# Evolving AI Risk Management: A Maturity Model based on the NIST AI Risk Management Framework

RAVIT DOTAN, TechBetter, USA

BORHANE BLILI-HAMELIN, AI Risk and Vulnerability Alliance, USA

RAVI MADHAVAN, The University of Pittsburgh, USA

JEANNA MATTHEWS, Clarkson University, USA

JOSHUA SCARPINO, TrustEngine Assessed.Intelligence, USA

Researchers, government bodies, and organizations have been repeatedly calling for a shift in the responsible AI community from general principles to tangible and operationalizable practices in mitigating the potential sociotechnical harms of AI. Frameworks like the NIST AI RMF embody an emerging consensus on recommended practices in operationalizing sociotechnical harm mitigation. However, private sector organizations currently lag far behind this emerging consensus. Implementation is sporadic and selective at best. At worst, it is ineffective and can risk serving as a misleading veneer of trustworthy processes, providing an appearance of legitimacy to substantively harmful practices. In this paper, we provide a foundation for a framework for evaluating where organizations sit relative to the emerging consensus on sociotechnical harm mitigation best practices: a flexible maturity model based on the NIST AI RMF.

CCS Concepts: • **Software and its engineering** → **Software creation and management**.

Additional Key Words and Phrases: Maturity model, AI Risk Management, NIST AI RMF

## 1 INTRODUCTION

In recent years, increasingly more professionals in the AI ethics space have been calling for "operationalizing AI ethics" or "translating principles into practice" – meaning moving away from articulating general priorities and principles, which has been prominent in the last decade, into establishing processes that rigorously anticipate, evaluate, mitigate, and provide redress for AI harm [2, 12, 26–28, 48]. On that backdrop, practitioners, researchers, and government bodies have developed recommendations and practices to bridge the gap [14, 17, 24, 30].

Despite the existence of a host of AI ethics frameworks and tools, organizations have been lagging behind the recommended best practices [11, 18, 23]. For example, McKinsey [23] shows that in 2022 only 17% of companies reported they worked to mitigate fairness and bias issues, merely a small increase from the 13% who did so in 2019, despite the many tools and frameworks that were developed in between. Dotan et al. [11] show that, while 76% of large companies reported making AI ethics commitments, employing AI ethics personnel, or engaging in thought leadership in 2022, only a small fraction reported engaging in AI risk management best practices such as data and model documentation (9%) and maintaining an incident log (2%).

This paper presents a foundation for a maturity model that is intended to help companies decrease this gap. Maturity models guide companies by laying a sequence of stages for progress and are widely used in many areas, from cybersecurity to software development best practices [39]. A maturity model grounded in AI ethics could help organizations evaluate their existing AI risk management practices and plan how to do better [44].

We chose to base the maturity model we describe in this paper on the NIST AI Risk Management Framework (AI RMF)[30] for practical reasons and because of an alignment with the values embodied in the AI RMF. We argue for

---

practical advantages of basing AI ethics maturity models in widely accepted frameworks, and we explain why we favored the RMF over the EU AI Act. Moreover, we highlight the fact that designing a maturity model requires making substantive value decisions because of the inherent evaluative aspect [4]. We chose the NIST AI RMF for its focus on sociotechnical harm and its elevation of marginalized voices through numerous design choices.

The structure of the paper is as follows. In Section 2, we provide a brief background of maturity models and their purpose, highlight the current maturity model landscape within AI, and provide a more detailed overview of the NIST AI RMF and our reasons for aligning our maturity model with it. In Section 3, we describe the questionnaire leveraged within our maturity model, and in Section 4, we explain the scoring approach we use, as well as a detailed description of the components of the scoring evaluation process. Section 5 is focused on explaining maturity model score aggregation. Section 6 covers research limitations and the next steps. The full questionnaire and an example of how it works in practice are in the appendices.

## 2 BACKGROUND

### 2.1 About maturity models

A staple of current technology management toolkits, maturity models are "conceptual multistage models that describe typical patterns in the development of organizational capabilities" [39]. They have been characterized as a Crawl/Walk/Run-style set of factors depicting the progression of capabilities while also serving as a tool to benchmark current capabilities and help set goals and priorities for improvement [43]. The practical utility of maturity models stems from their simplicity, conceptual power, and evolutionary orientation, which result in effective managerial guidance on where to invest attention, effort, and other resources in order to build capability in successive stages (see Poeppelbus et al. [39] and [47] for overviews of the large body of literature on maturity models, and Poeppelbuss & Röglinger [39] for a discussion of design principles for maturity models).

The use of maturity models in technology management dates back to the 1980s, with predecessors dating back to the 1960s [10, 21, 36, 37]. As the popularity of maturity models has increased, their application has spread to many capability arenas. Well-known examples include the Software Capability Maturity Model [1] and the risk management maturity model [37]. NIST-related maturity models include the NIST cybersecurity maturity model (National Cybersecurity Center of Excellence [31]), and the simple Privacy maturity model in NIST's privacy framework in the form of "Ready, Set, Go" labels [29]. Moreover, NIST worked with the US Department of Defense (DOD) to help create Cybersecurity Maturity Model Certification (CMMC)[7] and with the Department of Energy (DOE) to produce mappings between the NIST 800-53 and 800-171, the Cybersecurity Framework (CSF) and Cybersecurity Capability Maturity Model (C2M2) [43] [31].

Maturity models have great potential in AI ethics, too. They can help companies understand where they stand with regard to standards, analyze gaps, and plan for improvement [44]. This support can help overcome the low implementation of AI ethics, which, as discussed above, is prevalent and can result in a great deal of harm to individuals, society, and the companies themselves.

### 2.2 Maturity models in AI ethics

A handful of maturity models already exist in AI ethics, most of which were created by private sector companies, especially big tech and consulting firms. Those offered by Salesforce [3] and Microsoft [45] may be the most well-known. Other maturity models dedicated to AI ethics include ODI [32], Ethical Intelligence [13], and Krijger et al. [20]. Moreover,

some maturity models for the development of AI capabilities, such as IBM [19], MITRE [25], and PwC [38], include maturity in AI ethics as an aspect of maturity in AI adoption.

The existing AI ethics maturity models are limited. To start, they assume specific trajectories toward good AI governance, while in reality, there may be multiple legitimate trajectories, especially in different contexts. Rigid trajectory expectations may render a model unhelpful or even misleading when used in the wrong context.

One trajectory expectation, for example, is that buy-in from senior management comes at the last stage (e.g. Microsoft [45]). This assumption sounds reasonable in large corporations such as Microsoft. However, in small to medium-sized enterprises (SMEs), research shows that the motivation for responsible innovation typically comes from the founder, so C-suite buy-in comes first ([9], [5]). Another expectation that AI governance starts with policy making that later develops into implementation (e.g., Salesforce [3], ODI [32]). Again, this expectation may sound reasonable in large corporations such as Salesforce. However, in SMEs, which tend to be more informal, company-wide policies may only develop later in the process. A startup we got feedback from raised this concern explicitly, worrying that the fact that their AI ethics processes have not been formalized into written company-wide policies will overshadow the fact that they are implementing AI ethics processes regardless. Other models assume a trajectory in which teams start with localized activities that are later generalized into company-wide policies [20]). While this trajectory may be more amenable to smaller companies that allow leeway to individual teams, it may be a poor fit for a large corporation.

All of these different trajectories may be fruitful in the right context, and AI ethics maturity models should be inclusive of multiple trajectories fitting to a variety of contexts. The existing models lack this feature. As we will show, the maturity model presented in this paper is open to multiple maturity trajectories.

Another limitation of the existing AI ethics maturity models is that they are based on bespoke conceptual frameworks that they develop, including their own dimensions of progress and maturity stages along the dimensions. Microsoft's [45] model, for example, has five maturity stages (Latent, Emerging, Developing, Realizing, and Leading) and five dimensions of progress (Organizational foundations, Team approach, Cross-discipline collaboration, and Responsible AI practices). For comparison, Ethical Intelligence's [13] model presents three maturity stages (Below, Average, and Exemplary) and five dimensions (Accountability, Social impact, Intentional design, Trust and transparency, and Fairness).

The bespoke approach reinvents the wheel and fails to take advantage of the conclusions, expertise, terminology, or authority of widely accepted documents and initiatives. This creates practical challenges. For one thing, reinventing the wheel with new terminology or a custom framework increases the complexity of adoption, and it makes it difficult for companies to communicate their results to others and benchmark their results relative to other organizations. Moreover, bespoke approaches are difficult to relate to industry standards and regulatory proposals. Therefore, even when an organization meets the expectations of a bespoke maturity model, there will likely still be confusion about whether it meets the expectations of the industry or regulators.

Further, the reliance on bespoke frameworks is also problematic due to the political undertone of AI ethics maturity models, especially in light of the fact that many of the existing models were created by technology companies themselves. By nature, maturity models guide organizations to achieve some "good," in this case, good governance of AI. Defining what counts as "good," especially in the case of AI governance, is not only thorny but also deeply political. In a narrow sense, defining what counts as good AI governance is political because various organizations, in particular big tech, are lobbying regulators to shape regulation about AI governance in ways that are favorable to them [35]. As illustrated in the previous paragraph, what is favorable to big tech may be unfavorable to SMEs, not to mention other stakeholders such as the public. In a wider sense, AI governance is political because it involves decision-making about topics that are political, such as expectations around the protection of civil and human rights. Due to this political nature, it matters

who gets to determine what counts as "good" AI governance. Big tech and other private sector organizations may fail to represent the public good in their AI ethics maturity models due to conflicts of interest or a more narrow view of what risks to prioritize in risk management.

## 2.3 NIST AI RMF

The NIST AI RMF, the basis of our maturity model, is a voluntary framework describing best practices for AI risk management, including concrete activities for the development and deployment of AI in a socially responsible way. It is one of the most well-respected documents on AI governance and is growing in influence, especially in light of the October 2023 Executive Order on the Safe, Secure, and Trustworthy Development and Use of Artificial Intelligence that specifically calls out the NIST AI RMF [42]. The many AI companies based in the United States may view the US-based policies as especially relevant.

The RMF's influence is one of the primary reasons we chose to use it as the basis of our maturity model, as we seek to avoid the challenges of bespoke frameworks discussed above.

In addition to the NIST AI RMF, we considered the creation of a maturity model with respect to the European Union's AI Act, but chose NIST's AI RMF for two primary reasons. First, many of the governance requirements in the EU AI Act are focused on high-risk systems, and companies that don't fall into this bucket may argue that it is not relevant to them. We wanted a maturity model that was not specifically limited to high-risk systems and could be leveraged for all AI systems. Second, we wanted to encourage viewing maturity as a continually evolving lifecycle and not simply a checklist for compliance with regulation. The voluntary nature of the NIST AI RMF allows organizations to view maturity of AI risk management in this way.

Another reason we chose the NIST AI RMF as the basis of our maturity model is its emphasis on a sociotechnical perspective on AI harm. Over the past 10 years, researchers and organizations have converged on the importance of relying on a plurality of sociotechnical methods to map, measure, disclose, and manage the interlocking technical and social causes of algorithmic harms. [40, 46] Moreover, grappling seriously with the interlocking technical and social causes of AI harm requires centering the perspectives of stakeholders impacted by AI systems, rather than merely of stakeholders involved in building and deploying those systems [22]. The AI RMF explicitly and repeatedly calls upon organizations to embrace this kind of sociotechnical perspective on AI risk management [30].

Last, the RMF naturally lends itself to a maturity model, and a maturity model based on it complements it. First, the RMF lends itself to a maturity model because the RMF touches directly on many high-level activities that are relevant to maturity, such as measuring risk and internal policy making. Their focus on activities throughout the lifecycle of an AI system makes it an especially good candidate for a maturity model. Second, our maturity model complements the RMF because the RMF is not intended to provide guidance on how organizations might evolve towards the best practices being recommended or on how to evaluate the extent to which organizations are aligned with those best practices. Therefore, organizations may struggle with how to prioritize plans for improvement. Similarly, external actors, such as investors and consumers, may struggle to use the NIST AI RMF to evaluate organizations and their level of maturity with respect to implementing the RMF. In fact, one could wonder why the NIST AI RMF does not itself contain a maturity model, especially given that NIST has developed explicit maturity model language for other related areas such as cybersecurity and privacy. The creation of a maturity model will assist in the transition from ad hoc implementation of responsible AI to more mature processes and programs [44].

## 2.4 Design process

We apply the maturity model approach to translate the AI risk considerations identified in the RMF to an evolutionary description. Our work unfolds in two steps: an inductive step to lay out the foundation of the maturity model, which we do in this paper, and a confirmatory step with case studies and empirical refinement, which we do in ongoing and future work, described in the last section, Section 6.

This approach differs from current models within AI ethics and addresses challenges encountered in the field of maturity model design. Maturity models have been critiqued for oversimplification and inadequate empirical grounding. Accordingly, an overarching current concern in the field is how to develop maturity models that are both theoretically robust and empirically grounded [39]. In our case, the conceptual strength of our framework comes from the conceptual strength of the RMF and the approach developed in this paper. The empirical strength will come from the empirical work based on this foundation.

In developing the conceptual framework in this paper, the inputs we used were the NIST AI RMF and the broader literature on the development of maturity models, leading to the initial specification of the phases and content of the maturity model. The model was iteratively refined by the authors to reflect the appropriate stages of, and linkages between, the RMF elements over the development of AI capability. The refined model was then pilot-tested through consultations and dry runs with a convenience sample of AI stakeholders, including startup leaders, active investors and AI ethics experts, including NIST team members involved in developing the RMF. Throughout this process, the model was successively refined via the incorporation of feedback with regard to accuracy, clarity, relevance and ease of use. The model presented in this paper represents the current stage of its development based on that process.

## 3 FLEXIBLE QUESTIONNAIRE

Our maturity model includes a flexible questionnaire and scoring guidelines. The questionnaire consists of a list of statements, and evaluators are asked to rank them using the scoring guidelines discussed below. The statements in the questionnaire center on concrete and verifiable actions, such as conducting certain processes and documenting the outcomes. For example:

> "We regularly evaluate and document bias and fairness issues caused by our AI systems".

The questionnaire avoids general and abstract statements such as "Our AI systems are fair". Further, the statements use the plural first pronoun "we" and the active present tense, e.g., "we document." This is an intentional choice made to emphasize the responsibilities of the companies and people who manage AI.

The statements cover the content of the RMF's governance recommendations, which are divided into four pillars: MAP - Learning about AI risks and opportunities; MEASURE - Measuring risks and impacts; MANAGE - Implementing practices to mitigate risks and maximize benefits; and GOVERN - Systematizing and organizing activities across the organization. Each of the pillars includes a list of categories and subcategories. For example, one of the categories in the MEASURE pillar is "MEASURE 2: AI systems are evaluated for trustworthy characteristics." One of the subcategories in this category is "MEASURE 2.11: Fairness and bias – as identified in the MAP function – are evaluated and results are documented." [30]. In isolation, each statement in the questionnaire covers one or more of the NIST AI RMF subcategories. For example, the statement above covers the subcategory MEASURE 2.11. Jointly, the statements in the questionnaire cover all RMF subcategories.

The questionnaire is flexible in that evaluators are not required to evaluate all statements. The questionnaire allows the evaluator to adjust the evaluation to the organization's specific context in three key ways: 1) level of granularity, 2)

life-cycle stage of the AI system, and 3) multiplicity of AI systems within the organization. We elaborate on each of these in the subsections that follow.

### 3.1   Flexibility 1: Granularity

Evaluators who are interested in a fine-grained evaluation can rank each of the 60 statements in the questionnaire. However, the feedback we received from practitioners indicates that many are interested in a more coarse-grain evaluation. Therefore, the 60 statements are divided into 9 topics. Each topic is represented by a sentence that describes the statements in that topic. For example, one of the topics is

- Topic 4 - "*Measuring risk:* We measure our potential negative impacts".

Under this topic, there are finer grain individual statements, including for example:

- Statement 4e: "We regularly evaluate and document bias and fairness issues related to our AI systems".
- Statement 4i: "We regularly evaluate and document security issues related to our AI systems".

Those interested in a coarse-grained evaluation can score only the topic statement. However, the individual statements are still taken into account because, as discussed in more detail below, the scoring guidelines instruct evaluators to give higher scores the better the coverage of the individual statements.

### 3.2   Flexibility 2: Life-cycle stage

A second aspect of flexibility comes from observing that a subset of RMF subcategories only becomes relevant once the AI system has reached a particular development stage. For example, RMF subcategory MANAGE 4.1 is only relevant after the system has been deployed:

> MANAGE 4.1: Post-deployment AI system monitoring plans are implemented, including mechanisms for capturing and evaluating input from users and other relevant AI actors, appeal and override, decommissioning, incident response, recovery, and change management. [31]

For this reason, the questionnaire is divided into phases of the development lifecycle, based on the AI life cycle described in the RMF [30]. We grouped the life cycle into three stages: (1) Planning and design; (2) Data collection and model building, including verifying and validating the system; and (3) Deployment - including deploying, using, operating, and monitoring the system. Each life-cycle stage contains topics and statements appropriate for that stage from multiple RMF pillars. The evaluator only uses the statements suitable for the relevant life-cycle stage. This flexibility explicitly guides evaluators to avoid questions that are not yet relevant to a particular AI system.

### 3.3   Flexibility 3: Multiplicity of AI systems

Organizations may have multiple AI systems, and the questionnaire allows for flexibility in approaching this multiplicity. Evaluators may score each AI system to be scored separately and aggregate those to get scores for the organization as a whole. Those interested in a more coarse-grain evaluation may instead score the organization holistically without delving into the details of each individual system.

### 3.4   Putting it together

Putting together all three aspects of flexibility, those interested in the most fine-grained version of the evaluation will score each AI system using the individual statements appropriate to that system's life-cycle stage. Those interested in

| Life-cycle stage | Topic | Statement | Score | | | Evidence |
|---|---|---|---|---|---|---|
| | | | Overall | AI 1 | AI 2 | |
| Plannign and on | | | | | | |
| Data collection, model building, and on | | Measuring risk: We measure the potential negative impacts of our AI | | | | |
| | | We evaluate and document bias and fairness issues caused by our AI system(s) | | | | |
| | | We evaluate and document security issues caused by our AI system(s) | | | | |
| Deployment and on | | | | | | |

Fig. 1. The structure of the questionnaire

the most coarse-grained version of the evaluation will score the organization as a whole using only topic statements appropriate to the life-cycle stage of the most advanced AI system that the organization manages.

You can see an illustration of the overall structure of the questionnaire in Figure 1.

## 4  SCORING GUIDELINES

### 4.1  Scoring Metrics

In addition to a questionnaire, the maturity model also includes scoring guidelines. Scoring AI responsibility involves a great degree of personal judgment. The goal is not to produce "objective" scores but rather to communicate evaluations and reasons for those evaluations that can help organizations understand where they are and how they can improve. Our scoring guidelines are designed to facilitate this kind of evaluation by guiding evaluators to provide evidence-based, well-reasoned evaluations that are based on expectations drawn from NIST's work.

When scoring, each evaluated statement should be ranked on a scale of 1-5, where 1 is the lowest and 5 is the highest, based on how well it satisfies three metrics (explained in more detail below): coverage of the RMF categories, robustness based on NIST's implementation tiers, and input diversity. For each of these metrics, evaluators should determine the degree to which it is satisfied – low, medium, or high. We explain how to put it all together below, after presenting the three metrics. We illustrate using these metrics on a case study in appendix A.2.

*4.1.1  Coverage of RMF Subcategories.* As discussed above, the questionnaire allows evaluators to evaluate topic statements only, rather than all of the individual statements. For example, the evaluator can score the statement "*Measuring Risk*: we measure the potential negative impacts of our AI", but not all the statements it contains, such as "We evaluate and document bias and fairness issues caused by our AI system(s)" and "We evaluate and document security issues caused by our AI system(s)." When the evaluator does so, the scoring of the topic statement should reflect coverage of all the individual statements included in that topic. For example, companies that evaluate and document security but not fairness risks satisfy this metric to a degree lower than companies that address both.

*4.1.2  Robustness.* We use the name "robustness" to refer to the ideals expressed through NIST's "implementation tiers." The implementation tiers are distinctions NIST uses to describe degrees of risk management activities in areas such as

privacy and cyber-security ([29], [31]). These tiers represent an increasing degree of rigor, and showcases how well an organization has implemented the component under evaluation. There are four tiers:

(1) PARTIAL- Activities are ad-hoc, reactive, occasional, or isolated from key organizational activities. For example, in organizations with a "partial" level of privacy risk management, "[o]rganizational privacy risk management practices are not formalized, and risk is managed in an ad hoc and sometimes reactive manner." [29]

(2) RISK INFORMED - Activities occur but they are informal and irregular. For example, in organizations with an "informed" level of privacy risk management, "[p]rivacy risk assessment occurs, but is not typically repeatable or reoccurring." [29]

(3) REPEATABLE- Activities are formalized into organization-wide policies or systematic practices. For example, in organizations with a "repeatable" level of cybersecurity risk management "[r]isk-informed policies, processes, and procedures are defined, implemented as intended, and reviewed." [31]

(4) ADAPTIVE - Risk management activities can adapt to changes in the landscape and product, including regular reviews and effective contingency processes to respond to failure. For example, in organizations with an "adaptive" level of cybersecurity risk "The organization uses real-time or near real-time information to understand and consistently act upon cybersecurity risks associated with the products and services it provides and uses." [31]

The implementation tiers are meant to be tools for internal communication to help organizations set priorities. Organizations are not expected to treat these tiers as targets, but should evaluate the desired tier that is appropriate based on their organizational goals, to both reduce risk to an acceptable level and that is feasible to implement and manage. However, the ideals they express can be used to evaluate maturity: The more an organization embodies the ideals, the more mature it is. We have extracted six interrelated ideals for the purposes of maturity evaluation. For convenience, we refer to them collectively as "robustness":

*Robustness* - The risk management activities are

(1) Regular - Performed in a routine manner
(2) Systematic - Follow policies that are well-defined and span company-wide
(3) Trained Personnel - Performed by people who are properly trained and whose roles in the activities are clearly defined
(4) Sufficiently Resourced - Supported by sufficient resources, including budget, time, compute power, and cutting-edge tools
(5) Adaptive - Adapting to changes in the landscape and product, including regular reviews and effective contingency processes to respond to failure
(6) Cross-functional - Involve all core business units and senior management. They are informed of the outcomes and contribute to decision-making, strategy, and resource allocation related to the activities (core business units include finance, customer support, HR, marketing, sales, etc)

*4.1.3  Input diversity.* Input diversity means that risk management activities receive input from diverse internal and external stakeholders. A low level of input diversity means that the relevant activities receive input from relatively few kinds of stakeholders. High levels of input diversity mean that the activities receive input from diverse internal and external stakeholders. For example, suppose that a company chooses its fairness metrics in consultation with civil

society organizations, surveys of diverse customers administered by the customer success team, and conversations with diverse employees in the company. In that case, the company demonstrates a high level of input diversity with regard to the statement "We evaluate and document bias and fairness issues related to this AI system".

The input diversity ideal is not highlighted in the NIST implementation tiers for privacy and cybersecurity. However, it is a key aspect of the AI RMF and is important in AI ethics. AI systems often impact masses of end-users and data subjects as well as society at large. Properly understanding, measuring, and managing AI risks requires an in-depth understanding of the potential impacts which, in turn, requires input from a wide range of perspectives.

One way to add input diversity to the maturity model is to include statements that specifically target activities that solicit feedback and incorporate it into the design of the AI system. There are even two RMF subcategories that articulate this content (GOVERN 5.1 and 5.2). However, we choose to include input diversity as a scoring metric because it is a topic that intersects most, or even all, the subcategories. Including input diversity as a scoring guideline allows the evaluator to reflect, and look for evidence for, how well the company solicits and uses feedback on its AI systems and the risks associated with them.

### 4.2    Scores and Evidence

The score of each statement depends on how well the three metrics are satisfied, and evaluators are asked to provide the evidence and rationale for the scoring.

Scores ranges between 1-5, where 1 is the lowest and 5 is the highest. We developed the following as a rule of thumb for determining scores:

- 5: HHH

    All three metrics are satisfied to a high degree
- 4: HHM

    Two of the metrics are satisfied to a high degree and one to a medium degree
- 3: HMM, HHL, HML, or MMM

    One of the following is the case: (1) Two of the metrics are satisfied to a medium degree and one to a high degree; (2) Two of the metrics are satisfied to a high degree and one to a low degree; (3) One metric is satisfied to a high degree, one to a medium degree, and one to a low degree; or (4) all metrics are satisfied to a medium degree.
- 2: MML, MLL, or HLL

    One of the following is the case: (1) Two of the metrics are satisfied to a medium degree and one to a low degree; (2) One metric is satisfied to a medium degree and two to a low degree; (3) One of the metrics is satisfied to a high degree and two to a low degree.
- 1: LLL

    All metrics are satisfied to a low degree
- N/A

    The statement is not applicable

This rule of thumb is based on thresholds which we determined in the following way. The satisfaction of one metric to a low degree counts as one point, medium counts as two points, and high counts as three points. Score 1 signifies a total of 3 points, Score 2 signifies 4-5 points, score 3 signifies 6-7 points, score 4 signifies 8 points, and score 5 signifies 9

points. For simplicity, we do not ask evaluators to keep track of these points. Instead, they are asked to use the list above.

Evidence includes information about what organizations do, about what they don't do, and reports of lack of evidence. For example, evidence may include describing artifacts that indicate that the company is engaged in the relevant activities or the evaluator's first-hand experience in the company. E.g., they may describe which company documents contain the relevant information and how detailed that information is, the evaluator's first-hand knowledge about the execution of the relevant tasks, and so on. Evidence may also include indications that certain activities are not performed, which may happen, for example, when company documents imply that these activities are outside of the company's current scope. Further, evidence discussions may also include pointing out a lack of evidence. We ask evaluators to note in their comments a distinction between lack of any evidence and presence of evidence to the contrary. Last, evaluators provide evidence that a statement is not applicable, which may happen for example, due to the life-cycle stage of the evaluated AI system.

We chose to base the scoring on evidence and ask the evaluators to describe or provide it for accountability and to increase the usefulness of the evaluation. Providing evidence encourages accountability in the evaluation process because it requires the evaluator to base the scoring on information that others can assess, too. Moreover, requiring evaluators to provide evidence also encourages accountability on the part of the evaluated companies, because it encourages them to ensure that such evidence is available. Companies can do so, for example, by documenting key processes and their outcomes.

Providing evidence for scoring improves the usefulness of the evaluation because it contextualizes and explains the reason for the score. Numbers on their own don't offer much information about the company, what they currently do, what is missing, and how they can improve. The evidence an evaluator cites helps others understand how the evaluator interprets the scoring guidelines and what a given score means to that evaluator. This can help companies understand what they are doing right and how to do better.

### 4.3 Applicability of the Scoring Guidelines

Inevitably, there is going to be some divergence in the scores when completed by different evaluators. This will be true for any set of guidelines, as no set of guidelines can cover all the details relevant to the wide range of contexts and circumstances evaluators may encounter. No matter how detailed the guidelines may be, evaluators will always need to exercise some judgment, deciding what satisfies the metrics and to what degree, deciding what counts as evidence, deciding which contextual factors matter most, etc. These are some of the reasons why scoring AI responsibility is an inherently subjective activity.

This subjectivity is a feature, not a bug, and the goal of the scoring guidelines is to help evaluators express these subjective judgments in a structured and helpful way. To further support the process, we articulate our own judgment about two examples in appendix A.2. In the last section, Section 6, we discuss ongoing and future work to provide further information for evaluators.

### 5 SCORE AGGREGATION

After scoring the individual statements or topics, evaluators can aggregate the scoring to get a unified score. Our maturity model offers two modes of aggregation: By NIST pillars and by responsibility dimensions (see Figure 2 for an illustration).
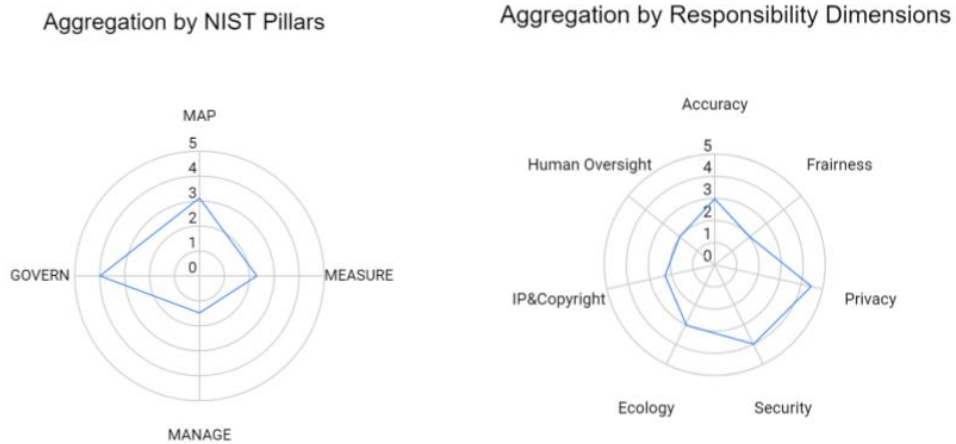
## Aggregation Illustration



Fig. 2. Illustration of aggregation modes in radio charts: To the left, aggregation by NIST Pillar. To the right, aggregation by responsibility dimension

In aggregating by NIST pillars, the output is a score for each of the NIST pillars, MAP, MEASURE, MANAGE, and GOVERN, based on the scores of the statements that belong to it. For example, the MAP score is the average of all the statements that are based on recommendations included in the MAP pillar.

Aggregation by NIST pillar can help discover organizations' strengths and weaknesses in different kinds of activities. In particular, this mode of aggregation can expose systematic failures in organizations' approaches to AI responsibility. For example, when organizations show strength in GOVERN activities but weakness in all other pillars, they may be engaged in ethics washing. For example, they may be establishing policies that are largely not implemented. Other organizations may show strength in GOVERN and MANAGE but weakness in MAP and MEASURE. These organizations' risk management activities may be ill-informed, as the low level of MAP and MEASURE may indicate that their understanding of the risks is lacking.

Another option is to aggregate based on some or all the dimensions of AI responsibility the RMF identifies, e.g., fairness, privacy, and security. In this aggregation mode, the score of each dimension is an average of all the statements relevant to that dimension and is possible only when the relevant individual statements get their own score.

Aggregation by responsibility dimensions can help discover when organizations ignore certain issues. For example, some organizations boast AI responsibility based on their activity in a handful of risk areas, such as privacy and security. Focus on each dimension can highlight the other areas, in which the company is lacking.

Aggregation can help track companies' progress over time. Our maturity model isn't prescriptive about the trajectory of the progress. It allows tracking progress which may take place in many different ways. For example, In large corporations, for example, we may see a top-down progress trajectory, where the company starts with strong GOVERN

## Maturity Trajectories Illustration

— Maturity at time 1   — Maturity at time 2

### Bottom-up trajectory
(Aggregation by NIST Pillar)

### Top-down trajectory
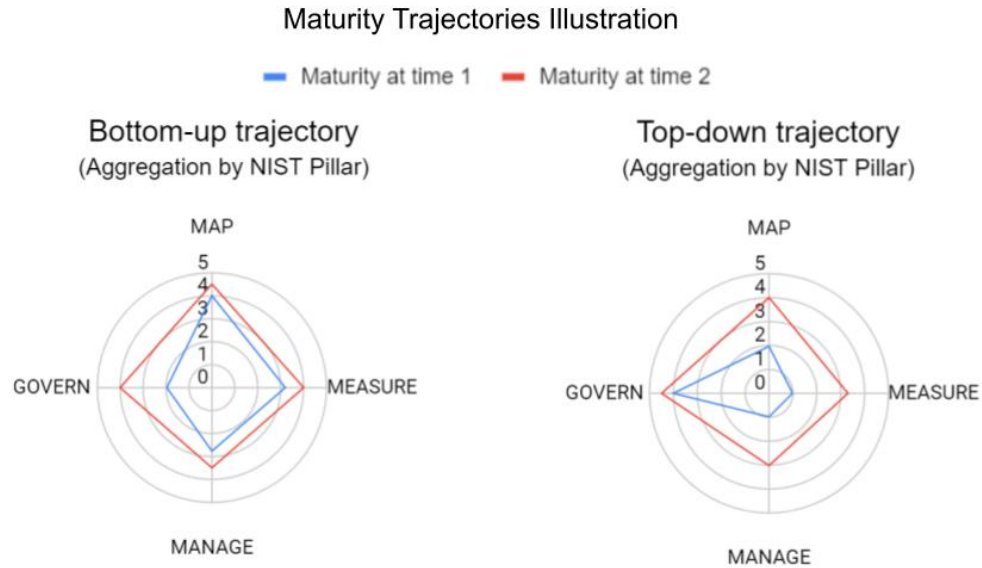(Aggregation by NIST Pillar)

Fig. 3.  Illustration of maturity progress trajectories. To the left, a bottom-up trajectory. To the right, a top-town trajectory

activities and advances to stronger MEASURE and MANAGE activities. In smaller companies, we might see a bottom-up progress trajectory, where the company starts with strong MEASURE, MAP, and MANAGE activities and progresses to stronger GOVERN activities (see Figure 3 for an illustration).

## 6   ADVANTAGES, LIMITATIONS, AND FUTURE WORK

This paper lays the foundation for a maturity model for responsible AI governance. It is based on the NIST AI RMF, a widely accepted voluntary framework, which increases its ability to facilitate comparability between organizations, gets away from questions of legal compliance, and highlights the mitigation of sociotechnical harms. Further, the model is flexible in that evaluators can adjust the questionnaire to the context and in that the evaluation accommodates multiple maturity trajectories.

We see this paper as laying the groundwork for future, practice-informed iterations of the maturity model. A key area for future work is the scoring guidelines. In their current form, the scoring guidelines are likely to lead to vast variations in evaluator behaviors, e.g., evaluators may interpret what counts as "evidence" differently or they might treat the same evidence as supporting different levels of satisfaction of the metrics. We anticipate that evaluators might face uncertainty in how to best apply the current guidelines. Moreover, we anticipate uncertainty in interpreting aggregate scores. Organizations might not know what to make of having a certain aggregated score.

Due to the inevitable subjectivity of evaluations, some degree of uncertainty will always remain. Having said that, in ongoing work based on the foundation laid out in this paper, we seek to better understand the perspective of people who aren't us about what it's like to use this model, and especially scoring. In that work, we are facilitating group experimentation with the model: We are recording how evaluators use it and their impressions of it, including what

they count as "evidence" and what merits satisfying the metrics to different degrees. Our analysis will reveal trends in scoring and evidence-giving that will help set expectations for scoring, such as which activities typically merit various scores and which kinds of evidence generally are given to support them. This practice will allow us to embed more diverse perspectives into the model and the scoring guidelines in particular. Moreover, in future work, we will use the model to evaluate a wide range of companies. This will help create scoring expectations by different company characteristics such as size and sector.

As with other evaluative frameworks, our maturity model carries the risk of ethics washing. Our maturity model is strongly process-oriented. The model doesn't evaluate companies on the success of their risk mitigation efforts, e.g. how "fair" their model is. This choice is intentional. First, we follow the NIST AI RMF which makes the same choice. Second, concepts related to socio-technical harm, such as fairness, are deeply context-sensitive. Allowing organizations to determine for themselves how to interpret and measure key concepts is crucial for pluralism. The downside of this choice is that it creates opportunities for ethics washing. Per Goodhart's law [8, 15, 16], "[w]hen a measure becomes a target, it ceases to be a good measure" [41]. Our framework is no exception. Organizations may find ways to appear to have high responsible AI maturity in accordance with the model, while still having poor performance in terms of their sociotechnical harm. However, the same is true for metrics that measure harms directly. Our model decreases gaming risks by moving away from strict compliance by using a voluntary framework as a foundation, which decreases compliance pressures.

## 7 CONCLUSION

This paper lays out a foundation for a maturity model to evaluate the responsibility of AI governance in organizations that develop and manage AI systems. This foundation includes a flexible questionnaire and scoring guidelines, both based on industry standards set out by NIST. The strengths of this model include a rigorous conceptual framework that is drawn from industry standards, a focus on the mitigation of sociotechnical harm and inclusivity, flexibility in the questionnaire and aggregation options to accommodate the needs of different organizations, compatibility with multiple maturity trajectories, and the facilitation of evidence-based evaluations that flesh out subjective judgments and the reasoning to support them. All these are intended to make this model practical and helpful in the hopes of supporting organizations in improving their AI risk management and supporting the field in enhancing the overall levels of AI ethics implementation, which are currently dangerously low.

## 8 RESEARCH ETHICS AND SOCIAL IMPACT

### 8.1 Ethical Considerations

Our work involved consultation and feedback from practitioners. Therefore, we have consulted with an Institutional Review Board, and they have determined that this project is not Human Subject Research. To ensure transparency, we have conveyed to all our interlocutors that the conversation is being held in the context of a research paper and that insights based on the conversation may appear in the paper. We have obtained consent before proceeding with the conversation.

### 8.2 Adverse Impacts

As we discussed in the body of the paper, when used inappropriately, the framework we developed in this paper could be used to provide a false assurance of compliance. The framework should be used to guide processes of improvement of AI

governance. The explanations and evidence evaluators are asked to provide are aimed to facilitate that. Usages that ignore that, such as treating the framework as a checklist, are misguided and could especially lead to the misrepresentation of the evaluated company and lack of progress.

## 9 ACKNOWLEDGMENTS

## REFERENCES

[1] [n. d.]. CMMI Institute - Home. https://cmmiinstitute.com/

[2] Jacqui Ayling and Adriane Chapman. 2022. Putting AI ethics to work: are the tools fit for purpose? *AI and Ethics* 2, 3 (Aug. 2022), 405–429. https://doi.org/10.1007/s43681-021-00084-x

[3] Baxter, K. 2021. *AI Ethics Maturity Model*. Technical Report. Salesforce. https://www.salesforceairesearch.com/static/ethics/EthicalAIMaturityModel.pdf

[4] Rishi Bommasani. 2022. Evaluation for Change. (2022). https://doi.org/10.48550/ARXIV.2212.11670 Publisher: arXiv Version Number: 1.

[5] Hilke Elke Jacke Bos-Brouwers. 2009. Corporate sustainability and innovation in SMEs: Evidence of themes and activities in practice. *Business Strategy and the Environment* 19 (June 2009), 417–435. https://doi.org/10.1002/bse.652

[6] Burstein, Jill. 2023. *Duolingo English Test Responsible AI Standards*. Technical Report. Duolingo. https://duolingo-papers.s3.amazonaws.com/other/DET+Responsible+AI+033123.pdf

[7] CCMC. 2021. *CMMC PROGRAM PROPOSED RULE PUBLISHED - PUBLIC COMMENT PERIOD BEGINS*. Technical Report. US Department of Defense. https://dodcio.defense.gov/CMMC/about/

[8] Alec Chrystal and Paul Mizen. 2003. Goodhart's Law: its origins, meaning and implications for monetary policy. In *Central Banking, Monetary Theory and Practice*. Edward Elgar Publishing, 2329. https://doi.org/10.4337/9781781950777.00022

[9] Christina Covello and Konstantinos Iatridis. 2020. On the challenges and drivers of implementing responsible innovation in foodpreneurial SMEs. In *Assessment of Responsible Innovation* (1 ed.), Emad Yaghmaei and Ibo van de Poel (Eds.). Routledge. https://doi.org/10.4324/9780429298998

[10] Philip B. Crosby. 1979. *Quality is free: the art of making quality certain*. McGraw-Hill, New York.

[11] Ravit Dotan, Gil Rosenthal, Tess Buckley, Josh Scarpino, Luke Patterson, and Thorin Bristow. 2024. *Evaluating AI Governance: Insights from Public Disclosures*. Technical Report. https://www.ravitdotan.com/_files/ugd/f83391_b853450bcc274e9ba9454d618ee41a94.pdf

[12] Cat Drew. 2018. Design for data ethics: using service design approaches to operationalize ethical principles on four projects. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences* 376, 2128 (Sept. 2018), 20170353. https://doi.org/10.1098/rsta.2017.0353

[13] Ethical Intelligence, BCV, and EAIGG. 2022. *ETHICS MATURITY CONTINUUM*. Technical Report. https://static1.squarespace.com/static/5f6dbf464a8eec79c3d177c0/t/61e8821d53b74041072d556d/1642627614838/Ethics+Maturity+Continuum+Report.pdf

[14] Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé Iii, and Kate Crawford. 2021. Datasheets for datasets. *Commun. ACM* 64, 12 (Dec. 2021), 86–92. https://doi.org/10.1145/3458723

[15] C. A. E. Goodhart. 1984. Problems of Monetary Management: The UK Experience. In *Monetary Theory and Practice*. Macmillan Education UK, London, 91–121. https://doi.org/10.1007/978-1-349-17295-5_4

[16] Christopher Hennessy and Charles A.E. Goodhart. 2020. Goodhart's Law and Machine Learning. *SSRN Electronic Journal* (2020). https://doi.org/10.2139/ssrn.3639508

[17] Sarah Holland, Ahmed Hosny, Sarah Newman, Joshua Joseph, and Kasia Chmielinski. 2018. The Dataset Nutrition Label: A Framework To Drive Higher Data Quality Standards. https://doi.org/10.48550/arXiv.1805.03677 arXiv:1805.03677 [cs].

[18] IBM. [n. d.]. IBM Global AI Adoption Index 2022 | IBM. https://www.ibm.com/watson/resources/ai-adoption

[19] IBM. 2021. *AI maturity framework for enterprise applications*. Technical Report. IBM. https://www.ibm.com/watson/supply-chain/resources/ai-maturity

[20] J. Krijger, T. Thuis, M. de Ruiter, E. Ligthart, and I. Broekman. 2023. The AI ethics maturity model: a holistic approach to advancing ethical data science in organizations. *AI and Ethics* 3, 2 (May 2023), 355–367. https://doi.org/10.1007/s43681-022-00228-7

[21] Simon Smith Kuznets. 1965. *Economic growth and structure: selected essays*. Norton, New York, N.Y.

[22] Seth Lazar and Alondra Nelson. 2023. AI safety on whose terms? *Science* 381, 6654 (July 2023), 138–138. https://doi.org/10.1126/science.adi8982

[23] McKinsey. [n. d.]. The state of AI in 2022—and a half decade in review | McKinsey. https://www.mckinsey.com/capabilities/quantumblack/our-insights/the-state-of-ai-in-2022-and-a-half-decade-in-review#/download//~/media/mckinsey/business%20functions/quantumblack/our%20insights/the%20state%20of%20ai%20in%202022%20and%20a%20half%20decade%20in%20review/the-state-of-ai-in-2022-and-a-half-decade-in-review.pdf?cid=soc-web

[24] Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. 2019. Model Cards for Model Reporting. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*. 220–229. https://doi.org/10.1145/3287560.3287596 arXiv:1810.03993 [cs].

[25] MITRE. 2023. *The MITRE AI Maturity Model and Organizational Assessment Tool Guide: A Path to Successful AI Adoption*. Technical Report. MITRE.

[26] Jessica Morley, Luciano Floridi, Libby Kinsey, and Anat Elhalal. 2020. From What to How: An Initial Review of Publicly Available AI Ethics Tools, Methods and Research to Translate Principles into Practices. *Science and Engineering Ethics* 26, 4 (Aug. 2020), 2141–2168. https://doi.org/10.1007/s11948-019-00165-5

[27] Jessica Morley, Libby Kinsey, Anat Elhalal, Francesca Garcia, Marta Ziosi, and Luciano Floridi. 2023. Operationalising AI ethics: barriers, enablers and next steps. *AI & SOCIETY* 38, 1 (Feb. 2023), 411–423. https://doi.org/10.1007/s00146-021-01308-8

[28] Luke Munn. 2023. The uselessness of AI ethics. *AI and Ethics* 3, 3 (Aug. 2023), 869–877. https://doi.org/10.1007/s43681-022-00209-w

[29] National Institute of Standards and Technology. 2020. *NIST PRIVACY FRAMEWORK:: A TOOL FOR IMPROVING PRIVACY THROUGH ENTERPRISE RISK MANAGEMENT, VERSION 1.0*. Technical Report NIST CSWP 01162020. National Institute of Standards and Technology, Gaithersburg, MD. NIST CSWP 01162020 pages. https://doi.org/10.6028/NIST.CSWP.01162020

[30] NIST. 2023. *AI Risk Management Framework: AI RMF (1.0)*. Technical Report NIST AI 100-1. National Institute of Standards and Technology, Gaithersburg, MD. error: NIST AI 100−1 pages. https://doi.org/10.6028/NIST.AI.100-1

[31] NIST, NCCoE, DOE, and CESER. 2023. Cybersecurity Capability Maturity Model to NIST Cybersecurity Framework Mapping | NCCoE. https://www.nccoe.nist.gov/news-insights/cybersecurity-capability-maturity-model-nist-cybersecurity-framework-mapping

[32] Open Data Institute. 2022. *Data Ethics Maturity Model: benchmarking your approach to data ethics*. Technical Report. https://theodi.org/insights/tools/data-ethics-maturity-model-benchmarking-your-approach-to-data-ethics/

[33] OpenAI. [n. d.]. Preparedness. https://openai.com/safety/preparedness

[34] OpenAI. 2023. Our approach to AI safety. https://openai.com/blog/our-approach-to-ai-safety

[35] Caitlin Oprysko. 2023. OpenAI registers to lobby. *Politico* (Nov. 2023). https://www.politico.com/newsletters/politico-influence/2023/11/17/openai-registers-to-lobby-00127874

[36] Jean Piaget. 1964. Cognitive development in children: Piaget development and learning. *Journal of Research in Science Teaching* 2, 3 (Sept. 1964), 176–186. https://doi.org/10.1002/tea.3660020306

[37] Diogo Proenca, Joao Estevens, Ricardo Vieira, and Jose Borbinha. 2017. Risk Management: A Maturity Model Based on ISO 31000. In *2017 IEEE 19th Conference on Business Informatics (CBI)*. IEEE, Thessaloniki, Greece, 99–108. https://doi.org/10.1109/CBI.2017.40

[38] PwC. 2021. *Responsible AI - Maturing from theory to practice*. Technical Report. PwC.

[39] Jens Pöppelbuß and Maximilian Röglinger. 2011. WHAT MAKES A USEFUL MATURITY MODEL? A FRAMEWORK OF GENERAL DESIGN PRINCIPLES FOR MATURITY MODELS AND ITS DEMONSTRATION IN BUSINESS PROCESS MANAGEMENT. *ECIS 2011 Proceedings* (Oct. 2011). https://aisel.aisnet.org/ecis2011/28

[40] Andrew D. Selbst, Danah Boyd, Sorelle A. Friedler, Suresh Venkatasubramanian, and Janet Vertesi. 2019. Fairness and Abstraction in Sociotechnical Systems. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*. ACM, Atlanta GA USA, 59–68. https://doi.org/10.1145/3287560.3287598

[41] Marilyn Strathern. 1997. 'Improving ratings': audit in the British University system. *European Review* 5, 3 (July 1997), 305–321. https://doi.org/10.1002/(SICI)1234-981X(199707)5:3<305::AID-EURO184>3.0.CO;2-4

[42] The White House. 2023. [Executive Order 14110 of October 30, 2023] Safe, Secure, and Trustworthy Development and Use of Artificial Intelligence. https://www.federalregister.gov/documents/2023/11/01/2023-24283/safe-secure-and-trustworthy-development-and-use-of-artificial-intelligence

[43] US Department of Energy. 2022. *Cybersecurity Capability Maturity Model (C2M2)*. Technical Report 2.1. US Department of Energy. https://www.energy.gov/ceser/cybersecurity-capability-maturity-model-c2m2

[44] Ville Vakkuri, Marianna Jantunen, Erika Halme, Kai-Kristian Kemell, Anh Nguyen-Duc, Tommi Mikkonen, and Pekka Abrahamsson. 2021. Time for AI (Ethics) Maturity Model Is Now. http://arxiv.org/abs/2101.12701 arXiv:2101.12701 [cs].

[45] Mihaela Vorvoreanu, Amy Heger, Samir Passi, Shipi Dhanorkar, Zoe Kahn, and Ruotong Wang. 2023. *Responsible AI Maturity Model*. Technical Report MSR-TR-2023-26. Microsoft. https://www.microsoft.com/en-us/research/publication/responsible-ai-maturity-model/

[46] Laura Weidinger, Maribeth Rauh, Nahema Marchal, Arianna Manzini, Lisa Anne Hendricks, Juan Mateos-Garcia, Stevie Bergman, Jackie Kay, Conor Griffin, Ben Bariach, Iason Gabriel, Verena Rieser, and William Isaac. 2023. Sociotechnical Safety Evaluation of Generative AI Systems. http://arxiv.org/abs/2310.11986 arXiv:2310.11986 [cs].

[47] Roy Wendler. 2012. The maturity of maturity model research: A systematic mapping study. *Information and Software Technology* 54, 12 (Dec. 2012), 1317–1339. https://doi.org/10.1016/j.infsof.2012.07.007

[48] Liming Zhu, Xiwei Xu, Qinghua Lu, Guido Governatori, and Jon Whittle. 2021. AI and Ethics – Operationalising Responsible AI. https://doi.org/10.48550/arXiv.2105.08867 arXiv:2105.08867 [cs].

## A   APPENDICES

### A.1   The Questionnaire

In this appendix, we present the maturity model's questionnaire. As discussed in the body of the paper, the questionnaire is divided by stages of the development life-cycle. Each stage contains a number of topic statements, which contain individual statements. Evaluators can choose to evaluate topic statements only or all statements and are asked to provide both a score and evidence to support that score. Figure 1 provides a visual illustration of the questionnaire's structure, and below is a list of all the statements.

**For AI at or after the planning and design stage:**

(1)  Topic 1 - *Mapping impacts:* We document what the AI will do and its potential impacts.
   Substatements:
   (a)  We document the *goals, scope, and methods* of this AI system. (MAP 1.3, 2.1, 2.4, 3.3)
   (b)  We document the *benefits and potential positive impacts* of this AI system, including the likelihood and magnitude. (MAP 1.1, 3.1, 5.1; GOV 4.2)
   (c)  We document the *business value* of this AI system. (MAP 1.4, 3.1)
   (d)  We document the *possible negative impacts* of this AI system, including the likelihood and magnitude. (GOV 1.1, 4.2, 5.1)
   (e)  We document the *potential costs* of malfunctions of this AI system, including non-monetary costs such as decreased trustworthiness. (MAP 3.2)
   (f)  We implement processes to integrate input about *unexpected impacts.* (MAP 5.2)
   (g)  We document the *methods and tools* we use for mapping impacts. (MAP 2.3, 4.1)
(2)  Topic 2 - *Documenting requirements:* We document basic requirements the system must meet
   Substatements:
   (a)  We document the *human oversight* processes the system needs. (MAP 3.5)
   (b)  We document the *technical standards and certifications* the system will need to satisfy. (MAP 3.4)
   (c)  We document AI *legal requirements* that apply to this AI system. (GOV 1.1)
(3)  Topic 3 - *Culture:* We cultivate AI ethics mindsets
   Substatements
   (a)  We write *policies and guidelines* about AI ethics. (GOV 1.2, 1.4)
   (b)  We document *roles, responsibilities, and lines of communication* related to AI risk management. (GOV 2.1)
   (c)  We provide *training* about AI ethics to relevant personnel. (GOV 2.2)
   (d)  We implement practices to foster *critical thinking* about AI risks. (GOV 4.1)

**For AI at or after the model building and data collection stage:**

(4)  Topic 4 - *Measuring risk:* We measure our potential negative impacts.
   (a)  We document and periodically re-evaluate our *strategy for measuring the impacts* of this AI system. It includes choosing which impacts we measure. It also includes how we will approach monitoring unexpected impacts and impacts that can't be captured with existing metrics. (MEA 1.1)

(b) We document the *methods and tools* we use to measure the impacts of this AI system. It includes which metrics and datasets we use. MEA (2.1, 3.1, 3.2)

(c) We document the *effectiveness of our measurement processes*. (MEA 2.13)

(d) We regularly evaluate and document the *performance* of this AI system in conditions similar to deployment. (MEA 2.3)

(e) We regularly evaluate and document *bias and fairness* issues related to this AI system. (MEA 2.11)

(f) We regularly evaluate and document *privacy* issues related to this AI system. (MEA 2.10)

(g) We regularly evaluate and document *environmental* impacts related to this AI system. (MEA 2.12)

(h) We regularly evaluate and document *transparency and accountability* issues related to this AI system. (MEA 2.8)

(i) We regularly evaluate and document *security and resilience* issues related to this AI system. (MEA 2.7)

(j) We regularly evaluate and document *explainability* issues related to this AI system. (MEA 2.9)

(k) We regularly evaluate and document *third-party issues, such as IP infringement*, related to this AI system. (MEA GOV 6.1)

(l) We regularly evaluate and document *other impacts* related to this AI system. (We added on top of the RMF)

(m) If evaluations use *human subjects*, they are representative and meet appropriate requirements. (MEA 2.2)

(5) Topic 5 - *Transparency:* We document information about the system's limitations and risk control

  Substatements:

(a) We document information about the system's *limitations and options for human oversight* related to this AI system. The documentation is good enough to assist those who need to make decisions based on the system's outputs. (MAP 2.2)

(b) We document the system *risk controls*, including in third-party components. (MAP 4.2)

(c) We *explain the model* to ensure responsible use. (MEA 2.9)

(d) We inventory information about this AI system in a *repository* of our AI systems. (GOV 1.6)

(6) Topic 6 - *Management plan:* We plan how we will respond to risks

  Substatements:

(a) We *plan* and document how we will respond to the risks caused by this AI system. The response options can include mitigating, transferring, avoiding, or accepting risks. (MAN 1.3)

(b) We *prioritize the responses* to the risks of this AI system based on impact, likelihood, available resources or methods, and the organization's risk tolerance. (MAN 1.2)

(c) We document the *residual risks* of this AI system (the risks that we do not mitigate). The documentation includes risks to buyers and users of the system. (MAN 1.4)

(d) We have a plan for addressing *unexpected risks* related to this AI system as they come up. (MAN 2.3)

(7) Topic 7 - *Risk mitigation:* We act to minimize the risks we identify

  Substatements:

(a) We proactively evaluate whether this system *meets its stated objectives* and whether its development or deployment should proceed. (MAN 1.1)

(b) We ensure this AI's *bias and fairness* performance stays meets our standards. (MAN 4.2)

(c) We ensure this AI's *privacy* performance meets our standards. (MAN 4.2)

(d) We ensure this AI's *environmental* performance meets our standards. (MAN 4.2)

(e)  We ensure this AI's *transparency and accountability* meets our standards. (MAN 4.2)

(f)  We ensure this AI's *security and resilience* meets our standards. (MAN 4.2)

(g)  We ensure this AI's *explainability* performance meets our standards. (MAN 4.2)

(h)  We ensure this AI's *third-party impacts, such as IP infringement*, meet our standards. (GOV 6.1)

(i)  We implement processes for *human oversight* related to this AI system. (MAN 3.5)

(j)  We implement processes for *appeal* related to this AI system. (MAN 4.1)

(k)  We maintain *end-of-life* mechanisms to supersede, disengage, or deactivate this AI system if its performance or outcomes are inconsistent with the intended use. (MAN 2.4, GOV 1.6)

(l)  We address all *other risks* prioritized in our plans related to this system by conducting measurable activities. (We added on top of the RMF)

(m)  We address *unexpected risks* related to this system by conducting measurable activities. (MAN 2.3)

(n)  We track and respond to *errors and incidents* related to this system by conducting measurable activities. (MAN 4.3)

**For AI at or after the deployment stage:**

(8)  Topic 8 - *Pre-deployment checks:* We only release versions that meet our AI ethics standards
   Substatments:

   (a)  We demonstrate that this system is *valid, reliable, and meets our standards*. We document the conditions under which it falls short. (MEA 2.5; MAN 1.1)

(9)  Topic 9 - *Monitoring:* We monitor and resolve issues as they arise
   Substatments:

   (a)  We *plan* how to monitor risks related to this system post-deployment. (MAN 4.1)

   (b)  We monitor this system's *functionality* and behavior post-deployment. (MEA 2.4)

   (c)  We apply mechanisms to *sustain the value* of this AI system post-deployment. (MAN 2.2)

   (d)  We capture and evaluate *input from users* about this system post-deployment. (MAN 4.1)

   (e)  We monitor *appeal and override* processes related to this system post-deployment. (MAN 4.1)

   (f)  We monitor *incidents* related to this system and responses to them post-deployment. (MAN 4.1)

   (g)  We monitor incidents related to *high-risk third-party* components and respond to them. (GOV 6.2)

   (h)  We implement *all other* components of our post-deployment monitoring plan for this system. (MAN 4.1)

   (i)  We monitor issues that would trigger our *end-of-life* mechanisms for this system, and we take the system offline if issues come up. (MAN 2.4)

### A.2   Scoring Examples

In this appendix, we illustrate working with our maturity model by using it to evaluate governance based on the public disclosures of two companies: OpenAI and Duolingo. OpenAI develops generative AI chatbots. They detailed their "Approach to AI Safety" in [34], which was published in April 2023, at the release of GPT 4. As the name suggests, the document presents the company's AI ethics approach. Dulingo's most prominent product is an app for learning languages. The Duolingo English Test is an English proficiency test intended to be used as a university entry requirement for non-native English speakers, similar to the TOEFL exam. Dulingo uses generative AI in the question-writing process

and in the evaluation process. The "English Test Responsible AI Standards" [6] is intended to guide their research and documentation.

For the sake of the example, we will use these documents as the sole source of evidence for scoring one topic: "Measuring risk - We measure our potential negative impacts". The "Measuring Risk" topic requires the company to evaluate and document the impact in the following areas regularly: performance, bias and fairness, privacy, environmental impact, transparency and accountability, security and resilience, explainability, and third-party issues such as IP infringement. We add to these NIST requirements that the company regularly evaluate and document other relevant impacts, even if they were not listed by NIST. In addition, "Measuring Risk" requirements also include evaluation of the measurement process itself, which means documenting and regularly re-evaluating the strategy for impact measurement, the methods and tools used for measurement, and the effectiveness of the measurement process. Last, "Measuring Risk" also requires that if the evaluations use human subjects, they are representative and meet appropriate requirements.

We discuss evidence relevant to the scoring of "Measuring Risk" topic and assess the degree to which the three metrics are satisfied. This discussion reflects the results of our deliberations and it is intended to illustrate our way of thinking.

*A.2.1    Coverage of NIST's recommended activities.*  OpenAI's "Approach of AI Safety" [34] is a 1200-word document that covers the following topics: Building increasingly safe AI systems; learning from real-world use to improve safeguards; Protecting children; Respecting privacy; and Improving factual accuracy. Each of these topics is addressed in a few paragraphs that describe the company's approach efforts related to that topic.

The document presents a partial coverage of the risk area NIST specifies. Only performance and privacy are mentioned. Several of the other risk areas that NIST mentions are omitted despite their centrality to OpenAI's product. For example, the lack of attention to IP infringement and fairness issues is especially striking. However, OpenAI's document adds a topic not specified by NIST: protecting children. Therefore, the document satisfies the "coverage" metric to a low degree.

Duolingo's The "English Test Responsible AI Standards" [6] is an 18-page document intended to guide their research and documentation. It covers the following topics: Validity and reliability; Fairness; Privacy and security; and Accountability and transparency. Each of these topics includes a statement of the company's goals related to that topic and details on the processes they implement to achieve these goals.

Duolingo's document covers most of the risk areas NIST specifies, as it addresses performance, fairness, privacy, transparency and accountability, explainability (as part of other sections), and security and resilience are covered. The topics that are missing are environmental impacts and third-party issues (e.g., copyright issues stemming from the use of pre-trained models). Therefore, it satisfies the "coverage" metric to a high degree.

*A.2.2    Robustness.*  As discussed above, robustness covers six ideals for the AI risk management activities: (i) Regularity - Performed in a routine manner; (ii) Systematicity - Follow policies that are well-defined and span company-wide; (iii) Trained Personnel - Performed by people who are properly trained and whose roles in the activities are clearly defined; (iv) Sufficient Resources - Supported by sufficient resources, including budget, time, compute power, and cutting-edge tools; (v) Adaptivity - Adapting to changes in the landscape and product, including regular reviews and effective contingency processes to respond to failure; and (iv) Cross-functionality - Involve all core business units and senior management. They are informed of the outcomes and contribute to decision-making, strategy, and resource allocation related to the activities.

OpenAI's document satisfies the robustness metric to a low degree because the description of the relevant processes lacks detail that would illustrate satisfaction of the robustness ideals. For example, the section about accuracy describes the efforts in one sentence only, this sentence:

> By leveraging user feedback on ChatGPT outputs that were flagged as incorrect as a main source of data—we have improved the factual accuracy of GPT-4. GPT-4 is 40% more likely to produce factual content than GPT-3.5.

This sentence describes an end result of decreasing inaccuracy by 40%. However, the document doesn't address topics such as what the level of inaccuracy is now, what it used to be, how exactly inaccuracy is quantified, or what practices the company implements to regularly monitor its accuracy (if at all). These details are crucial for determining whether accuracy-related processes satisfy the robustness ideals.

At the top of the document, OpenAI states that [A]fter our latest model, GPT-4, finished training, we spent more than 6 months working across the organization to make it safer and more aligned prior to releasing it publicly. This sentence suggests some robustness. However, the document provides no further descriptions or explanations that would indicate that this sentence is more than a platitude.

The lack of detail in OpenAI's "AI Safety Approach" is especially striking in comparison to their "Preparedness Framework (Beta)", which came out in December 2023 [33]. This document describes their approach to risk management with regards to what they call "catastrophic risk", which they define to be "any risk which could result in hundreds of billions of dollars in economic damage or lead to the severe harm or death of many individuals —this includes, but is not limited to, existential risk" (p. 2). The "Preparedness" document is very detailed. Over its 27 pages, it describes topics such as how they will use evaluations to track catastrophic risks, continuously identify unknown categories of catastrophic risks, what the "Preparedness" team will do, and the cross-functionality of the advisory board. This document suggests plans for high degree of robustness, but only regarding one risk area, catastrophic risk. It highlights the fact that the "AI Safety" document lacks similar evidence for robustness with regard to the risk areas recommended by the NIST AI RMF, such as fairness.

In comparison, Duolingo's document presents much stronger evidence for robustness. For example, the section about validity and reliability states that the rationale of these standards is "to ensure that the test is suitable for its intended purpose" and then goes on to specify goals their work processes aim to achieve as well as recommended activities.

The level of detail in this description makes the document stronger evidence that the company's efforts to ensure good AI performance are robust. They are well thought out, systematic, and lend themselves to implementation. However, the evidence would be stronger if the document included explanations about how these processes are integrated into the organization's day-to-day work or other indications that they are indeed implemented. Therefore, we see the document as providing medium evidence for the robustness of the "Mitigating Risk" topic.

*A.2.3 Input Diversity.* The input metric is about the extent to which the company integrates input from diverse stakeholders in its processes. OpenAI's document mentions user feedback in one sentence in the "Learning from real-world examples" section: We cautiously and gradually release new AI systems—with substantial safeguards in place—to a steadily broadening group of people and make continuous improvements based on the lessons we learn.

However, it is unclear to what extent the feedback is diverse and to what extent it is adopted. Therefore, the document satisfies the "input diversity" metric to a low degree.

Duolingo's document mentions details that are relevant to input diversity. For example, the document requires that employees "Evaluate and document demographic representation in data sets used to build AI" [6]. However, key

sections lack an indication that the company collects and uses feedback from diverse internal and external stakeholders. For example, Goal 4.1 is to "Assess how AI processes impact stakeholders" [6]. The requirements include processes to document impacts and risks but not processes for direct engagement with stakeholders such as test takers and universities. Therefore, Duolingo's document also satisfies the "input diversity" metric to a low degree.

*A.2.4   Overall score for the* Measuring Risk *topic.* Recall that the scoring question is whether "There is evidence that the company performs the relevant activities in a way that satisfies all metrics to a high degree", and evaluators are asked to indicate whether they strongly agree (4), somewhat agree (3), somewhat disagree (2), or strongly disagree (1).

OpenAI's document exhibits low satisfaction of all three metrics, coverage, robustenss, and input diversity, with regard to the "Measuring Risk" topic. Therefore, we determine that they deserve the score of "1" for this topic (based on this document).

Duloingo's document exhibits variability: high coverage, medium robustness, and low input diversity. Therefore, we determine that they deserve a score of "3" for the risk measurement topic (based on this document).