

# RAI Guidelines: Method for Generating Responsible AI Guidelines Grounded in Regulations and Usable by (Non-)Technical Roles

ANONYMOUS AUTHOR(S)\*

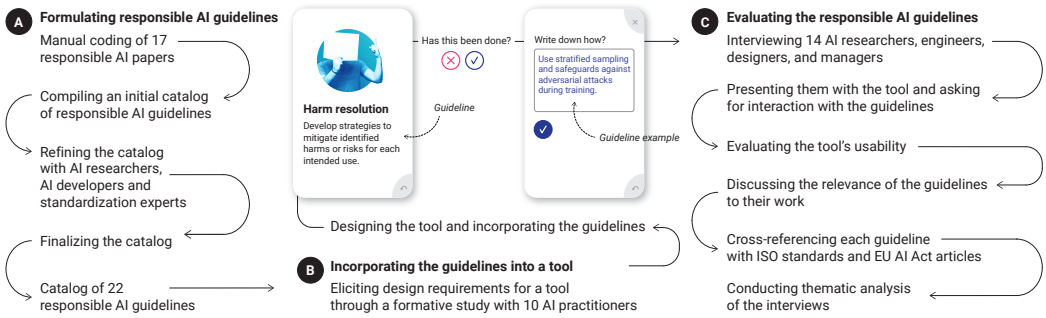


Fig. 1. Overview of our method for generating responsible AI guidelines and evaluating them: (A) formulating responsible AI guidelines that are grounded in regulations and are usable by different roles; (B) incorporating the guidelines into a tool; and (C) evaluating them.

Many guidelines for responsible AI have been suggested to help AI practitioners in the development of ethical and responsible AI systems. However, these guidelines are often neither grounded in regulation nor usable by different roles. To bridge this gap, we developed a four-step method to generate a list of responsible AI guidelines; these steps are: (1) manual coding of 17 papers on responsible AI; (2) compiling an initial catalog of responsible AI guidelines; (3) refining the catalog through interviews and expert panels; and (4) finalizing the catalog. To evaluate the resulting 22 guidelines, we incorporated them into an interactive tool and assessed them in a user study with 14 AI researchers, engineers, designers, and managers from a large technology company. Through interviews with these practitioners, we found that the guidelines were grounded in current regulations and usable across roles, encouraging self-reflection on ethical considerations at early stages of development. This significantly contributes to the concept of 'Responsible AI by Design'—a design-first approach that embeds responsible AI values throughout the development lifecycle and across various business roles.

CCS Concepts: • **Human-centered computing** → **Empirical studies in collaborative and social computing; Interactive systems and tools**; • **Computing methodologies** → **Machine learning; Artificial intelligence**.

Additional Key Words and Phrases: responsible AI, AI ethics, AI guidelines, system development, co-design

## ACM Reference Format:

Anonymous Author(s). 2024. RAI Guidelines: Method for Generating Responsible AI Guidelines Grounded in Regulations and Usable by (Non-)Technical Roles. In *Conference on Computer-Supported Cooperative Work & Social Computing (CSCW '24)*. ACM, New York, NY, USA, 29 pages.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

CSCW '24, November 09–13, San José, Costa Rica

© 2024 Association for Computing Machinery.

## 1 INTRODUCTION

The development of responsible AI systems [50, 87, 92] has become a significant concern as AI technologies continue to permeate various aspects of society [88]. While AI holds the potential to benefit humanity, concerns regarding biases [7, 15, 19] and the lack of transparency and accountability [66, 79] hinder its ability to unlock human capabilities on a large scale. In response, AI practitioners<sup>1</sup> are actively exploring ways to enhance responsible AI development and deployment. One popular approach is the use of tools such as checklists [59] or guideline cards [4, 26, 54] that are designed to promote AI fairness, transparency, and sustainability. These tools provide practical frameworks that enable practitioners to systematically assess and address ethical considerations throughout the AI development lifecycle. By incorporating checklists and guideline cards into their workflows, practitioners can evaluate key aspects such as data sources, model training, and decision-making processes to mitigate potential biases, ensure transparency, and promote the long-term sustainability of AI. However, these tools face two main challenges, creating a mismatch between their potential to support ethical AI development and their current design.

The first challenge is that these tools often exhibit a static nature, lacking the ability to dynamically incorporate the latest advancements in responsible AI literature and international standards [30, 73]. In the rapidly evolving field of responsible AI, new ethical considerations and regulatory guidelines constantly emerge (e.g., the EU AI Act [2]). It is therefore crucial for AI practitioners to stay updated of these developments to ensure their AI systems align with the current ethical and responsible AI practices. While checklists and guideline cards are increasingly used to assist and enhance the development of responsible AI systems, they are rarely grounded in current regulations. For example, Vakkuri et al. [94] proposed the ECCOLA cards that are based on AI ethics guidelines (e.g., IEEE Ethically Aligned Design and EU Trustworthy AI), which are not meant to be grounded on any specific regulations. Additionally, guidelines can quickly become outdated (e.g., the AI Blindspots deck has undergone several iterations [51, 52]), limiting their effectiveness in addressing evolving concerns related to fairness, transparency, and accountability.

The second challenge is that, while these tools emphasize the importance of involving stakeholders from diverse roles and backgrounds, they are often designed for specific AI practitioners (e.g., ML engineers), neglecting a broader spectrum of stakeholders (e.g., non-technical roles). Balayn et al. [8] found that less experienced practitioners in machine learning tend to use a limited set of metrics and methods from toolkits. Similarly, Deng et al. [20] stressed the lack of standardized guidelines in toolkits like AIF360 for introducing fairness issues to non-technical collaborators. Therefore, it is important that toolkits enhance communication, provide comprehensive guidance and support for cross-functional collaboration [101].

To overcome these challenges, we developed a four-step method to generate a list of responsible AI guidelines which we then incorporated in a tool to evaluate them (Figure 1). With this method, we aim to equip different roles with actionable guidelines that are grounded in regulations. To achieve this, we focused on answering this main research question: *How to generate responsible AI guidelines that are grounded in regulations and are usable by different roles?* In addressing this question, we made two main contributions:

- (1) We proposed a four-step method for generating Responsible AI guidelines; these steps are: (1) manual coding of 17 papers on responsible AI; (2) compiling an initial catalog of responsible AI guidelines; (3) refining the catalog through interviews with 10 AI researchers and engineers, and workshops with 4 standardization experts; and (4) finalizing the catalog. This procedure resulted into a set of 22 Responsible AI guidelines (§4).

<sup>1</sup>We use the term practitioners to cover a wide range of stakeholders including AI engineers, developers, researchers, designers, ethics experts.

- 99 (2) We evaluated the 22 guidelines in a user study with 14 AI researchers, engineers, designers,  
100 product managers from a large technology company (§5) by designing and deploying a  
101 tool incorporating the guidelines. To develop the tool, we conducted a formative study  
102 with 10 AI practitioners to determine key design requirements. Using these requirements,  
103 we populated the tool with the guidelines and conducted the case study. Interviews with  
104 the 14 AI researchers, engineers, designers, and managers revealed that the guidelines  
105 were grounded to current regulations and were effectively usable across different roles,  
106 promoting self-reflection on ethical considerations in early development stages.

107 In light of these findings, we discuss how our method contributes to the idea of “Responsible AI by  
108 Design” by contextualizing the guidelines, informing existing or new theories, and offering practical  
109 recommendations for incorporating responsible AI guidelines into toolkits and recommendations  
110 for technical and non-technical roles in enabling organizational accountability (§6).  
111

## 112 2 RELATED WORK

113 We surveyed various lines of research that our work draws upon, and grouped them into two main  
114 areas: *i)* AI regulation and governance (§2.1), and *ii)* responsible AI practices and toolkits (§2.2).  
115

### 116 2.1 AI Regulation and Governance

117 The landscape of AI regulation and governance is constantly evolving [47, 68]. At the time of writing,  
118 the European Union (EU) has endorsed new transparency and risk-management rules for AI systems  
119 known as the EU AI Act [2], which is expected to become law in 2024. Similarly, the United States  
120 (US) has recently passed a blueprint of the AI Bill of Rights in late 2022 [44]. This bill comprises  
121 “*five principles and associated practices to help guide the design, use, and deployment of automated*  
122 *systems to protect the rights of the American public in the age of AI.*” While both the EU and US share  
123 a conceptual alignment on key principles of responsible AI, such as fairness and explainability, as  
124 well as the importance of international standards (e.g., ISO 24028 for Trustworthiness), the specific  
125 AI risk management regimes they are developing are potentially diverging, creating an “artificial  
126 divide” [32]. The EU aims to become the leading regulator for AI globally, while the US takes the  
127 view that excessive regulation may impede innovation.  
128

129 Notable predecessors to AI regulations include the EU GDPR law on data protection and privacy [28],  
130 the US Anti-discrimination Act [27], and the UK Equality Act 2010 [38]. GDPR’s Article 25  
131 mandates that data controllers must implement appropriate technical and organizational measures  
132 during the design and implementation stages of data processing to safeguard the rights of data  
133 subjects. The Anti-discrimination Act prohibits employment decisions based on an individual’s race,  
134 color, religion, sex (including gender identity, sexual orientation, and pregnancy), national origin,  
135 age (40 or older), disability, or genetic information. This legislation ensures fairness in AI-assisted  
136 hiring systems. Similarly, the UK Equality Act provides legal protection against discrimination in  
137 the workplace and wider society.

138 The National Institute of Standards and Technology (NIST), a renowned organization for de-  
139 veloping frameworks and standards, recently published an AI risk management framework [73].  
140 According to the NIST framework, an AI system is defined as “*an engineered or machine-based*  
141 *system capable of generating outputs such as predictions, recommendations, or decisions that influence*  
142 *real or virtual environments, based on a given set of objectives. These systems are designed to operate*  
143 *with varying levels of autonomy.*” Similarly, the Principled Artificial Intelligence white paper from  
144 the Berkman Klein Center [30] highlights eight key thematic trends that represent a growing  
145 consensus on responsible AI. These themes include privacy, accountability, safety and security,  
146 transparency and explainability, fairness and non-discrimination, human control of technology,  
147

148 professional responsibility, and the promotion of human values. Building on these themes, previous  
149 works have proposed a set of guidelines involving specific groups of AI practitioners. Saleema et  
150 al. [4] proposed 168 guidelines on how to design AI tailored to HCI practitioners. Similarly, Vakkuri  
151 et al. [94] formulated AI ethics guidelines tailored to researchers and technologists. No subsequent  
152 work has associated these guidelines with current international standards or regulations.

153 **Research Gaps.** As AI regulation and governance continue to evolve, AI practitioners are faced with  
154 the challenge of staying updated not only with the changing guidelines, but also with the regulations,  
155 requiring significant time and effort. Because prior guidelines lacked alignment with regulations,  
156 standards, and the input of experts in those fields, this work aims to create a methodology for  
157 crafting responsible AI guidelines that adhere to regulations and standards.  
158

## 159 2.2 Responsible AI Practices and Toolkits

160 **Responsible AI Toolkits.** At the time of writing, the OECD’s website lists 613 toolkits dedicated  
161 to fostering the development and deployment of responsible AI systems [74]. These toolkits are  
162 essential for operationalizing guidelines and regulations to assist AI practitioners such as engineers  
163 and researchers in addressing algorithmic bias [11, 34], explaining algorithmic decisions [6], and  
164 ensuring privacy in AI systems [30]. For addressing algorithmic bias, Google’s Fairness Indicators  
165 toolkit allows developers to assess data distribution and model performance across user-defined  
166 groups [37]. IBM’s AI Fairness 360 offers fairness metrics for bias mitigation [45]. Microsoft’s  
167 Fairlearn assesses model impact on specific groups (e.g., under-represented populations) in terms  
168 of fairness and accuracy [29]. For explaining algorithmic decisions, IBM’s AI Explainability 360  
169 provides metrics and guidance for explainability, and new visualization techniques to enhance  
170 transparency [12, 36, 70]. Finally, for ensuring privacy in AI systems, IBM’s AI Privacy 360 helps  
171 assess and mitigate privacy risks through data anonymization and minimization [17, 30, 84].  
172

173 **Toolkits Used in Practice.** Developing toolkits specialized for certain audiences such as AI devel-  
174 opers can lead to techno-solutionism, focusing exclusively on technical fixes. However, responsible  
175 AI entails broader socio-technical challenges (e.g., diversity and inclusion in decision-making)  
176 that require involvement of different roles with diverse expertise and background [85], typically  
177 discussed in venues with a long-standing commitment to human-centered design such as CHI,  
178 CSCW, AIES, and FAccT.

179 Different roles (e.g., data scientists, ML engineers and developers, UX designers) use toolkits in  
180 various ways. Data scientists often struggle to fully grasp visualizations of interpretable tools (e.g.,  
181 InterpretML [69] and SHAP [58]), hindering their ability to understand datasets and underlying  
182 models [48]. Experienced ML developers and engineers often go beyond what fairness toolkits offer  
183 to tackle algorithmic unfairness, while those with less experience typically use only a few metrics  
184 and methods from these toolkits [8]. UX designers often rely on custom prototypes and their own  
185 past experiences to help contextualize responsible AI issues for non-technical colleagues [20] due  
186 to communication gaps [101].

187 Major communication gaps between technical and non-technical roles typically arise because  
188 these roles are involved in different stages of a project, which is likely to create fragmentation in  
189 communication [77]. By exploring how data science teams collaborate, Zhang et al. [102] found  
190 that non-technical roles play more prominent roles in the early and late stages of projects, while  
191 technical roles primarily handle the core data and modeling tasks. However, this disparity in  
192 involvement at various project stages is likely to create fragmentation. In fact, Organizational  
193 Science research reinforces the notion that effective communication and collaboration is crucial for  
194 overcoming the “silo mentality” [35]. Consequently, due to this fragmentation and a lack of robust  
195 organizational support, practitioners often take on “bridging” roles to help the communication  
196

197 between the technical and non-technical project members [21]. One way of doing so is through  
198 “leaky abstractions” [91]. These are representations that are meant to communicate the inner work-  
199 ings and technical aspects of an AI system to these roles. Similarly, Nahar et al. [71] highlighted  
200 the extreme difficulty faced by non-technical practitioners in eliciting requirements due to the  
201 absence of suitable tools and the involvement of diverse stakeholders, highlighting the need for  
202 integrating communication features into toolkits. The design of such features was explored by  
203 Elsayed-Ali et al. [26] who developed question cards to facilitate stakeholder group discussions.  
204 These cards included built-in mechanisms for the automatic and cyclical assignment of cards to  
205 different participants, ensuring that everyone had the opportunity to share their opinions during  
206 the discussion.

207  
208 **Research Gaps.** While many toolkits emphasize the importance of involving stakeholders from  
209 diverse roles and backgrounds, it becomes evident that they are frequently designed for specific  
210 stakeholders (e.g., ML engineers), thereby neglecting a broader spectrum of roles, especially non-  
211 technical roles. To address this gap, we aim to develop a set of actionable guidelines that are usable  
212 by a diverse range of stakeholders.

### 213 214 3 AUTHOR POSITIONALITY STATEMENT

215 Understanding researcher positionality is crucial for transparently examining our perspectives on  
216 methodology, data collection, and analyses [33, 42]. In this paper, we situate ourselves in a Western  
217 country<sup>2</sup> during the 21<sup>st</sup> century, writing as authors primarily engaged in academic and industry  
218 research. Our team comprises three males and two females from Southern, Eastern, and North  
219 Europe, and Middle East with diverse ethnic and religious backgrounds. Our collective expertise  
220 spans various fields, including human-computer interaction (HCI), ubiquitous computing, software  
221 engineering, artificial intelligence, data visualization, and digital humanities.

222 It is important to recognize that our backgrounds and experiences have shaped our positionality.  
223 As HCI researchers affiliated with a predominantly Western organization, we acknowledge the  
224 need to expand the understanding of the research questions and methodology presented in this  
225 paper. Consequently, our positionality may have influenced the subjectivity inherent in framing our  
226 research questions, selecting our methodology, designing our study, and interpreting and analyzing  
227 our data.

### 228 229 4 METHOD FOR GENERATING RESPONSIBLE AI GUIDELINES

230 To generate a list of responsible AI guidelines, we followed a four-step process (Figure 2), based on  
231 the methodology proposed by Michie et al. [63]. This process allowed us to identify the essential  
232 element of a guideline, referred to as the “active ingredient,” focusing on the “what” rather than the  
233 “how” [62]. A similar parallel can be drawn in software engineering, where the “what” represents the  
234 software requirements and the “how” represents the software design, both of which are important  
235 for a successful software product [3]. However, by shifting the focus to the “what,” AI practitioners  
236 can develop a clearer understanding of the objectives and goals they need to achieve, fostering  
237 a deeper comprehension of complex underlying ethical concepts. Throughout this process, we  
238 actively engaged a diverse group of stakeholders, including AI engineers, researchers, designers,  
239 product managers, and experts in law and standardization. As a result, we were able to formulate a  
240 total of 22 responsible AI guidelines (Panel A of Figure 1).

241  
242  
243  
244 <sup>2</sup>REDACTED FOR BLIND REVIEW.

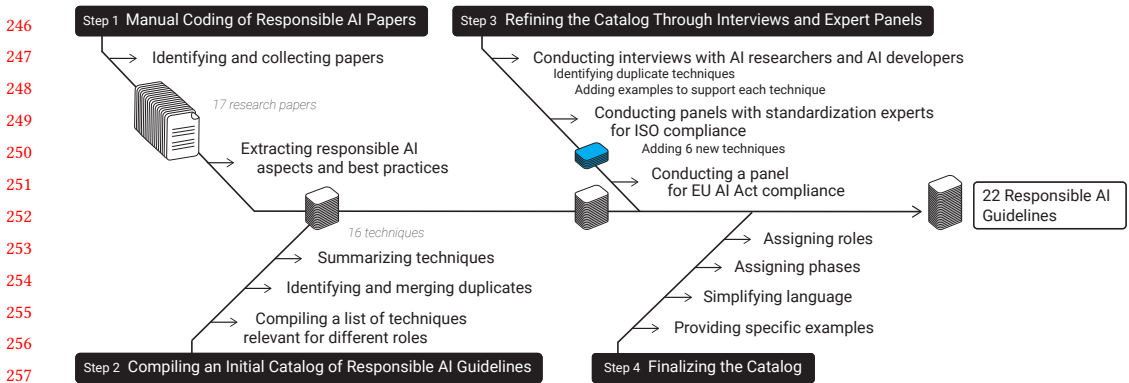


Fig. 2. Four-step method for generating responsible AI guidelines. These guidelines were derived from research papers, and are in line with ISO standards and the EU AI Act [2].

#### 4.1 Manual Coding of Responsible AI papers

In the first step, we compiled a list of key scientific articles focusing on responsible AI guidelines applicable to a diverse set of roles, and manually coded them. We created this list by targeting papers published in renowned computer science conferences, such as the ACM CHI, CSCW, FAccT, AAAI/ACM Conference on AI, Ethics, and Society (AIES), and scientific literature from the medical domain (e.g., the *Annals of Internal Medicine*). Note that we did not conduct a systematic literature. Instead, we identified 17 papers that served as a starting point to compile an initial catalogue of techniques covering a broad range of responsible AI aspects, including fairness, explainability, sustainability, and best practices for data and model documentation and evaluation. These are foundational papers in responsible AI, and, as we shall see in a subsequent step of our methodology (§4.3), we refined the techniques identified from those papers through interviews and expert panels as well as cross-referencing them with the EU AI Act and ISO standards.

These papers encompass a growing body of research focusing on the work practices (e.g., ensuring fairness or models' explainable outputs) of AI practitioners in addressing responsible AI issues. This strand of research covers various aspects of responsible AI, including fairness, explainability, sustainability, and best practices for data and model documentation and evaluation. Fairness is a fundamental value in responsible AI, but its definition is complex and multifaceted [72]. To assess bias in classification outputs, various research efforts have introduced quantitative metrics such as disparate impact and equalized odds, as discussed by Dixon et al. [23]. Another concept explored in the literature is "equality of opportunity," advocated by Hardt et al. [75], which ensures that predictive models are equally accurate across different groups defined by protected attributes like race or gender. Equally important is the development of dedicated checklists for fairness [59]. Explainable AI (XAI) is another aspect of responsible AI. XAI involves tools and frameworks that assist end users and stakeholders in understanding and interpreting predictions made by machine learning models [5, 25, 40, 53, 56, 70]. Furthermore, the environmental impact of training AI models should also be considered. Numerous reports have highlighted the significant carbon footprint associated with deep learning and large language models [41, 86, 90]. Best practices for data documentation and model evaluation have also been developed to promote fairness in AI systems. Gebru et al. [34] proposed "Datasheets for Datasets" as a comprehensive means of providing information about a dataset, including data provenance, key characteristics, relevant regulations, test results, and potential biases. Similarly, Bender et al. [10] introduced "data statements" as qualitative summaries

295 that offer crucial context about a dataset’s population, aiding in identifying biases and understanding  
296 generalizability. For model evaluation, Mitchell et al. [66] suggested the use of model cards, which  
297 provide standardized information about machine learning models, including their intended use,  
298 performance metrics, potential biases, and data limitations. Transparent reporting practices, such  
299 as the TRIPOD statement by Collins et al. [18] in the medical domain, emphasize standardized and  
300 comprehensive reporting to enhance credibility and reproducibility of AI prediction models.

## 302 4.2 Compiling an Initial Catalog of Responsible AI Guidelines

304 For each source, we compiled a list of techniques that could be employed to create responsible AI  
305 guidelines, focusing on the actions different roles (i.e., designers, researchers, developers, product  
306 managers) should consider during AI development. Following the methodology proposed by  
307 Michie et al. [63] (which was also used to identify community engagement techniques by Dittus  
308 et al. [22]), we sought techniques that describe the “active ingredient” of what needs to be done.  
309 This means that the phrasing of the technique should focus on *what* needs to do be done, rather  
310 than the specific implementation details of *how* it should be done. For example, a recommended  
311 practice for ensuring fairness involves evaluating an AI system across different demographic  
312 groups [23, 59, 75]. In this case, the technique specifies “what” needs to be done rather than  
313 “how” it should be implemented (e.g., using common fairness metrics such as demographic parity  
314 or equalized odds). In total, we formulated a set of 16 techniques based on relevant literature  
315 sources [5, 10, 18, 23, 30, 34, 41, 43, 53, 59, 66, 67, 75, 86, 95, 96].

316 We then conducted an iterative review of the collection of techniques to identify duplicates,  
317 which were instances where multiple sources referred to the same technique. For example, four  
318 sources indicated that data biases could affect the model [10, 34, 43, 66], emphasizing the need  
319 to report the characteristics of training and testing datasets. We consolidated such instances by  
320 retaining the specific actions to be taken (e.g., reporting dataset characteristics). This process  
321 resulted in an initial list of 16 distinct techniques. We provided a concise summary sentence for  
322 each technique, utilizing active verbs to emphasize the recommended actions.

## 325 4.3 Refining the Catalog Through Interviews and Expert Panels

326 The catalog of techniques underwent eleven iterations to ensure clarity and comprehensive thematic  
327 coverage. The iterations were carried out by two authors, with the first author conducting interviews  
328 with five AI researchers and developers. During the interviews, the participants were asked to  
329 consider their current AI projects and provide insights on the implementation of each technique,  
330 focusing on the “how” aspect. This served two purposes: firstly, to identify any statements that  
331 were unclear or vague, prompting suggestions for alternative phrasing; and secondly, to expand the  
332 catalog further. The interviews yielded two main recommendations for improvement: (1) mapping  
333 duplicate techniques to the same underlying action(s); and (2) adding examples to support each  
334 technique.

335 In addition to the interviews, the two authors who developed the initial catalog conducted a  
336 series of six 1-hour expert panels with two standardization experts from a large organization. The  
337 purpose of these panels was to review the initial catalog for ISO compliance. The standardization  
338 experts examined six AI-related ISOs, including ISO 38507, ISO 23894, ISO 5338, ISO 24028, ISO  
339 24027, and ISO 24368, which were developed at the time of writing. Then experts provided input  
340 on any missing techniques and mapped each technique in the initial catalog to the corresponding  
341 ISO that covers it. As a result of this exercise, six new techniques (#2, #7, #12, #13, #14, #21) were  
342  
343

344 added to the catalog, resulting in a total of 22 guidelines. Next, we provide we provide a high-level  
 345 summary of each ISO.<sup>3</sup>

346 **ISO 38507 (Governance, 28 pages).** It offers guidance on responsible AI use (e.g., identify potential  
 347 harms and risks for each intended use(s) of the systems) and recommendations about current and  
 348 future AI uses to governing bodies and various stakeholders such as managers and auditors.

349 **ISO 23894 (Risk Management, 26 pages).** It provides guidelines for managing AI-related risks  
 350 (e.g., mechanisms for incentivizing reporting of system harms) in developing, producing, deploying,  
 351 or using AI products and systems, including recommendations for integrating risk management  
 352 into AI processes.

353 **ISO 5338 (AI Lifecycle Process, 27 pages).** It provides a framework for the life cycle of AI  
 354 systems, detailing processes for managing and enhancing these systems from development to  
 355 implementation (e.g., through reporting of harms and risks, obtaining approval of intended uses).

356 **ISO 24028 (Trustworthiness, 43 pages).** It offers guidance on trustworthiness in AI systems,  
 357 focusing on transparency, explainability, controllability, and addressing potential risks with mit-  
 358 igation techniques. It also covers AI systems' availability, resiliency, reliability, accuracy, safety,  
 359 security, and privacy.

360 **ISO 24027 (Bias, 39 pages).** It discusses bias in AI systems related to protected attributes such  
 361 as age and gender, especially in AI-aided decision-making, providing techniques to measure and  
 362 assess bias throughout the AI system lifecycle.

363 **ISO 24368 (Ethical and Societal Concerns, 48 pages).** It provides an introduction to ethical and  
 364 societal concerns related to AI (e.g., principles, processes, and methods), targeting technologists,  
 365 regulators, interest groups, and society as a whole.

366 As the final step of refining the catalog, the two authors reviewed the 85 articles of the EU  
 367 AI Act [2] to map each of the 22 guidelines with the most relevant article(s), as shown in the  
 368 last column of Table 1. They began with Article 3 of the Act, which defines the key concepts of  
 369 an AI system, including its definition, intended purpose, performance, training, validation, and  
 370 post-deployment monitoring. After reading all the articles and annotating them, they identified 22  
 371 unique articles corresponding to the guidelines. Articles 9, 10, and 17 were mapped to multiple  
 372 guidelines. For example, Article 9 (*Risk management system*) states that “a risk management system  
 373 shall be established, implemented, documented and maintained throughout the entire lifecycle  
 374 of a high-risk AI system”. This article aligns with guidelines #1, #3-5, and #13 as it is about  
 375 the identification of harms and risks of the AI system's intended use. Article 10 (*Data and data  
 376 governance*) states that “*training, validation and testing data sets shall be subject to appropriate  
 377 data governance and management practices*”. This article aligns with guidelines #8 and #15-18 as it  
 378 discusses the management and quality of data for training, validation, and testing, including aspects  
 379 of diversity and minimizing biases. Finally, Article 17 (*Quality management system*) states that “*an  
 380 AI system shall be documented in a systematic and orderly manner in the form of written policies,  
 381 procedures and instructions*”. This article aligns with guidelines #6, #7, #10, and #14-18 because it is  
 382 about documentation of all system components, including AI models and testing and validation  
 383 procedures. The full mapping along with justifications is provided in Appendix B.

#### 385 4.4 Finalizing the Catalog

386 In response to the interviews with AI developers and standardization experts, we incorporated an  
 387 example for each guideline. For instance, under the guideline on system interpretability (guideline  
 388 #9), the example provided reads: “output feature importance and provide human-understandable  
 389 explanations.” Furthermore, we simplified the language by avoiding domain-specific or technical  
 390

391 <sup>3</sup>Note that the summary provided is a brief and simplified description due to a paywall restriction.  
 392



Table 1. Responsible AI guidelines are actionable items that a developer should consider during the 3 phases of AI development lifecycle. These guidelines are grounded in the scientific literature (main sources are reported), and were checked for ISO compliance: ISO 38507 (Governance); ISO 23894 (Risk management); ISO 5338 (AI lifecycle processes); ISO 24028 (Trustworthiness); ISO 24027 (Bias); ISO 24368 (Ethical considerations). They were also cross-referenced with the EU AI Act’s articles [2]. They are marked with the ‘Phase’ during which a guideline can be applied. There are three phases: development ( $P_1$ ), deployment ( $P_2$ ), and use ( $P_3$ ). Guidelines are also marked with the job ‘Role’ that should consider them. There are three roles: designer ( $R_D$ ), engineer or researcher ( $R_E$ ), and manager or executive ( $R_M$ ). Each guideline is followed by an example, and the guidelines are categorized thematically into six categories, concerning the *intended uses, harms, system, data, oversight, and team*.

Number	Guideline	Phase	Role	Source(s)	ISO	AI Act
<b>INTENDED USES</b>						
1	Work with relevant parties to identify intended uses. (e.g., identify the system’s usage, deployment, and contextual conditions)	$P_{1-3}$	$R_{D,EM}$	[66]	5338, 38507, 23894, 24027, 24368	Art. 6, 9
2	Obtain approval from an Ethics Committee or similar body for intended uses. (e.g., Obtain Ethics Committee approval for the intended use, aligned with sustainability goals)	$P_{1-3}$	$R_{D,EM}$	–	38507, 5338, 23894	Art. 11, 69
<b>HARMS</b>						
3	Identify potential harms and risks associated with the intended uses. (e.g., prevent privacy violation, discrimination, and adversarial attacks, provide interpretable output)	$P_{1-3}$	$R_{D,EM}$	[59]	23894, 24028, 38507, 24368	Art. 9, 65
4	Provide mechanism(s) for incentivizing reporting of system harms. (e.g., provide contact emails and feedback form to raise concerns)	$P_1$	$R_{D,E}$	[59]	38507, 23894	Art. 9, 60-63
5	Develop strategies to mitigate identified harms or risks for each intended use. (e.g., use stratified sampling and safeguards against adversarial attacks during training)	$P_{1-3}$	$R_{D,M}$	[66]	24368, 23894	Art. 9, 67
<b>SYSTEM</b>						
6	Document all system components, including the AI models, to enable reproducibility and scrutiny. (e.g., create UML diagrams, flowcharts, and specify model types, versions, hardware architecture)	$P_{1-3}$	$R_{D,E}$	[59]	5338, 23894, 24027	Art. 11, 12, 16-18, 50
7	Review the code for reliability (e.g., manage version control using software.)	$P_{1-3}$	$R_{D,E}$	–	5338	Art. 17
8	Report evaluation metrics for various groups based on factors such as age, gender, and ethnicity. (e.g., evaluate false positive/negative, AUC, and feature importance across protected attributes)	$P_{1-3}$	$R_{D,EM}$	[23, 59, 75] [67, 95]	23894, 5338, 24028, 24027	Art. 10, 13
9	Provide mechanisms for interpretable outputs and auditing. (e.g., output feature importance and provide human-understandable explanations)	$P_{1-3}$	$R_{D,EM}$	[5, 53]	38507, 24028	Art. 12-14
10	Document the security of all system components in consultation with experts. (e.g., guard against adversarial attacks and unauthorized access)	$P_{1-3}$	$R_{EM}$	[30]	24028, 24368	Art. 12, 13, 15, 17
11	Provide an environmental assessment of the system. (e.g., report the number of GPU hours used in training and deployment)	$P_{1-3}$	$R_E$	[41, 86]	38507, 23894, 5338, 24368	Art. 69
12	Develop feedback mechanisms to update the system. (e.g., provide contact email, feedback form, and notification of new knowledge extracted)	$P_{1-3}$	$R_{D,E}$	–	24028	Art. 61
13	Ensure safe system decommissioning. (e.g., ensure decommissioned data is either deleted or restricted to authorized personnel.)	$P_3$	$R_E$	–	38507, 24368	Art. 9
14	Redocument model information and contractual requirements at every system update. (e.g., update the model information when re-training the system or using datasets with new contractual requirements)	$P_3$	$R_E$	–	23894, 5338, 24368	Art. 11, 12, 17, 61
<b>DATA</b>						
15	Ensure compliance with agreements and legal requirements when handling data. (e.g., create data sharing and non-disclosure agreements and secure servers)	$P_{1-3}$	$R_{D,EM}$	–	38507, 23894, 5338	Art. 10, 17, 61
16	Compare the quality, representativeness, and fit of training and testing datasets with the intended uses. (e.g., report dataset details such as public/private, personal information, demographics, and data provenance)	$P_{1-3}$	$R_E$	[10, 34, 43, 96] [59, 67, 95]	38507, 5338, 24028, 24027	Art. 10, 13, 17, 64
17	Identify any measurement errors in input data and their associated assumptions. (e.g., account for potential input errors in the input device, text data, audio, and video)	$P_{1-3}$	$R_E$	[18]	38507	Art. 10, 13, 17, 64
18	Protect sensitive variables in training/testing datasets. (e.g., protect sensitive data and use techniques such as k-anonymity and differential privacy)	$P_{1-3}$	$R_{D,EM}$	[24]	38507, 24028	Art. 10, 13, 17
<b>OVERSIGHT</b>						
19	Continuously monitor metrics and utilize guardrails or rollbacks to ensure the system’s output stays within a desired range. (e.g., validate against concept drift and test with diverse testers and compliance and adversarial cases)	$P_{1-3}$	$R_{D,E}$	[30]	38507, 5338, 24028, 24027, 24368	Art. 12, 20, 29, 61
20	Ensure human control over the system, particularly for designers, developers, and end-users. (e.g., include human in the loop with the ability to inspect data, models, and training methods)	$P_{1-3}$	$R_{D,EM}$	–	38507, 5338, 24028, 24368	Art. 13, 14
<b>TEAM</b>						
21	Ensure team diversity. (e.g., consider diversity in gender, neurotypes, personality traits, and thinking styles)	$P_{1-3}$	$R_{D,EM}$	–	38507, 5338, 24028, 24368	Art. 69
22	Train team members on ethical values and regulations. (e.g., train on privacy regulations, ethical issues, and raising concerns)	$P_{1-3}$	$R_{D,EM}$	[30]	38507, 24368	Art. 69

jargon. We also categorized each guideline into six thematically distinct categories, namely *intended uses, harms, system, data, oversight, and team*.

Recognizing that certain guidelines may only be applicable at specific stages (e.g., monitoring AI after deployment) by specific roles (e.g., developers, managers), we went through two steps. First, we assigned the guidelines to three phases based on previous research (e.g., [59, 66]). These phases are development (designing and coding the system), deployment (transferring the system into the production stage), and use (actual usage of the system). For example, guidelines like identifying the

system's intended uses (guideline #1) are relevant to all three phases, while those related to system updates (guideline #14) or decommissioning (guideline #13) are applicable during the use phase.

Second, we assigned the guidelines to three types of role (Table 1) based on previous literature. Wang et al. [97] interviewed UX practitioners and responsible AI experts to understand their work practices. UX practitioners included designers, researchers, and engineers, while responsible AI experts included ethics advisors and specialists. Wong et al. [100] analyzed 27 ethics toolkits to identify the intended audience of these toolkits, specifically those who are expected to engage in AI ethics work. The intended audience roles identified included software engineers, data scientists, designers, members of cross-functional or cross-disciplinary teams, risk or internal governance teams, C-level executives, and board members. Additionally, Madaio et al. [59] co-designed a fairness checklist with a diverse set of stakeholders, including product managers, data scientists and AI/ML engineers, designers, software engineers, researchers, and consultants. Following guidance therefore from these studies [59, 97, 100], we formulated three roles as follows:

- (1) Designer: This role includes interaction designers and UX designers.
- (2) Engineer or Researcher: This role includes AI/ML engineers, AI/ML researchers, data scientists, software engineers, UX engineers, and UX researchers.
- (3) Manager or Executive: This role includes product managers, C-suite executives, ethics advisors/responsible AI consultants, and ethical board members.

The revised and final catalog, consisting of 22 unique guidelines, is presented in Table 1.

## 5 EVALUATION OF THE 22 RESPONSIBLE AI GUIDELINES

We first conducted a formative study with 10 AI practitioners from a large technology company to elicit design requirements for a tool that would allow us to evaluate the guidelines and then implemented this tool (Panel B in Figure 1 and §5.1). To evaluate the guidelines, we conducted a user study with 14 AI researchers, engineers, designers, and product managers from the same company (Panel C in Figure 1 and §5.2).

### 5.1 Incorporating the guidelines into a tool

**Eliciting design requirements for a tool through a formative study.** We conducted a formative study that included semi-structured interviews with 10 participants. These participants, comprising 6 males and 4 females, were AI practitioners in their 30s and 40s employed at a large technology company. The participants had a range of work experience, spanning from 1 to 8 years, and were skilled in areas such as data science, data visualization, UX design, natural language processing, and machine learning. The interview study took place online and consisted of three parts. In the first part, we encouraged participants to share information about their ongoing AI projects. In the second part, we presented them with the table containing the 22 guidelines and asked them to think about how each guideline could apply to their projects. Finally, in the third part, we conducted semi-structured interviews to discuss how these guidelines could be incorporated into an interactive responsible AI tool.

Each study lasted about half an hour. Two authors took notes during the interviews, and afterward, they analyzed the interview transcripts using inductive thematic analysis [13, 61, 64, 80]. This analysis then resulted in the following four design requirements (participant quotes are marked with FP):

*R1: Simplify the guidelines by breaking them into smaller visual components.* Participants found it challenging to reflect on guidelines and examples because of their quantity. According to FP5, *"the sheer number of the guidelines is the main difficulty [...] they should be separated in bite-sized*

491 questions”. Additionally, participants requested to visually separate the guidelines from the examples.

492  
493 *R2: Implement clear navigation features to systematically guide users through the guidelines.* Partici-  
494 pants were unsure about the best way to navigate through the guidelines. FP9 suggested that “*the*  
495 *system should provide clear navigation [...] for example, using a progress bar*”. FP5 further emphasized  
496 that the design of the progress bar could facilitate “*gaining insights while engaging with the 22*  
497 *guidelines*”.

498  
499 *R3: Track how guidelines are applied and share progress among team members.* Participants faced  
500 difficulty in tracking their responses on how to apply the guidelines to their projects and share  
501 progress among team members. To address this challenge, FP5 suggested implementing a feature  
502 that would save user responses as they progress through the guidelines: “*there should be some*  
503 *functionality there that captures the answers I gave, so it'd allow me to track progress and share it among*  
504 *team members*”. These responses would then be transformed into comprehensive documentation  
505 and made accessible to users for download.

506 *R4: Develop a mechanism for post-hoc reflections on how the project aligns with responsible AI*  
507 *guidelines.* Participants found it challenging to envision how well their AI systems aligned with  
508 the guidelines. Therefore, FP8 suggested developing “*visual feedback or a score that shows how*  
509 *responsible [their] AI system is.*”. However, FP2 cautioned that this mechanism “*should not make me*  
510 *anxious and feel like I have not done enough*”. Instead, it should create a positive learning experience  
511 and encourage users to generate ideas for improving their AI systems.

512  
513 **Designing the tool and incorporating the guidelines.** From these requirements, we designed  
514 an interactive web-based tool<sup>4</sup> (Figure 3) and populated it with the 22 guidelines in Table 1.

515 To meet the design requirement R1 (*Simplify the guidelines*), each guideline is presented as a  
516 digital card [57] with interactive boxes on both the front and back sides. The front side includes a  
517 symbolic graphic collage representing the guideline, followed by its name and full text. The back  
518 side includes an input box for users to write their thoughts on implementing each guideline in their  
519 project [82]. We also used this box to showcase an example for each guideline (refer to Figure 1).  
520 Initially, the example in the box is visible, but it disappears once the user inputs their specific  
521 implementation details. Users can view the guideline from both sides by using the flip buttons at  
522 the bottom-left corner of each side.

523 Each guideline is paired with two guiding questions [101] that help users think about the  
524 relevance of the guideline to their specific AI system and context (Figure 4). The first question asks  
525 the user whether the guideline has been successfully implemented in their AI system. For example,  
526 for an engineer addressing fairness, the question asks if they have reported evaluation metrics for  
527 various groups based on factors like age, gender, and ethnicity (technique #8 in Table 1). If the  
528 engineer answers “yes,” they are then prompted to provide specific details on how fairness was  
529 implemented in the input box on the card’s back. After sharing this information, the tool moves  
530 the guideline to the “successfully implemented” stack. In contrast, if the engineer answers “no,” the  
531 tool asks a second follow-up question regarding whether the guideline should be implemented in a  
532 future iteration of the AI system. If the engineer answers “yes,” they are prompted to provide specific  
533 details on how to implement it. The tool then moves the guideline to the “should be considered”  
534 stack. However, if the engineer answers “no” to both questions, indicating that the guideline is not  
535 applicable to their AI system, the tool moves the guideline to the “inapplicable” stack.

536  
537  
538 <sup>4</sup>anonymous-url

540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588



Fig. 3. Interactive Responsible AI Tool with 22 guidelines. The first part (A) allows for entering information about the developed AI system and (B) selecting the applicable user role. The second part (C) enables interaction with the guidelines. The third part (D) presents a summary of user responses for post-hoc reflections. Guidelines for other project phase can be viewed through the phase selectors (E/A).

To meet the design requirement R2 (*Implement clear navigation*), we explored different layout options and considered previous research that involved swiping [99], scrolling, or organizing guidelines into different groups [22]. Due to the limited screen size and the repetition of guidelines for each phase and role, we chose to organize the guidelines into nine groups. These groups were derived from three phases of the AI system: *i*) development (designing and coding), *ii*) deployment (transitioning into production), and *iii*) use (actual usage of the system), as well as from three user roles: *i*) designer, *ii*) engineer or researcher, and *iii*) manager or executive. The number of guidelines in each group varied and accommodated the specific requirements of each phase and

589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637

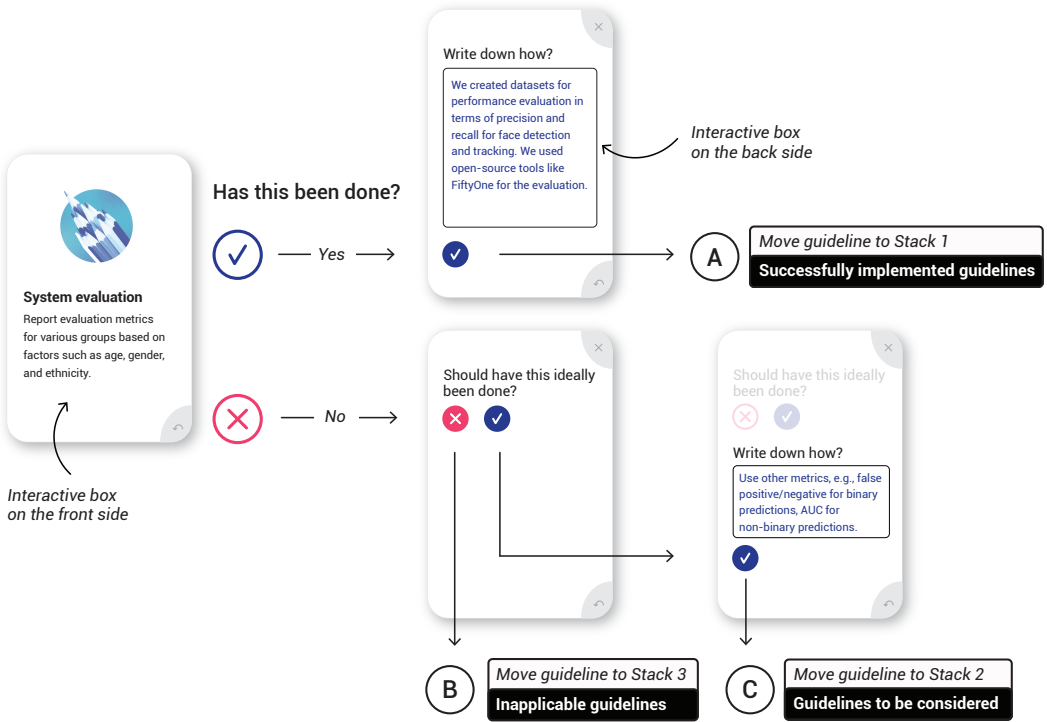


Fig. 4. Guideline sorting procedure. Users can place a guideline in any of the three stacks (i.e., successfully implemented, should be considered, inapplicable) by: (1) considering two guiding questions, and (2) using the Yes/No buttons located next to the card and on the back side of it.

role. For example, engineers or researchers needed to go through 20 guidelines for development, 18 for deployment, and 20 for use (§4, Step 4).

To facilitate the browsing of guidelines and addressing them in the user’s preferred order, we introduced two arrow buttons on both the left and right side of the card group. We then introduced a graphical progress bar next to the group that not only displayed the number of remaining guidelines but also color-coded them to indicate their assignment to the three stacks. Blue leaves in the bar represented successfully implemented guidelines, magenta leaves represented guidelines for future consideration, and empty leaves represented inapplicable guidelines. Finally, we added an “Exit” button that becomes available as soon as the user goes through minimum two guidelines. In that way user can quit the experience at a preferable moment.

To meet the design requirement R3 (*Track how guidelines are applied and share progress among team members*), we added a feature to store user responses locally in the browser session. Users can download their responses as a structured PDF report at any time using the floating download button.

To fulfill the last requirement, R4 (*Develop a mechanism for post-hoc reflections*), after completing the sorting process, we display to the user a summary page. The summary is divided into three sections, one for each stack of cards (i.e., successfully implemented, should be considered, and inapplicable), with in-text counters indicating the number of guidelines in each stack. To read the responses for each guideline, hover-over functionality is provided.

Table 2. User study participants' demographics, including their job 'Role' (designer ( $R_D$ ), engineer or researcher ( $R_E$ ), and manager or executive ( $R_M$ )).

ID	Gender	Yrs of expr. in AI	Education	Current continent	Expertise	Role
1	Male	6	Ph.D.	EU	Deep learning, computer vision	$R_M$
2	Male	+10	Ph.D.	North America	Machine learning, computer vision	$R_E$
3	Male	8	Ph.D.	EU	Machine learning	$R_E$
4	Male	4	Ph.D.	North America	Deep learning, IoT, computer vision	$R_E$
5	Female	5	Ph.D.	EU	Machine learning	$R_D$
6	Female	8	Ph.D.	EU	Computer vision	$R_D$
7	Male	2	Ph.D.	North America	Computer vision	$R_E$
8	Male	10	Ph.D.	EU	Machine learning	$R_M$
9	Male	4	Ph.D.	North America	Computer vision	$R_E$
10	Male	+10	M.S.	EU	Machine learning, natural language processing	$R_E$
11	Male	+10	Ph.D.	EU	Machine learning	$R_M$
12	Male	6	Ph.D.	EU	Machine learning	$R_E$
13	Male	4	Ph.D.	EU	Reinforcement learning, decision making	$R_E$
14	Male	8	Ph.D.	EU	Computer vision, robotics	$R_D$

Figure 3 shows the tool with its three parts that meet the four design requirements. The first part enables users to enter the name of the developed AI system (Figure 3A), select the phase it belongs to and specify user's role (Figure 3B). Once the phase and role are selected, the second part displays the guidelines (Figure 3C) and allows to address them one by one. The third part presents the user with the summary for post-hoc reflections (Figure 3D). If desired, the user can repeat the experience and generate documentation for other phases (Figure 3E).

## 5.2 Evaluating the Guidelines Through a User Study

To evaluate whether our guidelines are usable by different roles and whether they match the EU AI Act articles and ISO standards, we conducted a user study with 14 AI researchers, engineers, designers, and managers from a large technology company (Panel C in Figure 1).

**Participants.** We recruited 14 participants from the same large technology company.<sup>5</sup> The recruitment process took place in October and November 2022. We aimed for a balanced sample of participants, including a variety of roles such as researchers (5), designers (3), engineers (3), and managers (3). All participants had significant expertise in AI, including areas such as machine learning, deep learning, and computer vision. Additionally, each participant was actively involved in at least one ongoing AI project during the time of the interviews. Table 2 summarizes participants' demographics.

**Procedure.** Ahead of the interviews, we sent an email to all participants, providing a concise explanation of the study along with a brief demographics survey. The survey consisted of questions regarding participants' age, domain of expertise, role, and years of experience in AI system development. The survey is available in Appendix A. It is important to note that our organization<sup>6</sup> approved the study, and we adhered to established guidelines for user studies, ensuring that no personal identifiers were collected, personal information was removed, and the data remained accessible solely to the research team.

During the interview session, we presented to the participants two systems: (a) our tool with the 22 guidelines and (b) a web page with the checklist items from Microsoft's Fairness Checklist. We

<sup>5</sup>Participants who took part in the formative study were not eligible to participate in this evaluation study.

<sup>6</sup>REDACTED FOR BLIND REVIEW

687 used the Microsoft’s AI Fairness Checklist as a baseline alternative because it is a published work  
688 in a human-computer interaction conference (CHI 2020), is freely available, and has a rigorous,  
689 transparent creation process.<sup>7</sup> We asked participants to interact with each system for 20 minutes (or  
690 less if finished sooner), alternating between them to avoid any learning effect. To make the scenario  
691 as realistic as possible, we encouraged participants to reflect on their ongoing AI projects and  
692 consider how the guidelines could be applied in their roles. We also presented them with excerpts  
693 from the EU AI Act articles [2] and summaries of each ISO standard (§4.3), and asked them whether  
694 the guidelines link to these articles and summaries. We further engaged participants by asking  
695 about their preferences, dislikes, and the relevance of the guidelines to their work. Subsequently,  
696 we administered the System Usability Scale (SUS) [14] to assess the usability of the guidelines and  
697 the checklist items.

698 We piloted our study with two researchers (1 female, 1 male), which helped us make minor  
699 changes to the study guide (e.g., clarifying question-wording and changing the order of questions  
700 for a better interview flow). These pilot interviews were not included in the analysis.

701 **Analysis.** First, we compared the two usability scores after using each system (i.e., the guidelines  
702 and the checklist items). Second, two authors conducted an inductive thematic analysis (bottom-up)  
703 of the interview transcripts, following established coding methodologies [61, 64, 80]. The transcripts  
704 included how the guidelines could be applied in the ongoing AI projects, how they link to the EU  
705 AI Act articles and ISO standards, and any other preferences or dislikes. The authors used sticky  
706 notes on the Miro platform [65] to capture the participants’ answers, and collaboratively created  
707 affinity diagrams based on these notes. They held seven meetings, totaling 14 hours, to discuss and  
708 resolve any disagreements that arose during the analysis process. Feedback from the last author  
709 was sought during these meetings. In some cases, a single note was relevant to multiple themes,  
710 leading to overlap between themes. All themes included quotes from at least two participants,  
711 indicating that data saturation had been achieved [39]. As a result, participant recruitment was  
712 concluded after the 14<sup>th</sup> interview.

713 **Results.** Participants, on average, rated the guidelines’ usability with a score of 66 out of 100 in SUS,  
714 with a standard deviation of 16.01 (Figure 5). This indicates a generally positive user experience [83].  
715 The moderately high usability score can likely be attributed to factors such as familiarity and  
716 efficiency [55] in interactions with the guidelines, further supporting their usability across different  
717 roles. In contrast, participants, on average, rated the checklist items’ usability with a score of 44  
718 out of 100 in SUS, with a standard deviation of 21.16. Despite the comparative lower SUS score,  
719 checklist items were seen as relevant for audit, formal processes, and certification purposes—acting  
720 as a ‘safeguard’. As for the thematic analysis, the resulting themes are provided in Table 3 in the  
721 Appendix. These themes pertain to how our participants saw the application of guidelines, what  
722 worked well and what could have been improved.

723  
724 Guidelines were generally well-received by the participants. The majority of them (12 out of  
725 14 participants) considered the guidelines valuable for raising awareness and facilitating self-  
726 learning about responsible AI, though to different extents. Participants found the set of guidelines  
727 to be comprehensive and aligned with their roles (10 out of 14 participants), as evidenced by P8’s  
728 observation that *“There are some aspects of responsible AI in the project that I knew about, but I never*  
729 *faced them in such an organized manner”*. Similarly, P4 *“felt that the guidelines were concrete and*  
730 *well-scoped, instead of the lengthy documents of current regulations.”* Participants also stated that the

731  
732 <sup>7</sup>We decided not to compare our tool with existing card-based systems for responsible AI as they serve different purposes.  
733 Card-based systems such as the IDEO AI Ethics and the Feminist Tech card aim at providing thought-provoking activities [46]  
734 and stimulating ethical conversations [54]. However, they cannot be used as tools ensuring compliance with internal ethical  
735 procedures (like Microsoft’s AI Fairness Checklist) or ISO standards.

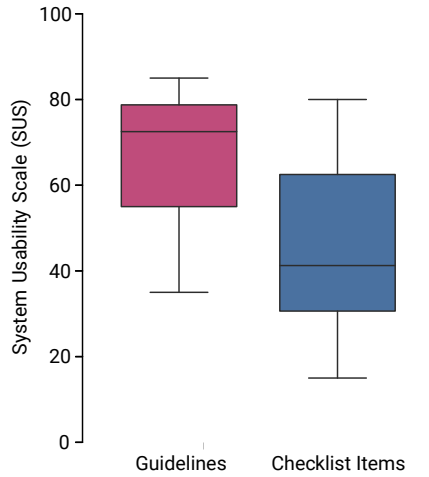


Fig. 5. SUS (usability) results. Guidelines are more usable than checklist items.

guidelines align with current regulations (10 out of 14 participants). P7 mentioned that “*he could understand the guidelines relevancy to the ISO standards and their applicability to his work.*” Similarly, P11 found the excerpts from the EU AI Act “*relevant and the guidelines helped him to reflect how the current regulations will affect his project.*” Additionally, seven participants acknowledged the usefulness of the provided examples, which helped them think about potential scenarios and make the guidelines more actionable. One participant expressed that “*the guidelines made me reflect on my previous choices and how I would describe my decisions when I had to develop the system (P3).*” P3 specifically mentioned four guidelines about data (listed 15-18 in the Table 1 in the manuscript) and thought of the following improvements: “*I would collect more annotated data from diverse populations and incentivize underrepresented groups to participate in data annotation.*” An action plan was also devised by P1, who recognized that, “*I need an expert in different areas of assessments, because I am probably not in the right position to do that,*” and they intended to consult system security and sustainability experts. Finally, after becoming familiar with the guidelines, P2 felt more empowered to introduce the topic of responsible system development during group discussions with his team, stating that “*I can at least raise a few questions during team discussions—these are some additional aspects we may need to consider.*”

## 6 DISCUSSION

To assist AI practitioners in navigating the rapidly evolving landscape of AI ethics, governance, and regulations, we have developed a method for generating responsible AI guidelines that are grounded in regulation and usable by different roles. We validated our method in a user study at a large technology company, where we designed and evaluated a tool that incorporates our responsible AI guidelines. We conducted a formative study involving 10 AI practitioners to design the tool, and evaluated our guidelines in a user study with an additional 14 AI practitioners. The results indicate that the guidelines were perceived as practical and actionable, promoting self-reflection and enhancing understanding of the ethical considerations associated with AI during the early stages of development. In light of these results, we discuss how our method contributes to the idea of “Responsible AI by Design”, that is, a design-first approach that considers responsible AI values throughout the development lifecycle and across business roles. We now discuss the inherent problem of decontextualization in responsible AI toolkits, the concept of meta-responsibility, and



785 provide practical recommendations for incorporating responsible AI guidelines into toolkits and  
 786 recommendations for technical and non-technical roles in enabling organizational accountability.  
 787

## 788 6.1 Theoretical Implications

789 **Decontextualization.** The inherent challenge in responsible AI toolkits lies in their attempt to  
 790 reconcile the tension between scalability and context specificity [100]. Traditional approaches  
 791 to toolkit development have often favored a universal, top-down approach that assumes a one-  
 792 size-fits-all solution [49, 60]. However, participatory development, such as the methodology we  
 793 followed in designing and populating a responsible AI tool with our guidelines, emphasizes the  
 794 importance of tailoring responsible AI guidelines to specific contexts and job roles needs. It is crucial  
 795 therefore to recognize that different AI practitioners, such as designers, developers, engineers, and  
 796 executives, have distinct requirements and considerations that cannot be treated as identical. This  
 797 highlights the complexity of developing toolkits that cater to a diverse range of practitioners while  
 798 accounting for their unique roles and settings—the problem of decontextualization in responsible  
 799 AI toolkits [100].  
 800

801 To tackle the problem of decontextualization, our proposed method incorporates two key ele-  
 802 ments: *guidelines usable by different roles* and *guiding questions*. Firstly, the integration of guidelines,  
 803 tailored to different roles and projects, provides practical steps and recommendations that technical  
 804 practitioners can easily implement, or C-level executives can make informed decisions upon. These  
 805 guidelines serve as a starting point for ethical decision-making throughout the AI lifecycle, con-  
 806 tributing to the vision of responsible AI by design (borrowing from the idea of ‘privacy by design’<sup>8</sup>).  
 807 Secondly, the inclusion of the two guiding questions (§5.1) enhances our toolkit’s ability to capture  
 808 the complexities of different social and organizational contexts. Expanding upon the concept that  
 809 guiding questions are an effective means of communication [98], as they help in gaining deeper  
 810 insights, clarifying responses, and uncovering underlying meanings, AI practitioners can engage  
 811 with these questions to explore the ethical considerations and challenges that are unique to their  
 812 deployment context.

813 **Meta-responsibility** Scholars have long recognized the need for a socio-technical approach that  
 814 considers the contextual factors governing the use of AI systems, including social, organizational,  
 815 and cultural factors [93]. In fact, Ackerman [1] introduced the concept of socio-technical gap to  
 816 highlight the disparity between human requirements in technology deployment contexts (socio-  
 817 requirements) and the technical solutions. This gap arises due to the flexible and nuanced nature of  
 818 human activity compared to the rigid and brittle nature of computational mechanisms, resulting  
 819 from necessary formalization and abstraction. Along these lines, Stahl [89] introduced the concept  
 820 of meta-responsibility to stress that AI systems should be viewed as systems of systems (ecosystems)  
 821 rather than single entities. To establish a regime of meta-responsibility, Stahl argued for an adaptive  
 822 governance structure to effectively respond to new insights and external influences (e.g., upcoming  
 823 AI regulation), and for a knowledge base that equips AI stakeholders with technical, ethical,  
 824 legal, and social understanding. By integrating ethical, legal, and social knowledge into the AI  
 825 development process—what Stahl referred to as adaptive governance structure, our work contributes  
 826 to this line of research by providing empirical evidence to it and pushing the theoretical boundaries  
 827 further.  
 828  
 829  
 830

831 <sup>8</sup>“Privacy by design” is a standard practice for incorporating data protection into the design of technology. In other words,  
 832 data protection is achieved when it is already integrated into the technology during its design and development [17].  
 833

## 6.2 Practical Implications

**Recommendations for incorporating responsible AI guidelines into toolkits.** Our work identified four essential design requirements for incorporating guidelines into tools. These requirements are simplifying guidelines into smaller visual components; implementing clear navigation through the guidelines; tracking and sharing progress how guidelines are applied; and developing mechanisms for reflection on the guidelines.

For simplifying guidelines, we displayed each guideline as a digital card and accompanied it with two guiding questions. Future work could explore how to further divide guidelines into additional visual elements on the cards and how to refine the guiding questions. For example, guideline #15—*ensuring compliance with agreements and legal requirements when handling data*—could be further divided into step-by-step processes, with each one marked by a visual element like a card tab or a link to a specific ISO, or excerpts from the EU AI Act. Regarding the guiding questions, we observed that formulating them is a delicate task, requiring a balance between directness and respect for the user’s autonomy. For example, a question formulated as “How did you consider the potential impact of your AI system on different user groups?” employs a proactive stance, avoiding any direct accusation or presumption of oversight. This method resonates with the experiences of our participants (e.g., P14), who found value in open-ended questions. However, guiding questions can be refined in various ways. For example, the *confront* type of question can incorporate the elements of “reminding consequences” and “providing multiple viewpoints,” encouraging users to consider alternative directions and diverse perspectives [16].

For ensuring clear navigation, we organized the guidelines into a one-page layout and incorporated multiple buttons along with a counter for easy navigation. Future work could explore how to develop alternative layouts and include different navigation mechanisms. For example, complementary guidelines with related content, such as guideline #3 *identify potential harms and risks associated with the intended use* and guideline #5 *develop strategies to mitigate identified harms or risks for each intended use* can be paired side by side to improve the quality of responses. Additionally, new navigation mechanisms might include a chart to illustrate the relationships between guidelines and a search bar to enable users to quickly locate specific guidelines.

For tracking how guidelines are applied and sharing progress among team members, we introduced a feature to store user responses locally within the browser session and dynamically generate a PDF report from these responses. Future work could explore how to structure user responses in formats suitable for automated analysis and integration with other tools. For instance, using JSON format as input for machine learning algorithms and Large Language Models (LLMs) can enable the analysis of user responses and the generation of automated insights and recommendations within the PDF report.

For enabling post-hoc reflections, we created a summary page where users can view the number of guidelines they have considered and their responses to each guideline. Future work could explore how to improve this summary page, for example by adding visual elements for recognizing responsible AI champions (e.g., responsible AI badges) and fostering empathy (e.g., animations presenting the environmental impact of an AI system) or by implementing a collaborative aspect where users can share and discuss their summary pages with peers or mentors.

**Recommendations for enabling organizational accountability.** While individual adoption of responsible AI best practices is crucial, fostering effective communication between technical and non-technical roles is equally important. Many existing responsible AI toolkits prioritize individual usage [100]. However, addressing complex ethical and societal challenges associated with AI systems requires diverse perspectives. Our interactive tool populated with guidelines addresses this need by offering features (e.g., adjusting number of guidelines on the fly to different

883 roles and system phases) that make the guidelines usable by different roles. However, our tool can  
884 further improve communication between roles by creating a knowledge base of responses. Such a  
885 knowledge base, according to Stahl [89], empowers team members to fulfill their responsibilities and  
886 supports distributed teams in constructing a shared understanding of their AI system. Furthermore,  
887 we suggest a mechanism for keeping this knowledge base up to date and enriched with diverse  
888 perspectives. This includes regularly revisiting the guidelines through our tool and providing  
889 responses at key project milestones, such as when the AI system enters a new phase or when  
890 the team communicates with external members. This approach ensures that the knowledge base  
891 remains dynamic and reflect the evolving insights and perspectives within the team.

892 Our guidelines and the tool that incorporates them can also be used to enable organizational  
893 accountability. Similar to Google’s five-stage internal algorithmic auditing framework [78], our  
894 guidelines serve as a practical tool for closing the AI accountability gap. The automatically generated  
895 report plays a crucial role in this process by providing a summary of the guidelines that were  
896 effectively implemented, those that should be considered for future development, and the non-  
897 applicable ones. These reports establish an additional chain of accountability that can be shared with  
898 stakeholders at various levels, including managers, senior leadership, and AI engineers. By offering  
899 more oversight and the ability to troubleshoot if needed, these reports help mitigate unintentional  
900 harm. However, it is important to note that when an organization adopts our guidelines, it should  
901 establish clear ethical guidelines for their intended uses. Our tool is not intended to discourage  
902 AI developers, engineers, designers, and managers from using it due to the fear of being held  
903 accountable for their responses. On the contrary, their responses, as documented in the report,  
904 provide an opportunity to identify potential ethical issues and address them early in the design  
905 stages. This proactive approach prevents the need for post-hoc fixes and repairs, aligning with the  
906 principle of addressing ethical considerations during the development process rather than as an  
907 afterthought [81]—the idea of *Responsible AI by Design*.

908  
909

### 6.3 Limitations and Future Work

910 Our work has four main limitations that highlight the need for future research efforts. Firstly,  
911 although we followed a rigorous four-step process involving multiple stakeholders, the list of  
912 22 guidelines may not be exhaustive. The rapidly evolving nature of AI ethics, governance, and  
913 regulations necessitates an ongoing effort to stay abreast of emerging developments. However, one  
914 of the strengths of our method lies in its modular design, which allows for ongoing refinement and  
915 expansion of the set of guidelines. This ensures that our guidelines maintain their relevance and  
916 stay up to date in the ever-evolving landscape of AI ethics, governance, regulations, and ISOs (e.g.,  
917 functional safety (ISO 5469), data quality (ISO 5259), explainability (ISO 6254), AI management  
918 system (ISO 42001)). However, we acknowledge that there may be limitations in ensuring that all  
919 standards are accessible to everyone and that experts may not always be available to evaluate them.  
920 A partial solution would be to create forums or discussion groups where individuals can share  
921 their experiences and insights about regulations and standards. At the same time, future research  
922 could also investigate the frequency with which our method should be updated as new literature  
923 emerges. One possibility would be to create an automated system that regularly collects research  
924 articles on responsible AI best practices, pairing them with current and upcoming regulations, to  
925 extract new guidelines.

926 Secondly, it is important to consider the qualitative nature of our user study, which involved  
927 in-depth interviews and analysis of participants’ responses. The findings from this study should be  
928 interpreted with caution though, understanding that the reported frequency of themes should be  
929 viewed in a comparative context rather than taken at face value [31]. This approach helps to avoid  
930 potential misinterpretation or overgeneralization of the results.

931

932 Thirdly, we need to acknowledge the limitations associated with the sample size and demo-  
933 graphics of our user study. The study was conducted with a specific group of participants, and  
934 therefore, the findings may not fully represent the practices and perspectives of all AI practi-  
935 tioners. Our sample predominantly consisted of male participants, which aligns with the gender  
936 distribution reported in Stack Overflow’s 2022 Developer Survey, where 92.85% of professional  
937 developer respondents identified as male [76]. Additionally, our participants were drawn from a  
938 large research-focused technology company. While the results may offer insights into practices  
939 within certain companies, they serve as a case study for future research.

940 Lastly, our qualitative data suggests indicators of ease of use for AI practitioners but does not  
941 provide direct information on the actual effectiveness of the guidelines. Understanding the impact  
942 of guidelines (or other AI toolkits [100]) requires long-term studies that consider multiple projects,  
943 with some utilizing the toolkit and others not. One potential avenue, as suggested by clinical  
944 researchers developing deep learning tools for patient care [9], is to conduct observational studies  
945 with users of the AI system to assess its performance. Another approach is to use proxies, such as  
946 measuring users’ attitudes, beliefs, and mindset regarding ethical values before and after utilizing  
947 the guidelines. We intend to explore these directions in future research.

## 948 7 CONCLUSION

949 We proposed a method for generating a list of responsible AI guidelines that are grounded in  
950 regulations and are usable by different roles. The resulting 22 guidelines were integrated into an  
951 interactive tool and evaluated through a user study with 14 AI researchers, engineers, designers,  
952 and managers from a large technology company. Our participants found the guidelines well-aligned  
953 with their roles, enabling them to communicate complex ethical concepts in a structured manner.  
954 The guidelines are also grounded in ISOs and the EU AI Act articles, receiving positive feedback  
955 for being comprehensive. The usefulness of examples in guidelines was particularly noted as they  
956 enabled participants to reflect on their choices concerning ethical issues. As these guidelines are  
957 likely to become part of future responsible AI toolkits, it is important to implement features that  
958 provide users with time and space for reflection. Additionally, these toolkits should take users’  
959 reflections and roles into account to offer actionable recommendations tailored to a specific project,  
960 using, for example, large language models.

## 961 REFERENCES

- 962 [1] Mark S. Ackerman. 2000. The Intellectual Challenge of CSCW: The Gap Between Social Requirements and Technical  
963 Feasibility. *Human-Computer Interaction* 15, 2-3 (2000), 179–203. [https://doi.org/10.1207/S15327051HCI1523\\_5](https://doi.org/10.1207/S15327051HCI1523_5)
- 964 [2] EU AI Act. 2021. *Proposal for a Regulation of the European Parliament and of the Council laying down harmonised rules  
965 on Artificial Intelligence (Artificial Intelligence Act) and amending certain Union legislative acts*. Retrieved December  
966 2023 from <https://artificialintelligenceact.eu/the-act/>
- 967 [3] K.K. Aggarwal and Yogesh Singh. 2008. *Software engineering*. New Age International.
- 968 [4] Saleema Amershi, Dan Weld, Mihaela Vorvoreanu, Adam Fourney, Besmira Nushi, Penny Collisson, Jina Suh, Shamsi  
969 Iqbal, Paul N. Bennett, Kori Inkpen, Jaime Teevan, Ruth Kikin-Gil, and Eric Horvitz. 2019. Guidelines for Human-AI  
970 Interaction. In *Proceedings of the ACM CHI Conference on Human Factors in Computing Systems (CHI)*. ACM, 1–13.  
971 <https://doi.org/10.1145/3290605.3300233>
- 972 [5] Alejandro B. Arrieta, Natalia Díaz-Rodríguez, Javier Del Ser, Adrien Bannetot, Siham Tabik, Alberto Barbado, Salvador  
973 García, Sergio Gil-López, Daniel Molina, Richard Benjamins, et al. 2020. Explainable Artificial Intelligence (XAI):  
974 Concepts, taxonomies, opportunities and challenges toward responsible AI. *Information Fusion* 58 (2020), 82–115.  
975 <https://doi.org/10.1016/j.inffus.2019.12.012>
- 976 [6] Vijay Arya, Rachel K.E. Bellamy, Pin-Yu Chen, Amit Dhurandhar, Michael Hind, Samuel C. Hoffman, Stephanie  
977 Houde, Q. Vera Liao, Ronny Luss, Aleksandra Mojsilović, et al. 2019. One Explanation Does Not Fit All: A Toolkit  
978 And Taxonomy Of AI Explainability Techniques. arXiv:1909.03012
- 979 [7] Ricardo Baeza-Yates. 2018. Bias on the Web. *Commun. ACM* 61, 6 (2018), 54–61. <https://doi.org/10.1145/3209581>

- 981 [8] Agathe Balayn, Mireia Yurrita, Jie Yang, and Ujwal Gadiraju. 2023. “Fairness Toolkits, A Checkbox Culture?” On  
 982 the Factors that Fragment Developer Practices in Handling Algorithmic Harms. In *Proceedings of the AAAI/ACM*  
 983 *Conference on AI, Ethics, and Society (AIES)*. 482–495. <https://doi.org/10.1145/3600211.3604674>
- 984 [9] Emma Beede, Elizabeth Baylor, Fred Hersch, Anna Iurchenko, Lauren Wilcox, Paisan Ruamviboonsuk, and Laura M.  
 985 Vardoulakis. 2020. A Human-Centered Evaluation of a Deep Learning System Deployed in Clinics for the Detection  
 986 of Diabetic Retinopathy. In *Proceedings of the ACM CHI Conference on Human Factors in Computing Systems (CHI)*.  
 987 ACM, 1–12. <https://doi.org/10.1145/3313831.3376718>
- 988 [10] Emily M. Bender and Batya Friedman. 2018. Data Statements for Natural Language Processing: Toward Mitigating  
 989 System Bias and Enabling Better Science. *Transactions of the Association for Computational Linguistics* 6 (2018),  
 990 587–604. [https://doi.org/10.1162/tacl\\_a\\_00041](https://doi.org/10.1162/tacl_a_00041)
- 991 [11] Sarah Bird, Miro Dudík, Richard Edgar, Brandon Horn, Roman Lutz, Vanessa Milan, Mehrnoosh Sameki, Hanna  
 992 Wallach, and Kathleen Walker. 2020. *Fairlearn: A toolkit for assessing and improving fairness in AI*. Technical  
 993 Report. Microsoft. Retrieved December 2022 from [https://www.microsoft.com/en-us/research/publication/fairlearn-](https://www.microsoft.com/en-us/research/publication/fairlearn-a-toolkit-for-assessing-and-improving-fairness-in-ai/)  
 994 [a-toolkit-for-assessing-and-improving-fairness-in-ai/](https://www.microsoft.com/en-us/research/publication/fairlearn-a-toolkit-for-assessing-and-improving-fairness-in-ai/)
- 995 [12] Michael Boone, Nikki Pope, Chaowei Xiao, and Anima Anandkumar. 2022. *Enhancing AI Transparency and Ethical*  
 996 *Considerations with Model Card++*. Nvidia. Retrieved November 2022 from [https://developer.nvidia.com/blog/](https://developer.nvidia.com/blog/enhancing-ai-transparency-and-ethical-considerations-with-model-card/)  
 997 [enhancing-ai-transparency-and-ethical-considerations-with-model-card/](https://developer.nvidia.com/blog/enhancing-ai-transparency-and-ethical-considerations-with-model-card/)
- 998 [13] Virginia Braun and Victoria Clarke. 2006. Using thematic analysis in psychology. *Qualitative Research in Psychology*  
 999 3, 2 (2006), 77–101. <https://doi.org/10.1191/1478088706qp063oa>
- 1000 [14] John Brooke. 1996. SUS: A “Quick and Dirty” Usability Scale. In *Usability Evaluation In Industry*, Patrick W. Jordan,  
 1001 B. Thomas, Ian L. McClelland, and Bernard Weerdmeester (Eds.). CRC Press, Chapter 12, 107–114.
- 1002 [15] Joy Buolamwini and Timnit Gebru. 2018. Gender Shades: Intersectional Accuracy Disparities in Commercial Gender  
 1003 Classification. In *Proceedings of the ACM Conference on Fairness, Accountability, and Transparency (FAccT)*. Proceedings  
 1004 of Machine Learning Research, 77–91. <https://proceedings.mlr.press/v81/buolamwini18a.html>
- 1005 [16] Ana Caraban, Evangelos Karapanos, Daniel Gonçalves, and Pedro Campos. 2019. 23 Ways to Nudge: A Review  
 1006 of Technology-Mediated Nudging in Human-Computer Interaction. In *Proceedings of the ACM CHI Conference on*  
 1007 *Human Factors in Computing Systems (CHI)*. ACM, 1–15. <https://doi.org/10.1145/3290605.3300733>
- 1008 [17] Ann Cavoukian. 2009. *Privacy by Design: The 7 Foundational Principles*. Information & Privacy Commissioner of  
 1009 Ontario, Canada. Retrieved November 2022 from [https://iab.org/wp-content/IAB-uploads/2011/03/fred\\_carter.pdf](https://iab.org/wp-content/IAB-uploads/2011/03/fred_carter.pdf)
- 1010 [18] Gary S. Collins, Johannes B. Reitsma, Douglas G. Altman, and Karel G.M. Moons. 2015. Transparent reporting of a  
 1011 multivariable prediction model for individual prognosis or diagnosis (TRIPOD): the TRIPOD statement. *Journal of*  
 1012 *British Surgery* 102, 3 (2015), 148–158. <https://doi.org/10.1161/CIRCULATIONAHA.114.014508>
- 1013 [19] Henriette Cramer, Jean Garcia-Gathright, Sravana Reddy, Aaron Springer, and Romain Takeo Bouyer. 2019. Transla-  
 1014 tion, tracks & data: an algorithmic bias effort in practice. In *Extended Abstracts of the ACM CHI Conference on Human*  
 1015 *Factors in Computing Systems (CHI)*. 1–8. <https://doi.org/10.1145/3290607.3299057>
- 1016 [20] Wesley H. Deng, Manish Nagireddy, Michelle Seng Ah Lee, Jatinder Singh, Zhiwei S. Wu, Kenneth Holstein, and Haiyi  
 1017 Zhu. 2022. Exploring how machine learning practitioners (try to) use fairness toolkits. In *Proceedings of the ACM*  
 1018 *Conference on Fairness, Accountability, and Transparency (FAccT)*. 473–484. <https://doi.org/10.1145/3531146.3533113>
- 1019 [21] Wesley H. Deng, Nur Yildirim, Monica Chang, Motahhare Eslami, Kenneth Holstein, and Michael Madaio. 2023.  
 1020 Investigating Practices and Opportunities for Cross-functional Collaboration around AI Fairness in Industry Practice.  
 1021 In *Proceedings of the ACM Conference on Fairness, Accountability, and Transparency (FAccT)*. 705–716. <https://doi.org/10.1145/3593013.3594037>
- 1022 [22] Martin Dittus, Luca M. Aiello, and Daniele Quercia. 2017. Community Engagement Triage: Lightweight Prompts  
 1023 for Systematic Reviews. *Proceedings of the ACM on Human-Computer Interaction* 1, CSCW (2017), 1–22. <https://doi.org/10.1145/3134674>
- 1024 [23] Lucas Dixon, John Li, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. 2018. Measuring and mitigating  
 1025 unintended bias in text classification. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society (AIES)*.  
 1026 67–73. <https://doi.org/10.1145/3278721.3278729>
- 1027 [24] Cynthia Dwork. 2008. Differential Privacy: A Survey of Results. In *Theory and Applications of Models of Computation*.  
 1028 Springer, 1–19. [https://doi.org/10.1007/978-3-540-79228-4\\_1](https://doi.org/10.1007/978-3-540-79228-4_1)
- 1029 [25] Upol Ehsan and Mark O. Riedl. 2020. Human-Centered Explainable AI: Towards a Reflective Sociotechnical Approach.  
 In *HCI International 2020 - Late Breaking Papers: Multimodality and Intelligence*, Constantine Stephanidis, Masaaki  
 Kurosu, Helmut Degen, and Lauren Reinerman-Jones (Eds.). Springer International Publishing, 449–466.
- [26] Salma Elsayed-Ali, Sara E. Berger, Vagner F. De Santana, and Juana Catalina Becerra Sandoval. 2023. Responsible  
& Inclusive Cards: An Online Card Tool to Promote Critical Reflection in Technology Industry Work Practices. In  
*Proceedings of the ACM CHI Conference on Human Factors in Computing Systems (CHI)*. 1–14. <https://doi.org/10.1145/3544548.3580771>

- 1030 [27] Equal Employment Opportunity Commission. 1977. *Prohibited Employment Policies/Practices*. Retrieved May 2023  
 1031 from <https://www.eeoc.gov/prohibited-employment-policiespractices>
- 1032 [28] European Union. 2018. *General Data Protection Regulation*. Retrieved May 2023 from <https://gdpr-info.eu/>
- 1033 [29] Fairlearn. 2022. *Improve fairness of AI systems*. Retrieved November 2022 from <https://fairlearn.org>
- 1034 [30] Jessica Fjeld, Nele Achten, Hannah Hilligoss, Adam Nagy, and Madhulika Srikumar. 2020. Principled Artificial  
 1035 Intelligence: Mapping Consensus in Ethical and Rights-Based Approaches to Principles for AI. *Berkman Klein Center  
 1036 Research Publication* 2020-1 (2020). <https://doi.org/10.2139/ssrn.3518482>
- 1037 [31] Ellie Fossey, Carol Harvey, Fiona McDermott, and Larry Davidson. 2002. Understanding and Evaluating Qualitative  
 1038 Research. *Australian & New Zealand Journal of Psychiatry* 36, 6 (2002), 717–732. <https://doi.org/10.1046/j.1440-1614.2002.01100.x>
- 1039 [32] Ulrike Franke. 2021. *Artificial divide: How Europe and America could clash over AI*. European Council on Foreign  
 1040 Relations. Retrieved May 2023 from <https://ecfr.eu/publication/artificial-divide-how-europe-and-america-could-clash-over-ai/>
- 1041 [33] Hana Frluckaj, Laura Dabbish, David G. Widder, Huilian Sophie Qiu, and James Herbsleb. 2022. Gender and  
 1042 Participation in Open Source Software Development. *Proceedings of the ACM on Human-Computer Interaction* 6,  
 1043 CSCW2 (2022), 1–31. <https://doi.org/10.1145/3555190>
- 1044 [34] Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé III, and  
 1045 Kate Crawford. 2021. Datasheets for Datasets. *Commun. ACM* 64, 12 (2021), 86–92. <https://doi.org/10.1145/3458723>
- 1046 [35] Brent Gleeson. 2013. *The Silo Mentality: How To Break Down The Barriers*. Retrieved December 2023 from <https://www.forbes.com/sites/brentgleeson/2013/10/02/the-silo-mentality-how-to-break-down-the-barriers/>
- 1047 [36] Google. 2022. *AI Explorables*. Retrieved November 2022 from <https://pair.withgoogle.com/explorables/>
- 1048 [37] Google. 2022. *Fairness Indicators*. Retrieved November 2022 from <https://github.com/tensorflow/fairness-indicators>
- 1049 [38] Government Equalities Office and Equality and Human Rights Commission. 2010. *Equality Act 2010: guidance*.  
 1050 Retrieved May 2023 from <https://www.gov.uk/guidance/equality-act-2010-guidance>
- 1051 [39] Greg Guest, Arwen Bunce, and Laura Johnson. 2006. How many interviews are enough? An experiment with data  
 1052 saturation and variability. *Field Methods* 18, 1 (2006), 59–82. <https://doi.org/10.1177/1525822X05279903>
- 1053 [40] David Gunning, Mark Stefik, Jaesik Choi, Timothy Miller, Simone Stumpf, and Guang-Zhong Yang. 2019. XAI—  
 1054 Explainable artificial intelligence. *Science Robotics* 4, 37 (2019). <https://doi.org/10.1126/scirobotics.aay7120>
- 1055 [41] Karen Hao. 2019. Training a single AI model can emit as much carbon as five cars in their lifetimes. *MIT technology  
 1056 Review* (2019). <https://www.technologyreview.com/2019/06/06/239031/training-a-single-ai-model-can-emit-as-much-carbon-as-five-cars-in-their-lifetimes/>
- 1057 [42] Lucy Havens, Melissa Terras, Benjamin Bach, and Beatrice Alex. 2020. Situated Data, Situated Systems: A Methodology  
 1058 to Engage with Power Relations in Natural Language Processing Research. In *Proceedings of the Second Workshop on  
 1059 Gender Bias in Natural Language Processing*. Association for Computational Linguistics, 107–124. <https://aclanthology.org/2020.gebnlp-1.10>
- 1060 [43] Sarah Holland, Ahmed Hosny, Sarah Newman, Joshua Joseph, and Kasia Chmielinski. 2020. The Dataset Nutrition  
 1061 Label: A Framework to Drive Higher Data Quality Standards. In *Data Protection and Privacy*, Dara Hallinan,  
 1062 Ronald Leenes, Serge Gutwirth, and Paul De Hert (Eds.). Hart Publishing, Chapter 1, 1–26. <https://doi.org/10.5040/9781509932771.ch-001>
- 1063 [44] The White House. 2023. *Blueprint for an AI Bill of Rights*. Retrieved May 2023 from <https://www.whitehouse.gov/ostp/ai-bill-of-rights/>
- 1064 [45] IBM. 2022. *AI Fairness 360*. Retrieved November 2022 from <https://aif360.mybluemix.net>
- 1065 [46] IDEO. 2019. *AI needs ethical compass. This tool can help*. Retrieved November 2022 from <https://www.ideo.com/blog/ai-needs-an-ethical-compass-this-tool-can-help>
- 1066 [47] Anna Jobin, Marcello Ienca, and Effy Vayena. 2019. The global landscape of AI ethics guidelines. *Nature Machine  
 1067 Intelligence* 1, 9 (2019), 389–399. <https://doi.org/10.1038/s42256-019-0088-2>
- 1068 [48] Harmanpreet Kaur, Harsha Nori, Samuel Jenkins, Rich Caruana, Hanna Wallach, and Jennifer Wortman Vaughan.  
 1069 2020. Interpreting Interpretability: Understanding Data Scientists' Use of Interpretability Tools for Machine Learning.  
 1070 In *Proceedings of the ACM CHI Conference on Human Factors in Computing Systems (CHI)*. 1–14. <https://doi.org/10.1145/3313831.3376219>
- 1071 [49] Christopher M. Kelty. 2018. *The Participatory Development Toolkit*. Retrieved December 2023 from <https://limn.it/articles/the-participatory-development-toolkit/>
- 1072 [50] Henry Kissinger, Eric Schmidt, and Daniel P. Huttenlocher. 2021. *The age of AI: And our human future*. John Murray  
 1073 London.
- 1074 [51] Knowledge Centre Data and Society. 2019. AI Blindspots Card Set 1.0. Retrieved February 2023 from <https://data-en-maatschappij.ai/en/tools/ai-blindspot>

- 1079 [52] Knowledge Centre Data and Society. 2019. AI Blindspots Card Set 2.0. Retrieved February 2023 from [https://data-](https://data-en-maatschappij.ai/en/tools/ai-blindspots-2.0)  
1080 [en-maatschappij.ai/en/tools/ai-blindspots-2.0](https://data-en-maatschappij.ai/en/tools/ai-blindspots-2.0)
- 1081 [53] Todd Kulesza, Margaret Burnett, Weng-Keen Wong, and Simone Stumpf. 2015. Principles of Explanatory Debugging  
1082 to Personalize Interactive Machine Learning. In *Proceedings of the ACM Conference on Intelligent User Interfaces (IUI)*.  
126–137. <https://doi.org/10.1145/2678025.2701399>
- 1083 [54] Superrrr Lab. 2022. *The Feminist Tech Card Deck*. Retrieved November 2022 from <https://superrrr.net/feministtech/deck>
- 1084 [55] James R. Lewis, Brian S. Utesch, and Deborah E. Maher. 2015. Measuring Perceived Usability: The SUS, UMUX-LITE,  
1085 and AltUsability. *International Journal of Human-Computer Interaction* 31, 8 (2015), 496–505. [https://doi.org/10.1080/](https://doi.org/10.1080/10447318.2015.1064654)  
1086 [10447318.2015.1064654](https://doi.org/10.1080/10447318.2015.1064654)
- 1087 [56] Q. Vera Liao and Kush R. Varshney. 2021. Human-Centered Explainable AI (XAI): From Algorithms to User Experiences.  
arXiv:2110.10790
- 1088 [57] Qinghua Lu, Liming Zhu, Xiwei Xu, Jon Whittle, Didar Zowghi, and Aurelie Jacquet. 2023. Responsible AI Pattern  
1089 Catalogue: A Collection of Best Practices for AI Governance and Engineering. *Comput. Surveys* (2023). <https://doi.org/10.1145/3626234>
- 1090 [58] Scott M. Lundberg and Su-In Lee. 2017. A Unified Approach to Interpreting Model Predictions. In *Proceedings of the*  
1091 *Conference on Neural Information Processing Systems (NeurIPS)*. Curran Associates Inc., 4768–4777. Retrieved August  
1092 2022 from [https://proceedings.neurips.cc/paper\\_files/paper/2017/file/8a20a8621978632d76c43dfd28b67767-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/8a20a8621978632d76c43dfd28b67767-Paper.pdf)
- 1093 [59] Michael A. Madaio, Luke Stark, Jennifer Wortman Vaughan, and Hanna Wallach. 2020. Co-designing checklists  
1094 to understand organizational challenges and opportunities around fairness in AI. In *Proceedings of the ACM CHI*  
1095 *Conference on Human Factors in Computing Systems (CHI)*. 1–14. <https://doi.org/10.1145/3313831.3376445>
- 1096 [60] Shannon Mattern. 2021. *Unboxing the Toolkit*. Retrieved December 2023 from [https://tool-shed.org/unboxing-the-](https://tool-shed.org/unboxing-the-toolkit/)  
1097 [toolkit/](https://tool-shed.org/unboxing-the-toolkit/)
- 1098 [61] Nora McDonald, Sarita Schoenebeck, and Andrea Forte. 2019. Reliability and Inter-Rater Reliability in Qualitative  
1099 Research: Norms and Guidelines for CSCW and HCI Practice. *Proceedings of the ACM on Human-Computer Interaction*  
3, CSCW (2019). <https://doi.org/10.1145/3359174>
- 1100 [62] Susan Michie, Charles Abraham, Martin P. Eccles, Jill J. Francis, Wendy Hardeman, and Marie Johnston. 2011.  
1101 Strengthening evaluation and implementation by specifying components of behaviour change interventions: a study  
1102 protocol. *Implementation Science* 6, 1 (2011), 1–8. <https://doi.org/10.1186/1748-5908-6-10>
- 1103 [63] Susan Michie, Michelle Richardson, Marie Johnston, Charles Abraham, Jill J. Francis, Wendy Hardeman, Martin P.  
1104 Eccles, James Cane, and Caroline Wood. 2013. The behavior change technique taxonomy (v1) of 93 hierarchically  
1105 clustered techniques. *Annals of Behavioral Medicine* 46, 1 (2013). <https://doi.org/10.1007/s12160-013-9486-6>
- 1106 [64] Matthew Miles and Michael Huberman. 1994. *Qualitative Data Analysis: A Methods Sourcebook*. Sage.
- 1107 [65] Miro. 2022. *Miro | Online Whiteboard for Visual Collaboration*. Retrieved August 2022 from <https://miro.com/>
- 1108 [66] Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer,  
1109 Ioliuwa D. Raji, and Timnit Gebru. 2019. Model Cards for Model Reporting. In *Proceedings of the ACM Conference on*  
1110 *Fairness, Accountability, and Transparency (FAccT)*. ACM, 220–229. <https://doi.org/10.1145/3287560.3287596>
- 1111 [67] Shira Mitchell, Eric Potash, Solon Barocas, Alexander D’Amour, and Kristian Lum. 2018. Prediction-based decisions  
1112 and fairness: A catalogue of choices, assumptions, and definitions. arXiv:1811.07867
- 1113 [68] Brent D. Mittelstadt, Patrick Allo, Mariarosaria Taddeo, Sandra Wachter, and Luciano Floridi. 2016. The ethics of  
1114 algorithms: Mapping the debate. *Big Data & Society* 3, 2 (2016). <https://doi.org/10.1177/2053951716679679>
- 1115 [69] Interpret ML. 2019. *Interpret ML*. Retrieved August 2022 from <https://interpret.ml/>
- 1116 [70] Aleksandra Mojsilovic. 2019. *Introducing AI Explainability 360*. IBM. Retrieved November 2022 from <https://www.ibm.com/blogs/research/2019/08/ai-explainability-360/>
- 1117 [71] Nadia Nahar, Shurui Zhou, Grace Lewis, and Christian Kästner. 2022. Collaboration challenges in building ML-enabled  
1118 systems: communication, documentation, engineering, and process. In *Proceedings of the ACM/IEEE International*  
1119 *Conference on Software Engineering (ICSE)*. 413–425. <https://doi.org/10.1145/3510003.3510209>
- 1120 [72] Arvind Narayanan. 2018. 21 Fairness definitions and their politics. In *Tutorial presented at the ACM Conference on*  
1121 *Fairness, Accountability, and Transparency (FAccT)*.
- 1122 [73] National Institute of Standards and Technology. 2023. *AI Risk Management Framework*. Retrieved May 2023 from  
1123 <https://www.nist.gov/itl/ai-risk-management-framework>
- 1124 [74] OECD. 2023. *Catalogue of Tools & Metrics for Trustworthy AI*. Retrieved December 2023 from [https://oecd.ai/en/](https://oecd.ai/en/catalogue/tools)  
1125 [catalogue/tools](https://oecd.ai/en/catalogue/tools)
- 1126 [75] Equality of opportunity in supervised learning. 2016. Hardt, Moritz and Price, Eric and Srebro, Nati. In *Proceedings of the*  
1127 *Conference on Neural Information Processing Systems (NeurIPS)*. Curran Associates Inc., 3323–3331. Retrieved August  
2022 from [https://proceedings.neurips.cc/paper\\_files/paper/2016/file/9d2682367c3935defcb1f9e247a97c0d-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2016/file/9d2682367c3935defcb1f9e247a97c0d-Paper.pdf)
- [76] Stack Overflow. 2022. *Stack Overflow Developer Survey 2022*. Retrieved November 2022 from [https://survey.](https://survey.stackoverflow.co/2022/)  
[stackoverflow.co/2022/](https://survey.stackoverflow.co/2022/)

- 1128 [77] David Piorkowski, Soya Park, April Yi Wang, Dakuo Wang, Michael Muller, and Felix Portnoy. 2021. How AI  
 1129 developers overcome communication challenges in a multidisciplinary team: A case study. *Proceedings of the ACM on*  
 1130 *Human-Computer Interaction* 5, CSCW1 (2021), 1–25. <https://doi.org/10.1145/3449205>
- 1131 [78] Inioluwa D. Raji, Andrew Smart, Rebecca N. White, Margaret Mitchell, Timnit Gebru, Ben Hutchinson, Jamila  
 1132 Smith-Loud, Daniel Theron, and Parker Barnes. 2020. Closing the AI accountability gap: Defining an end-to-end  
 1133 *Transparency (FAccT)*. 33–44.
- 1134 [79] Bogdana Rakova, Jingying Yang, Henriette Cramer, and Rumman Chowdhury. 2021. Where responsible AI meets  
 1135 reality: Practitioner perspectives on enablers for shifting organizational practices. *Proceedings of the ACM on*  
 1136 *Human-Computer Interaction* 5 (2021), 1–23. <https://doi.org/10.1145/3449081>
- 1137 [80] Johnny Saldaña. 2015. *The Coding Manual for Qualitative Researchers*. Sage.
- 1138 [81] Nithya Sambasivan and Jess Holbrook. 2018. Toward responsible AI for the next billion users. *Interactions* 26, 1  
 1139 (2018), 68–71. <https://doi.org/10.1145/3298735>
- 1140 [82] Conrad Sanderson, David Douglas, Qinghua Lu, Emma Schleiger, Jon Whittle, Justine Lacey, Glenn Newnham, Stefan  
 1141 Hajkowicz, Cathy Robinson, and David Hansen. 2023. AI ethics principles in practice: Perspectives of designers and  
 1142 developers. *IEEE Transactions on Technology and Society* (2023), 171–187. <https://doi.org/10.1109/TTS.2023.3257303>
- 1143 [83] Jeff Sauro. 2011. *A practical guide to the system usability scale: Background, benchmarks & best practices*. Measuring  
 1144 Usability LLC.
- 1145 [84] Peter Schaar. 2010. Privacy by Design. *Identity in the Information Society* 3, 2 (April 2010), 267–274. <https://doi.org/10.1007/s12394-010-0055-x>
- 1146 [85] Andrew D. Selbst, Danah Boyd, Sorelle A. Friedler, Suresh Venkatasubramanian, and Janet Vertesi. 2019. Fairness  
 1147 and Abstraction in Sociotechnical Systems. In *Proceedings of the ACM Conference on Fairness, Accountability, and*  
 1148 *Transparency (FAccT)*. ACM, 59–68. <https://doi.org/10.1145/3287560.3287598>
- 1149 [86] Or Sharir, Barak Peleg, and Yoav Shoham. 2020. The Cost of Training NLP Models: A Concise Overview.  
 1150 arXiv:2004.08900
- 1151 [87] Ben Shneiderman. 2021. Responsible AI: Bridging from ethics to practice. *Commun. ACM* 64, 8 (2021), 32–35.  
 1152 <https://doi.org/10.1145/3445973>
- 1153 [88] Ben Shneiderman. 2022. *Human-centered AI*. Oxford University Press.
- 1154 [89] Bernd C. Stahl. 2023. Embedding responsibility in intelligent systems: from AI ethics to responsible AI ecosystems.  
 1155 *Scientific Reports* 13, 1 (2023), 7586. <https://doi.org/10.1038/s41598-023-34622-w>
- 1156 [90] Emma Strubell, Ananya Ganesh, and Andrew McCallum. 2020. Energy and Policy Considerations for Modern  
 1157 Deep Learning Research. *Proceedings of the AAAI Conference on Artificial Intelligence* 34, 09 (2020), 13693–13696.  
 1158 <https://doi.org/10.1609/aaai.v34i09.7123>
- 1159 [91] Hariharan Subramonyam, Jane Im, Colleen Seifert, and Eytan Adar. 2022. Solving separation-of-concerns problems  
 1160 in collaborative design of human-AI systems through leaky abstractions. In *Proceedings of the ACM Conference on*  
 1161 *Human Factors in Computing Systems (CHI)*. 1–21. <https://doi.org/10.1145/3491102.3517537>
- 1162 [92] Mohammad Tahaei, Marios Constantinides, Daniele Quercia, Sean Kennedy, Michael Muller, Simone Stumpf, Q. Vera  
 1163 Liao, Ricardo Baeza-Yates, Lora Aroyo, Jess Holbrook, et al. 2023. Human-Centered Responsible Artificial Intelligence:  
 1164 Current & Future Trends. In *Extended Abstracts of the ACM Conference on Human Factors in Computing Systems (CHI)*.  
 1165 1–4. <https://doi.org/10.1145/3544549.3583178>
- 1166 [93] Mohammad Tahaei, Marios Constantinides, Daniele Quercia, and Michael Muller. 2023. A Systematic Literature  
 1167 Review of Human-Centered, Ethical, and Responsible AI. arXiv:2302.05284
- 1168 [94] Ville Vakkuri, Kai-Kristian Kemell, Marianna Jantunen, Erika Halme, and Pekka Abrahamsson. 2021. ECCOLA  
 1169 – A method for implementing ethically aligned AI systems. *Journal of Systems and Software* 182 (2021), 111067.  
 1170 <https://doi.org/10.1016/j.jss.2021.111067>
- 1171 [95] Sahil Verma and Julia Rubin. 2018. Fairness definitions explained. In *Proceedings of the IEEE/ACM International*  
 1172 *Workshop on Software Fairness (FairWare)*. Association for Computing Machinery, 1–7.
- 1173 [96] Angelina Wang, Alexander Liu, Ryan Zhang, Anat Kleiman, Leslie Kim, Dora Zhao, Iroha Shirai, Arvind Narayanan,  
 1174 and Olga Russakovsky. 2022. REVISE: A Tool for Measuring and Mitigating Bias in Visual Datasets. *International*  
 1175 *Journal of Computer Vision* 130, 7, 1790–1810. <https://doi.org/10.1007/s11263-022-01625-5>
- 1176 [97] Qiaosi Wang, Michael Madaio, Shaun Kane, Shivani Kapania, Michael Terry, and Lauren Wilcox. 2023. Designing  
 1177 Responsible AI: Adaptations of UX Practice to Meet Responsible AI Challenges. In *Proceedings of the ACM Conference*  
 1178 *on Human Factors in Computing Systems (CHI)*. 1–16. <https://doi.org/10.1145/3544548.3581278>
- 1179 [98] Harry Weger Jr, Gina Castle Bell, Elizabeth M. Minei, and Melissa C. Robinson. 2014. The relative effectiveness of  
 1180 active listening in initial interactions. *International Journal of Listening* 28, 1 (2014), 13–31. <https://doi.org/10.1080/10904018.2013.813234>



- 1177 [99] Stefan Werning. 2020. Making data playable: a game co-creation method to promote creative data literacy. *Journal of*  
1178 *Media Literacy Education* 12, 3 (2020), 88–101. <https://doi.org/10.23860/JMLE-2020-12-3-8>
- 1179 [100] Richmond Y. Wong, Michael A. Madaio, and Nick Merrill. 2023. Seeing like a toolkit: How toolkits envision  
1180 the work of AI ethics. *Proceedings of the ACM on Human-Computer Interaction* 7, CSCW1 (2023), 1–27. <https://doi.org/10.1145/3579621>
- 1181 [101] Nur Yildirim, Mahima Pushkarna, Nitesh Goyal, Martin Wattenberg, and Fernanda Viégas. 2023. Investigating How  
1182 Practitioners Use Human-AI Guidelines: A Case Study on the People+ AI Guidebook. In *Proceedings of the ACM*  
1183 *Conference on Human Factors in Computing Systems (CHI)*. 1–13. <https://doi.org/10.1145/3544548.3580900>
- 1184 [102] Amy X. Zhang, Michael Muller, and Dakuo Wang. 2020. How do data science workers collaborate? roles, workflows,  
1185 and tools. *Proceedings of the ACM on Human-Computer Interaction* 4, CSCW1 (2020), 1–23. <https://doi.org/10.1145/3392826>
- 1186
- 1187
- 1188
- 1189
- 1190
- 1191
- 1192
- 1193
- 1194
- 1195
- 1196
- 1197
- 1198
- 1199
- 1200
- 1201
- 1202
- 1203
- 1204
- 1205
- 1206
- 1207
- 1208
- 1209
- 1210
- 1211
- 1212
- 1213
- 1214
- 1215
- 1216
- 1217
- 1218
- 1219
- 1220
- 1221
- 1222
- 1223
- 1224
- 1225

**A ADDITIONAL MATERIALS FOR THE USER STUDY**

- How old are you?
- What is your gender? [Male, Female, Non-binary, Prefer not to say, Open-ended option]
- How many years of experience do you have in AI systems?
- What’s your educational background?
- In which country do you currently reside?
- What is domain or sector of your work? (e.g., health, energy, education, finance, technology, food)
- What is your current role?
- What kinds of AI systems do you work on? (e.g., machine learning, computer vision, NLP, game theory, robotics)

Table 3. Constructed themes for the user study based on how our participants saw the application of guidelines, what worked well and what could have been improved.

Theme	Participants
Raising awareness, facilitating self-learning	12
Aligning with roles	10
Aligning with regulations	10
Providing helpful examples	7
Engaging team members and external experts	5
Maintaining the visual simplicity of the guidelines	3
Documenting guidelines in a concise summary PDF	3
Providing a systematic flow of information and guidelines	2

**B MAPPING GUIDELINES WITH EU AI ACT ARTICLES**

1275 **Article 6 (Classification rules for high-risk AI systems):** It states that an AI system shall be  
1276 considered high-risk when “it [the AI system] is intended to be used as a safety component of a  
1277 product, or is itself a product”. This article aligns with **guideline #1** as it mandates the identification  
1278 of an AI system’s intended use to determine whether its use poses a low or high risk.  
1279

1280 **Article 9 (Risk management system):** It states that “a risk management system shall be established,  
1281 implemented, documented and maintained throughout the entire lifecycle of a high-risk AI system”.  
1282 This article aligns with **guidelines #1, #3-5, and #13** as it is about the identification of harms and  
1283 risks of the AI system’s intended use.  
1284

1285 **Article 10 (Data and data governance):** It states that “training, validation and testing data  
1286 sets shall be subject to appropriate data governance and management practices”. This article aligns  
1287 with **guidelines #8 and #15-18** as it discusses the management and quality of data for training,  
1288 validation, and testing, including aspects of diversity and minimizing biases.

1289 **Article 11 (Technical documentation):** It states that the technical documentation of a high-risk  
1290 AI system shall “be drawn up before that system is placed on the market or put into service and shall be  
1291 kept up-to date”, and “provide national competent authorities and notified bodies with all the necessary  
1292 information to assess the compliance of the AI system”. This article aligns with **guidelines #2, #6,  
1293 #14** as it about documentation of the system and its contractual requirements, which may also be  
1294 needed for obtaining ethical approvals.

1295 **Article 12 (Record-keeping):** It states that high-risk AI systems shall include “logging capabilities  
1296 to enable the monitoring of the operation of the high-risk AI system with respect to the occurrence of  
1297 situations that may result in the AI system presenting a risk”. This article aligns with **guidelines #6,  
1298 #9, #10, and #14** as it is about providing mechanisms for interpretable outputs and auditing, and  
1299 improving the security of the system.  
1300

1301 **Article 13 (Transparency and provision of information to users):** It states that “high-risk AI  
1302 systems shall be designed and developed in such a way to ensure that their operation is sufficiently  
1303 transparent to enable users to interpret the system’s output and use it appropriately”. This article  
1304 aligns with **guidelines #8-10, #16-18, and #20** as it is about quality, representativeness, and fit of  
1305 training and testing datasets with the intended use.

1306 **Article 14 (Human oversight):** It states that “high-risk AI systems shall be designed and developed  
1307 in such a way, including with appropriate human-machine interface tools, that they can be effectively  
1308 overseen by natural persons during the period in which the AI system is in use”. This article aligns  
1309 with **guidelines #9 and #20** as it about ensuring human control over the system.

1310 **Article 15 (Accuracy, robustness and cybersecurity):** It states that “high-risk AI systems shall  
1311 be designed and developed in such a way that they achieve, in the light of their intended purpose, an  
1312 appropriate level of accuracy, robustness and cybersecurity, and perform consistently in those respects  
1313 throughout their lifecycle”. This article aligns with **guideline #10** as it is about documenting the  
1314 security of all system components.  
1315

1316 **Article 16 (Obligations of providers of high-risk AI systems):** It states that “providers of high-risk  
1317 AI systems shall draw-up the technical documentation of the high-risk AI system”. This article aligns  
1318 with **guideline #6** as it is about system documentation.

1319 **Article 17 (Quality management system):** It states that “an AI system shall be documented in a  
1320 systematic and orderly manner in the form of written policies, procedures and instructions”. This  
1321 article aligns with **guidelines #6, #7, #10, and #14-18** because it is about documentation of all  
1322 system components, including AI models and testing and validation procedures.  
1323

1324 **Article 18 (Obligation to draw up technical documentation):** It states that “providers of high-risk  
 1325 AI systems shall draw up the technical documentation”. This article aligns with **guideline #6** as it is  
 1326 about system documentation.

1327 **Article 20 (Automatically generated logs):** It states that “providers of high-risk AI systems shall  
 1328 keep the logs automatically generated by their high-risk AI systems, to the extent such logs are under  
 1329 their control by virtue of a contractual arrangement with the user or otherwise by law”. This article  
 1330 aligns with **guideline #19** as it is about monitoring of the system.

1331 **Article 29 (Obligations of users of high-risk AI systems):** It states that users shall “monitor the  
 1332 operation of the high-risk AI system on the basis of the instructions of use.”, and “inform the provider  
 1333 or distributor when they have identified any serious incident or any malfunctioning and interrupt the  
 1334 use of the AI system”. This article aligns with **guideline #19** as it about monitoring of the system  
 1335 and utilizing guardrails or rollbacks.

1336 **Article 50 (Document retention):** It states that “the provider shall, for a period ending 10 years after  
 1337 the AI system has been placed on the market or put into service, keep at the disposal of the national  
 1338 competent authorities the technical documentation”. This article aligns with **guideline #6** as it about  
 1339 system documentation.

1340 **Article 60 (EU database for stand-alone high-risk AI systems):** It states that information  
 1341 contained in the EU database shall “be accessible to the public” and “include the names and contact  
 1342 details of natural persons who are responsible for registering the system and have the legal authority  
 1343 to represent the provider”. This article aligns with **guideline #4** as it is about providing mechanisms  
 1344 for reporting system harms.

1345 **Article 61 (Post-market monitoring by providers and post-market monitoring plan for high-  
 1346 risk AI systems):** It states that “the post-market monitoring system shall actively and systematically  
 1347 collect, document and analyse relevant data provided by users or collected through other sources on the  
 1348 performance of high-risk AI systems throughout their lifetime”. This article aligns with **guidelines**  
 1349 **#12, #14, #15, #19** as it is about data handling and model updates when the AI system is in use.

1350 **Article 62 (Reporting of serious incidents and of malfunctioning):** It states that “providers of  
 1351 high-risk AI systems placed on the Union market shall report any serious incident or any malfunction-  
 1352 ing of those systems which constitutes a breach of obligations under Union law intended to protect  
 1353 fundamental rights to the market surveillance authorities of the Member States where that incident or  
 1354 breach occurred”. This article aligns with **guideline #4** as it is about incentivizing the reporting of  
 1355 system harms.

1356 **Article 63 (Market surveillance and control of AI systems in the Union market):** It states  
 1357 that “the national supervisory authority shall report to the Commission on a regular basis the out-  
 1358 comes of relevant market surveillance activities.”. This article aligns with **guideline #4** as it about  
 1359 incentivizing the reporting of system harms.

1360 **Article 64 (Access to data and documentation):** It states that “access to data and documentation  
 1361 in the context of their activities, the market surveillance authorities shall be granted full access to  
 1362 the training, validation and testing datasets used by the provider, including through application  
 1363 programming interfaces (‘API’) or other appropriate technical means and tools enabling remote access”.  
 1364 This article aligns with **guidelines #16 and #17** as it is about data documentation.

1365 **Article 65 (Procedure for dealing with AI systems presenting a risk at national level):** It  
 1366 states that “AI systems presenting a risk shall be understood as a product presenting a risk defined  
 1367 in Article 3, point 19 of Regulation (EU) 2019/1020 insofar as risks to the health or safety or to the  
 1368 protection of fundamental rights of persons are concerned”. This article aligns with **guideline #3** as  
 1369 it is about harms and risks identification.

1370  
 1371  
 1372

1373 **Article 67 (Compliant AI systems which present a risk):** It states that if the AI system is  
1374 compliant with the EU AI Act but still presents a risk to the health or safety of persons, the market  
1375 surveillance authority “shall require the relevant operator to take all appropriate measures to ensure  
1376 that the AI system concerned, when placed on the market or put into service, no longer presents that risk,  
1377 to withdraw the AI system from the market or to recall it within a reasonable period, commensurate  
1378 with the nature of the risk, as it may prescribe”. This article aligns with **guideline #5** as it is about  
1379 mitigation strategies about the identified harms and risks.

1380 **Article 69 (Codes of conduct):** It states that “the Commission and the Board shall encourage and  
1381 facilitate the drawing up of codes of conduct intended to foster the voluntary application to AI systems  
1382 of requirements related for example to environmental sustainability, accessibility for persons with a  
1383 disability, stakeholders participation in the design and development of the AI systems and diversity  
1384 of development teams on the basis of clear objectives and key performance indicators to measure the  
1385 achievement of those objectives”. This article aligns with **guidelines #2, #11, #21, #22** as it is about  
1386 the environmental assessment of the system, the ethical approvals obtained from ethics committees  
1387 and boards, and the characteristics of the development team.

1388

1389

1390

1391

1392

1393

1394

1395

1396

1397

1398

1399

1400

1401

1402

1403

1404

1405

1406

1407

1408

1409

1410

1411

1412

1413

1414

1415

1416

1417

1418

1419

1420

1421