to establish guiding principles, objectives, and management direction for all AI-related activities according to business requirements. A standalone policy for AI can also centralize and highlight an organization's responsible AI strategy to accelerate internal adoption and external awareness. Developing an AI Policy is a requirement of an AI management system under leading standard ISO/IEC 42001, and it is also a main policy recommendation under the Responsible AI Institute's framework for organizational maturity.

This document serves as a template for an AI Policy, which can be used to establish a comprehensive framework for an organization's development, procurement, supply, and use of AI technologies. The template provides one interpretation of how global and regional guidance, including the NIST AI Risk Management Framework (RMF), and ISO/IEC 42001, can be operationalized through corporate policy, informed by RAI Institute's deep expertise in accelerating member organizations' maturity for AI. Organizations, depending on their size, existing policies, and industry requirements, may find that some schemes detailed in the template may need to be adapted or substituted to better fit their context.

Organizations are encouraged to customize this Policy Template to address the specific needs and risks of their AI use cases. In alignment with ISO/IEC 42001, organizations are recommended to use a suite of factors to inform their AI Policy, including but not limited to business strategy; organizational values and culture; organizational risk environment and tolerance; statutory, regulatory and contractual requirements; and possible AI risks and impacts of its use cases.<sup>1</sup>

The following template is RAI Institute's first draft version. The Institute will be accepting feedback into July 2024 for an updated version, to be released in the following months.

© Responsible Al Institute 2024 | All Rights Reserved | Do Not Use Without Permission

<sup>&</sup>lt;sup>1</sup> Adapted from ISO/IEC 42001 B.2.2.

<ol> <li>Purpose a</li> </ol>	and Scope
-------------------------------	-----------

A.	inte our to e	corganization name], we recognize the transformative potential of artificial elligence (AI) to enhance our operations, products, and services. This Policy outlines commitment to responsible AI [e.g., development, deployment, supply, use]ensure ethical considerations are upheld, AI risks are managed, and compliance to erging regulation is achieved.
B.	guid all [ rele	e following [policy document name e.g. "Al Policy"], provides a framework to de all activities at [organization name] related to Al. This Policy applies to e.g., employees, contractors, and third-party suppliers] involved in [all evant Al activities e.g. "the development, procurement, and use of Al"] within ganization name]2
objectives, generate outputs such as predictions, recommendations, or decisions		uencing real or virtual environments. Al systems are designed to operate with varying
	1.	We <b>define an AI model as</b> [e.g., "a component of an information system that implements AI technology and uses computational, statistical, or machine-learning techniques to produce outputs from a given set of inputs" ———.
	2.	We <b>define an AI system as</b> [e.g., "any data system, software, hardware, application, tool, or utility that operates in whole or in part using AI" <sup>5</sup> ]
	3.	We have developed these definitions in alignment with leading industry and ecosystems definitions, including [e.g., the OECD, EU AI Act, U.S. EO], to define the bounds of AI life cycle management and inventorying and of compliance requirements, as appropriate for our context and use cases.
II.	ΑI	Principles
A.	dive cor	lignment with the established enterprise values, including [e.g., sustainability, ersity, equity, and inclusion, corporate social responsibility], we are mmitted to upholding ethical principles in the [AI activities e.g., procurement, relopment, deployment, supply, and/or use] of AI technologies.6

<sup>&</sup>lt;sup>2</sup> Aligned with ISO/IEC 42001 4.3.

<sup>&</sup>lt;sup>3</sup> Definition of AI provided by the National Institute of Standards and Technology AI Risk Management Framework (NIST AI RMF). Adapted from: OECD Recommendation on AI:2019; ISO/IEC 22989:2022.

<sup>&</sup>lt;sup>4</sup> Definition of AI model provided by United States Executive Order No. 14110 on Safe, Secure, and Trustworthy Development and Use of Artificial Intelligence.

<sup>&</sup>lt;sup>5</sup> Definition of AI system provided by United States Executive Order No. 14110 on Safe, Secure, and Trustworthy Development and Use of Artificial Intelligence.

<sup>&</sup>lt;sup>6</sup> Aligned with ISO/IEC 42001 A.2.2.

- B. Our Al principles encompass the following<sup>7</sup>:
  - System Trustworthiness: Every AI system that is [bought, built, used, or sold]
     \_\_\_\_\_\_ by [organization name]\_\_\_\_\_\_ shall strive to achieve appropriate levels
     of all trustworthy characteristics, as defined by the National Institute of Standards
     and Technology AI Risk Management Framework (NIST AI RMF).8 The achievement
     of each of these characteristics will depend on the use case and may require
     tradeoffs between characteristics, which will be justified and documented.
    - a) <u>Validity and Reliability</u>: All Al systems shall consistently provide accurate outputs
      or otherwise behave within a defined range of acceptability when subject to
      expected conditions of use.
    - b) <u>Safety</u>: No AI system shall endanger human life, health, property, or the environment.
    - c) <u>Security and Resiliency</u>: All Al systems shall withstand unexpected adverse events or unexpected changes in their environment or use, maintaining confidentiality, integrity, and availability in the event of adversarial or unauthorized actions.
    - d) <u>Accountability and Transparency</u>: Meaningful and timely information about every AI system shall be provided to all relevant stakeholders, tailored to the expected knowledge and accessibility needs of each audience. An accountability structure governs each decision made related to an AI system.
    - e) Explainability and Interpretability: All Al systems shall be designed and documented to answer how and why a decision was made by the system, to the fullest extent possible.
    - f) <u>Privacy-Enhanced</u>: All Al systems shall safeguard human autonomy, identity, and dignity with respect to privacy to the fullest possible extent.
    - g) <u>Fairness with Harmful Bias Managed</u>: All Al systems shall meet a defined metric of fairness appropriate for its context and shall manage all forms of harmful bias, including system bias, computational and statistical bias, and human-cognitive bias.<sup>9</sup>
  - Human Oversight and Accountability: We commit to ensuring that all processes and
    materials related to AI in our organization are subject to proper oversight
    mechanisms to enable responsible development and use of AI. We also embrace all
    sources of external accountability, including seeking independent audits and

<sup>8</sup> The following statements by characteristic are aligned with the definitions of each characteristic provided in the NIST AI RMF. Additionally aligned with GOVERN 1.2.

<sup>&</sup>lt;sup>7</sup> Aligned with ISO/IEC 42001 B.6.1.2.

<sup>&</sup>lt;sup>9</sup> NIST has identified these three as major categories of AI bias to be considered and managed. Each of these can occur in the absence of prejudice, partiality, or discriminatory intent.

certifications	s, monitoring by $\mathfrak g$	governmental	l organizations,	and meaningful
transparency	with interested	public parties	S.	

- 3. **Beneficence, Equity, and Ethics:** We shall align our AI strategy with the broader interest of our organization to preserve and promote societal well-being. This includes our commitments to ethical frameworks related to [e.g., sustainability, equity, human rights] \_\_\_\_\_\_.
- 4. **Continual Learning:** As the technological and regulatory environment of AI rapidly develops, we are committed to a culture of open-mindedness, flexibility, and dialogue. We shall engage with partners, peers, stakeholders, and the public to invest in shared knowledge and a shared vision for responsible AI.

## III. Al Objectives and Strategy

A.	Uni ent sel	a [entity in a context e.g. "leading company in the financial services industry in the ited States"], AI has the power to [benefits], adding significant terprise value and bolstering our competitive advantage. We will [build, buy, and/or applications], and aim to oly AI in [list of functions or use cases] <sup>10</sup>	
B.	However, we also foresee challenges to responsible [procurement, development, deployment, supply, and/or use] of AI, given the context of our organization and of our use cases. This includes [examples of regulatory requirements, areas of deficient organizational capacity or governance, industry-specific data challenges]11		
C.	. With these potential benefits and challenges in mind, a long-term AI business strategy shall be developed, documented, and implemented. The <b>key objectives</b> that shall guide this strategy are as follows <sup>12</sup> :		
	1.	[Key objective with an explanation of how progress along it will be measured]	
	2.	[Key objective with an explanation of how progress along it will be measured]	
	3.	[Key objective with an explanation of how progress along it will be measured]	
	4.	This strategy shall align with [organization name]'s broader objectives and commitments, including [e.g., achieving net-zero by 2030, investing in leadership by historically marginalized groups within and beyond the organization]	

<sup>&</sup>lt;sup>10</sup> Aligned with NIST AI RMF MAP 1.3 and ISO/IEC 42001.

<sup>&</sup>lt;sup>11</sup> Aligned with ISO/IEC 42001 4.1.

<sup>&</sup>lt;sup>12</sup> Aligned with NIST AI RMF MAP 1.3 and ISO/IEC 42001 6.2 and B.6.1.2.

D.	Recognizing the need for bespoke, flexible, and unobtrusive organizational adaptation for AI, the contents of this AI Policy shall be used as an augmenting layer on top of [organization name]'s existing governance, policies, and processes. <sup>13</sup>		
E. While [organization name]'s AI strategy will evolve into a complete pract over time, the urgency of AI risks requires that [organization name] ident areas of highest priority as a [developer, procurer, and/or supplier] of AI:			
		For bought systems <sup>14</sup> :  a) Responsibly buying AI will require significant investment in the following capabilities:  (1) A rigorous and principles-driven procurement process that sufficiently weighs marketplace options and comprehensively assesses the risks attached to potential suppliers and their product or service;  (2) A legal means to clarify liability and other requirements with suppliers, to whatever extent is not articulated by regulation;  (3) Role- and system-specific training and educational materials for employees such that procured AI can be used safely and responsibly; and  (4) [additional areas of strategic focus]  For built systems:  a) Responsibly building AI will require significant investment in the following capabilities:  (1) A product management program that balances the forward inertia of innovation with the necessary governance gates and other processes, such as AI Impact Assessments, to sufficiently manage risk;  (2) Processes and tools to enable systematic and comprehensive documentation of AI systems throughout their life cycle, in accordance with regulatory requirements;  (3) Responsible AI training for all technical and nontechnical roles involved in an AI system's life cycle across design, development, deployment, and operation; and  (4) [additional areas of strategic focus]	
	3.	For sold systems <sup>15</sup> :  a) Responsibly selling AI will require significant investment in the following capabilities:  (1) Development of audience-specific guidance and documentation for each AI system, such as for potential buyers, users, or the general public;  (2) A legal means to clarify liability and other requirements with buyers, to whatever extent is not articulated by regulation;	

<sup>13</sup> Aligned with ISO/IEC 42001 A.2.3 and NIST AI RMF GOVERN 1.2.
14 Aligned with ISO/IEC 42001 A.10.3.
15 Aligned with ISO/IEC 42001 A.10.4.

	<ul> <li>(3) Regular assessment of downstream impacts of AI products, enabled by transparency and shared learning agreements with buyers while also protecting users' privacy and other rights; and</li> <li>(4) [additional areas of strategic focus]</li> </ul>
F.	[organization name] maintains a commitment to continually improve the suitability, adequacy, and effectiveness of this AI Policy and of its AI management system. <sup>16</sup>
IV.	Governance
A.	The following executive or senior management positions are designated as the <code>[owners, champions, and/or sponsors]</code> of <code>[organization name]</code> 's responsible Al approach. Executive <code>[owners, champions, and/or sponsors]</code> bear responsibility for ensuring that <code>[organization name]</code> 's responsible Al strategy is developed and executed effectively and are ultimately accountable for its success.\frac{17}{2}
	<ol> <li>The [position title e.g., CIO, CTO, COO] is accountable for [responsibilities and key results e.g., leading the Steering Committee, tracking, reporting, and owning overall progress]</li> </ol>
	2. The [position title e.g., VP of Legal, VP of Compliance] is accountable for [responsibilities and key results e.g. ensuring that AI systems do not violate legal or regulatory requirements]
	3. [Additional positions and fields of accountability]
В.	[organization name]'s responsible AI strategy shall be led by two major governance bodies <sup>18</sup> :
	1. A High-Level Board/Steering Committee shall provide executive leadership and oversight, including direction, mandates, and resourcing for responsible AI efforts, in a timely manner. The Steering Committee has representation across senior-level management, including [e.g., CIO, CDO, Board members, Chief AI Officer]
scope establis develop top-lev	ned with ISO/IEC 42001 10.1. An organization's AI management system, in accordance with the of ISO/IEC 42001, refers to the collective body of management structures and processes shed for an organization to responsibly perform their role with respect to AI systems (e.g. to use, p, monitor or provide products or services that utilize AI). This AI Policy has been designed as a sel framework for an organization's AI management system.

<sup>&</sup>lt;sup>18</sup> ISO/IEC 42001 3.22 notes that "not all organizations, particularly small organizations, will have a governing body separate from top management." The creation of separate steering and operational groups may not be necessary.

<sup>&</sup>lt;sup>19</sup> ISO/IEC 42001 5.1 provides a list of activities for top management to demonstrate leadership and commitment in this capacity.

<sup>©</sup> Responsible Al Institute 2024 | All Rights Reserved | Do Not Use Without Permission

	2.	An Operational Committee shall direct the implementation of Steering Committee objectives, oversee the life cycle progression and impact of AI systems, and act as a convening power for responsible AI efforts. The Operational Committee has representation across functions and/or divisions of [organization name], including [e.g., Head of Information Security, Head of Procurement, Legal, Head of HR]
C.		e Operational Committee <sup>20</sup> shall convene every [e.g. two weeks] Its sponsibilities include:
	1.	Approving progression of an AI system's development to the next stage;
	2.	Directing the development and updating of both cross-functional and function- specific AI guidance and tools alongside function or division leaders;
	3.	Developing and managing Al-related inventories, such as of Al systems and Al incidents;
	4.	Developing responsible AI training;
	5.	Determining and implementing changes to the AI management system in a planned manner $^{21}$ ; and
	6.	[additional responsibilities]
D.		e Steering Committee shall convene every [e.g. two months] Its sponsibilities include <sup>22</sup> :
	1.	Developing an organizational-level AI strategy with a clear timeline and measures of success (KPIs);
	2.	Articulating how its AI strategy amplifies or creates trade-offs with other organizational objectives or commitments (highest-level SWOT or ROI analyses);
	3.	Determining organizational AI risk tolerances;
	4.	Determining and allocating the resources needed for the establishment, implementation, maintenance, and continual improvement of the AI management system <sup>23</sup> ;
	5.	Aligning workforce planning with responsible AI human capital needs; and
	6.	[additional responsibilities]
·ho /	000	rational Committee is often named differently in practice, such as simply "Responsible AL (RAL)

 $<sup>^{20}</sup>$  The Operational Committee is often named differently in practice, such as simply "Responsible AI (RAI) Team."

<sup>&</sup>lt;sup>21</sup> Aligned with ISO/IEC 42001 6.3.

<sup>&</sup>lt;sup>22</sup> See ISO/IEC 42001 5.1 for further guidance on responsibilities of top management.

<sup>&</sup>lt;sup>23</sup> Aligned with ISO/IEC 42001 7.1.

E.		ditional individual roles shall be created to direct and support [organization name]'s Al strategy, including [e.g., Chief Al Officer, Responsible Al Chair] ich shall be responsible for [description of responsibilities]
dat		ore specific areas related to AI, such as [e.g., cybersecurity, supplier relationships, and ta management], can be managed by specific [e.g. functions], in llaboration with or with broad oversight by the Operational Committee. <sup>24</sup>
	1.	Leaders of [organization name]'s [functions, segments, and/or divisions] shall direct the development of Al-specific processes and tools tailored for their operations, informed by Operational Committee directives, similar efforts by peer [functions, segments, and/or divisions], and consultations with domain experts and internal or external stakeholders.
G.	ac	ganization name] shall additionally determine <b>roles, responsibilities, and</b> countability of internal AI actors with respect to AI systems and the systems' external actors. <sup>25</sup> The position of each role in a chain of accountability shall be made clear.
	1.	<b>Internal AI actors</b> include all employees who contribute to AI design, development, deployment, operation and monitoring; TEVV (test, evaluation, verification, and validation) tasks; risk management and impact assessment tasks <sup>26</sup> ; procurement tasks; governance and oversight tasks; and [additional task areas]
		a) This can include developers; data scientists; data engineers; product managers; system integrators; system operators; domain experts; socio-cultural analysts and experts (e.g. DEI, accessibility, governance); human factors experts (e.g., UX/UI design); procurers; AI governance and oversight professionals; legal and privacy officers; and [additional internal AI actor profiles]
	2.	Internal AI actors shall be assigned responsibilities to appropriately manage relationships with <b>external AI actors</b> , which include all individuals, groups, and society members that are involved in AI systems' life cycle or may be impacted by the system.
		a) This can include suppliers and partners (of data or Al platforms, products, or services); third-party assessors or evaluators (of data, algorithms, models, and/or systems); clients; data subjects; end users; members of impacted communities; and [additional external Al actor profiles]
H.		<b>mmunication and feedback channels</b> shall be augmented to enable proper ormation sharing and issue management between different Al actors. <sup>27</sup> These

<sup>&</sup>lt;sup>24</sup> ISO/IEC 42001 B.3.2 provides examples of areas where roles and responsibilities can be defined.
<sup>25</sup> Adapted from NIST AI RMF Appendix A.
<sup>26</sup> Aligned with NIST AI RMF GOVERN 2.1.
<sup>27</sup> Aligned with ISO/IEC 42001 7.4.

		annels shall facilitate various processes percolating upwards or downwards ganization name]'s structure, including <sup>28</sup> :
	1.	A mechanism for employees to report their concerns about any of [organization name]
	2.	Readily accessible means for employees to contact a human representative within the organization for support and guidance on Al-related inquiries, including [e.g., on where to find responsible Al upskilling material, to identify executive RAI champions];
	3.	Formalized and regular performance and progress reporting between organization [functions, segments, or divisions] to AI governance bodies to top management; <sup>31</sup> and
	4.	A formalized and regular process to proactively solicit feedback from internal AI actors on the suitability, adequacy, and effectiveness of [organization name]'s AI management system.
l <b>.</b>	sys	vernance gates shall be established to facilitate the standardized approval of Al stems' life cycle progression and the collection of documentation at each life cycle age. Governance gates shall accomplish the following <sup>32</sup> :
	1.	Be performed and managed by [structure e.g. "the Operational Committee with high-level oversight by the Steering Committee"] through [process e.g. "asynchronous and digital voting, bimonthly meetings to confirm approval decisions, and triggers and mechanisms for escalation of decisions"];
	2.	Pause or terminate an AI proof of concept (PoC) or project once a risk has been identified that measures beyond [organization name]'s risk tolerance (see Risk Management Section G);
		<ul> <li>a) Approval requirements are tied to measures of AI risk criteria and enable exercise of AI risk triage processes in alignment with risk priorities and tolerance.<sup>33</sup></li> </ul>

<sup>&</sup>lt;sup>28</sup> Aligned with NIST AI RMF GOVERN 2.1.
<sup>29</sup> Aligned with ISO/IEC 42001 B.3.3.
<sup>30</sup> Aligned with ISO/IEC 42001 B.3.3.
<sup>31</sup> Aligned with ISO/IEC 42001 5.3.
<sup>32</sup> Aligned with ISO/IEC 42001 B.6.1.3.
<sup>33</sup> Aligned with ISO/IEC 42001 6.1.1.

- b) Approval requirements are tied to continued alignment with organizational RAI principles and objectives, determined via established measures.
  - (1) Approval is also dependent on demonstrated achievement or projected achievement of the system's intended purpose and stated objectives.<sup>34</sup>
  - (2) Approval is also dependent on the completion of all required responsible Al tasks at the current life cycle stage, including the documentation of all relevant decisions made and Al actors involved.
- c) Depending on the catalyst for pause/termination, projects may proceed at a later date once concerns have been remediated, or they may need to be resubmitted at the beginning of the life cycle pipeline as a new PoC.
- 3. Reprioritize an Al project in relation to its counterparts depending on the risk level and the human and technological resource requirements identified, in the context of overall Al objectives, as necessary;

4.	Utilize AI Impact Assessments (see Risk Management Section I) as a tool to inform
	approval decisions and set conditions for progression, such as [e.g. specific design
	or use requirements], <sup>35</sup> performed by [details on roles e.g. "mainly developer
	and procurement teams with support of and consultation with compliance, external
	stakeholders, etc."]36 at the appropriate level of comprehensiveness and at
	the appropriate life cycle stages; and

5. Be responsive to risk triggers, both established and unexpected, and outline a review or recourse process to manage the cause of the trigger. Risk triggers include [e.g., feedback or AI performance requiring immediate attention, switching vendors, changes to applicable regulation, reports of a serious incident from a similar use case]

## V. Data Management

A.	Existing data management at [organization name]	is implemented by
	[governance and policies] and [tools and	processes] The following
	guidance shall be incorporated into the existing da absent. <sup>37</sup>	ta management framework wherever

В.	Additionally, [organization name]	$\_$ identifies and aligns with existing regulations
	and guidelines for data management and	reporting, including [e.g., appropriate use and

<sup>&</sup>lt;sup>34</sup> Aligned with NIST AI RMF MANAGE 1.1.

<sup>35</sup> Aligned with ISO/IEC 42001 B.5.2.

<sup>&</sup>lt;sup>36</sup> Aligned with ISO/IEC 42001 B.5.2.

<sup>&</sup>lt;sup>37</sup> Aligned with ISO/IEC 42001 A.2.3 and NIST AI RMF GOVERN 1.2.

	disclosure of IP-protected content, proper and industry-specific consent mechanisms for data subjects]38
C.	[organization name] shall be committed to meaningful <b>data transparency</b> , characterized by the ongoing delivery and testing of concise and audience-specific data information and mechanisms for recourse, informed by engagement with all interested parties, including [e.g., users, impacted groups, third-party auditors, researchers]
D.	During the early stages of a project, teams shall identify, characterize, and justify the type and quantity of data needed for an Al system.
E.	Rigorous exploratory data analysis shall be conducted for each potential data set to gain insights into underlying structure and statistical properties and to assess alignment with data fit-for-purpose and quality standards.
F.	General information about every selected and utilized data set shall be logged in an <b>enterprise-wide data inventory system</b> , including [e.g., in which projects it is used, internal cross-project usage permissions, data source and contact information, data type and features] <sup>39</sup>
G.	Systems and standards shall be developed by [e.g. functions such as Information Security] to manage the secure and organized storage of data sets enterprisewide.
H.	For all data sets from third-party sources, project teams shall follow a standard protocol set forth by the procurement team to communicate procurement needs during project scoping and to allow the procurement team to manage the procurement process.
l.	For each data set used in a project, the project team shall document as part of project-level documentation <sup>40</sup> :
	1. <u>Data source</u> <sup>41</sup> : internal or external, including how data is acquired and/or accessed from the source;

b) If external, from which supplier and their contact information, and how the data is created, if known;

a) If internal, how the data is created and contact information of the data owner/manager;

 $<sup>^{38}</sup>$  Aligned with NIST AI RMF 1.2.2 and GOVERN 1.1.

<sup>&</sup>lt;sup>39</sup> Aligned with ISO/IEC 42001 A.4.3. Documenting resources at the organizational level in addition to the project level enables more efficient tracking and handling, such as for estimating data storage needs and reusing data sets, when appropriate.

<sup>&</sup>lt;sup>40</sup> Aligned with ISO/IEC 42001 B.4.3.

<sup>&</sup>lt;sup>41</sup> Aligned with ISO/IEC 42001 A.7.3.

<sup>©</sup> Responsible Al Institute 2024 | All Rights Reserved | Do Not Use Without Permission

- Data types and features<sup>42</sup>: whether the data is directly observable, collected from subjects and their demographics/characteristics, and/or indirectly inferred/derived from other data; includes sensitive personal identifiable information (PII); associated metadata; structured or unstructured; continuous or discrete; or multimodal, time series, etc.;
- 3. <u>Data consent process and results</u>: any legal requirements for consent in the use case; whether data subjects provided free, prior, and informed consent; how data can be excluded/removed in response to retracted consent;
- 4. <u>Data provenance process and results</u><sup>43</sup>: information about the creation, update, transcription, abstraction, validation and transferring of the control of data within; data sharing and data transformation;
- Internal data access and usage: permissions and restrictions on the access (permission to view, edit, or share) and use of the data set, including whether it can be used in other projects based on collection conditions (with supplier, with respect to consent);
- 6. External data access and cybersecurity process and results<sup>44</sup>: terms of data sharing with external parties, such as model suppliers or researchers, and measures to prevent unauthorized access to or control over the data;
- 7. <u>Data privacy process and results</u><sup>45</sup>: the required level of privacy based on the use case and privacy-enhancing technologies and techniques applied; sensitive data or PII anonymization, etc.;
- 8. <u>Data proxy risks and fairness</u>: identifying potential proxies for protected class characteristics, and mitigating bias risks to align with fairness standards;
- 9. <u>Data collection and preparation process and results</u><sup>46</sup>: the collection, preparation, and data transformation methods used and a justification for each based on intended use(s) and system context, and including categories of data for machine learning (e.g. training, validation, test and production data);
- 10. <u>Data quality process and results</u><sup>47</sup>: assessment of data set quality along defined metrics, such as fit-for-purpose, representativeness, completeness, consistency, and accuracy, and any methods used to increase data quality and resulting measured improvements;

<sup>&</sup>lt;sup>42</sup> Aligned with ISO/IEC 42001 B.7.3.

<sup>&</sup>lt;sup>43</sup> Aligned with ISO/IEC 42001 A.7.5 and B.7.3.

<sup>&</sup>lt;sup>44</sup> Aligned with ISO/IEC 42001 B.7.2.

<sup>&</sup>lt;sup>45</sup> Aligned with ISO/IEC 42001 B.7.2.

<sup>&</sup>lt;sup>46</sup> Aligned with ISO/IEC 42001 B.7.3 and B.7.6, which lists common methods.

<sup>&</sup>lt;sup>47</sup> Aligned with ISO/IEC 42001 A.7.4 and B.7.2 and NIST AI RMF MAP 2.3.

- 11. <u>Data retention and disposal</u>: protocols for retaining or disposing of data used for training, operating, and maintaining the AI system in a timely manner, if not prohibited by regulations and other transparency requirements;
- 12. <u>Data drift and versioning</u>: a process to ensure data do not become outdated as the culture and context around the system change, tracked by a data set version control system.

#### J. For bought systems<sup>48</sup>:

1.	In partnership with the procurement team, project teams shall request information on all data used to develop the system from the supplier. Based on the data documentation requirements listed in Section I, suppliers of systems shall be required to disclose [e.g., data types and features, sources, provenance, consent policy, privacy policy, collection and preparation process, quality assessment, disposal policy], to advance in the procurement process.
2.	Suppliers shall be required to disclose additional details related to data used to train, validate, and test the system, including [e.g., IP compliance, consent and privacy processes specific to use in training], to advance in the procurement process.
3.	In partnership with the procurement team, project teams shall clarify and establish data control and sharing requirements with the supplier, including [e.g., what data received or created by the system will be shared with the supplier, what enterprise (buyer) data will be visible to the supplier, what data related or resulting from use of the system will be retained by the supplier and/or used to train the system, what data related or resulting from use of the system will be controlled by the enterprise (buyer)]
4.	Project teams shall document all additional data used to support the operation of the system, including [e.g., enterprise reference databases, data sets for fine-tuning or prompt engineering, customer data], and shall document details on how such data is collected, utilized, stored, and/or disposed.
5.	Project teams shall document post-deployment and/or post-decommission data protocols to manage data generated through the operation of the system, including [e.g., input and output logs, event and incident logs, feedback from users], and shall document details on how such data is collected, utilized, stored, and/or disposed.

#### K. For built systems:

1. If any components (e.g., data sets, AI models, platforms) of the built system are procured, project teams shall request information on all data used to develop them, in partnership with the procurement team. Based on the data documentation

<sup>&</sup>lt;sup>48</sup> Aligned with ISO/IEC 42001 B.10.3.

<sup>©</sup> Responsible AI Institute 2024 | All Rights Reserved | Do Not Use Without Permission

	requirements listed in Section I, suppliers shall be required to disclose [e.g., data types and features, sources, provenance, consent policy, privacy policy, collection and preparation process, quality assessment, disposal policy], to advance in the procurement process.
2.	Project teams shall document how data sets, regardless of origin, are used in the development of the system. Teams shall disclose additional details related to the data used to train, validate, and test the system, including [e.g., IP compliance, consent and privacy processes specific to use in training]
	a) If an AI model is procured, suppliers shall be required to disclose additional details related to data used to train, validate, and test the model, including e.g., IP compliance, consent and privacy processes specific to use in training], to advance in the procurement process.
3.	Project teams shall document how data sets are used to support the operation of the system, including [e.g., enterprise reference databases, data sets for iterative retraining, fine-tuning, or prompt engineering, customer data], and shall document details on how such data is collected, utilized, stored, and/or disposed.
	a) If an AI model is procured, project teams shall, in partnership with the procurement team, clarify and establish data control and sharing requirements with the supplier, including [e.g., what data received or created by a model will be shared with the supplier, what enterprise (buyer) data will be visible to the supplier, what data related or resulting from use of the model will be retained by the supplier and/or used to train the system, what data related or resulting from use of the model will be controlled by the enterprise (buyer)]
4.	Project teams shall document post-deployment and/or post-decommission data protocols to manage data generated through the operation of the system, including [e.g., input and output logs, event and incident logs, feedback from users], and shall document details on how such data is collected, utilized, stored, and/or disposed.
Fo	r sold systems <sup>49</sup> :
1.	In partnership with compliance, legal, and sales teams, project teams shall determine and compile data documentation that is relevant, required, or requested by buyers of systems.
2.	In partnership with compliance, legal, and sales teams, project teams shall clarify and establish data control and sharing requirements with the buyer, including [e.g., what data received or created by the system will be shared with the enterprise (supplier), what buyer data will be visible to the enterprise (supplier), what data related or resulting from use of the model will be retained by the enterprise (supplier) and/or

L.

<sup>&</sup>lt;sup>49</sup> Aligned with ISO/IEC 42001 B.10.4.

	controlled by the buyer]
	VI. Risk Management
A.	Existing risk management at [organization name] is implemented by [governance and policies] and [tools and processes] The following guidance shall be incorporated into the existing risk management framework wherever absent. <sup>50</sup>
B.	[Organization name] shall determine to what extent discipline-specific risk management processes sufficiently integrate AI considerations for those specific aspects, such as for [e.g., information security, safety, or privacy] 51 AI-specific risk management processes and tools shall be implemented once existing processes are exhausted.
C.	Additionally, [organization name] identifies and aligns with existing regulations and guidelines for risk management, including established documentation, reporting, and disclosure requirements and industry or use case-specific requirements, such as [e.g. "reporting results of safety tests of high-risk models, as required by the U.S. Executive Order on the Safe, Secure, and Trustworthy Development and Use of Artificial Intelligence"] 52
D.	[Organization name] defines an Al impact as [e.g. "a negative effect experienced by the organization, individuals, groups of individuals, or societies as a result of the organization's development, use, or otherwise interaction with Al technologies"] This definition is developed to be operationalized through our Al Impact Assessment, and therefore [e.g. "considers negative impacts exclusively but additionally considers impacts to the organization"]
	1. This definition was developed in alignment with leading industry and ecosystem definitions, including [e.g., industry definitions of harm, the scope of AI impacts as detailed by regulators requiring an AI Impact Assessment]
	2. Using this definition, [organization name] has also developed an <b>impact taxonomy</b> that enables the identification and categorization of specific impacts and assigns a severity score to each impact, based on a justified methodology of measuring magnitudes of harm.
E.	[Organization name] defines an Al incident as [e.g. "an event precipitated by the organization's mismanagement of an Al system at any point of its life cycle which has the potential to lead to an Al impact"] This definition is developed to align with

used to train the system, what data related or resulting from use of the model will be

<sup>&</sup>lt;sup>50</sup> Aligned with ISO/IEC 42001 A.2.3 and NIST AI RMF GOVERN 1.2. <sup>51</sup> Aligned with ISO/IEC 42001 B.5.2.

<sup>&</sup>lt;sup>52</sup> Aligned with NIST AI RMF 1.2.2 and GOVERN 1.1.

<sup>©</sup> Responsible Al Institute 2024 | All Rights Reserved | Do Not Use Without Permission

	dis	tinc	finition and taxonomy for Al impacts, and therefore [e.g. "does not make a tion between 'incidents', 'accidents', 'hazards', and 'near-misses' nor limits 'incidents' unexpected events, unlike in other disciplines"]
F.	cor	ent's res <sub>l</sub>	ization name] defines an AI risk as [e.g. "the composite measure of an As probability of occurring (likelihood) and the magnitude of the consequences of the bonding AI impact" [53] This definition is developed in alignment with g industry and ecosystem definitions, including [e.g. the NIST AI RMF]
	1.	of Al	systems and individual AI risks identified for a system can be each assigned one the following risk levels: [e.g. "minimal (1), moderate (2), high (3), and very high (4)"] A system's risk level is calculated as the maximum risk level of all known risks for that system; individual AI risks are calculated using a risk matrix enstructed from the two independently measured dimensions of likelihood and verity. 54
G.	pu ris	rsue ks, v	anization name], we recognize that tolerating risk is necessary to innovation and remain competitive. As we balance the opportunities of AI with its we define the following boundaries of [organization name]'s risk nce <sup>55</sup> with respect to AI:
	1.	Un [e.g	e cases that carry unacceptable risks are prohibited without condition. <sup>56</sup> acceptable risks are aligned with leading and emerging requirements, such as g. the EU AI Act], and include [e.g., the use of emotion recognition stems, manipulative behavior-change techniques, and social scoring algorithms], among others.
		a)	Based on our own values and risk appetite, we additionally prohibit systems that introduce the following unacceptable risks: [e.g., unfair and unexplainable promotion or termination decisions for employees based on an automated decision, use of a specific Large Language Model (LLM)]
	2.	Th	e following <b>risk thresholds</b> apply to all systems immediately prior to deployment:
		a)	Systems that remain very high-risk or high-risk post-treatment shall not be allowed to deploy. <sup>57</sup>
		b)	Systems deemed moderate-risk post-treatment shall be subject to more rigorous and frequent testing and monitoring requirements, including impact assessments, than those deemed minimal-risk.
<sup>54</sup> An o	rgan	izati	NIST AI RMF 1.1. ion's approach to characterizing and calculating risk of an AI system can vary in structure method. The provided approach is only one example. NIST AI RMF MEASURE 1.1

recommends that the most significant Al risks are determined, measured, and addressed first.

Aligned with NIST AI RMF GOVERN 1.3 and MAP 1.5.
 Aligned with ISO/IEC 42001 6.1.1.

<sup>&</sup>lt;sup>57</sup> Aligned with NIST AI RMF MEASURE 2.6.

H.		•	onents of [organization name]'s risk management process for Al as are as follows:
	1.	lmp	pact Identification with Al Impact Assessments
	2.	Ris	k Measurement
	3.	Ris	k Prioritization
	4.	Ris	k Treatment
	5.	Res	sidual Risk
	6.	Imp	pact and Risk Tracking, Inventorying, and Transparency
	7.	lmp	pact Contingency Planning and Recourse
	8.	lmp	pact Reassessment
Impact Assessment (AIIA) identifies potential impacts of an AI system			Assessment (AlIA) identifies potential impacts of an Al system through a not evaluation of current and anticipated system development, deployment, and on. <sup>58</sup>
	1.	res for	ant to guide teams to expand their impact considerations, improve their ponsible AI practices throughout the system life cycle, and determine next steps risk treatment, the AIIA can be used at different points in a system's life cycle at ying levels of specificity based on system maturity. <sup>59</sup>
		a)	A Low-Touch AllA is completed after the [e.g. Plan and Design] stage with the appropriate stakeholders, including [e.g., Al developers and representatives from the compliance team, procurement team, data engineering team, system user and/or operator groups, domain experts, socio-cultural analysts, potentially impacted individuals and groups of individuals, and responsible Al (RAI) Operational Committee]
		b)	A Medium-Touch AllA is completed during the [e.g. Verify and Validate]stage once the system has reached reliable performance, with the appropriate stakeholders, including [e.g., Al developers and representatives from the data science team, TEVV (Test, Evaluation, Verification, and Validation; e.g. red-teaming) team and experts, system user group, domain experts, socio-cultural analysts, potentially impacted individuals and groups of individuals, and responsible Al (RAI) Operational Committee] This AllA can be skipped for systems designated minimal-risk.

Aligned with NIST AI RMF MEASURE 4.1.
 Aligned with ISO/IEC 42001 B.6.1.3. ISO/IEC 42001 B.5.2 lists potential circumstances under which an Al system impact assessment should be performed, should an organization desire an Al impact assessment process independent of life cycle stage.

	c)	A High-Touch AllA is completed immediately before deployment with the appropriate stakeholders, including [e.g., Al developers and representatives from the compliance team, procurement team, MLOps team, TEVV (Test, Evaluation, Verification, and Validation; e.g. red-teaming) team and experts, system user group domain experts, socio-cultural analysts, potentially impacted individuals and group of individuals, and responsible Al (RAI) Operational Committee]				
2.	. As a tool for impact identification, the AIIA identifies potential AI incidents and the sources and outcomes (i.e. impacts). The AIIA considers the impacts on 60:					
	a)	Individuals and their [e.g., legal position, life opportunities, physical or psychological well-being]				
		(1) Relevant individuals include those [e.g. "using the AI system (users) or whose PII are processed by the AI system (data subjects)"];				
	b)	Groups of individuals and their [e.g., legal position, life opportunities, physical or psychological well-being]				
		(1) Relevant groups of individuals include [e.g., children, impaired persons, elderly persons, marginalized ethnic groups, and workers];				
	c)	Societies, which can include impacts to <sup>61</sup> :				
		(1) Environment sustainability, including [e.g., the impacts on natural resources and greenhouse gas emissions from compute, data storage, etc.]62;				
		(2) Economic opportunity and protections, including [e.g., access to financial services, employment opportunities, taxes, trade and commerce];				
		(3) Government, including [e.g., legislative processes, misinformation for political gain, national security, and criminal justice systems];				
		(4) Health and safety, including [e.g., access to healthcare, medical diagnosis and treatment, and potential physical and psychological harms];				
		(5) Norms, traditions, culture and values, including [e.g., human rights, accessibility rights, potential misinformation that leads to biases or harms to individuals]; and				
	d)	The financial and reputational health of [organization name]63				
3.	Th	e Al Impact Assessment shall:				

<sup>&</sup>lt;sup>60</sup> Aligned with ISO/IEC 42001 B.5.2 and B.5.4. <sup>61</sup> Aligned with ISO/IEC 42001 B.5.5. <sup>62</sup> Aligned with NIST AI RMF MEASURE 2.12. <sup>63</sup> Aligned with ISO/IEC 42001 6.1.2.

- a) Produce consistent, valid and comparable results across different systems and levels of specificity<sup>64</sup>;
- Assess the system according to trustworthiness characteristics, including validity and reliability; safety; security and resiliency; accountability and transparency; explainability and interpretability; privacy-enhanced; fairness with harmful bias managed<sup>65</sup>;
- c) Calibrate its assessment criteria based on various aspects of the system and its specific technical and societal context, including [e.g., expected scope of use, the data used for the development of the AI system, the AI technologies used, the thirdparty components or services used,<sup>66</sup> the functionality of the overall system, and use case context, including deployment and operating environments<sup>67</sup>] \_\_\_\_\_\_;
- d) Incorporate consultation insights and feedback on potential individual and societal impacts from external AI actors<sup>68</sup>;
- e) Continually improve its ability to identify potential AI impacts by tracking and applying historical and emerging data from [e.g., past uses of AI systems in similar contexts, public incident reports, external stakeholder feedback, results from other impact assessments<sup>69</sup>] \_\_\_\_\_\_; and
- f) Be documented, have its results retained for a defined period, and be made available to all relevant and impacted audiences.<sup>70</sup>
- 4. The AI Impact Assessment shall be compatible with, but distinct from, other existing impact or risk assessments conducted at [organization name] \_\_\_\_\_\_, including [e.g., financial impact assessments, security risk assessments, privacy risk assessments, business impact assessments] \_\_\_\_\_\_.
  - a) Where domain-specific assessments do not exist, the AIIA shall be comprehensive and domain representatives shall be involved in the design and completion of the AIIA.
  - b) Otherwise, to reduce redundancy, the AIIA can focus on AI-driven risk sources<sup>71</sup> and AI-specific risks.<sup>72</sup> Means shall be established to ensure that the assessments collectively cover the entire known AI risk landscape and that all AI-

<sup>&</sup>lt;sup>64</sup> Aligned with ISO/IEC 42001 6.1.2.

<sup>&</sup>lt;sup>65</sup> Aligned with ISO/IEC 42001 A.5.4 and NIST AI RMF 3 and MEASURE 2.5-2.11.

<sup>&</sup>lt;sup>66</sup> Aligned with NIST AI RMF GOVERN 6.1 and MAP 4.1.

<sup>&</sup>lt;sup>67</sup> Aligned with ISO/IEC 42001 6.1.4 and B.5.2. Also aligned with NIST AI RMF MEASURE 4.1.

<sup>&</sup>lt;sup>68</sup> Aligned with NIST AI RMF GOVERN 5.1, MEASURE 1.3, and MEASURE 4.1-4.2.

<sup>&</sup>lt;sup>69</sup> Aligned with NIST AI RMF MAP 5.1.

<sup>&</sup>lt;sup>70</sup> Aligned with ISO/IEC 42001 B.5.3.

<sup>&</sup>lt;sup>71</sup> Examples can be found in ISO/IEC 42001 Annex C and ISO/IEC 42001 B.5.2.

<sup>&</sup>lt;sup>72</sup> Examples can be found in NIST AI RMF Appendix B.

related risks across assessments can be efficiently aggregated for each Al system.

- J. <u>Risk Measurement</u>: Once potential impacts are identified through the AI Impact Assessment and assigned a severity score, AI risk is calculated through a likelihood analysis.
  - Likelihood and magnitude of each identified impact are determined in a standardized manner through an organizational risk and impact taxonomy delineated with common characteristics of AI systems and their use case context.<sup>73</sup>
    - a) Various factors can be used to determine the likelihood and magnitude of impacts,<sup>74</sup> including [e.g., whether the AI system is trained on large data sets composed of sensitive or protected data such as personally identifiable information; whether it uses any third-party resources or components<sup>75</sup>; whether it is designed or deployed to directly interact with humans; whether its outputs have direct impacts on humans; and how the organization itself measures risk based on its own values and context] \_\_\_\_\_\_.
    - b) [organization name] \_\_\_\_\_\_'s risk and impact taxonomy shall be reasonably aligned with applicable regulatory regimes, including [e.g. the EU AI Act] \_\_\_\_\_\_, to ease compliance efforts.
  - 2. The risk and impact taxonomy is informed and continually improved through the tracking of historical and emerging data in an **Al Incident, Impact, and Risk (IIR) database**, including from sources like [e.g., past uses of Al systems in similar contexts, public incident reports, external stakeholder feedback, results from other impact assessments<sup>76</sup>] \_\_\_\_\_\_.
  - 3. Risks or system trustworthiness characteristics that cannot be measured<sup>77</sup> and the specific reasons why<sup>78</sup> shall be properly documented.
- K. <u>Risk Prioritization</u>: Recognizing that risk management of AI systems must also be cost-effective to bolster a competitive AI strategy, [organization name] \_\_\_\_\_ also establishes the following risk priorities<sup>79</sup> and risk triage process<sup>80</sup> to enable the efficient use of resources to treat risks<sup>81</sup>:

<sup>&</sup>lt;sup>73</sup> Aligned with NIST AI RMF MAP 5.1.

<sup>&</sup>lt;sup>74</sup> Aligned with NIST AI RMF 1.2.3.

<sup>&</sup>lt;sup>75</sup> Aligned with NIST AI RMF GOVERN 6.1 and MAP 4.1.

<sup>&</sup>lt;sup>76</sup> Aligned with NIST AI RMF MAP 5.1.

<sup>&</sup>lt;sup>77</sup> Aligned with NIST AI RMF MEASURE 1.1.

<sup>&</sup>lt;sup>78</sup> NIST AI RMF 1.2.1 lists potential challenges to accurate measurement of risks.

<sup>&</sup>lt;sup>79</sup> Aligned with NIST AI RMF GOVERN 1.4.

<sup>80</sup> Aligned with NIST AI RMF GOVERN 1.4.

<sup>&</sup>lt;sup>81</sup> Aligned with NIST AI RMF 1.2.3.

- The order of risk treatment is decided first by a cost-benefit analysis at the systemlevel; systems with the highest upside potential and the lowest total cost of treatments that reduce the system risk level most effectively are placed at the top of the queue.<sup>82</sup>
  - a) Risk triage shall be performed at the level at which developer resources are selfcontained, for example, at the level of a function, segment, or division of the organization.
  - b) Risk treatment of a system must be comprehensive; while a system's risk-level is determined by its highest individual risk, all known risks of a system must be appropriately treated in one round of effort.
  - c) The order of prioritization for AI system risk treatment shall be mainly automated using metrics as follows<sup>83</sup>:
    - Highest net opportunity-to-treatment cost ratio
    - Estimated risk residual level is within deployment threshold (minimal or moderate risk)
    - Highest net opportunity/benefit-to-risk residual and monitoring cost ratio
    - Absolute highest opportunity/benefit measure
    - Absolute lowest risk residual level
- L. <u>Risk Treatment</u>: Risk treatment of a system shall consist of an appropriate and justified response to each identified risk and a system-level estimate of the resources needed to complete all risk responses collectively.<sup>84</sup>
  - Responses to a specific identified risk shall first consider possible options for avoidance of or safeguarding from the risk. Avoidance techniques intervene to reduce the likelihood that an AI incident occurs in the first place, while safeguarding techniques attempt to insulate individuals or society from harm after an AI incident has already occurred.<sup>85</sup>

a)	Other risk response options can include [e.g. "transferring the risk to another
	entity, such as a downstream enterprise, or accepting the risk if no other treatment
	action is possible"]

2.	Risk treatment shall aim to address weaknesses in a system's trustworthiness
	characteristics in a holistic, complete manner; a system is only as trustworthy as its
	weakest characteristic.86

<sup>&</sup>lt;sup>82</sup> Aligned with NIST AI RMF MANAGE 1.2.

<sup>83</sup> Aligned with NIST AI RMF MANAGE 1.2.

<sup>84</sup> Aligned with ISO/IEC 42001 6.1.3.

<sup>85</sup> Aligned with NIST AI RMF MANAGE 1.3.

<sup>&</sup>lt;sup>86</sup> Aligned with NIST AI RMF GOVERN 1.2.

- a) Progress toward any one characteristic may be dependent on or even in conflict with progress toward others. A risk treatment strategy shall balance the tradeoffs among trustworthiness characteristics for a system.
- - a) Cognizant of existing workflows and how to efficiently allocate resources, teams shall also be afforded the agency to prioritize risk treatment efforts across their own projects, as long as this is aligned with top-down risk triage directives.
- M. <u>Residual Risk</u>: Residual risk of a system, or the unmitigated risk remaining after risk treatment, shall be documented and made transparent, and resources needed to manage residual risk shall be documented.<sup>88</sup>
  - For each system, a residual risk level shall be assigned (which shall determine whether the system can be deployed) and individual residual risks shall be documented in a manner compatible for impact reassessment after potential deployment.
    - a) Residual risks shall measure the full scope of a system's potential impact, including to downstream acquirers of the system and to end users.<sup>89</sup>
  - 2. Residual risk of a deployed system shall be publicly reported to inform end users and society about potential negative impacts of the system.<sup>90</sup>
  - 3. Internal risk controls to manage residual risk of a system and of each of its components (including third-party technologies) for the duration of its operation shall be developed and documented.<sup>91</sup>
    - a) The resources needed to manage residual risk of a system for the duration of its operation shall be estimated and later measured, in the same manner as for the system's initial risk treatment.
- N. <u>Impact and Risk Tracking, Inventorying, and Transparency</u>: Materialized impacts from deployed systems shall be documented in an internal inventory and shall be broadly and publicly reported.<sup>92</sup>
  - Deployed systems shall be subject to continuous impact monitoring and regular impact measurement at a specificity-level and cadence commensurate with residual risk level and types. The process and tools (including AI actors involved, metrics and

<sup>&</sup>lt;sup>87</sup> Aligned with NIST AI RMF MANAGE 2.1.

<sup>&</sup>lt;sup>88</sup> Aligned with NIST AI RMF 1.2.3 and MANAGE 1.4.

<sup>&</sup>lt;sup>89</sup> Aligned with NIST AI RMF MANAGE 1.4.

<sup>90</sup> Aligned with NIST AI RMF 1.2.3.

<sup>&</sup>lt;sup>91</sup> Aligned with NIST AI RMF GOVERN 4.2 and MANAGE 3.1. Also aligned with ISO/IEC 42001 8.1.

<sup>92</sup> Aligned with NIST AI RMF GOVERN 4.2 and GOVERN 4.3.

<sup>©</sup> Responsible AI Institute 2024 | All Rights Reserved | Do Not Use Without Permission

	measurement of each system shall be documented and justified. Artifacts and insights shall be properly stored and made accessible.
2.	Al incidents and the impacts that may have resulted from the incident shall be entered into [organization name]'s Al Incident, Impact, and Risk (IIR) database.
	a) The database shall track existing, unanticipated, and emergent AI risks based on factors such as intended and actual performance in deployed contexts. <sup>93</sup>
	(1) Additional risk tracking approaches are considered for settings where AI risks are difficult to assess using currently available measurement techniques or where metrics are not yet available. <sup>94</sup>
	b) This database shall have a standardized entry card that aligns with [organization name]'s risk and impact taxonomy and that allows for quick reference by teams.
tir re	npact Contingency Planning and Recourse: [organization name] shall have nely and clear incident management policies, including contingency plans and course protocols, to respond to and recover from an Al impact or incident when it is entified or has occurred. <sup>95</sup>
1.	Contingency plans shall include [e.g., assigned and understood internal responsibilities, standardized protocols to quickly determine the proper response (automated when appropriate), processes to address issues related to third-party providers or buyers] <sup>96</sup> Plans are inclusive of processes to handle failures or incidents originating from third-party data or AI systems. <sup>97</sup>
2.	Al systems that display performance or outcomes inconsistent with intended use, as demonstrated through performance and safety metrics or impact reassessment results, shall be subject a proper recourse protocol, including to [e.g., supersede, disengage, or deactivate] the system. <sup>98</sup>
	a) In cases when an AI system presents unacceptable negative risk levels, such as when [e.g., significant negative impacts are imminent, harms are measurably occurring, or catastrophic risks are identified], the system shall be immediately taken offline in a safe manner. The system shall only be redeployed according to a formal protocol that determines whether the impact or risk has
<sup>94</sup> Aligned	with NIST AI RMF MEASURE 3.1. with NIST AI RMF MEASURE 3.2. with NIST AI RMF MANAGE 2.3.

technology used) established for both continuous monitoring and discrete

<sup>96</sup> Aligned with NIST AI RMF MANAGE 2.4.
<sup>97</sup> Aligned with NIST AI RMF GOVERN 6.2.
<sup>98</sup> Aligned with NIST AI RMF MANAGE 2.4.

<sup>©</sup> Responsible Al Institute 2024 | All Rights Reserved | Do Not Use Without Permission

been appropriately managed and how operation and monitoring of the system can be adjusted henceforth.99 3. In the aftermath of an AI incident, [organization name] \_\_\_\_\_ shall engage with affected users, operators, data subjects, and third parties according to the following requirements<sup>100</sup>: a) Conspicuous notifications with information relevant to the audience are delivered b) Individuals have clear and simple means to report any adverse experiences and expect a response within a reasonable time frame 101; and c) Updates are made to system characteristics or operations, and individuals are provided timely updates of any relevant changes. P. Impact Reassessment: In addition to continuous impact monitoring and regular impact measurement, an AI system shall undergo impact reassessment using a level of AI Impact Assessment (AIIA) commensurate to its risk-level in specific circumstances, including when [e.g., an incident has precipitated significant changes to system design or operations, the business or context scope of a system is significantly updated or broadened, technical components are switched out in a manner that changes third-party relationships with vendors] \_\_\_\_\_ 1. The results of impact reassessments shall be used to update system-level documentation and the Al Incident, Impact, and Risk (IIR) database. 2. If the reassessment is triggered by an Al incident, its results shall be incorporated into project retrospective reports and used to update organization-wide best practices.

#### VII. **Project Management**

A.	Existing project and/or product management at [organization name] is implemented by [governance and policies] and [tools and processes]
	The following guidance shall be incorporated into the existing project/product management framework wherever absent. 102
B.	Additionally, $[organization\ name]$ identifies and aligns with existing regulations and guidelines impacting or guiding project and product management, including $[e.g.$

<sup>&</sup>lt;sup>99</sup> Aligned with NIST AI RMF 1.2.3.

<sup>&</sup>lt;sup>100</sup> Aligned with ISO/IEC 42001 A.8.

<sup>&</sup>lt;sup>101</sup> Aligned with ISO/IEC 42001 A.8.3.

<sup>&</sup>lt;sup>102</sup> Aligned with ISO/IEC 42001 A.2.3 and NIST AI RMF GOVERN 1.2.

requirements to register the develop	oment	of foundation	models of	or other	models	with
national or global impact]	103					

- C. Project management of an AI system encompasses the processes of all or a subset of the following AI system life cycle stages, depending on the built or bought status of the system or of its components: [e.g., Plan and Design, Collect and Process Data, Build and Use Model, Verify and Validate, Deploy and Use, Operate and Monitor<sup>104</sup>] \_\_\_\_\_\_.<sup>105</sup>
  - 1. Project management leaders and [e.g. the RAI Operational Committee] \_\_\_\_\_\_ shall determine and build the tools and processes (including tech platforms or applications) necessary to operationalize cross-functional collaboration and to standardize documentation creation and collection across AI systems' life cycles.
  - Project leads/managers, with the support of relevant personnel, shall ensure that all requirements detailed in other sections of this Policy, including [e.g., Governance, Data Management, Risk Management, Stakeholder Management, Regulatory Compliance, Al Procurement] \_\_\_\_\_\_, are properly executed throughout each Al system's life cycle.

#### 3. For bought systems:

- a) Project leads and system integrator teams shall work closely with procurement teams to collect sufficient information to assess potential suppliers' responsible Al development practices and to evaluate and verify potential systems' performance across trustworthiness characteristics, in alignment with the processes outlined in Al Procurement.
- b) Project leads shall ensure that a selected (validated and deployed) system is used according to its intended uses and in alignment with documented objectives and processes for the responsible use of AI systems.<sup>106</sup>

#### 4. For built systems:

- a) Project leads shall ensure that a system is developed in accordance with documented objectives and processes for its responsible design and development.<sup>107</sup>
- b) Developer teams shall work closely with procurement teams to collect sufficient information to assess potential suppliers' responsible development practices and to evaluate and verify potential components' performance in the system

<sup>&</sup>lt;sup>103</sup> Aligned with NIST AI RMF 1.2.2.

<sup>&</sup>lt;sup>104</sup> Aligned with NIST AI RMF 2.

<sup>&</sup>lt;sup>105</sup> For more detailed guidance on management of a system across all lifecycle stages, refer to Responsible AI Institute's system-level resources.

<sup>&</sup>lt;sup>106</sup> Aligned with ISO/IEC 42001 A.9.

<sup>&</sup>lt;sup>107</sup> Aligned with ISO/IEC 42001 A.6.1.

- across trustworthiness characteristics, in alignment with the processes outlined in Al Procurement.
- c) Project leads shall ensure that the system is used according to its intended uses and in alignment with documented objectives and processes for the responsible use of AI systems.<sup>108</sup>

#### 5. For sold systems:

- a) Project leads shall integrate buyer expectations and needs into how a system is developed and can be used. 109
- D. <u>Human-AI Interaction and Configurations</u>: Project teams shall clearly define and differentiate the various human roles and responsibilities when using, interacting with, or managing AI systems.<sup>110</sup>
  - 1. Human oversight protocols shall be established and tested. 111
    - a) Oversight roles can depend on the system's degree of autonomy. More autonomous or low-risk systems may be largely maintained by system operators (responsible for maintaining infrastructure, monitoring automated performance logging, and executing contingency plans) while less autonomous or higher risk systems may require human intervention, interpretation, or manipulation as part of regular operation (i.e. a human-in-the-loop).
  - Evaluations of whether users (and/or end users), defined as [e.g. "relevant interested parties who make decisions or are subject to decisions based on the AI system outputs"] \_\_\_\_\_\_, can adequately interpret the AI system outputs shall be regularly conducted to inform the design of system output interfaces, communication methods, and user documentation.<sup>112</sup>
  - 3. Design of sites of human-Al system interaction (such as between the system and operator or between the system and users) shall be informed by research and consultation on human biases and individual preferences, traits, and skills that may influence interpretation of system outputs and generate or amplify harms.<sup>113</sup>
- E. <u>DEI and Stakeholder Engagement</u>: Activities across a system life cycle, including [e.g., AI system design and development<sup>114</sup>, evaluation<sup>115</sup>, performance monitoring<sup>116</sup>, and impact

<sup>108</sup> Aligned with ISO/IEC 42001 A.9.

<sup>&</sup>lt;sup>109</sup> Aligned with ISO/IEC 42001 A.10.4.

<sup>&</sup>lt;sup>110</sup> From NIST AI RMF Appendix C.

<sup>&</sup>lt;sup>111</sup> Aligned with ISO/IEC 42001 B.9.3 and NIST AI RMF GOVERN 3.2, MAP 2.2, and MAP 3.5.

<sup>&</sup>lt;sup>112</sup> Aligned with ISO/IEC 42001 B.6.2.4.

<sup>&</sup>lt;sup>113</sup> From NIST AI RMF Appendix C.

<sup>&</sup>lt;sup>114</sup> Aligned with NIST AI RMF GOVERN 5.2.

<sup>&</sup>lt;sup>115</sup> Aligned with NIST AI RMF MEASURE 3.3.

<sup>&</sup>lt;sup>116</sup> Aligned with NIST AI RMF MEASURE 4.3.

<sup>©</sup> Responsible AI Institute 2024 | All Rights Reserved | Do Not Use Without Permission

assessments<sup>117</sup>] \_\_\_\_\_\_, should all be **informed by a representatively diverse immediate team** (e.g., diversity of demographics, disciplines, experience, expertise, and backgrounds)<sup>118</sup> and by **consultation with and feedback from internal and external Al actors**.<sup>119</sup> Refer to Workforce Management for more on DEI guidance and to Stakeholder Management and Engagement on consultation and feedback mechanisms.

- F. <u>Managing Adaptation and Drift</u>: Project owners shall implement mechanisms to sustain the value of deployed AI systems against unintended scope drift. 120
  - 1. Systems and their components (e.g., data for fine-tuning, pre-trained models<sup>121</sup>) shall be continually monitored for drift with respect to characteristics, quality, suitability, and behavior. Undesired drift with respect to performance criteria and other documented systems requirements shall be corrected and measures shall be applied to prevent future drift. All other relevant types of drift, including increased capabilities from adaptive learning, shall be documented once identified.
  - 2. Project owners, with support of [e.g. the RAI Operational Committee] \_\_\_\_\_\_, shall proactively determine the current system's boundaries for business use case scope and technical scope. Change or expansion beyond these boundaries requires [e.g., an impact assessment reassessment, a new project proposal] \_\_\_\_\_\_.
- G. <u>System-level Documentation</u>: As artifacts of a system's progression across life cycle stages, system-level documentation shall be created and maintained, including:
  - 1. Entries into enterprise-wide inventories, including the **AI system inventory**<sup>122</sup> and [e.g., Data inventory, Project inventory, AI Incident, Impact, and Risk (IIR) database<sup>123</sup>]
    - a) Entries into inventories for resource documentation and management can include for [e.g., AI system components, data resources<sup>124</sup>, tooling resources<sup>125</sup>, system and computing resources<sup>126</sup>, human resources<sup>127</sup>] \_\_\_\_\_. <sup>128</sup> Resource

<sup>&</sup>lt;sup>117</sup> Aligned with NIST AI RMF MAP 5.2 and MEASURE 1.3.

<sup>&</sup>lt;sup>118</sup> Aligned with NIST AI RMF GOVERN 3.1.

<sup>&</sup>lt;sup>119</sup> Aligned with NIST AI RMF GOVERN 5.1.

<sup>&</sup>lt;sup>120</sup> Aligned with NIST AI RMF MANAGE 2.2.

<sup>&</sup>lt;sup>121</sup> Aligned with NIST AI RMF MANAGE 3.2.

<sup>&</sup>lt;sup>122</sup> Aligned with NIST AI RMF GOVERN 1.6.

<sup>&</sup>lt;sup>123</sup> See Appendix A for more examples.

<sup>&</sup>lt;sup>124</sup> See ISO/IEC 42001 B.4.3 for a list of potential documentation on data resources utilized for the AI system.

<sup>&</sup>lt;sup>125</sup> See ISO/IEC 42001 B.4.4 for a list of potential documentation on tooling resources utilized for the Al system.

<sup>&</sup>lt;sup>126</sup> See ISO/IEC 42001 B.4.5 for a list of potential documentation on system and computing resources utilized for the AI system.

<sup>&</sup>lt;sup>127</sup> See ISO/IEC 42001 B.4.5 for a list of potential documentation on human resources utilized for the AI system.

<sup>&</sup>lt;sup>128</sup> Aligned with ISO/IEC 42001 B.4.2.

documentation shall be incrementally updated across each life cycle stage and be used to inform future risk triage, project prioritization, and resource allocation efforts.

- 2. Business use case documentation that can describe, but is not limited to:
  - a) Intended purpose, business value, and context(s) of business use<sup>129</sup>;
    - (1) Potential benefits of intended AI system uses, capabilities, and performance. 130
    - (2) Laws, norms and expectations specific to the prospective settings in which the AI system will be deployed.<sup>131</sup>
  - b) Targeted application scope, based on the system's capability, established context, and AI system categorization<sup>132</sup>;
    - (1) Viable non-Al alternative systems, approaches, or methods based on the system's intended purpose and scope.<sup>133</sup>
  - c) [organization name] \_\_\_\_\_\_'s role and activities with respect to the system's intended purpose, in alignment with regulatory definitions such as [e.g., Al Provider, Deployer, Importer, Distributor] \_\_\_\_\_\_ and [e.g., making available on the market, putting into service] \_\_\_\_\_\_, respectively. 134
  - d) System requirements, including for new AI systems or material enhancements to existing systems<sup>135</sup>;
  - e) Assumptions and related limitations about AI system's purpose or knowledge and how system output may be utilized and overseen by humans<sup>136</sup>;
  - f) Potential risks across the development or product AI life cycle<sup>137</sup>; and
    - (1) Potential costs, including non-monetary costs, which result from expected or realized AI errors or system functionality and trustworthiness<sup>138</sup>; and
    - (2) Potential positive and negative impacts of system uses to individuals, communities, organizations, society, and the planet. 139

<sup>&</sup>lt;sup>129</sup> From NIST AI RMF MAP 1.4.

<sup>&</sup>lt;sup>130</sup> From NIST AI RMF MAP 3.1.

<sup>&</sup>lt;sup>131</sup> From NIST AI RMF MAP 1.1.

<sup>&</sup>lt;sup>132</sup> From NIST AI RMF MAP 3.3.

<sup>&</sup>lt;sup>133</sup> From NIST AI RMF MANAGE 2.1.

<sup>&</sup>lt;sup>134</sup> Aligned with ISO/IEC 42001 4.1.

<sup>&</sup>lt;sup>135</sup> Aligned with NIST AI RMF MAP 1.6 and ISO/IEC 42001 B.6.2.2.

<sup>&</sup>lt;sup>136</sup> From NIST AI RMF MAP 1.1 and MAP 2.2.

<sup>&</sup>lt;sup>137</sup> From NIST AI RMF MAP 1.1.

<sup>&</sup>lt;sup>138</sup> From NIST AI RMF MAP 3.2.

<sup>&</sup>lt;sup>139</sup> From NIST AI RMF MAP 1.1.

- g) Specific set or types of users and their expectations 140;
- 3. **System technical documentation** that is made appropriate for each relevant or interested party, including users, partners, and supervisory authorities.<sup>141</sup> The documentation can describe, but is not limited to<sup>142</sup>:
  - a) Technical assumptions and limitations (e.g. related to system interoperability, run-time environment, data quality, AI explainability)<sup>143</sup>;
  - b) Human-Al configurations, including operator, user, or human-in-the-loop capabilities and instructions<sup>144</sup>, and configuration evaluations that meet applicable requirements (including human subject protection) and are representative of the relevant population<sup>145</sup>;
  - c) The system's intended capabilities and the implementation methods and architecture (e.g., classifiers, generative models, recommenders)<sup>146</sup>;
  - d) Test, evaluation, validation, and verification (TEVV) and system metrics, test sets, and tools calibrated to specific scientific integrity considerations, including [e.g. those related to experimental design, data collection and selection (e.g., availability, representativeness, suitability), and construct validation.<sup>147</sup>] \_\_\_\_\_\_\_<sup>148</sup>;
  - e) Results from an internal audit or technical assurance process to enable impartial system evaluation by those outside of the developer team, if deemed necessary<sup>149</sup>;
  - Release criteria requirements, including acceptable ranges for operational factors and performance errors, and any acceptable factors that affect a system's ability to reach minimum release criteria<sup>150</sup>;
  - g) Qualitative or quantitative performance or assurance criteria, calibrated to deployment setting(s)<sup>151</sup>;
  - h) Roles and processes for the responsible operation of the AI system, including for monitoring and improvements<sup>152</sup> (including of pre-trained models)<sup>153</sup>, appeal and

<sup>&</sup>lt;sup>140</sup> From NIST AI RMF MAP 1.1.

<sup>&</sup>lt;sup>141</sup> See ISO/IEC 42001 B.6.2.7 for more details on potential technical documentation elements.

<sup>&</sup>lt;sup>142</sup> Aligned with ISO/IEC 42001 B.6.2.4 and B.6.2.6.

<sup>&</sup>lt;sup>143</sup> Aligned with ISO/IEC 42001 B.6.2.7.

<sup>&</sup>lt;sup>144</sup> Aligned with ISO/IEC 42001 B.6.2.7.

<sup>&</sup>lt;sup>145</sup> From NIST AI RMF MEASURE 2.2.

<sup>&</sup>lt;sup>146</sup> From NIST AI RMF MAP 2.1.

<sup>&</sup>lt;sup>147</sup> From NIST AI RMF MAP 2.3.

<sup>&</sup>lt;sup>148</sup> Aligned with NIST AI RMF GOVERN 4.3, MAP 1.1, and MEASURE 2.1.

<sup>&</sup>lt;sup>149</sup> Aligned with ISO/IEC 42001 9.2.2.

<sup>&</sup>lt;sup>150</sup> Aligned with ISO/IEC 42001 B.6.2.4.

<sup>&</sup>lt;sup>151</sup> From NIST AI RMF MEASURE 2.3.

<sup>&</sup>lt;sup>152</sup> Aligned with NIST AI RMF MANAGE 4.2.

<sup>&</sup>lt;sup>153</sup> Aligned with NIST AI RMF MANAGE 3.2.

override, decommissioning<sup>154</sup>, incident management<sup>155</sup> (and enabling systems to fail safely)<sup>156</sup>, recovery, and change management<sup>157</sup>:

		, , , , , , , , , , , , , , , , , , ,
4.	spe wit out sys	mplete, up-to-date, and accurate <b>user documentation</b> , including technical ecifications but also general notification and information about their interaction in an AI system, including [e.g., necessary information to properly interpret system to the system of system's accuracy and performance, disclosure of recent extern updates, incidents, or impacts, contact information, means to report feedback tharms, links for additional informational materials], presented in an dessible and understandable manner. 159
5.		cumentation of processes, decisions, and results across each life cycle stage, h justifications of each, including of:
	a)	<b>Design choices</b> made during [e.g. "the Plan and Design, Collect and Process Data, and Build and Use Model"] stage(s), including [e.g., the choice between models, machine learning architecture and methods, data sets, infrastructure options for compute and interoperability, response to security threats, user or output interface design, human-Al configuration]160;
	b)	Measures for <b>verification and validation</b> of the system during [e.g. "the Verify and Validate"] stage(s), including [e.g. the process and results of evaluating the system with respect to the environmental impact and sustainability of model training and management activities <sup>161</sup> and across each trustworthiness characteristic <sup>162</sup> , in support of but also beyond the AI Impact Assessment process] <sup>163</sup> ;
	c)	Design and implementation of a <b>deployment plan</b> for [e.g. "the Deploy and Use"]  stage(s) that [e.g. details roles and required approvals, a timeline, and a strategy for testing and feedback (e.g., phased roll-outs, pilots, beta or full release)  tailored to deployment contexts and requirements and to the risk profile of the system <sup>164</sup> ; and
	d)	Measures for post-deployment system <b>operation and monitoring</b> during [e.g. "the Operate and Monitor"] stage(s), including [e.g., event logs <sup>165</sup> , monitoring
าed าed	with with	NIST AI RMF GOVERN 1.7. NIST AI RMF MANAGE 4.3. NIST AI RMF MEASURE 2.6. AI RMF MANAGE 4.1.

<sup>154</sup> Align

<sup>155</sup> Align

<sup>156</sup> Align

<sup>&</sup>lt;sup>157</sup> From

<sup>&</sup>lt;sup>158</sup> Aligned with NIST AI RMF MAP 2.2.

<sup>&</sup>lt;sup>159</sup> Aligned with ISO/IEC 42001 B.8.2.

<sup>&</sup>lt;sup>160</sup> See ISO/IEC 42001 B.6.2.3 for a list of design choices to be documented.

<sup>&</sup>lt;sup>161</sup> Aligned with NIST AI RMF MEASURE 2.13.

 $<sup>^{162}</sup>$  See Al Principles. Aligned with NIST AI RMF MEASURE 2.5-2.12 and NIST AI RMF 3.

<sup>&</sup>lt;sup>163</sup> Aligned with ISO/IEC 42001 9.1 and A.6.2.4.

<sup>&</sup>lt;sup>164</sup> Aligned with ISO/IEC 42001 B.6.2.5.

<sup>&</sup>lt;sup>165</sup> Aligned with ISO/IEC 42001 B.6.2.8.

		logs of system behavior and user input <sup>166</sup> , incident reports, impact reassessment results] <sup>167</sup>
H.	Operation best post databaseduca	pective Learning and Integration: Project owners, with support of [e.g. the RAI tional Committee], shall direct efforts to integrate lessons learned and ractices from AI projects (attempted or deployed) into relevant enterprise-level ases (e.g., on unexpected risks, into a bank of red-teaming prompts) and tional resources, to update and augment the AI and responsible AI capabilities of ization name]'s workforce for future AI projects.
•	/III.	Stakeholder Management and Engagement
A.	impler	ig stakeholder management and engagement at [organization name] is nented by [governance and policies] and [tools and processes] The following guidance shall be incorporated into the existing stakeholder gement strategy wherever absent. 168
B.	throug their li	ors, both internal and external to [organization name], shall be engaged the regular and formalized processes to inform and improve AI systems throughout fe cycles, the knowledge of and responses to risks and impacts from systems, rganization name]'s AI management system as a whole.
C.		ging and engaging stakeholders (internal and external AI actors) consist of the ing activities:
	rep hu exp	nsultation and Feedback <sup>169</sup> : All relevant and interested parties, including [e.g., presentatives of internal functions like Legal or Procurement, the intended user and man-in-the-loop groups, potentially impacted groups; and domain and socio-cultural poerts], shall be regularly consulted and solicited for feedback roughout a system's life cycle.
	a)	Consultation and feedback processes shall be formal and compulsory but tailored to the context and risk level of the system. Those invited to participate in consultations or to provide feedback shall be identified through a balanced and accessible recruitment strategy that prioritizes demographic diversity and broad domain and user experience expertise. <sup>170</sup> Participants shall also be compensated for their efforts except in proper extenuating circumstances.
	b)	The needs, expectations, concerns, and experienced impacts of those consulted or providing feedback shall be used to inform system requirements and to guide
_		NIST AI RMF MEASURE 2.4 and MANAGE 4.1. ISO/IEC 42001 9.1 and B.6.2.6.

166

Aligned with ISO/IEC 42001 A.2.3 and NIST AI RMF GOVERN 1.2.

 $<sup>^{\</sup>rm 169}$  Aligned with NIST AI RMF GOVERN 5.1 and MAP 1.2.

<sup>&</sup>lt;sup>170</sup> Aligned with NIST AI RMF MAP 1.2.

development and improvement of the system, including related to [e.g., system design and implementation<sup>171</sup>; performance<sup>172</sup>; system trustworthiness<sup>173</sup>; assessments<sup>174</sup>; potential risks and impacts<sup>175</sup>; continual improvements<sup>176</sup>]

c) Voluntary feedback mechanisms shall also be made accessible for all interested parties who wish to provide feedback at any time. Such mechanisms shall be directly available from a system's output interface and provide options for different types of feedback (e.g., rating of output quality, corrective actions for RLHF, options to send a message, escalation to a report of concerns or impacts (see <u>Reporting and Response</u>)).

2.	Notification and Communication 178: [organization name] shall provide
	necessary, accessible, and appropriate information to all relevant and interested
	parties, including [e.g., representatives of internal functions like Legal or Procurement,
	of the intended user and human-in-the-loop groups, of potentially impacted groups;
	domain and socio-cultural experts; policymakers and other civil society actors; and
	third-party assessors]

- a) Plans for notifications and other communication of information about systems, including [e.g., that one is interacting with or will be subject to a decision made (in part) by a system; user documentation<sup>179</sup>; incidents, impacts, and risks<sup>180</sup>]

   shall be proactive, timely, and meet legal, regulatory, and stakeholder-driven transparency obligations.
- 3. Reporting and Response<sup>181</sup>: All relevant and interested parties, including [e.g., representatives of internal functions like Legal or Procurement, of the intended user and human-in-the-loop groups, of potentially impacted groups; domain and sociocultural experts; policymakers and other civil society actors; and third-party assessors] \_\_\_\_\_\_, shall be be provided mechanisms to report concerns and impacts related to a system and receive corrective action from [organization name] \_\_\_\_\_\_.
  - a) Such a mechanism shall [e.g., protect individuals from identification and reprisals; be accessible to all workforce members; have appropriate personnel and

<sup>&</sup>lt;sup>171</sup> Aligned with NIST AI RMF GOVERN 5.2.

<sup>&</sup>lt;sup>172</sup> Aligned with NIST AI RMF MEASURE 4.3.

<sup>&</sup>lt;sup>173</sup> Aligned with NIST AI RMF MEASURE 4.2.

<sup>&</sup>lt;sup>174</sup> Aligned with NIST AI RMF MEASURE 1.3.

<sup>&</sup>lt;sup>175</sup> Aligned with ISO/IEC 42001 B.5.4 and NIST AI RMF MAP 5.2 and MEASURE 4.1.

<sup>&</sup>lt;sup>176</sup> Aligned with NIST AI RMF MANAGE 4.2.

<sup>&</sup>lt;sup>177</sup> Aligned with ISO/IEC 42001 4.2.

<sup>&</sup>lt;sup>178</sup> Aligned with ISO/IEC 42001 7.4 and A.8 and NIST AI RMF GOVERN 4.3.

<sup>&</sup>lt;sup>179</sup> Aligned with ISO/IEC 42001 A.8.2.

<sup>&</sup>lt;sup>180</sup> Aligned with ISO/IEC 42001 B.8.5 and NIST AI RMF MANAGE 4.3.

<sup>&</sup>lt;sup>181</sup> Aligned with ISO/IEC 42001 B.8.4 and NIST AI RMF MEASURE 1.2 and MEASURE 3.3.

			capabilities (including investigation, resolution, escalation, and reporting powers); respond and act in a timely manner <sup>182</sup> ];
		b)	Reports shall be used to inform system requirements and to improve the system and [organization name]'s overall AI management system. Any individuals who find the response or remedy provided for their report insufficient have means to note their public dissatisfaction and escalate their issue.
	4.	be an	blic and Ecosystem Engagement: In the spirit of continual learning and neficence, [organization name] shall contribute to collaborative, shared, d open-source initiatives with the public and with the broader AI ecosystem to smote the responsible development, use, procurement, and supply of trustworthy
		a)	Activities can include [e.g., creating thought leadership; leveraging network and marketplace synergies to create stronger technological or informational resources (e.g., developing an industry-specific bias evaluation toolkit, sharing common lessons from the AI Incident, Impact, and Risk (IIR) database; providing insights to shared learning spaces like working groups]
	IX	ζ.	Workforce Management
A.	[go	overi idar	ig workforce management at [organization name] is implemented by nance and policies] and [tools and processes] The following ace shall be incorporated into the existing workforce management framework wer absent. 183
A.	[gd gu wh	idar idar nere <u>Div</u>	nance and policies] and [tools and processes] The following are shall be incorporated into the existing workforce management framework
A.	[gd gu wh	idar idar nere <u>Div</u>	nance and policies] and [tools and processes] The following note shall be incorporated into the existing workforce management framework over absent. 183  Versity, Equity, and Inclusion (DEI) shall be integrated into every component of
Α.	[go gu wh	idar idar nere <u>Div</u> [or a)	nance and policies] and [tools and processes] The following ace shall be incorporated into the existing workforce management framework wer absent.     Versity, Equity, and Inclusion (DEI)   Shall be integrated into every component of ganization name] 's AI strategy.     DEI personnel at [organization name], if present, shall work with [e.g. the RAI Operational Committee] to integrate DEI objectives and requirements into responsible AI efforts. DEI experts shall also be represented in cross-functional stakeholder processes for AI, including [e.g., AI impact]

<sup>182</sup> Aligned with ISO/IEC 42001 B.3.3.
183 Aligned with ISO/IEC 42001 A.2.3 and NIST AI RMF GOVERN 1.2.
184 Aligned with NIST AI RMF GOVERN 3.1 and MAP 1.2.

		c)	All interactions between Al actors shall be conducted with mutual respect and foster an environment of belonging that values unique and dissenting perspectives about [organization name]
		d)	Existing DEI policies, initiatives, and resources, including [e.g., Employee Resource Groups (ERGs), development programs, events for underrepresented employees] shall be updated or augmented as necessary to support employees in
			Al-specific upskilling or in Al-specific equity concerns.
		e)	The use of AI at [organization name] shall not violate existing DEI commitments nor regulatory and legal requirements (e.g., discrimination, harassment, and equal employment opportunity), for example, in the context of [e.g., employee performance evaluations, productivity tracking using computer vision]
В.	<u>Em</u>	plo	yee Awareness and Competence <sup>185</sup> : Employees at [organization name]
	res	por	evelop the necessary awareness of and competence for their AI role and asibilities through training, educational resources, and other upskilling initiatives nented by [e.g. the RAI Operational Committee]
	1.	res res	ployees and relevant partners of [organization name] shall receive ponsible AI training to enable them to understand and perform their duties with spect to the AI management system outlined in this Policy and to other related licies and processes. <sup>186</sup>
		a)	Employees shall demonstrate adherence to responsible use policies and processes, including only using approved systems (to avoid "shadow AI" risks) and only using systems as intended by system documentation. <sup>187</sup>
	2.	to de:	ucational and RAI outreach efforts within [organization name] shall help cultivate a RAI culture, including a critical thinking and safety-first mindset in the sign, development, deployment, and use of AI systems, 188 and can include [e.g., mmunities of Practice for function- or role-specific RAI collaboration, a Center of cellence to centralize RAI guidance and resources]
	3.	ed	ployees shall be provided with role- and/or system-specific training and ucational tools to enable and evaluate their understanding and fulfillment of RAI uirements for a system. 189
		a)	Employees interacting with a system (e.g., as a user, operator, human-in-the-loop) shall demonstrate proficiency with respect to knowledge of relevant

 $<sup>^{\</sup>rm 185}$  Aligned with ISO/IEC 42001 7.2 and 7.3.

<sup>&</sup>lt;sup>186</sup> Aligned with NIST AI RMF GOVERN 2.2.

187 Aligned with ISO/IEC 42001 B.9.2 and B.9.4.

<sup>&</sup>lt;sup>188</sup> From NIST AI RMF GOVERN 4.1.

<sup>&</sup>lt;sup>189</sup> Aligned with NIST AI RMF MAP 1.6.

system information and to the ability to responsibly execute and complete tasks. 190

C. <u>Hiring for Responsible AI</u>: [organization name] \_\_\_\_\_\_ shall identify deficiencies in AI and RAI capacity and skills, identify which gaps can be addressed by training and which are relevant to hiring, and seek new talent as necessary. RAI workforce planning shall be informed by human resource inventorying activities and documentation.<sup>191</sup>

## X. Regulatory Compliance

A.	Compliance	with existing data, analytics, and technology regulation	at [organization
	name]	is implemented by [governance and policies]	and [tools and
	processes] _	The following guidance for compliance for Al	systems shall be
	incorporated	into existing compliance activities wherever absent. 192	

- B. The compliance team shall actively monitor and understand all existing and emerging legal and regulatory requirements relevant to all AI activities.<sup>193</sup>
- C. The compliance team shall map the compliance requirements of each AI system, or otherwise develop tools and processes, such as checklists or compliance meetings, to enable technical teams to determine requirements at each stage of a system's life cycle.
- D. The compliance team shall be ultimately accountable for the perceived or realized compliance of every system at each stage of the life cycle, from initial project scoping to post-deployment or post-sale to potential decommission, and shall be provided with timely and sufficient transparency to enable this responsibility.
- E. Compliance team members shall receive appropriate foundational legal and AI training, developed and distributed by [e.g. the RAI Operational Committee] \_\_\_\_\_\_, to enable successful implementation of compliance organization-wide. 194
- F. The compliance team shall develop role- or system-specific compliance training and other guidance materials for the workforce, with the support of function leadership or system owners, respectively.

<sup>&</sup>lt;sup>190</sup> Aligned with NIST AI RMF MAP 3.4.

<sup>&</sup>lt;sup>191</sup> Aligned with ISO/IEC 42001 B.4.6. Refer to Project Management (System-level Documentation) for more on resource documentation.

<sup>&</sup>lt;sup>192</sup> Aligned with ISO/IEC 42001 A.2.3 and NIST AI RMF GOVERN 1.2.

<sup>&</sup>lt;sup>193</sup> Aligned with NIST AI RMF GOVERN 1.1.

<sup>&</sup>lt;sup>194</sup> Aligned with NIST AI RMF GOVERN 2.2.

#### XI. Al Procurement

A.	Existing procure	ment at [organization name]	is implemented by [governance
	and policies]	and [tools and processes]	The following guidance
	shall be incorpor	rated into the existing procurement str	ategy wherever absent.195

- B. The procurement team shall engage with product teams and intended users of procured products/services in the earliest stage of use case scoping, and shall develop standardized processes to initiate and manage such engagements, such as [e.g., ticket submissions for procurement requests] \_\_\_\_\_\_.
- C. The procurement team shall identify and document all external parties involved during the life cycle of an AI system, including all partners, suppliers, and buyers. 196
- D. Procurement team members shall receive appropriate foundational compliance and Al training and maintain continual and informed use of risk management tools, including [e.g., all levels of the Al Impact Assessment, the Al Incident, Impact, and Risk (IIR) database, the Al resources inventory] \_\_\_\_\_\_\_.<sup>197</sup>
- E. Procurement teams shall support the creation of analyses for third-party resources in the cost-risk/benefit repository, and shall use the repository as a resource to inform future procurement efforts.<sup>198</sup>
- F. In the case of products or services that do not employ AI directly but may be delivered in part by AI used internally by the supplier, the procurement team shall subject every potential supplier to a **Responsible Supplier Assessment** to assess the responsible AI maturity of the supplier organization. Results from the Responsible Supplier Assessment shall inform the decision to accept the supplier.
  - If an existing procured product or service that does not employ AI directly but may be newly delivered in part by AI used internally by the suppliers, the procurement team shall initiate a Responsible Supplier Assessment with the supplier. Results from the Responsible Supplier Assessment shall inform the decision to maintain a relationship with the supplier.

#### G. For built systems<sup>199</sup>:

 The procurement team shall clarify with the compliance team, product team, and intended users (if internal) the purpose, technical needs, and legal requirements of an AI component (e.g., data sets, AI models, platforms) in its use case to develop a must-have list for potential suppliers.

<sup>&</sup>lt;sup>195</sup> Aligned with ISO/IEC 42001 A.2.3 and NIST AI RMF GOVERN 1.2.

<sup>&</sup>lt;sup>196</sup> Aligned with ISO/IEC 42001 A.10.2.

<sup>&</sup>lt;sup>197</sup> Aligned with NIST AI RMF GOVERN 2.2.

<sup>&</sup>lt;sup>198</sup> Aligned with NIST AI RMF MANAGE 3.1.

<sup>&</sup>lt;sup>199</sup> Aligned with ISO/IEC 42001 B.10.3.

<sup>©</sup> Responsible AI Institute 2024 | All Rights Reserved | Do Not Use Without Permission

- 2. The procurement team shall consider multiple possible suppliers during a process, while understanding that it is possible that none of the suppliers will fulfill all technical and responsible AI needs, and therefore, none may be accepted.
- The procurement team shall subject every potential supplier to a Responsible Supplier Assessment (for AI products like models) or to a standard supplier assessment (for other components), depending on the development or use of AI to deliver the product.
- 4. For every potential component, the procurement team shall request sufficient information from the potential supplier to conduct a Low-Touch Al Impact Assessment of the proposed system, including the component. If the information is refused, the procurement team shall not move forward with the supplier.
- 5. The procurement team shall determine what documentation the potential supplier must provide for the responsible procurement and integration of the component. This includes [e.g., data documentation, training process, change policies] \_\_\_\_\_\_.
- The procurement team shall establish absolute thresholds for the level of responsible AI maturity required, level of AI risk allowed, and type of documentation required for all suppliers, which shall determine whether an individual prospective supplier is accepted.
- 7. If the potential supplier is accepted, the procurement team shall negotiate and document terms of use with the supplier, including [e.g., delineation of liability and responsibility for corrective action in downstream impacts, data sharing, deletion, and privacy agreements, force majeure clause, termination clause] \_\_\_\_\_\_.
- 8. Upon acceptance of the terms and conditions, the procurement team shall clarify and document ongoing communication and transparency processes with the supplier, including [e.g., timely notification from suppliers when the component is changed or an unmitigated risk has emerged, high-level performance monitoring reports to the supplier in the case of Al models, a feedback or dispute mechanism with the supplier] \_\_\_\_\_\_.
- 9. The procurement team shall share information about the component and supplier with all relevant functions and audiences, including [e.g., compliance teams, product teams, the RAI Operational Committee, the user group] \_\_\_\_\_\_ to ensure responsible use and integration of the component into built systems.

#### H. For bought systems<sup>200</sup>:

1. The procurement team shall clarify with the compliance team, product team, and intended users (if internal) the purpose, technical needs, and legal requirements of the use case to develop a must-have list for potential suppliers.

<sup>&</sup>lt;sup>200</sup> Aligned with ISO/IEC 42001 B.10.3.

<sup>©</sup> Responsible AI Institute 2024 | All Rights Reserved | Do Not Use Without Permission

- The procurement team shall consider multiple possible suppliers during a process, while understanding that it is possible that none of the suppliers will fulfill all technical and responsible AI needs, and therefore, none may be accepted.
- The procurement team shall subject every potential supplier of an AI system to a Responsible Supplier Assessment to assess the responsible AI maturity of the supplier organization.
  - a) If an existing procured tool or application has rolled out new AI capabilities, the procurement team shall initiate a Responsible Supplier Assessment with the supplier.
- 4. The procurement team shall request sufficient information from the potential supplier to conduct a Medium-Touch AI Impact of the system in the context of its intended use, and if refused, shall not move forward with the supplier.
  - a) If an existing procured tool or application has rolled out new Al capabilities, the procurement team shall conduct its own High-Touch Al Impact Assessment, requesting information from the supplier as needed.
- 5. The procurement team shall determine what system documentation the potential supplier must provide for the responsible procurement and use of the system. This includes [e.g., the supplier's Al risk or impact assessment, data documentation, summary of incidents, change policies] \_\_\_\_\_\_.
  - a) If an existing procured tool or application has rolled out new AI capabilities, the procurement team shall request additional system documentation, including [e.g., the supplier's AI risk or impact assessment, data documentation, training process, summary of incidents, change policies] \_\_\_\_\_\_.
- The procurement team shall establish absolute thresholds for the level of responsible AI maturity required, level of AI risk allowed, and type of documentation required for all suppliers, which shall determine whether an individual prospective supplier is accepted.
  - a) If an existing procured tool or application has rolled out new AI capabilities, the procurement team shall establish absolute thresholds for the level of responsible AI maturity required, level of AI risk allowed, and type of documentation required to retain the product. All relationships with suppliers who do not meet the required thresholds shall be terminated, within reason.
- 7. If the potential supplier is accepted, the procurement team shall negotiate and document terms of use with the supplier, including [e.g., delineation of liability and responsibility for corrective action in downstream impacts; data sharing, deletion, and

		privacy agreements; PII processor or controller role <sup>201</sup> ; force majeure clause; termination clause]
		a) If an existing procured tool or application has rolled out new AI capabilities and the procurement team has determined that the product will be retained, the procurement team shall review existing terms of use with the supplier and negotiate and document any updates needed for responsible procurement and use.
	8.	If the terms and conditions are agreed upon, the procurement team shall clarify and document ongoing communication and transparency processes with the supplier, including [e.g., prompt notification from suppliers when the system is changed or an unmitigated risk has emerged, high-level performance monitoring reports to the supplier, a feedback or dispute mechanism with the supplier]
		a) If an existing procured tool or application has rolled out new AI capabilities and the procurement team has determined that the product will be retained, the procurement team shall review ongoing communication and transparency processes with the supplier and negotiate and document any updates needed for responsible procurement and use.
	9.	The procurement team shall share information about the system and supplier with all relevant functions and audiences, including [e.g., compliance teams, product teams, the RAI Operational Committee, the user group] to ensure responsible use and integration of the system and enable the development of system-specific guidance.
I.	Fo	r sold systems <sup>202</sup> :
	1.	If the system is intended to be sold, the procurement team shall additionally determine what information shall be requested from suppliers to enable downstream transparency with buyers.
XII.	Do	ocumentation Management
A.	[go	isting documentation processes at [organization name] are implemented by overnance and policies] and [tools and processes] The following idance shall be incorporated into existing documentation practices wherever sent. <sup>203</sup>
В.		ditionally, [organization name] identifies and aligns with existing regulations d guidelines for documentation and reporting for AI systems, including [e.g. "reporting
<sup>202</sup> Alig	ned	with ISO/IEC 42001 A.10.2. with ISO/IEC 42001 B.10.3 and B.10.4. with ISO/IEC 42001 A.2.3 and NIST AI RMF GOVERN 1.2.

<sup>©</sup> Responsible Al Institute 2024 | All Rights Reserved | Do Not Use Without Permission

	Sai	ults of safety tests of high-risk models, as required by the U.S. Executive Order on the fe, Secure, and Trustworthy Development and Use of Artificial Intelligence"]   204
C.	ma	ganization name] identifies the following roles as responsible for proper intenance and control of documentation relevant to its management of AI systems d its broader AI strategy <sup>205</sup> :
	1.	The [e.g. RAI Operational Committee] shall maintain [e.g. "all enterprise-level AI-specific policy documents and artifacts"] This includes [e.g., charter documents for new AI bodies such as the Steering/Operational Committees and working groups, AI component inventories, and databases for AI impacts], among others;
	2.	The [e.g. "business lead, with support of the technical lead"] shall maintain [e.g. "documentation for each AI project and/or system they oversee"] This includes [e.g., business use case documentation, system technical documentation, AI impact assessment results, and approval gate documentation], among others;
	3.	The [e.g., function leaders, owners of specific risk areas] shall maintain [e.g. "all related documentation, including policies and system-specific artifacts"] This includes [e.g., data management policies and standards, governance for AI procurement, and privacy risk assessment results], among others; and
	4.	The [role or group] shall maintain [documentation category]  This includes [examples of documents], among others.
D.		intenance and control of documentation encompasses the following ponsibilities <sup>206</sup> :
	1.	Enablement and review of the creation or update (i.e. version control) of documentation by the proper AI actors;
	2.	Storage and preservation of documents' legibility and suitability for use, including through [features e.g. "identification and description, format, and media" 207];
	3.	Enforcement of access controls for permissions to view, change, retrieve, use, or distribute documentation, to enable traceability and meet reporting or audit trail requirements; and

<sup>&</sup>lt;sup>204</sup> Aligned with NIST AI RMF 1.2.2. <sup>205</sup> See Appendix A for a list of documents under each identified documentation category. In accordance with ISO/IEC 42001 7.5.3, organizations shall also identify and control necessary external documentation.  $^{206}$  Aligned with ISO/IEC 42001 7.5.3.  $^{207}$  See ISO/IEC 42001 7.5.2 for examples of each listed feature.

4. Retention and disposition, in accordance with regulatory requirements, for as long as other organizational policies, and with adequate protection of documentation in line with confidentiality and usage requirements.

### XIII. Review and Enforcement of the AI Policy

A.	this Policy	ole or governance body e.g. RAI Operational Committee] shall maintain olicy and is responsible for <b>monitoring</b> and <b>reviewing</b> the effectiveness of the and of the AI management system described therein [e.g. every 3 months], with ongoing input from cross-functional stakeholders and periodic review by tive [owners, champions, and/or sponsors] <sup>208</sup>
	a)	<b>Corrective updates</b> and other changes to the Policy shall be planned, and their effects are reviewed with any adverse effects mitigated. <sup>209</sup>
	b)	Each component of the AI Policy shall undergo periodic review and ongoing monitoring, with the [e.g. RAI Operational Committee] assigning subroles and responsibilities. <sup>210</sup>
В.		ole or governance body e.g. RAI Operational Committee] shall be a sible for procedures to <b>detect</b> and <b>respond to deviations and violations</b> of the

<sup>&</sup>lt;sup>208</sup> Aligned with ISO/IEC 42001 5.2, 9.2.1, 9.3.1 and A.2.4. ISO/IEC 42001 9.3.2-9.3.3 provides guidance on the desired inputs and results of a review of the AI Policy and the AI management system described therein.

 $<sup>^{209}</sup>$  Aligned with ISO/IEC 42001 6.3, 8.1, and 10.1-10.2. ISO/IEC 42001 10.2 provides guidance on actions to take when a nonconformity occurs.

<sup>&</sup>lt;sup>210</sup> Aligned with NIST AI RMF GOVERN 1.5.

 $<sup>^{211}</sup>$  Aligned with ISO/IEC 42001 5.3. ISO/IEC 42001 10.2 provides guidance on actions to take when a nonconformity occurs.

# XIV. Conclusion/Acknowledgement

[Name, Role]	[Date]
[Name, Role]	[Date]
 [Name, Role]	[Date]

## **Appendix A**

Appendix A provides a potential list of documents and artifacts that can be linked to the Al Policy to provide a complete ecosystem view of the organization's Al strategy.

#### **Enterprise-level Policy Documents:**

- Governance structure policies, including charter documents for new Al bodies or groups
- Existing ethics, DEI, ESG, social responsibility, corporate human rights policies
- Stated and achieved voluntary commitments related to AI (e.g., Frontier AI Safety Commitments, 2023 White House AI commitments)
- Data management policies and standards
- Risk management policies and standards
- Product development or project management policies and standards
- Stakeholder engagement policies
- Procurement framework and policies
- Infosecurity and/or cybersecurity policies and standards
- Workforce policies, including usage policies, planning and hiring protocols
- Any other existing policies that have been augmented with Al-specific guidance

#### **Enterprise-level Artifacts:**

- Risk/impact taxonomy
- Complete and detailed list of prohibited use cases
- Al system inventory with model cards
- Al regulatory tracker
- Al Incident, Impact, and Risk (IIR) database
- Resource utilization tracker
- Data inventory
- Project inventory
- Best practices for AI database
- Workforce RAI development, education, and training resources

#### **System-level Artifacts (per AI system):**

- Al technical documentation, including role-specific documentation (e.g. for users, operators)
- Al business use case documentation
- Contracts, with suppliers or buyers
- Al impact assessment results
- Documentation across life cycle stages (e.g., data set quality test results, model evaluation results, performance and event logs, governance gate approval forms)
- Relevant external documented information, such as from partners or suppliers