

Towards a multi-stakeholder value-based assessment framework for algorithmic systems

ANONYMOUS AUTHOR(S)

In an effort to regulate Machine Learning-driven (ML) systems, current auditing processes mostly focus on detecting harmful algorithmic biases. While these strategies have proven to be impactful, some values outlined in documents dealing with ethics in ML-driven systems are still underrepresented in auditing processes. Such *unaddressed* values mainly deal with contextual factors that cannot be easily quantified. In this paper, we develop a value-based assessment framework that is not limited to bias auditing and that covers prominent ethical principles for algorithmic systems. Our framework presents a circular arrangement of values with two bipolar dimensions that make common motivations and potential tensions explicit. In order to operationalize these high-level principles, values are then broken down into specific criteria and their manifestations. However, some of these value-specific criteria are mutually exclusive and require negotiation. As opposed to some other auditing frameworks that merely rely on ML researchers' and practitioners' input, we argue that it is necessary to include stakeholders that present diverse standpoints to systematically negotiate and consolidate value and criteria tensions. To that end, we map stakeholders with different insight needs, and assign tailored means for communicating value manifestations to them. We, therefore, contribute to current ML auditing practices with an assessment framework that visualizes closeness and tensions between values and we give guidelines on how to operationalize them, while opening up the evaluation and deliberation process to a wide range of stakeholders.

CCS Concepts: • **General and reference** → **Evaluation**; • **Human-centered computing** → **Human computer interaction (HCI)**; • **Social and professional topics** → **User characteristics**.

Additional Key Words and Phrases: values, ML development and deployment pipeline, algorithm assessment, multi-stakeholder

ACM Reference Format:

Anonymous Author(s). 2018. Towards a multi-stakeholder value-based assessment framework for algorithmic systems. In *Woodstock '18: ACM Symposium on Neural Gaze Detection, June 03–05, 2018, Woodstock, NY*. ACM, New York, NY, USA, 37 pages. <https://doi.org/10.1145/1122445.1122456>

1 INTRODUCTION

In recent years, it has become clear that algorithmic systems might encode harmful biases and might lead to unfair outcomes [133, 158]. The dangers of using Machine Learning (ML) in Computer Vision (CV) [31] or Natural Language Processing (NLP) [8, 21, 44, 111, 153], for assessing recidivism [170], for candidate screening [146] and for recommending content on social media platforms [99, 142, 148, 192] have been pinpointed. The origins of harmful algorithmic bias¹ might be diverse [158, 165]. Just to mention a few, representativeness issues, play a key role in disparate algorithmic performance [3, 34]. The way in which data is collected [21, 144] and labelled [39, 51, 144] is a major menace to data soundness. Beyond the data generation process, aggregation, learning, evaluation and deployment biases have been identified throughout

¹Following the approach adopted by Shen et al. [158], we will distinguish between harmful algorithmic biases and harmful algorithmic behaviors, since not all harmful algorithmic behaviors originate from biases and not all algorithmic biases are necessarily harmful [28].

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.

© 2018 Association for Computing Machinery.

Manuscript submitted to ACM

the ML pipeline [165]. In response to harmful algorithmic bias, current auditing processes² [2, 151, 181] have provided numerous useful bias detection techniques [10, 17, 19, 55, 71, 91, 184, 185, 188].

However, harmful algorithmic behavior is not limited to biases encoded in the ML life cycle [158]. The lack of social and cultural context in the mathematical representation of socio-technical systems [124, 158] or the omission of changing practices and long-term effects of the deployed systems [29, 42, 95, 108] are also some problematic aspects that are hardly considered in current auditing processes. Such processes mostly consist of quantitative analysis for assessing the conformance of those systems to applicable standards [94], rather than additionally gaining insights into their contextual implications [147, 158]. Furthermore, these auditing approaches solely rely on ML researchers, and practitioners, who can fail to detect issues that arise from context-dependent unanticipated circumstances during usage time [158].

In this paper, we argue that: firstly, assessment processes for algorithmic systems should go beyond bias auditing and take into account additional high-level values³ that are outlined in Artificial Intelligence (AI) ethics documents [11, 36, 57, 60, 75, 93, 96, 127, 139, 168, 172]. Contestability, for example, has been identified as a key value of algorithmic systems, but there is still little guidance on what contestability requires [121]. In order to provide a good coverage of values that deal with principled algorithmic behavior, we develop a value-based assessment framework, where contextual conditions are considered along with quantifiable measurements. We organize such values in a circular layout with two bipolar dimensions. As claimed by Friedman et al. [65], values do not exist in isolation. They are situated in a delicate balance and touching one value might have implications in another value [65]. This means that value interactions need to be taken into account when making choices about value prioritization and situating algorithmic systems in a space of trade-offs [15]. The circularity of our framework makes such interactions explicit and facilitates the identification of common motivations and tensions among values.

Secondly, an assessment process should give tangible guidelines for the operationalization⁴ of values, so as to eventually put ethics into practice following a context-aware approach [159]. To this end, each value in our framework is broken down into criteria manifested through quantifiable indicators, process-oriented practices or signifiers⁵. These value-specific criteria and their manifestations can be used either as a checklist if our framework is applied for evaluating a system that is already developed, or for promoting such values if it is being used during design time.

Thirdly, assessment processes should allow critical reflection on algorithmic systems and engage in conflictual plurality⁶. Inevitable value tensions inherent in the nature of socio-technical systems [78] require spaces for ethical discussions [159], that can benefit from the insights of multiple stakeholders beyond ML practitioners [15, 158]. To enable fruitful multi-stakeholder discussions [115], we map and match value-specific communication means with different stakeholders. We, therefore, contribute with:

- A review of prominent high-level values in AI ethics and translation into specific criteria through the:

²We will use the term *auditing* processes to refer to external audits, where third parties only have access to model outputs [152]. We will use the term *assessment* processes to refer to an evaluation process that is applied “throughout the development process and that enables proactive ethical intervention methods” [147]. We will not use the term *Internal Audit* defined by Raji and Smart [147] to avoid erroneous inferences that would limit the stakeholders of our framework to the employees of an organization.

³We will adopt the definition of *values* used in philosophy of science, following Bihrane et al. [26]. Values of an entity are, thus, defined as properties that are desirable for that kind of entity.

⁴Our strategy follows the definition by Shahin et al. [157], where “operationalizing values” is defined as the process of identifying values and translating them into concrete system specifications that can be implemented.

⁵We adopt the definition given by Don Norman in his 2013 edition of “The Design of Everyday Things”. Signifiers are perceivable cues of an affordance, affordances being “the relationship between the properties of an object and the capabilities of the agent that determine how the object could be possibly used”. In this paper, the “object” in question is the ML-driven system.

⁶We understand *conflictuality* as a solution for dealing with the “figure of alterity”. Unlike *conflict*, it represents a method for linking opposing views and opening out onto the exercise of thinking [66]

- Design of an assessment framework that facilitates the identification of common motivations and tensions among values encoded in ML-driven systems.
- Definition of guidelines to deal with the complex middle ground between abstract values and concrete system specifications.
- Translation of value-specific criteria into manifestations that are understandable for diverse stakeholders through the:
 - Review of available means to communicate value manifestations to different stakeholders based on their insight needs and nature of knowledge.
 - Definition of steps to introduce those communication means into multi-stakeholder deliberation dynamics.

The remainder of the paper is organized as follows: in section 2, we analyze related work for documenting and auditing ML systems. We also introduce the theoretical basis of our framework. Section 3 describes and justifies the selected values, criteria and manifestations and their arrangement in our framework. Section 4 maps the stakeholders involved in the algorithm evaluation process and reviews the available means for communicating system-specific information to them. In sections 3 and 4, we illustrate the necessary steps for navigating our framework through an example in the area of life insurance application. We discuss our approach, its implications, and future lines of work in section 5, and we conclude this paper in section 6.

2 BACKGROUND AND RELATED WORK

In this section, we survey current practices for documenting and auditing technical specifications of algorithmic systems. We also provide the theoretical basis of our framework.

2.1 Background

2.1.1 Standardized documentation. In order to facilitate the audit of ML-driven systems, it is important that technical specifications are collected and documented in a standardized way. So far, ML system documentation practices are limited to datasets and models.

Documenting datasets. Recent studies in documentation practices for ML datasets claim the need for greater data transparency [92]. Since the quality of the prediction made by the ML system highly depends on the way the data has been collected, the need for setting rigorous practices (as it is the case in other areas of knowledge, such as social sciences or humanities [70]) has been highlighted [144]. Likewise, the choice of what data to collect and how to collect this data is in itself a value-laden decision [46, 155]. To standardize documentation for ML datasets and make data-related decisions more transparent for other practitioners, various methodologies have been suggested in the last years, “Datasheets” [69] and “Dataset Nutrition Labels” [88], for instance. For NLP techniques, “Data Statements” are regarded as a dataset characterization approach that helps developers anticipate biases in language technology and understand how these can be better deployed [20].

Documenting models. In addition to documenting datasets, the importance of disclosing the technical characteristics of ML models has also been emphasized. A good example of model documentation practices are the “Model Cards” [130].

2.1.2 Auditing techniques. Various methodologies and tools for incorporating auditing tasks into the Machine Learning workflow have been suggested. Aequitas [151] is an open source toolkit to detect traces of bias in models. The toolkit designed by Saleiro et al. [151] facilitates the creation of equitable algorithmic decision-making systems where data scientists and policymakers can easily use Aequitas for model selection, evaluation and approval. Wilson et al. [181] described a framework that helps ensure fairness in socio-technical systems, and used it for auditing the model of the startup *pymetrics*. Adler et al. [2] studied auditing techniques for black-box models to discover whether proxy variables linked to sensitive

157 attributes indirectly influence the predictions of the system. The end-to-end “Internal Audit Framework” suggested by Raji
158 and Smart [147] is of special interest for justifying the need of setting specific guidelines to enable multi-stakeholder delib-
159 eration in assessment processes. It consists of five main stages where the need for stakeholder diversity is highlighted, e.g.
160 the scoping stage calls for covering a “critical range of viewpoints” to review the ethical implications of the system use case.
161

162
163 **2.1.3 Motivation.** While standardized documentation practices [20, 69, 88, 130] and audits [2, 151, 181] have been
164 influential methodologies for dealing with harmful algorithmic bias, their scope is limited to performing quantitative
165 analysis over data and model outputs so as to ensure compliance with applicable standards [94]. Such an approach does
166 not deal with additional ethical values which cannot be easily quantified [115] and that are essential for ensuring desirable
167 algorithmic behavior. One could argue that “Datasheets” [69] and “ModelCards” [130] already devote a section to the
168 description of ethical considerations of datasets and models. Yet, there are no specific guidelines on how to identify
169 ethical issues. As Shklovski et al. [159] discovered, technical people both in industry and academia struggle to identify
170 what an ethical issue entails. To address this caveat, as part of our value-based framework, we give tangible guidelines
171 for putting ethics into practice [133, 159]. We operationalize each high-level value into actionable value criteria and
172 their manifestations. One could also argue that Raji and Smart [147] already included an Ethics Review as part of their
173 end-to-end internal audit framework. Indeed, they exemplified such a review by describing ethical considerations and
174 potential mitigation strategies against bias and privacy threats for a smile detection system. However, this review does
175 not address most of the values that are referred in AI ethics documents. We fill in this gap by offering a good coverage
176 of values to examine, including those that normally go unnoticed in current documenting and auditing practices.
177
178
179
180

181 **2.2 Accounting for human values in the assessment of algorithmic systems**

182
183 Our ML assessment framework identifies and arranges values encoded in algorithmic systems by covering prominent
184 principles in AI ethics and organizing them in a circular structure.
185

186 **2.2.1 Addressing human values in technology.** For the definition of our value-based framework, we followed other
187 theoretically grounded approaches, such as Value Sensitive Design (VSD) [65]. VSD represents a pioneering endeavour
188 where human values are proactively considered throughout the process of technology design [45]. Just as VSD does
189 with interactive systems, we address the need to account for human values during the design, implementation, use, and
190 evaluation [45] of algorithmic systems. To this end, we select and define values involved in ML-driven systems, and
191 we identify stakeholders that will be in contact with such systems and whose standpoints need to be considered. Our
192 approach resonates with conceptual investigations described in VSD literature [45].
193
194

195 The circular nature of our framework is inspired by Schwartz’s Theory of Basic Human Values [156]. This theory
196 identifies individual value priorities based on ten basic personal values. Values are arranged in a circular form and
197 categorized in four quadrants. These quadrants are located in two bipolar dimensions, which visualize “oppositions
198 between competing values”. In addition, adjacency between values denotes a common motivation, which results in
199 these values forming a circular continuum. The advantage of adopting a circular arrangement, like the one suggested
200 by Schwartz, for ML-driven systems is that value commonalities and trade-offs can be easily identified thanks to their
201 positioning. Considering the struggles of technical people when addressing ethical issues [159], an explicit representation
202 of value interactions will facilitate the analysis of trade-offs and decision-making about value prioritization.
203
204

205 **2.2.2 Ethical principles for ML-driven systems.** The values considered in our assessment framework cover prominent
206 principles outlined in AI ethics. In the last five years, many institutions have studied and defined high-level principles that
207
208

209 AI systems should follow [60]. As a matter of fact, documents that aim at guiding the “ethical development, deployment
210 and governance of AI” are converging into a common set of principles [132, 133]. However, high-level principles are far
211 from being actionable [133] and it is necessary to provide answers on how to proceed [4]. Efforts for going from “what” to
212 “how”⁷ include the review carried out by Morley et al. [133], where available tools for operationalizing ethical principles
213 were examined. Similarly, the AI Ethics Impact (AIEI) Group designed a framework for rating the presence of ethical
214 principles in AI systems, getting inspiration from energy efficiency labels [3].

215 Our value-based framework differs from previous applied ethics frameworks [3, 133] in various ways. Firstly, we
216 arrange values in a circular form, which makes it easier to navigate common motivations and trade-offs between values.
217 Although such common motivations and trade-offs can be inferred from current AI ethics documents, we make them
218 explicit by arranging values in a geometrically meaningful way. This is especially useful for identifying overlaps between
219 values that are adjacent to each other and for detecting potential value tensions that need to be negotiated and consolidated.
220 Secondly, we do not limit our ethics framework to a mere checklist. We follow Shklovski et al. [159] and combine the
221 enumeration of tangible and actionable value manifestations with the generation of an open space for ethical debate.
222 As opposed to the deterministic approach adopted by the AIEI group [3], we map communication means for facilitating
223 ethical reflections of algorithmic systems and for addressing ethical issues in practice [159]. Thirdly, as opposed to
224 previous applied ethics frameworks [3, 133], we embrace diversity in ethical reflections and deal with the complexities
225 that arise from plurality. In order to facilitate multi-stakeholder discussions, we match available communication means
226 for addressing different value manifestations with stakeholders that present different insight needs.
227
228
229
230
231

232 3 DESIGN OF OUR VALUE-BASED FRAMEWORK

233
234 In this section, we describe the composition of our value-based framework and justify its arrangement. We provide the
235 definition of each of the selected values and the derived criteria and manifestations.
236

237 3.1 Methodology for reviewing values, criteria and manifestations in ML-driven systems

238
239 To design our framework, we analysed documents outlining high-level ethical principles that ML systems should follow.
240 Our starting point was the review performed by Fjeld et al. [60], where principles coming from governments, inter-
241 governmental organizations, multiple stakeholders, the private sector, and the civil society were examined. In their review,
242 Fjeld et al. identify nine key themes, some of which overlap with the values outlined in our framework. The identification
243 of prominent high-level values was also complemented with other reviews [3, 26, 43, 86, 133]. To identify the criteria
244 that define the fulfilment of prominent high-level values, we navigated the visual representation provided by Fjeld et
245 al. [60] and accessed the documents that offer a higher coverage of the value in question. For instance, for the value of
246 *privacy*, one of our main references has been the GDPR [56].
247

248
249 We went from criteria to value manifestations through an extensive exploration of available value-specific reviews
250 that identify such manifestations. For instance, for the value of *security* Xiong et al. [183] presented a thorough study
251 of mechanisms used for securing the ML pipeline against external threats. For *explainability*, Barredo-Arrieta et al. [18]
252 put together more than four hundred references and mapped strategies in the field of Explainable Artificial Intelligence
253 [18]. We partly rely on such reviews for identifying value manifestations because our contribution lies in covering and
254 putting together a set of values and their manifestations in ML-driven systems to end up with a “health-check” for
255 assessing algorithmic systems, rather than rediscovering such value manifestations ourselves. Similarly, for the values
256
257

258 ⁷Expression used by Morley et al. [133] to refer to the operationalization of ethical principles in AI. The ‘what’ refers to the ethical principles themselves,
259 whereas the ‘how’ refers to the act of putting such principles into practice.

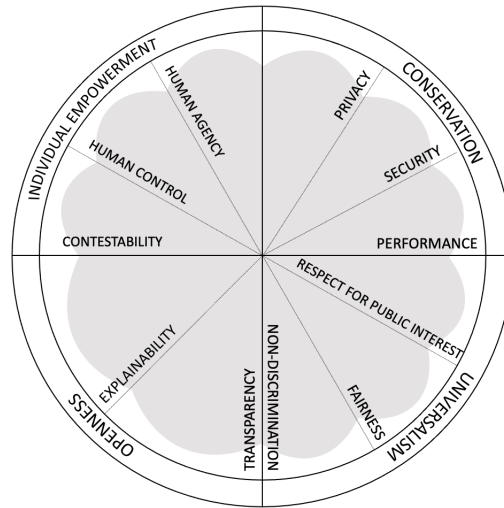


Fig. 1. Graphic representation of our circular value-based assessment framework. Oppositions between competing values are illustrated through the arrangement of those values in bipolar dimensions and common motivations through adjacency between values, which form a circular continuum.

of *performance* and *fairness*, we only included the main value manifestations that represent the basis for any other derived metrics. That is to say, just as Verma et al. [177] did, we outline the main quantifiable indicators (false positives, false negatives etc) used for measuring *performance* and *fairness*, but we are aware that many other metrics that derive from these ones can be insightful for specific contexts. Dealing with such compound metrics is out of the scope of this work.

3.2 Assessment of algorithmic systems through a circular value-based framework

Our resulting ML assessment framework arranges values in a circular form (figure 1). Adjacency between values denotes a common motivation and oppositions between competing values are represented through two bipolar dimensions. For instance, adjacency between *privacy* and *security* denotes a common objective towards the protection of sensitive information [60, 150] and resilience to external threats [127]. The trade-off between *privacy* and *explainability*, on the other hand, is made explicit by their opposing positioning in our circular framework. High-level values are then broken down into specific criteria and their manifestations, as indicated in figure 2. Criteria defining a specific value ultimately represent a set of questions to be asked as part of the assessment process to ensure the fulfillment of the value in question—if the framework is being applied before deployment—or the promotion of a specific value—if the framework is being applied during design time—. These sets of criteria are not unique and exclusive to one value. For instance, when defining the criteria for *privacy* we refer to “data protection”, which is also involved in *security* in the form of “resilience to attacks”. These overlaps are precisely what we want to highlight and make explicit thanks to the circularity of our framework and adjacency between values.

Manifestations are classified in three groups depending on their nature: (1) *Quantifiable indicators* are specific measurable parameters that numerically manifest the (lack of) adequacy in the standards set for a criterion (magenta). (2) *Process-oriented practices* are actions and mechanisms implemented during the ML development or deployment process that advocate for a certain value (olive). (3) *Signifiers*⁸ are files and reports that describe the relationship between the

⁸Check footnote 5

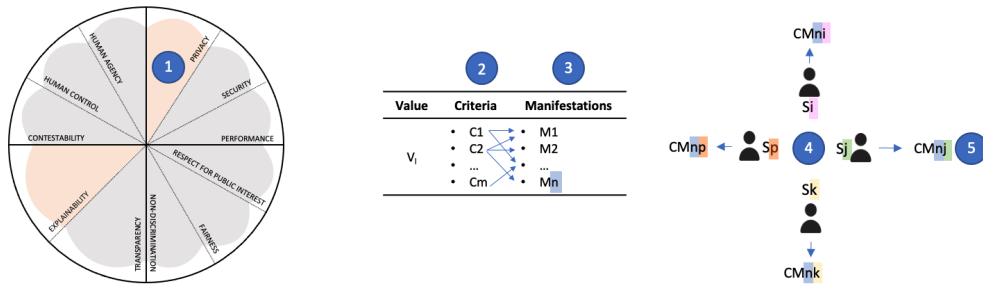


Fig. 2. Workflow for operationalizing high-level values and for enabling multi-stakeholder assessment of algorithmic systems. This workflow represents the methodology that we followed for structuring our framework and the steps that researchers and practitioners should take to make use of it. (1) Select and discuss project-specific values (V), (2) Decide on criteria (C) for embodying those values, (3) Select the manifestations (M) that enact value-specific criteria, (4) Map relevant stakeholders (S) to enable ethical reflection of value and criteria tensions, (5) Match adequate communication means (CM) to stakeholders.

about individuals [86, 178]. Furthermore, the need to provide humans with agency over their data is emphasized [60]. Based on these definitions, we identified six main criteria for the fulfillment of *privacy* within ML systems (table 1). (1) Consent for data usage [3, 56, 60]: data subjects should be appropriately informed when their data is being used and their explicit approval is needed. (2) Implementation of data protection mechanisms [3, 60, 61]: during the development of algorithmic systems, resources should be devoted to making user data management secure and confidential. (3) Users having control over their data and ability to restrict its processing [56, 60]: users should be able to limit the way their personal data is being used. (4) Users having the right to rectify [3, 56, 60]: users should be able to modify their data at any time. (5) Users having the right to erase their data [3, 56, 60]: this criterion refers to the right that users have to be forgotten. (6) Users having right to access their data [56, 167]: this right empowers users to have agency over their data. These criteria manifest in various ways. Signifiers include: a written declaration of consent [56], detailed descriptions of the collected data, how data is handled, how long it will be kept and the purpose of collecting that data [125]. These signifiers are necessary for users to fully understand what sharing their data entails. Process-oriented practices include the obfuscation of data [3] and forms and submission mechanisms to object data collection and make complaints [27].

Value	Criteria	Manifestations
Privacy	(1) Consent for data usage [3, 56, 60]	• Written declaration of consent [56]
	(2) Data protection [3, 60, 61]	• Description of what data is collected [125]
	(3) Control over data / ability to restrict processing [56, 60]	• Description of how data is handled [125]
	(4) Right to rectification [3, 56, 60]	• Purpose statement of data collection [125]
	(5) Right to erase the data [3, 56, 60]	• Statement of how long the data is kept [125]
	(6) Right of access by data subject, data agency [56, 167]	• Form and submission mechanisms to object data collection and to make complaints [27]
		• Obfuscation of data [3]

Table 1. Illustration of how to move from values, to criteria and their manifestations with an example for *privacy*. The rest of the values, criteria and manifestations are detailed in appendix A.

Security. Definitions characterizing *security* highlight the need for ML systems to be (1) resilient to potential maleficent attacks [60, 133] and to present a (2) predictable [3, 57, 60] and (3) robust [3] behavior at any time. This includes implementing mechanisms to protect user privacy, such as strategies that ensure that inferences about an individual cannot be made by interrogating the model [86, 127, 178]. Following the survey performed by Xiong et al. [183], different methodologies that aim at protecting algorithmic systems against external threats (process-oriented practices) have been classified into two main groups. The first group consists of defence methods against integrity threats at two different stages of the ML pipeline: during training time [23, 40, 73] and during prediction time [23, 74, 122, 141]. The second group aim at defending the ML system against privacy threats, namely membership inference attacks [53, 97, 136, 161, 187].

Performance. The value of *performance* is defined by the (1) correctness of predictions [57, 60], along with the (2-5) resources necessary to reach such predictions [3, 26, 109]. The conditions under which systems are evaluated will have a direct impact on the “appropriateness score” that these systems will obtain in the form of a quantifiable indicator [50]. In other words, if the level of performance is solely measured in terms of accuracy, regardless of the needed data, prerequisites will be inherently favoring big “data-hungry” [112] models. As far as the measurement of performance is concerned, this is mainly done through quantifiable indicators, either referring to the preciseness of the results [130, 180] or to the estimated consumption of environmental resources [14, 41, 67, 68, 123].

3.3.2 *Openness*. Transparency and explainability advocate for making system components and specifications accessible.

417 *Transparency.* Documents providing high-level principles for AI define *transparency* as the property that enables
418 traceability and monitoring of algorithmic systems [60, 133]. *Transparency* relates to the right to information [60] and
419 requires that data or algorithms present some level of accessibility [168]. That is to say, data and models should present
420 some level of (1) interpretability [26, 168], so as to (2) enable human oversight [60, 133]. Those data and models should also
421 be (3) accessible [3, 60, 168], as a step towards achieving (4) traceability [133] and (5) reproducibility [26]. Manifestations of
422 such criteria emerge mostly in the form of documentation detailing technical aspects of the algorithmic system (considered
423 signifiers in our framework) [3, 20, 34, 69, 70, 130, 133, 168]. Process-oriented practices mostly focus on giving open access
424 to data and algorithms [3, 26, 60, 168], regularly reporting key information about the system [60] and notifying users
425 whenever they are being subject to or interacting with an algorithmic system [60].
426
427

428 *Explainability.* Explainable Artificial Intelligence (XAI) is formed by a set of techniques that allow a wide range of
429 stakeholders to understand why or how a decision was reached by an algorithmic system [61, 168]. *Explainability* is, thus,
430 conceived as an interface that translates reasoning mechanisms of the system into formats that are (1) comprehensible
431 [18, 26, 57, 60–62, 139, 168]. In addition, strategies for making black-box algorithms more interpretable facilitate their
432 (2) monitoring [133] and, therefore, make them (3) suitable for evaluation [60, 133]. XAI techniques (process-oriented
433 practices) are very diverse in nature. As claimed by Vera Liao et al. [117] and Barredo-Arrieta et al. [18], *explainability*
434 methodologies are usually classified by the scope of the explanation, complexity of the model, model specificity and the
435 stage of the ML pipeline where such a strategy is to be used. For our framework, we will consider that explainable models
436 can be either (a) interpretable by design or they can be (b) explained by additional *post-hoc* explanations [18].
437
438
439
440

441 3.4 Universalism vs Individual Empowerment

442 The second dimension captures the conflict between *universalism* and *individual empowerment*. Values included within
443 the *individual empowerment* category emphasize the defense of the decision subjects' interests. These principles advocate
444 for giving decision subjects the means to oppose to the conclusion reached and uphold the need for putting humans in
445 the loop. Values within the *universalism* category emphasize the need to equalize system behavior to *all* and to ensure
446 that such a system adheres to the interests of society as a whole, beyond the interests of a few individuals.
447
448

449 3.4.1 *Universalism.* Respect for public interest, fairness and non-discrimination uphold the need to ensure equitable
450 and socially acceptable system behavior for *all*.
451

452 *Respect for public interest.* The value of *respect for public interest* deals with the (1) appropriateness of developing algo-
453 rithmic systems for a certain purpose within a specific context. As Keyes et al. [104] claimed, making ML-driven systems
454 fairer, more transparent and more accountable is insufficient if we ignore the purpose of developing and implementing
455 these systems in a certain context in the very first place [116, 171]. Algorithmic systems should, therefore, (2) be beneficial
456 to society and humanity as a whole [60–62, 133], respect law [26] and be aligned with human norms [60]. This involves
457 giving a clear justification of the purpose and benefits of building such a system [1, 34, 104, 133], so that the deployment of
458 the system in question upholds public-spirited goals [60]. Universalism aims at protecting the welfare of *all*, both people
459 and nature [109]. AI systems' (3) negative impacts on environment should, therefore, be considered and valued [3, 21]. To
460 this end, process-oriented practices include the creation of diverse and inclusive forums for discussion [60, 129], whereas
461 signifiers include the qualitative measurement of social and environmental impact [21, 133, 147].
462
463
464

465 *Fairness.* The value of *fairness* represents a complex concept that accepts multiple definitions [15, 103], some of which
466 cannot be satisfied simultaneously [80, 86, 103]. Overall, we will understand *fairness* in terms of parity in output [49]
467
468

469 and equal treatment [3] among individuals. When addressing more specific definitions of *fairness* (1-8), we will adopt the
470 approach followed by Verma et al. [177], which was also echoed by Mehrabi et al. [126] (for a detailed enumeration and
471 explanation of each of the definitions, the reader is encouraged to check appendix A). ML techniques generally conceive
472 fairness in terms of statistical metrics [76] and observe whether specific quantifiable indicators are above or below the
473 thresholds set for a certain application. Even if error rates were equal across groups for a certain application, if those rates
474 are too high, the system could still be considered unfair [80]. This means that for our value-based framework we outline
475 the quantifiable indicators that are normally used for manifesting fairness-related criteria, but we do not determine
476 the threshold for these indicators to be considered good enough for a specific application. Similarly, the quantifiable
477 indicators relate to the output of the system, rather than the outcome that these outputs lead to.
478
479

480
481 *Non-discrimination.* The value of *non-discrimination*, as defined in our framework, deals with algorithmic systems
482 not being socially biased [26] and ensuring that equal accessibility is provided to all individuals [133]. This means that (1)
483 quality and integrity of data should be evaluated and ensured [60, 70, 86, 133, 144] in order to prevent “socially constructed
484 biases, inaccuracies, errors, and mistakes” [133] from being present in the data. Processes that safeguard inclusive data
485 generation [3, 34, 70, 133] and analysis procedures for identifying potential biases in data and for assessing its quality
486 [60, 70, 86, 133, 144] are strategies that avoid social stereotypes being codified, maintained and amplified [86]. Furthermore,
487 non-discriminatory systems should (2) ensure diversity and inclusiveness in the design process [57, 60, 133]. From a
488 process-oriented perspective, participants involved in the development process should, thus, present diverse profiles
489 [3, 60, 114, 189]. Finally, giving (3) equal access to the technology [3, 26, 60, 133] avoids the growth of inequalities as a
490 consequence of deploying AI systems [60].
491
492
493

494 3.4.2 *Individual empowerment.* Contestability, human control and human agency address the politics behind algorithmic
495 systems [13, 182] and deal with the issues caused by power imbalances [26, 34, 100, 121, 174].
496
497

498 *Contestability.* The value of *contestability* is defined as the value that ensures that users have the necessary information to
499 (1) enable argumentation against conclusions reached by algorithmic systems [6, 16, 57, 60, 100, 113, 121, 168]. This involves
500 (2) empowering citizens [16, 57, 100] to investigate and influence AI [100], as part of a broader regulatory approach [121].
501 As a matter of fact, *contestability* has been identified as a “critical aspect of future public decision-making systems” [6]. This
502 implies that, from a documentation perspective (signifiers), users should be made aware of who determines what constitutes
503 a contestable decision, who is accountable for it and who can contest a decision. This last point is particularly necessary to
504 determine whether (legal) representatives of decision subjects can act on their behalf. The review mechanism in place and
505 the workflow of contestations [121] are policy-related details that users should also be informed about. From a process-
506 oriented standpoint, mechanisms for users to ask questions and to record disagreements should also be put in place [87, 131].
507
508
509

510 *Human Control.* The value of *Human control* addresses the influence that data-driven technologies have over humans
511 and that leads to a reduction of human agency, power and control [143]. Algorithmic systems should be controllable [26]
512 and (1) subject to user and collective influence [26, 113]. They should also be (2) subject to human review [60]. Governance
513 mechanisms that ensure human oversight of automated decisions are, thus, necessary to maintain control and influence
514 over such systems [133]. It should be possible to (3) choose how and even whether (in the very first place) to delegate a
515 decision to an automated system [60]. From a development perspective, levels of human discretion should be established
516 [57, 127] and the ability to override decisions made by a system [57] ought to be set up by design. Once the system is
517 deployed, it should be continuously monitored to enable adequate intervention when necessary [57, 60, 166].
518
519
520

521 *Human Agency.* The value of *human agency* deals with the risks of algorithmic systems displacing human autonomy
522 [57, 60]. As claimed by Cila et al. [38], algorithmic systems may displace human agency in governance processes and may
523 undermine human autonomy. ML-driven systems advocating for human agency should, therefore, (1) respect human
524 autonomy [57, 60, 133] and (2) citizens' power to decide [26, 57]. In addition, (3) decision subjects should be able to opt
525 out of an automated decision [57, 60]. The manifestations of such criteria involve giving knowledge and tools to users to
526 comprehend and interact with AI systems [57] (signifier) and, from a process-oriented perspective, providing strategies
527 for users to self-assess the systems [57].
528
529

530 *Selecting values, criteria and manifestations for our example use case.* Returning to the hypothetical insurance modelling
531 team from our motivating example (section 3.2), they decided to apply our value-based framework before launching their
532 system. They quickly realised that they need to consider more values than those outlined in current auditing processes.
533 For example transparency, non-discrimination, supporting human agency and the public good. They also discovered a
534 range of methods for enacting those values: from data handling processes that ensure anonymity and meaningful consent
535 around the model, to models of fairness appropriate to their case.
536
537

538 Although we cover prominent ethical principles in AI and the assessment of the algorithmic system might include all
539 of them, here we focus on a subset of those values for illustrative purposes. We imagine that the researchers developing
540 the algorithmic life insurance application system want to focus on *explainability* and *privacy* (fig 2). We assume that
541 they are dealing with a blackbox algorithm that is not interpretable by design. Checking appendix A, the team needs
542 to examine whether the algorithmic system and the decision reached are understandable. Additionally, the deployed
543 XAI methods should enable traceability and evaluation of the system. As far as the *explainability* manifestations are
544 concerned, since they are dealing with a blackbox algorithm, they need to deploy adequate post-hoc explanations. When
545 it comes to *privacy*, the data used for training and testing the algorithmic model should have been obtained through the
546 explicit approval of the decision subjects. These subjects should have been informed about the nature and purpose of
547 the data that is collected, the way this data is handled and stored. Decision subjects should also have agency and control
548 over their data. Additionally, data protection mechanisms should have been implemented to make sure that there is no
549 possibility of identifying sensitive (in this case medical) data about the subjects. These two values that the team needs
550 to advocate for, represent some trade-offs: XAI methods uphold interpretability of algorithmic systems and some of them
551 even rely on comparing data instances at inference time with those used for training the system. This would directly
552 violate the subjects' right to have their data protected and confidentiality ensured.
553
554
555
556
557

558 4 TOWARDS A MULTI-STAKEHOLDER CRITICAL REFLECTION OF ALGORITHMIC SYSTEMS

559 Since there are value trade-offs, like the one outlined in our example use case, and certain value-specific criteria are
560 mutually exclusive, we follow the claim made by Raji and Smart [147], and advocate for standpoint diversity. This implies
561 involving a wide range of stakeholders in the negotiation process [159] to discuss and critically reflect on the degree to
562 which each of the values should be promoted in detriment of the other one and how the prioritization process should
563 take place. These stakeholders will possess different types of knowledge and will present different insight needs. In this
564 section we map those stakeholders and match them with the most suitable communication means.
565
566
567

568 4.1 Methodology for identifying relevant stakeholders and communication means

569 To identify relevant stakeholders, we follow the stakeholder characterization of Suresh et al. [164]. They classified
570 stakeholders in a two dimensional matrix, where one dimension captured the nature of the knowledge of the stakeholders
571
572

(formal, instrumental or personal) and the second one identified the context in which that knowledge manifests (Machine Learning, data domain, and the general milieu). Formal knowledge entails a deep understanding of the theories of a certain domain. Instrumental knowledge refers to the capability of applying formal knowledge in one of the three contexts. Personal knowledge is acquired by the participation of the subject in a specific context. The two dimensional-matrix classification results in nine different stakeholder profiles. To facilitate the process of mapping the stakeholders to tailored communication means, we narrow those stakeholders down into four categories ⁹.

We then proceed to identify the means to communicate system-related information to different stakeholders. We searched such means using arXiv and Google search, so as to cover the state of the art in terms of research papers and open source toolkits. Each search referred to specific value criteria and manifestations, although many of the found means address more than one value. This review does not intend to be exhaustive. We expect novel research to address value manifestations that still present scarce resources in our framework. Hence our review is just a snapshot of some of the available communication means until January 2022, but we host the latest version on an online repository ¹⁰ and is open to anyone's contribution. We aim at creating a living document that will keep growing and that will address current research gaps as time goes by.

4.2 Mapping stakeholders

We characterize four main stakeholders in our framework: (1) The development team: they have the formal, instrumental and personal knowledge in the domain of ML [164]. They want to ensure and improve product efficiency and research new functionalities [18]. (2) Auditing team: they have the formal and instrumental knowledge of the general milieu, meaning that they are aware of the social theories behind AI, and are able to evaluate technical specifications of ML systems. They aim at verifying model compliance with legislation [18] (3) Data domain experts: they have the theoretical (formal) and instrumental knowledge of the application context (healthcare, economics etc.). They look forward to gaining scientific or domain-specific knowledge [18, 164], trust the model [18, 164] and act based on the model output [164]. And (4) Data subjects: they have the personal knowledge of the data domain in which the AI is being applied and the general milieu. They aim at understanding their situation [18], verifying that the decision is fair [18], contesting the decision (if needed) [164] and understanding how their data is being used [164].

Mapping stakeholders in our example use case. Going back to our example, once *explainability* and *privacy* have been broken down into specific criteria manifestations, the team needs to map the stakeholders who will take part in the assessment process (fig 2). Based on the mapping presented in appendix B, the development team represents the stakeholders who have the knowledge of the math behind the system. An external auditing team will join the discussion to make sure that the model is aligned with current legislation. Since the algorithmic life insurance application system deals with medical data, the data domain experts will be represented by a medical team and a life insurance expert. Decision subjects will be laypeople who seek to understand and verify their situation with regards to data usage and the decision reached by the system.

4.3 Mapping tailored communication means

We then examine each of the reviewed means and identify their typology, the value manifestations that they cover and the stakeholders that can make use of it, as illustrated in table 2 for privacy dashboards. The objective of mapping value manifestations, stakeholder profiles and communication means is that of enabling a fruitful and informed discussion among stakeholders. We classify these means in three categories: (1) *Descriptive documents* (red), (2) *Design strategies*

⁹This reduced classification is backed up by the framework employed by Barredo-Arrieta et al. [18] when identifying the explainability needs of various stakeholders.

¹⁰For the sake of anonymity, the link to the online repository has not been disclosed for the peer-review process

(blue), and (3) *Ready-to-use tools* (green). **Appendix C summarizes the rest of the communication means we identified and maps them to value manifestations and stakeholders for whom such methods are suitable.**

Means	Value	Manifestation(s)	Stakeholder				Application (model)	Approach	Visual elements
			DT	AT	DE	DS			
[B] Privacy dashboards [54, 58, 59, 85, 191]	Privacy	<ul style="list-style-type: none"> Description of what and why data is collected Description of how data is handled 				✓	Agnostic	<ul style="list-style-type: none"> Timelines Bar charts Maps Network graphs 	
	Human agency	<ul style="list-style-type: none"> Self-assessment of the system 							
	Transparency	<ul style="list-style-type: none"> Disclosure of properties of data 							

Table 2. Illustration of how we mapped communication means with values, manifestations and stakeholders (DT = Development Team; AT = Auditing Team; DE = Data Domain Experts; DS = Decision Subjects). Privacy dashboards are tools (green) that allow users to interactively assess the collection and usage of their data. The rest of the reviewed communication means are characterized in appendix C.

The stakeholders assigned to a specific communication means are based on the audience addressed by the original authors of such methodologies. In some cases, the characterization of the intended audience was not as granular as our stakeholder mapping and the authors merely differed experts in ML from non-experts. Based on the nature of knowledge that we assigned to each of the mentioned stakeholders in section 4.1, we considered that the development and auditing teams are able to understand technically formulated system details (experts) whereas data domain experts and decision subjects would require more accessible communication means (non-experts). Similarly, some of the communication means identified for *explainability* are suitable for any stakeholder, but the original authors formulated the post-hoc explanations with varying degrees of complexity, which should be taken into account when trying to deploy such strategies. If the target audience are data subjects, we echo van Berkel et al. [175] and Cheng et al. [35] and recommend to limit presentation complexity and to instruct participants throughout the session.

It should be noted that this mapping process represents a first step to making a wide range of stakeholders with different backgrounds understand each other. We are aware that communicating system-related information in a tailored way does not directly lead to the resolution of value trade-offs, and that design strategies are necessary for facilitating such conversations [81]. In any case, the exercise of resolving value tensions should be a communicative process, rather than a simple explanation [140]. However, the means used for communicating specifications of the system will play a key role in the dynamics that will take place in those sessions.

Assigning communication means to each stakeholder in our example use case. The life insurance researchers are now looking into appropriate methods for communicating values to different stakeholders (fig 2), so that they can develop a comprehensive plan that ensures both compliance and communication of values.

Based on the mapping presented in appendix C, tables 5 and 6, for the value of *explainability* and its manifestations in the form of post-hoc explanations, the team can use various design strategies and tools as part of their assessment process. To facilitate the navigation of table 6, they first examine table 5 to locate the type of means (tool, strategy, or documentation), values and stakeholders they are interested in. Once they select the codes associated to each communication means, they check table 6 to see whether the value manifestations in question are addressed and to explore the selected communication means. If the team working on the life insurance case prefers a ready-to-use tool over the description of design strategies

677 for assessing *explainability*, they can use InterpretML [137] and especially the DiCE [134] functionality, (code [AC]) with
678 the development and auditing teams to evaluate counterfactual examples. These counterfactual examples tell how input
679 features should change in order for the output of the system to be different. That is to say, how the individual applying
680 for life insurance should be different, physically, or when it comes to insurance or medical history, for them to accept
681 the application (if the original output was a refusal). However, this tool might not be suitable for non-experts who are
682 not familiar with ML-related concepts. In the life insurance use case, the medical and insurance team and the decision
683 subjects should receive a description of how the output changes if a feature is perturbed, absent or present adapted to
684 their insight needs. This can be done by describing the answers to the questions “Why, Why not and How to be that”
685 for a certain output [117] (code [P]). As for *privacy* manifestations, the development and auditing teams can examine
686 data collection and storage specifications through the Datasheet [69] associated to the dataset in question (code [K]).
687 Special attention should be paid to the “Collection” and “Preprocessing/cleaning/labelling” sections. For decision subjects,
688 iconsets [56, 89, 125, 149] (code [A]) and privacy dashboards [54, 58, 59, 85, 191] (code [B]) are means for them to explore
689 how their data is being used. It should be noted that the cell that intersects between data domain experts and *privacy*
690 in table 5 is blank. Based on the characterization of stakeholders that we provided, data privacy-related matters are
691 not directly linked to the purpose that data domain experts show when willing to explore algorithmic systems. This is
692 translated into scarcity of methodologies related to *privacy* manifestations that directly address data domain experts.

698 5 DISCUSSION AND FUTURE WORK

699 We discuss important aspects of our framework below.

700 *Design choices for creating a value-based framework.* We aim at examining values that characterize ML-driven systems
701 rather than the organizations responsible for these systems. Hence, we did not integrate accountability or responsibility
702 as a value *per se* in our framework. We are aware that algorithms cannot be held responsible for the potential harm that
703 they might cause [30, 86], and that in order to effectively deploy such systems, there is an urgent call for accountability
704 [6, 186]. Likewise, we are aware of the need for rigorous frameworks that support accountability [92] and we consider
705 that the act of conceiving an assessment framework itself answers to the need to evaluate and audit algorithmic systems
706 [60]. Nevertheless, we did not explicitly highlight the profiles of the people accountable for the system. We decided to
707 follow Zhu et al. [190] and considered accountability as a governance issue. We do, however, believe that entities up the
708 chain of command should be held accountable for the potential harm caused by algorithmic systems [3, 60, 86, 121, 154]. It
709 should also be noted that values and criteria presented in this paper might not be unique [159]. We acknowledge current
710 discussions in VSD about the shortcomings of pre-selecting values [45] and, hence, do not claim universality. Extension
711 and modification of values is possible in our framework, but are subject to respecting continuity and opposition between
712 values. Similarly, criteria and manifestations can be extended and subsets could be included to create situationally-specific
713 versions of the framework. Since the aim of our framework is to encourage critical reflection [63, 173] and we identified
714 some value manifestations that require additional communication means, we particularly encourage those context-specific
715 adaptations to happen. Under no circumstances should the scarcity of communication means for certain values identified
716 in our framework represent an excuse to justify inaction or to ignore such values.

717 *Context dependence and consistency.* As echoed by Liscio et al. [119], in order to translate values into system require-
718 ments [145, 176], to reason about conflicting values [5, 135] and to communicate them to different stakeholders [64], it
719 is necessary to situate these values within a context. The prioritization of values depends on the application context of
720 such systems [3]. In this paper, we showed an example of how the framework could be applied to a particular use case.
721 However, considering the differences between value alignments and tensions that may arise due to context dependence,
722

729 the validity and consistency of our framework is still to be tested. Future work needs to validate our framework across
730 scenarios [72, 175] through user studies or synthetic experiments [162].

731 *Need for standardization.* To systematically review and revisit value priorities and tensions among different stakeholders,
732 our framework should be part of a broader evaluation workflow [115], such as the one suggested by Raji and Smart
733 [147]. Besides, practices from software engineering such as the Values Dashboard [138] could be adopted [160, 169]. This
734 dashboard promotes awareness of values and aims at triggering discussions among stakeholders. It claims to be beneficial
735 in each phase of the software development process, from inception to release, and establishes strategies, such as Timelines
736 or Issues, that are already common practice on software development platforms like Github.

737 *Implications of our work.* Our multi-stakeholder value-based framework facilitates the unveiling of assumptions that
738 encode political and social values made by developers [147]. By bringing together a wide range of stakeholders to evaluate
739 and discuss value manifestations, one can anticipate and remedy harmful algorithmic behaviors before deploying a system.
740 Besides, we provide researchers and industry practitioners with a good coverage of values to evaluate their systems
741 and the association of such values to actionable value manifestations. This contributes greatly to the adoption of ethical
742 approaches by practically-minded people [133]. For researchers, we provide them with an easy-to-navigate mapping of
743 value manifestations, stakeholders and communication means. Our framework also visually illustrates research gaps that
744 need to be addressed. Blank spaces in appendix C or values with a scarce number of associated communication means
745 directly refer to valuable research opportunities. For instance, for the value of *fairness*, a great deal of effort has been
746 devoted to designing ready-to-use tools for stakeholders with a deep understanding of ML (developers and auditing
747 teams). However, means for addressing *fairness* manifestations and communicating them to decision subjects have not
748 received the same attention. For industry practitioners, we gathered ready-to-use open source toolkits (appendix C) that
749 can be directly applied to their own use cases. Moreover, since we host this mapping on an online repository ¹¹ open
750 to future contributions, we hope that the number of tools addressing each of the identified value manifestations will grow
751 and that the benefits of designing such a framework will be even more tangible in the future.

752 6 CONCLUSIONS

753 In this paper, we designed a value-based framework for assessing algorithmic systems from a multi-stakeholder perspec-
754 tive. This provides investigators of algorithmic systems with an actionable set of criteria manifestations to operationalize
755 high-level ethical principles. We arranged eleven prominent values of ML-driven systems in a circular composition, so
756 that common motivations and trade-offs can be easily identified.

757 We then broke down each of these values into a set of criteria and their correspondent manifestations in the form of
758 quantifiable indicators, process-oriented practices, and signifiers. In addition, we examined available tools for communi-
759 cating those value manifestations to different stakeholders based on the nature of their knowledge and their insight needs.
760 This should enable to bring a wide range of stakeholders together to systematically assess values encoded in a system
761 and facilitate value- and ethics-related discussions among them. This work completes previous studies that claim the
762 need for incorporating a diverse range of stakeholders and viewpoints in the ML workflow, so that conflicting priorities
763 and value tensions can be reviewed, negotiated and consolidated.

764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780

¹¹Check footnote 10

ACKNOWLEDGMENTS

REFERENCES

- [1] Rediet Abebe, Solon Barocas, Jon Kleinberg, Karen Levy, Manish Raghavan, and David G. Robinson. 2020. Roles for computing in social change. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*. ACM, New York, NY, USA, 252–260. <https://doi.org/10.1145/3351095.3372871>
- [2] Philip Adler, Casey Falk, Sorelle A Friedler, Tionney Nix, Gabriel Rybeck, Carlos Scheidegger, Brandon Smith, and Suresh Venkatasubramanian. 2018. Auditing black-box models for indirect influence. *Knowledge and Information Systems* 54, 1 (2018), 95–122. <https://doi.org/10.1007/s10115-017-1116-3>
- [3] AI Ethics Impact Group (AIEIG). 2020. From Principles to Practice An interdisciplinary framework to operationalise AI ethics. <https://www.ai-ethics-impact.org/resource/blob/1961130/c6db9894ee73aefa489d6249f5ee2b9f/aieig---report---download-hb-data.pdf>
- [4] Evgeni Aizenberg and Jeroen van den Hoven. 2020. Designing for human rights in AI. *Big Data & Society* 7, 2 (7 2020), 2053951720949566. <https://doi.org/10.1177/2053951720949566>
- [5] Nirav Ajmeri, Hui Guo, Pradeep K Murukannaiah, and Munindar P Singh. 2020. Elessar: Ethics in Norm-Aware Agents. In *Proceedings of the 19th International Conference on Autonomous Agents and MultiAgent Systems (AAMAS '20)*. International Foundation for Autonomous Agents and Multiagent Systems, Richland, SC, 16–24.
- [6] Kars Alfrink, T. Turel, A. I. Keller, N. Doorn, and G. W. Kortuem. 2020. Contestable City Algorithms. International Conference on Machine Learning Workshop.
- [7] Ahmed Alqaraawi, Martin Schuessler, Philipp Weiß, Enrico Costanza, and Nadia Berthouze. 2020. Evaluating saliency map explanations for convolutional neural networks. In *Proceedings of the 25th International Conference on Intelligent User Interfaces*. ACM, New York, NY, USA. <https://doi.org/10.1145/3377325.3377519>
- [8] Thayer Alshaabi, David Rushing Dewhurst, Joshua R Minot, Michael V Arnold, Jane L Adams, Christopher M Danforth, and Peter Sheridan Dodds. 2021. The growing amplification of social media: measuring temporal and social contagion dynamics for over 150 languages on Twitter for 2009–2020. *EPJ Data Science* 10, 1 (2021), 15. <https://doi.org/10.1140/epjds/s13688-021-00271-0>
- [9] Saleema Amershi, Max Chickering, Steven M. Drucker, Bongshin Lee, Patrice Simard, and Jina Suh. 2015. ModelTracker. In *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems*. ACM, New York, NY, USA. <https://doi.org/10.1145/2702123.2702509>
- [10] Alexander Amini, Ava P Soleimany, Wilko Schwarting, Sangeeta N Bhatia, and Daniela Rus. 2019. Uncovering and Mitigating Algorithmic Bias through Learned Latent Structure. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society (AIES '19)*. Association for Computing Machinery, New York, NY, USA, 289–295. <https://doi.org/10.1145/3306618.3314243>
- [11] Access Now Amnesty International. 2018. Toronto Declaration: Protecting the Right to Equality and Non-Discrimination in Machine Learning Systems. https://www.accessnow.org/cms/assets/uploads/2018/08/The-Toronto-Declaration_ENG_08-2018.pdf
- [12] Ariful Islam Anik and Andrea Bunt. 2021. Data-Centric Explanations: Explaining Training Data of Machine Learning Systems to Promote Transparency. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. ACM, New York, NY, USA. <https://doi.org/10.1145/3411764.3445736>
- [13] Sherry R Arnstein. 2019. A Ladder of Citizen Participation. *Journal of the American Planning Association* 85, 1 (1 2019), 24–34. <https://doi.org/10.1080/01944363.2018.1559388>
- [14] Mahmoud Assran, Joshua Romoff, Nicolas Ballas, Joelle Pineau, and Michael Rabbat. 2019. Gossip-based Actor-Learner Architectures for Deep Reinforcement Learning. (6 2019).
- [15] Chloé Bakalar, Renata Barreto, Stevie Bergman, Miranda Bogen, Bobbie Chern, Sam Corbett-Davies, Melissa Hall, Isabel Kloumann, Michelle Lam, Joaquin Quiñero Candela, Manish Raghavan, Joshua Simons, Jonathan Tannen, Edmund Tong, Kate Vredenburg, and Jiejing Zhao. 2021. Fairness On The Ground: Applying Algorithmic Fairness Approaches to Production Systems. (3 2021).
- [16] Agathe Balayn and Seda Gürses. 2021. Beyond Debiasing: Regulating AI and its inequalities. <https://edri.org/our-work/if-ai-is-the-problem-is-debiasing-the-solution/>
- [17] Ari Ball-Burack, Michelle Seng Ah Lee, Jennifer Cobbe, and Jatinder Singh. 2021. Differential Tweetment: Mitigating Racial Dialect Bias in Harmful Tweet Detection. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (FAccT '21)*. Association for Computing Machinery, New York, NY, USA, 116–128. <https://doi.org/10.1145/3442188.3445875>
- [18] Alejandro Barredo Arrieta, Natalia Díaz-Rodríguez, Javier Del Ser, Adrien Bannetot, Siham Tabik, Alberto Barbado, Salvador Garcia, Sergio Gil-Lopez, Daniel Molina, Richard Benjamins, Raja Chatila, and Francisco Herrera. 2020. Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. *Information Fusion* 58 (6 2020), 82–115. <https://doi.org/10.1016/J.INFFUS.2019.12.012>
- [19] R K E Bellamy, K Dey, M Hind, S C Hoffman, S Houde, K Kannan, P Lohia, J Martino, S Mehta, A Mojsilović, S Nagar, K Natesan Ramamurthy, J Richards, D Saha, P Sattigeri, M Singh, K R Varshney, and Y Zhang. 2019. AI Fairness 360: An extensible toolkit for detecting and mitigating algorithmic bias. *IBM Journal of Research and Development* 63, 4/5 (2019), 1–4. <https://doi.org/10.1147/JRD.2019.2942287>
- [20] Emily M. Bender and Batya Friedman. 2018. Data Statements for Natural Language Processing: Toward Mitigating System Bias and Enabling Better Science. *Transactions of the Association for Computational Linguistics* 6 (12 2018). https://doi.org/10.1162/tacl-1ja_00041
- [21] Emily M Bender, Timmit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (FAccT '21)*. Association for Computing Machinery, New York, NY, USA, 610–623. <https://doi.org/10.1145/3442188.3445922>

- [22] Richard Berk, Hoda Heidari, Shahin Jabbari, Michael Kearns, and Aaron Roth. 2017. Fairness in Criminal Justice Risk Assessments: The State of the Art. (3 2017).
- [23] Battista Biggio and Fabio Roli. 2018. Wild patterns: Ten years after the rise of adversarial machine learning. *Pattern Recognition* 84 (12 2018), 317–331. <https://doi.org/10.1016/j.patcog.2018.07.023>
- [24] Reuben Binns, Max Van Kleek, Michael Veale, Ulrik Lyngs, Jun Zhao, and Nigel Shadbolt. 2018. 'It's Reducing a Human Being to a Percentage'; Perceptions of Justice in Algorithmic Decisions. (1 2018). <https://doi.org/10.1145/3173574.3173951>
- [25] Sarah Bird, Miro Dudik, Richard Edgar, Brandon Horn, Roman Lutz, Vanessa Milan, Mehrnoosh Sameki, Hanna Wallach, and Kathleen Walker. 2020. *Fairlearn: A toolkit for assessing and improving fairness in AI*. Technical Report MSR-TR-2020-32. Microsoft. <https://www.microsoft.com/en-us/research/publication/fairlearn-a-toolkit-for-assessing-and-improving-fairness-in-ai/>
- [26] Abeba Birhane, Pratyusha Kalluri, Dallas Card, William Agnew, Ravit Dotan, and Michelle Bao. 2021. The Values Encoded in Machine Learning Research. (6 2021).
- [27] Alice Namuli Blazevic, Patrick Mugalula, and Andrew Wandera. 2021. Towards Operationalizing the Data Protection and Privacy Act 2020: Understanding the Draft Data Protection and Privacy Regulations, 2020. *SSRN Electronic Journal* (2021). <https://doi.org/10.2139/ssrn.3776353>
- [28] Su Lin Blodgett, Solon Barocas, Hal Daumé, and Hanna Wallach. 2020. Language (Technology) is Power: A Critical Survey of "Bias" in NLP. (5 2020).
- [29] Dimitrios Bountouridis, Jaron Harambam, Mykola Makhortykh, Mónica Marrero, Nava Tintarev, and Claudia Hauff. 2019. SIREN: A Simulation Framework for Understanding the Effects of Recommender Systems in Online News Environments. In *Proceedings of the Conference on Fairness, Accountability, and Transparency (FAT* '19)*. Association for Computing Machinery, New York, NY, USA, 150–159. <https://doi.org/10.1145/3287560.3287583>
- [30] Joanna J. Bryson, Mihailis E. Diamantis, and Thomas D. Grant. 2017. Of, for, and by the people: the legal lacuna of synthetic persons. *Artificial Intelligence and Law* 25, 3 (9 2017), 273–291. <https://doi.org/10.1007/s10506-017-9214-9>
- [31] Joy Buolamwini and Timnit Gebru. 2018. Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification. In *Proceedings of the 1st Conference on Fairness, Accountability and Transparency (Proceedings of Machine Learning Research, Vol. 81)*, Sorelle A Friedler and Christo Wilson (Eds.). PMLR, 77–91. <https://proceedings.mlr.press/v81/buolamwini18a.html>
- [32] Carrie J. Cai, Jonas Jongejan, and Jess Holbrook. 2019. The effects of example-based explanations in a machine learning interface. In *Proceedings of the 24th International Conference on Intelligent User Interfaces*. ACM, New York, NY, USA. <https://doi.org/10.1145/3301275.3302289>
- [33] Simeon C Calvert, Daniël D Heikoop, Giulio Mecacci, and Bart Van Areem. 2019. A human centric framework for the analysis of automated driving systems based on meaningful human control. *Theoretical Issues in Ergonomics Science* 21, 4 (2019), 478–506. <https://doi.org/10.1080/1463922X.2019.1697390>
- [34] Kyla Chasalow and Karen Levy. 2021. Representativeness in Statistics, Politics, and Machine Learning. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (FAccT '21)*. Association for Computing Machinery, New York, NY, USA, 77–89. <https://doi.org/10.1145/3442188.3445872>
- [35] Hao-Fei Cheng, Ruotong Wang, Zheng Zhang, Fiona O'Connell, Terrance Gray, F Maxwell Harper, and Haiyi Zhu. 2019. Explaining Decision-Making Algorithms through UI: Strategies to Help Non-Expert Stakeholders. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems (CHI '19)*. Association for Computing Machinery, New York, NY, USA, 1–12. <https://doi.org/10.1145/3290605.3300789>
- [36] China Electronics Standardization Institute. 2018. Original CSET Translation of "Artificial Intelligence Standardization White Paper". <https://cset.georgetown.edu/research/artificial-intelligence-standardization-white-paper/>
- [37] Alexandra Chouldechova. 2016. Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. (10 2016).
- [38] Nazli Cila, Gabriele Ferri, Martijn de Waal, Inte Gloerich, and Tara Karpinski. 2020. The Blockchain and the Commons: Dilemmas in the Design of Local Platforms. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*. Association for Computing Machinery, New York, NY, USA, 1–14. <https://doi-org.tudelft.idm.oclc.org/10.1145/3313831.3376660>
- [39] Kate Crawford and Trevor Paglen. 2019. Excavating AI: The Politics of Training Sets for Machine Learning.
- [40] Gabriela F. Cretu, Angelos Stavrou, Michael E. Locasto, Salvatore J. Stolfo, and Angelos D. Keromytis. 2008. Casting out Demons: Sanitizing Training Data for Anomaly Sensors. In *2008 IEEE Symposium on Security and Privacy (sp 2008)*. IEEE, 81–95. <https://doi.org/10.1109/SP.2008.11>
- [41] Steven Dalton, Iuri Frosio, and Michael Garland. 2019. Accelerating Reinforcement Learning through GPU Atari Emulation. (7 2019).
- [42] Alexander D'Amour, Hansa Srinivasan, James Atwood, Pallavi Baljekar, D Sculley, and Yoni Halpern. 2020. Fairness is Not Static: Deeper Understanding of Long Term Fairness via Simulation Studies. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency (FAT* '20)*. Association for Computing Machinery, New York, NY, USA, 525–534. <https://doi.org/10.1145/3351095.3372878>
- [43] Dasha Simons. 2019. *Design for fairness in AI: Cooking a fair AI Dish*. Technical Report. Delft University of Technology. Graduation project. MSc in Strategic Product Design. <http://resolver.tudelft.nl/uuid:5a116c17-ce0a-4236-b283-da6b8545628c>
- [44] Thomas Davidson, Debasmita Bhattacharya, and Ingmar Weber. 2019. Racial Bias in Hate Speech and Abusive Language Detection Datasets. In *Proceedings of the Third Workshop on Abusive Language Online*. Association for Computational Linguistics, Stroudsburg, PA, USA. <https://doi.org/10.18653/v1/W19-3504>
- [45] Janet Davis and Lisa P. Nathan. 2015. Value Sensitive Design: Applications, Adaptations, and Critiques. In *Handbook of Ethics, Values, and Technological Design*. Springer Netherlands, Dordrecht, 11–40. https://doi.org/10.1007/978-94-007-6970-0_3
- [46] Emily Denton, Alex Hanna, Razvan Amironesei, Andrew Smart, Hilary Nicole, and Morgan Klaus Scheuerman. 2020. Bringing the People Back In: Contesting Benchmark Machine Learning Datasets. (7 2020). <https://arxiv.org/abs/2007.07399>
- [47] Amit Dhurandhar, Pin-Yu Chen, Ronny Luss, Chun-Chen Tu, Paishun Ting, Karthikeyan Shanmugam, and Payel Das. 2018. Explanations based on the Missing: Towards Contrastive Explanations with Pertinent Negatives. (2 2018).

- 885 [48] Lucas Dixon, John Li, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. 2018. Measuring and Mitigating Unintended Bias in Text Classification.
886 In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*. ACM, New York, NY, USA, 67–73. <https://doi.org/10.1145/3278721.3278729>
- 887 [49] Jonathan Dodge, Q. Vera Liao, Yunfeng Zhang, Rachel K. E. Bellamy, and Casey Dugan. 2019. Explaining Models: An Empirical Study of How
888 Explanations Impact Fairness Judgment. (1 2019). <https://doi.org/10.1145/3301275.3302310>
- 889 [50] Ravit Dotan and Smitha Milli. 2019. Value-laden Disciplinary Shifts in Machine Learning. (12 2019).
- 890 [51] Tim Draws, Alisa Rieger, Oana Inel, Ujwal Gadiraju, and Nava Tintarev. 2021. A Checklist to Combat Cognitive Biases in Crowdsourcing. *Proceedings*
891 *of the AAAI Conference on Human Computation and Crowdsourcing* 9, 1 (10 2021), 48–59. <https://ojs.aaai.org/index.php/HCOMP/article/view/18939>
- 892 [52] Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Rich Zemel. 2011. Fairness Through Awareness. (4 2011).
- 893 [53] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. 2006. Calibrating Noise to Sensitivity in Private Data Analysis. 265–284.
894 https://doi.org/10.1007/11681878_14
- 895 [54] Julia Earp and Jessica Staddon. 2016. "I had no idea this was a thing". In *Proceedings of the 6th Workshop on Socio-Technical Aspects in Security and*
896 *Trust*. ACM, New York, NY, USA, 79–86. <https://doi.org/10.1145/3046055.3046062>
- 897 [55] Bora Edizel, Francesco Bonchi, Sara Hajian, André Panisson, and Tamir Tassa. 2020. FaiRecSys: mitigating algorithmic bias in recommender systems.
898 *International Journal of Data Science and Analytics* 9, 2 (2020), 197–213. <https://doi.org/10.1007/s41060-019-00181-5>
- 899 [56] European Commission. 2018. 2018 reform of EU data protection rules. https://ec.europa.eu/commission/sites/beta-political/files/data-protection-factsheet-changes_en.pdf
- 900 [57] European Commission. 2019. Ethics guidelines for trustworthy AI. <https://www.aepd.es/sites/default/files/2019-12/ai-ethics-guidelines.pdf>
- 901 [58] Florian M. Farke, David G. Balash, Maximilian Golla, Markus Dürmuth, and Adam J. Aviv. 2021. Are Privacy Dashboards Good for End Users?
902 Evaluating User Perceptions and Reactions to Google's My Activity (Extended Version). (5 2021).
- 903 [59] Simone Fischer-Hübner, Julio Angulo, Farzaneh Karegar, and Tobias Pulls. 2016. Transparency, Privacy and Trust – Technology for Tracking and
904 Controlling My Data Disclosures: Does This Work? 3–14. https://doi.org/10.1007/978-3-319-41354-9_1
- 905 [60] Jessica Fjeld, Nele Achten, Hannah Hilligoss, Adam Nagy, and Madhulika Srikumar. 2020. Principled Artificial Intelligence: Mapping Consensus
906 in Ethical and Rights-Based Approaches to Principles for AI. *SSRN Electronic Journal* (2020). <https://doi.org/10.2139/ssrn.3518482>
- 907 [61] Luciano Floridi. 2019. Translating Principles into Practices of Digital Ethics: Five Risks of Being Unethical. *Philosophy & Technology* 32, 2 (6 2019).
908 <https://doi.org/10.1007/s13347-019-00354-x>
- 909 [62] Luciano Floridi, Josh Cowsls, Monica Beltrametti, Raja Chatila, Patrice Chazerand, Virginia Dignum, Christoph Luetge, Robert Madelin, Ugo Pagallo,
910 Francesca Rossi, Burkhard Schafer, Peggy Valcke, and Effy Vayena. 2018. AI4People—An Ethical Framework for a Good AI Society: Opportunities,
911 Risks, Principles, and Recommendations. *Minds and Machines* 28, 4 (12 2018). <https://doi.org/10.1007/s11023-018-9482-5>
- 912 [63] Christopher Frauenberger, Marjo Rauhala, and Geraldine Fitzpatrick. 2016. In-Action Ethics: Table 1. *Interacting with Computers* (6 2016).
913 <https://doi.org/10.1093/iwc/iww024>
- 914 [64] W. Fred van Raaij and Theo M.M. Verhallen. 1994. Domain-specific Market Segmentation. *European Journal of Marketing* 28, 10 (10 1994), 49–66.
915 <https://doi.org/10.1108/03090569410075786>
- 916 [65] Batya Friedman, David G. Hendry, and Alan Borning. 2017. A Survey of Value Sensitive Design Methods. *Foundations and Trends® in*
917 *Human–Computer Interaction* 11, 2 (2017), 63–125. <https://doi.org/10.1561/11000000015>
- 918 [66] Georges Gaillard. 2016. La conflictualité : une modalité de lien où s'arrime la destructivité humaine. *Connexions* 106, 2 (2016), 71.
919 <https://doi.org/10.3917/cnx.106.0071>
- 920 [67] Yanjie Gao, Yu Liu, Hongyu Zhang, Zhengxian Li, Yonghao Zhu, Haoxiang Lin, and Mao Yang. 2020. Estimating GPU memory consumption of
921 deep learning models. In *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations*
922 *of Software Engineering*. ACM, New York, NY, USA, 1342–1352. <https://doi.org/10.1145/3368089.3417050>
- 923 [68] Eva García-Martín, Crefeda Faviola Rodrigues, Graham Riley, and Håkan Grah. 2019. Estimation of energy consumption in machine learning.
924 *J. Parallel and Distrib. Comput.* 134 (12 2019), 75–88. <https://doi.org/10.1016/j.jpdc.2019.07.007>
- 925 [69] Timmit Gebru, Google Jamie Morgenstern, Briana Vecchione, and Jennifer Wortman Vaughan. 2020. Datasheets for Datasets.
- 926 [70] R Stuart Geiger, Kevin Yu, Yanlai Yang, Mindy Dai, Jie Qiu, Rebekah Tang, and Jenny Huang. 2020. Garbage in, Garbage out? Do Machine
927 Learning Application Papers in Social Computing Report Where Human-Labeled Training Data Comes From?. In *Proceedings of the 2020*
928 *Conference on Fairness, Accountability, and Transparency (FAT* '20)*. Association for Computing Machinery, New York, NY, USA, 325–336.
929 <https://doi.org/10.1145/3351095.3372862>
- 930 [71] Bhavya Ghai, Q. Vera Liao, Yunfeng Zhang, and Klaus Mueller. 2020. Measuring Social Biases of Crowd Workers using Counterfactual Queries. (4 2020).
- 931 [72] Marzyeh Ghassemi, Luke Oakden-Rayner, and Andrew L Beam. 2021. The false hope of current approaches to explainable artificial intelligence
932 in health care. *The Lancet Digital Health* 3, 11 (11 2021), e745–e750. [https://doi.org/10.1016/S2589-7500\(21\)00208-9](https://doi.org/10.1016/S2589-7500(21)00208-9)
- 933 [73] Amir Globerson and Sam Roweis. 2006. Nightmare at test time. In *Proceedings of the 23rd international conference on Machine learning - ICML '06*.
934 ACM Press, New York, New York, USA, 353–360. <https://doi.org/10.1145/1143844.1143889>
- 935 [74] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. 2014. Explaining and Harnessing Adversarial Examples. (12 2014).
- 936 [75] Google. 2018. AI at Google: Our Principles. <https://www.blog.google/technology/ai/ai-principles/>
- [76] Ben Green and Lily Hu. 2018. The Myth in the Methodology: Towards a Recontextualization of Fairness in Machine Learning. In *Machine Learning: The Debates workshop at the 35th International Conference on Machine Learning (ICML)*. Stockholm, Sweden.

- [77] Nina Grgic-Hlaca, Muhammad Bilal Zafar, Krishna P Gummadi, and Adrian Weller. 2018. Beyond Distributive Fairness in Algorithmic Decision Making: Feature Selection for Procedurally Fair Learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 32.
- [78] Christopher Groves. 2015. Logic of Choice or Logic of Care? Uncertainty, Technological Mediation and Responsible Innovation. *NanoEthics* 9, 3 (12 2015), 321–333. <https://doi.org/10.1007/s11569-015-0238-x>
- [79] Moritz Hardt, Eric Price, and Nathan Srebro. 2016. Equality of Opportunity in Supervised Learning. In *Proceedings of the 30th International Conference on Neural Information Processing Systems (NIPS'16)*. Curran Associates Inc., Red Hook, NY, USA, 3323–3331.
- [80] Galen Harrison, Julia Hanson, Christine Jacinto, Julio Ramirez, and Blase Ur. 2020. An empirical study on the perceived fairness of realistic, imperfect machine learning models. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*. ACM, New York, NY, USA, 392–402. <https://doi.org/10.1145/3351095.3372831>
- [81] Katrina Heijne and Han van der Meer. 2019. *Road Map for Creative Problem Solving Techniques Organizing and facilitating group sessions*. Boom Uitgevers Amsterdam.
- [82] Drew Hemment, Ruth Aylett, Vaishak Belle, Dave Murray-Rust, Ewa Luger, Jane Hillston, Michael Rovatsos, and Frank Broz. 2019. Experiential AI. *AI Matters* 5, 1 (4 2019), 25–31. <https://doi.org/10.1145/3320254.3320264>
- [83] Peter Henderson, Jieru Hu, Joshua Romoff, Emma Brunskill, Dan Jurafsky, and Joelle Pineau. 2020. Towards the Systematic Reporting of the Energy and Carbon Footprints of Machine Learning. (1 2020).
- [84] Clément Henin and Daniel Le Métayer. 2021. Beyond explainability: justifiability and contestability of algorithmic decision systems. *AI & SOCIETY* (7 2021). <https://doi.org/10.1007/s00146-021-01251-8>
- [85] Eelco Herder and Olaf van Maaren. 2020. Privacy Dashboards: The Impact of the Type of Personal Data and User Control on Trust and Perceived Risk. In *Adjunct Publication of the 28th ACM Conference on User Modeling, Adaptation and Personalization*. ACM, New York, NY, USA, 169–174. <https://doi.org/10.1145/3386392.3399557>
- [86] César Hidalgo, Diana Orghian, Jordi Albo-Canals, Filipa de Almeida, and Natalia Martin. 2021. *How Humans Judge Machines*. MIT Press. <https://hal.archives-ouvertes.fr/hal-03058652>
- [87] Tad Hirsch, Kritzia Merced, Shrikanth Narayanan, Zac E. Imel, and David C. Atkins. 2017. Designing Contestability. In *Proceedings of the 2017 Conference on Designing Interactive Systems*. ACM, New York, NY, USA. <https://doi.org/10.1145/3064663.3064703>
- [88] Sarah Holland, Ahmed Hosny, Sarah Newman, Joshua Joseph, and Kasia Chmielinski. 2018. The Dataset Nutrition Label: A Framework To Drive Higher Data Quality Standards. (5 2018).
- [89] Leif-Erik Holtz, Katharina Nocun, and Marit Hansen. 2011. Towards Displaying Privacy Information with Icons. 338–348. https://doi.org/10.1007/978-3-642-20769-3_27
- [90] Matthew K. Hong, Adam Fournery, Derek DeBellis, and Saleema Amershi. 2021. Planning for Natural Language Failures with the AI Playbook. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. ACM, New York, NY, USA, 1–11. <https://doi.org/10.1145/3411764.3445735>
- [91] Christoph Hube, Besnik Fetahu, and Ujwal Gadiraju. 2019. Understanding and Mitigating Worker Biases in the Crowdsourced Collection of Subjective Judgments. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*. Association for Computing Machinery, New York, NY, USA, 1–12. <https://doi-org.tudelft.idm.oclc.org/10.1145/3290605.3300637>
- [92] Ben Hutchinson, Andrew Smart, Alex Hanna, Emily Denton, Christina Greer, Oddur Kjartansson, Parker Barnes, and Margaret Mitchell. 2021. Towards Accountability for Machine Learning Datasets: Practices from Software Engineering and Infrastructure. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (FAccT '21)*. Association for Computing Machinery, New York, NY, USA, 560–575. <https://doi.org/10.1145/3442188.3445918>
- [93] IBM. 2019. IBM Everyday Ethics for AI. <https://www.ibm.com/watson/assets/duo/pdf/everydayethics.pdf>
- [94] IEEE. 2008. IEEE Standard for Software Reviews and Audits. *IEEE Std 1028-2008* (2008), 1–53. <https://doi.org/10.1109/IEEESTD.2008.4601584>
- [95] Stefania Ionescu, Anikó Hannák, and Kenneth Joseph. 2021. An Agent-Based Model to Evaluate Interventions on Online Dating Platforms to Decrease Racial Homogamy. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (FAccT '21)*. Association for Computing Machinery, New York, NY, USA, 412–423. <https://doi.org/10.1145/3442188.3445904>
- [96] Technology Japanese Cabinet Office, Council for Science and Innovation. 2019. Social Principles of Human-Centric Artificial Intelligence. <https://www8.cao.go.jp/cstp/english/humancentricai.pdf>
- [97] Jinyuan Jia, Ahmed Salem, Michael Backes, Yang Zhang, and Neil Zhenqiang Gong. 2019. MemGuard: Defending against Black-Box Membership Inference Attacks via Adversarial Examples. (9 2019).
- [98] Weina Jin, Jianyu Fan, Diane Gromala, Philippe Pasquier, and Ghassan Hamarneh. 2021. EUCA: A Practical Prototyping Framework towards End-User-Centered Explainable Artificial Intelligence. (2 2021). <https://arxiv.org/abs/2102.02437>
- [99] Jonas Kaiser and Adrian Rauchfleisch. 2020. Birds of a Feather Get Recommended Together: Algorithmic Homophily in YouTube’s Channel Recommendations in the United States and Germany. *Social Media + Society* 6, 4 (10 2020), 2056305120969914. <https://doi.org/10.1177/2056305120969914>
- [100] Pratyusha Kalluri. 2020. Don’t ask if artificial intelligence is good or fair, ask how it shifts power. *Nature* 583, 7815 (2020). <https://doi.org/10.1038/d41586-020-02003-2>
- [101] Yigitcan Kaya, Sanghyun Hong, and Tudor Dumitras. 2020. On the Effectiveness of Regularization Against Membership Inference Attacks. (6 2020).
- [102] Michael Kearns, Seth Neel, Aaron Roth, and Zhiwei Steven Wu. 2018. Preventing Fairness Gerrymandering: Auditing and Learning for Subgroup Fairness. In *Proceedings of the 35th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 80)*, Jennifer Dy and Andreas Krause (Eds.). PMLR, 2564–2572. <https://proceedings.mlr.press/v80/kearns18a.html>

- 989 [103] Michael Kearns and Aaron Roth. 2019. *The Ethical Algorithm: The Science of Socially Aware Algorithm Design*. Oxford University Press, Inc., USA.
- 990 [104] Os Keyes, Jevan Hutson, and Meredith Durbin. 2019. A Mulching Proposal: Analysing and Improving an Algorithmic System for Turning the
991 Elderly into High-Nutrient Slurry. In *Extended Abstracts of the 2019 CHI Conference on Human Factors in Computing Systems (CHI EA '19)*. Association
992 for Computing Machinery, New York, NY, USA, 1–11. <https://doi.org/10.1145/3290607.3310433>
- 993 [105] Jon Kleinberg, Sendhil Mullainathan, and Manish Raghavan. 2016. Inherent Trade-Offs in the Fair Determination of Risk Scores. (9 2016).
- 994 [106] Daniel Kluttz, Nitin Kohli, and Deirdre K. Mulligan. 2018. Contestability and Professionals: From Explanations to Engagement with Algorithmic
995 Systems. *SSRN Electronic Journal* (2018). <https://doi.org/10.2139/ssrn.3311894>
- 996 [107] TD Krafft and K Zweig. 2019. Transparenz und Nachvollziehbarkeit algorithmenbasierter Entscheidungsprozesse. *Ein Regulierungsvorschlag* (2019).
- 997 [108] Lenneke Kuijjer and Elisa Giaccardi. 2018. Co-Performance: Conceptualizing the Role of Artificial Agency in the Design of Everyday Life. In
998 *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*. Association for Computing Machinery, New York, NY, USA, 1–13.
<https://doi-org.tudelft.idm.oclc.org/10.1145/3173574.3173699>
- 999 [109] Bogdan Kulynych, Rebekah Overdorf, Carmela Troncoso, and Seda Gürses. 2018. POTs: Protective Optimization Technologies. (6 2018).
1000 <https://doi.org/10.1145/3351095.3372853>
- 1001 [110] Matt J Kusner, Joshua Loftus, Chris Russell, and Ricardo Silva. 2017. Counterfactual Fairness. In *Advances in Neural Information Processing*
1002 *Systems*, I Guyon, U V Luxburg, S Bengio, H Wallach, R Fergus, S Vishwanathan, and R Garnett (Eds.), Vol. 30. Curran Associates, Inc.
1003 <https://proceedings.neurips.cc/paper/2017/file/a486cd07e4ac3d270571622f4f316ec5-Paper.pdf>
- 1004 [111] Claire Larsonneur. 2021. Intelligence artificielle ET/OU diversité linguistique : les paradoxes du traitement automatique des langues.
1005 <http://www.hybrid.univ-paris8.fr/lodel/index.php?id=1542>
- 1006 [112] Douglass B. Lee. 1973. Requiem for Large-Scale Models. *Journal of the American Institute of Planners* 39, 3 (5 1973), 163–178.
1007 <https://doi.org/10.1080/01944367308977851>
- 1008 [113] Min Kyung Lee and Su Baykal. 2017. Algorithmic Mediation in Group Decisions: Fairness Perceptions of Algorithmically Mediated vs.
1009 Discussion-Based Social Division. In *Proceedings of the 2017 ACM Conference on Computer Supported Cooperative Work and Social Computing (CSCW*
'17). Association for Computing Machinery, New York, NY, USA, 1035–1048. <https://doi.org/10.1145/2998181.2998230>
- 1010 [114] Min Kyung Lee, Daniel Kusbit, Anson Kahng, Ji Tae Kim, Xinran Yuan, Allissa Chan, Daniel See, Ritesh Noothigattu, Siheon Lee, Alexandros Psomas,
1011 and Ariel D Procaccia. 2019. WeBuildAI: Participatory Framework for Algorithmic Governance. *Proc. ACM Hum.-Comput. Interact.* 3, CSCW (11
1012 2019). <https://doi.org/10.1145/3359283>
- 1013 [115] Michelle Seng Ah Lee and Jatinder Singh. 2021. Risk Identification Questionnaire for Unintended Bias in Machine Learning Development Lifecycle.
1014 *SSRN Electronic Journal* (2021). <https://doi.org/10.2139/ssrn.3777093>
- 1015 [116] kobi leins, Jey Han Lau, and Timothy Baldwin. 2020. Give Me Convenience and Give Her Death: Who Should Decide What Uses of NLP are
1016 Appropriate, and on What Basis?. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for
1017 Computational Linguistics, Stroudsburg, PA, USA. <https://doi.org/10.18653/v1/2020.acl-main.261>
- 1018 [117] Q. Vera Liao, Daniel Gruen, and Sarah Miller. 2020. Questioning the AI: Informing Design Practices for Explainable AI User Experiences. (1 2020).
1019 <https://doi.org/10.1145/3313831.3376590>
- 1020 [118] Q. Vera Liao, Milena Pribić, Jaesik Han, Sarah Miller, and Daby Sow. 2021. Question-Driven Design Process for Explainable AI User Experiences. (4 2021).
- 1021 [119] Enrico Liscio, Michiel van der Meer, Luciano C Siebert, Catholijn M Jonker, Niek Mouter, and Pradeep K Murukannaiah. 2021. Axes: Identifying and
1022 Evaluating Context-Specific Values. In *Proceedings of the 20th International Conference on Autonomous Agents and MultiAgent Systems*. International
1023 Foundation for Autonomous Agents and Multiagent Systems, Richland, SC, 799–808.
- 1024 [120] Duri Long and Brian Magerko. 2020. What is AI Literacy? Competencies and Design Considerations. In *Proceedings of the 2020 CHI*
Conference on Human Factors in Computing Systems. Association for Computing Machinery, New York, NY, USA, 1–16. <https://doi-org.tudelft.idm.oclc.org/10.1145/3313831.3376727>
- 1025 [121] Henrietta Lyons, Eduardo Velloso, and Tim Miller. 2021. Conceptualising Contestability: Perspectives on Contesting Algorithmic Decisions. (2
1026 2021). <https://doi.org/10.1145/3449180>
- 1027 [122] Chunchuan Lyu, Kaizhu Huang, and Hai-Ning Liang. 2015. A Unified Gradient Regularization Family for Adversarial Examples. In *2015 IEEE*
1028 *International Conference on Data Mining*. IEEE, 301–309. <https://doi.org/10.1109/ICDM.2015.84>
- 1029 [123] N. Mahendran. 2021. Analysis of memory consumption by neural networks based on hyperparameters. (10 2021).
- 1030 [124] Donald Martin, Jr Google Vinodkumar Prabhakaran Google Jill Kuhlberg, and Andrew S Smart Google William Isaac DeepMind. 2020. Extending
1031 the Machine Learning Abstraction Boundary: A Complex Systems Approach to Incorporate Societal Context. (2020).
- 1032 [125] M. Mehltau. 2007. Iconset for data-privacy declarations v 0.1. <https://netzpolitik.org/wp-upload/data-privacy-icons-v01.pdf>
- 1033 [126] Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. 2021. A Survey on Bias and Fairness in Machine Learning.
1034 *ACM Comput. Surv.* 54, 6 (7 2021). <https://doi.org/10.1145/3457607>
- 1035 [127] Microsoft. 2018. AI Principles. <https://www.microsoft.com/en-us/ai/responsible-ai?activetab=pivot1%3aprimarary6>
- 1036 [128] Swati Mishra and Jeffrey M Rzeszotarski. 2021. Designing Interactive Transfer Learning Tools for ML Non-Experts. In *Proceedings of the 2021 CHI*
1037 *Conference on Human Factors in Computing Systems*. ACM, New York, NY, USA. <https://doi.org/10.1145/3411764.3445096>
- 1038 [129] Mission assigned by the French Prime Minister. 2019. For a Meaningful Artificial Intelligence: Toward a French and European Strategy.
1039 https://www.aiforhumanity.fr/pdfs/MissionVillani_Report_ENG-VF.pdf

- 1041 [130] Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and
1042 Timnit Gebru. 2018. Model Cards for Model Reporting. (10 2018). <https://doi.org/10.1145/3287560.3287596>
- 1043 [131] Tanushree Mitra. 2021. Provocation: Contestability in Large-Scale Interactive {NLP} Systems. In *Proceedings of the First Workshop on Bridging*
1044 *Human{-}Computer Interaction and Natural Language Processing*. Association for Computational Linguistics, 96–100.
- 1045 [132] Brent Mittelstadt. 2019. Principles alone cannot guarantee ethical AI. *Nature Machine Intelligence* 1, 11 (2019), 501–507. [https://doi.org/10.1038/s42256-](https://doi.org/10.1038/s42256-019-0114-4)
1046 [019-0114-4](https://doi.org/10.1038/s42256-019-0114-4)
- 1047 [133] Jessica Morley, Luciano Floridi, Libby Kinsey, and Anat Elhalal. 2020. From What to How: An Initial Review of Publicly Available
1048 AI Ethics Tools, Methods and Research to Translate Principles into Practices. *Science and Engineering Ethics* 26 (2020), 2141–2168.
1049 <https://doi.org/10.1007/s11948-019-00165-5>
- 1050 [134] Ramaravind Kommiya Mothilal, Amit Sharma, and Chenhao Tan. 2019. Explaining Machine Learning Classifiers through Diverse Counterfactual
1051 Explanations. (5 2019). <https://doi.org/10.1145/3351095.3372850>
- 1052 [135] Pradeep K Murukannaiah and Munindar P Singh. 2014. Xipho: Extending Tropos to Engineer Context-Aware Personal Agents. In *Proceedings*
1053 *of the 2014 International Conference on Autonomous Agents and Multi-Agent Systems (AAMAS '14)*. International Foundation for Autonomous Agents
1054 and Multiagent Systems, Richland, SC, 309–316.
- 1055 [136] Milad Nasr, Reza Shokri, and Amir Houmansadr. 2018. Machine Learning with Membership Privacy using Adversarial Regularization. (7 2018).
- 1056 [137] Harsha Nori, Samuel Jenkins, Paul Koch, and Rich Caruana. 2019. InterpretML: A Unified Framework for Machine Learning Interpretability. (9 2019).
- 1057 [138] Arif Nurwidyanoro, Mojtaba Shahin, Michel Chaudron, Waqar Hussain, Harsha Perera, Rifat Ara Shams, and Jon Whittle. 2021. Towards a Human
1058 Values Dashboard for Software Development: An Exploratory Study. (7 2021).
- 1059 [139] OECD. 2019. Recommendation of the Council on Artificial Intelligence. <https://legalinstruments.oecd.org/en/instruments/OECD-LEGAL-0406>
- 1060 [140] Kieron O'Hara. 2020. Explainable AI and the philosophy and practice of explanation. *Computer Law & Security Review* 39 (11 2020), 105474.
1061 <https://doi.org/10.1016/j.clsr.2020.105474>
- 1062 [141] Nicolas Papernot, Patrick McDaniel, Xi Wu, Somesh Jha, and Ananthram Swami. 2016. Distillation as a Defense to Adversarial Perturbations Against
1063 Deep Neural Networks. In *2016 IEEE Symposium on Security and Privacy (SP)*. IEEE, 582–597. <https://doi.org/10.1109/SP.2016.41>
- 1064 [142] Lorenza Parisi and Francesca Comunello. 2020. Dating in the time of “relational filter bubbles”: exploring imaginaries, perceptions and tactics
1065 of Italian dating app users. *The Communication Review* 23, 1 (2020), 66–89. <https://doi.org/10.1080/10714421.2019.1704111>
- 1066 [143] Reema Patel. 2021. Reboot AI with human values. *Nature* 598, 7879 (10 2021). <https://doi.org/10.1038/d41586-021-02693-2>
- 1067 [144] Amandalynne Paullada, Inioluwa Deborah Raji, Emily M. Bender, Emily Denton, and Alex Hanna. 2020. Data and its (dis)contents: A survey of
1068 dataset development and use in machine learning research. (12 2020). <http://arxiv.org/abs/2012.05345>
- 1069 [145] Alina Pommeranz, Christian Detweiler, Pascal Wiggers, and Catholijn Jonker. 2012. Elicitation of situated values: need for tools to help stakeholders
1070 and designers to reflect and communicate. *Ethics and Information Technology* 14, 4 (12 2012), 285–303. <https://doi.org/10.1007/s10676-011-9282-6>
- 1071 [146] Manish Raghavan, Solon Barocas, Jon Kleinberg, and Karen Levy. 2020. Mitigating Bias in Algorithmic Hiring: Evaluating Claims and Practices.
1072 In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency (FAT* '20)*. Association for Computing Machinery, New York,
1073 NY, USA, 469–481. <https://doi.org/10.1145/3351095.3372828>
- 1074 [147] Inioluwa Deborah Raji, Andrew Smart, Rebecca N. White, Margaret Mitchell, Timnit Gebru, Ben Hutchinson, Jamila Smith-Loud, Daniel Theron,
1075 and Parker Barnes. 2020. Closing the AI accountability gap. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*.
1076 ACM, New York, NY, USA, 33–44. <https://doi.org/10.1145/3351095.3372873>
- 1077 [148] Manoel Horta Ribeiro, Raphael Ottoni, Robert West, Virgilio A F Almeida, and Wagner Meira. 2020. Auditing Radicalization Pathways on YouTube.
1078 In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency (FAT* '20)*. Association for Computing Machinery, New York,
1079 NY, USA, 131–141. <https://doi.org/10.1145/3351095.3372879>
- 1080 [149] Arianna Rossi and Monica Palmirani. 2017. A Visualization Approach for Adaptive Consent in the European Data Protection Framework. In *2017*
1081 *Conference for E-Democracy and Open Government (CeDEM)*. IEEE, 159–170. <https://doi.org/10.1109/CeDEM.2017.23>
- 1082 [150] Stuart Russell, Daniel Dewey, and Max Tegmark. 2015. Research Priorities for Robust and Beneficial Artificial Intelligence. *AI Magazine* 36, 4 (12
1083 2015). <https://doi.org/10.1609/aimag.v36i4.2577>
- 1084 [151] Pedro Saleiro, Benedict Kuester, Loren Hinkson, Jesse London, Abby Stevens, Ari Anisfeld, Kit T Rodolfa, and Rayid Ghani. 2018. Aequitas: A
1085 Bias and Fairness Audit Toolkit.
- 1086 [152] Christian Sandvig, Kevin Hamilton, Karrie Karahalios, and Cedric Langbort. 2014. Auditing algorithms: Research methods for detecting
1087 discrimination on internet platforms. In *Data and discrimination: converting critical concerns into productive inquiry 22*.
- 1088 [153] Maarten Sap, Dallas Card, Saadia Gabriel, Yejin Choi, and Noah A. Smith. 2019. The Risk of Racial Bias in Hate Speech Detection. *ACL 2019 - 57th*
1089 *Annual Meeting of the Association for Computational Linguistics, Proceedings of the Conference (2019)*, 1668–1678. <https://doi.org/10.18653/v1/P19-1163>
- 1090 [154] Enrique Schaefer, Richard Kelley, and Monica Nicolescu. 2009. Robots as animals: A framework for liability and responsibility in human-
1091 robot interactions. In *RO-MAN 2009 - The 18th IEEE International Symposium on Robot and Human Interactive Communication*. IEEE.
1092 <https://doi.org/10.1109/ROMAN.2009.5326244>
- 1093 [155] Morgan Klaus Scheuerman, Emily Denton, and Alex Hanna. 2021. Do Datasets Have Politics? Disciplinary Values in Computer Vision Dataset
1094 Development. *Proc. ACM Hum.-Comput. Interact.* 5, CSCW2, Article 317 (2021). <https://doi.org/10.1145/3476058>
- 1095 [156] Shalom H. Schwartz. 2012. An Overview of the Schwartz Theory of Basic Values. *Online Readings in Psychology and Culture* 2, 1 (12 2012).
1096 <https://doi.org/10.9707/2307-0919.1116>

- 1093 [157] Mojtaba Shahin, Waqar Hussain, Arif Nurwidyanoro, Harsha Perera, Rifat Shams, John Grundy, and Jon Whittle. 2021. Operationalizing Human
1094 Values in Software Engineering: A Survey. (8 2021).
- 1095 [158] Hong Shen, Alicia DeVos, Motahhare Eslami, and Kenneth Holstein. 2021. Everyday algorithm auditing: Understanding the power of everyday
1096 users in surfacing harmful algorithmic behaviors. (5 2021). <https://doi.org/10.1145/3479577>
- 1097 [159] Irina Shklovski and Carolina Némethy. 2022. Nodes of certainty and spaces for doubt in AI ethics for engineers. *Information, Communication &*
1098 *Society* (1 2022), 1–17. <https://doi.org/10.1080/1369118X.2021.2014547>
- 1099 [160] Ben Shneiderman. 2020. Bridging the Gap Between Ethics and Practice: Guidelines for Reliable, Safe, and Trustworthy Human-Centered AI Systems.
1100 *ACM Trans. Interact. Intell. Syst.* 10, 4 (10 2020). <https://doi.org/10.1145/3419764>
- 1101 [161] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. 2016. Membership Inference Attacks against Machine Learning Models. (10 2016).
- 1102 [162] Kacper Sokol and Peter Flach. 2020. Explainability fact sheets. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*.
1103 ACM, New York, NY, USA. <https://doi.org/10.1145/3351095.3372870>
- 1104 [163] Megha Srivastava, Hoda Heidari, and Andreas Krause. 2019. Mathematical Notions vs. Human Perception of Fairness. In *Proceed-*
1105 *ings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. ACM, New York, NY, USA, 2459–2468.
1106 <https://doi.org/10.1145/3292500.3330664>
- 1107 [164] Harini Suresh, Steven R. Gomez, Kevin K. Nam, and Arvind Satyanarayan. 2021. Beyond Expertise and Roles: A Framework to Characterize the
1108 Stakeholders of Interpretable Machine Learning and their Needs. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*.
1109 ACM, New York, NY, USA, 1–16. <https://doi.org/10.1145/3411764.3445088>
- 1110 [165] Harini Suresh and John Gutttag. 2021. A Framework for Understanding Sources of Harm throughout the Machine Learning Life Cycle. In *Equity*
1111 *and Access in Algorithms, Mechanisms, and Optimization*. ACM, New York, NY, USA, 1–9. <https://doi.org/10.1145/3465416.3483305>
- 1112 [166] Telia Company. 2019. Guiding Principles on Trusted AI Ethics. [https://www.teliacompany.com/globalassets/telia-company/documents/about-](https://www.teliacompany.com/globalassets/telia-company/documents/about-telia-company/public-policy/2018/guiding-principles-on-trusted-ai-ethics.pdf)
1113 [telia-company/public-policy/2018/guiding-principles-on-trusted-ai-ethics.pdf](https://www.teliacompany.com/globalassets/telia-company/documents/about-telia-company/public-policy/2018/guiding-principles-on-trusted-ai-ethics.pdf)
- 1114 [167] The IEEE Global Initiative on Ethics of Autonomous and Intelligent Systems. 2019. *Ethically Aligned Design: A Vision for Prioritizing Human*
1115 *Well-being with Autonomous and Intelligent Systems* (first edition ed.). IEEE.
- 1116 [168] The Royal Society. 2019. Explainable AI: the basics . [https://royalsociety.org/tudelft.idm.oclc.org/-/media/policy/projects/explainable-ai/AI-](https://royalsociety.org/tudelft.idm.oclc.org/-/media/policy/projects/explainable-ai/AI-and-interpretability-policy-briefing.pdf)
1117 [and-interpretability-policy-briefing.pdf](https://royalsociety.org/tudelft.idm.oclc.org/-/media/policy/projects/explainable-ai/AI-and-interpretability-policy-briefing.pdf)
- 1118 [169] Sarah Thew and Alistair Sutcliffe. 2018. Value-based requirements engineering: method and experience. *Requirements Engineering* 23, 4 (11 2018).
1119 <https://doi.org/10.1007/s00766-017-0273-y>
- 1120 [170] Songül Tolan, Marius Miron, Emilia Gómez, and Carlos Castillo. 2019. Why Machine Learning May Lead to Unfairness: Evidence from Risk
1121 Assessment for Juvenile Justice in Catalonia. In *Proceedings of the Seventeenth International Conference on Artificial Intelligence and Law (ICAIL*
1122 *'19)*. Association for Computing Machinery, New York, NY, USA, 83–92. <https://doi.org/10.1145/3322640.3326705>
- 1123 [171] Dimitrios Tsarapatsanis and Nikolaos Aletras. 2021. On the Ethical Limits of Natural Language Processing on Legal Text. (5 2021).
- 1124 [172] National Science United States Executive Office of the President and Technology Council Committee on Technology. 2016. Preparing for the Future
1125 of Artificial Intelligence. [https://obamawhitehouse.archives.gov/sites/default/files/whitehouse_files/microsites/ostp/NSTC/preparing_for_the_](https://obamawhitehouse.archives.gov/sites/default/files/whitehouse_files/microsites/ostp/NSTC/preparing_for_the_future_of_ai.pdf)
1126 [future_of_ai.pdf](https://obamawhitehouse.archives.gov/sites/default/files/whitehouse_files/microsites/ostp/NSTC/preparing_for_the_future_of_ai.pdf)
- 1127 [173] Funda Ustek-Spilda, Alison Powell, and Selena Nemorin. 2019. Engaging with ethics in Internet of Things: Imaginaries in the social milieu of
1128 technology developers. *Big Data & Society* 6, 2 (7 2019), 205395171987946. <https://doi.org/10.1177/2053951719879468>
- 1129 [174] Kristen Vaccaro, Christian Sandvig, and Karrie Karahalios. 2020. "At the End of the Day Facebook Does What It Wants". *Proceedings of the ACM*
1130 *on Human-Computer Interaction* 4, CSCW2 (10 2020), 1–22. <https://doi.org/10.1145/3415238>
- 1131 [175] Niels van Berkel, Jorge Goncalves, Daniel Russo, Simo Hosio, and Mikael B. Skov. 2021. Effect of Information Presentation on Fairness Perceptions
1132 of Machine Learning Predictors. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. ACM, New York, NY, USA,
1133 1–13. <https://doi.org/10.1145/3411764.3445365>
- 1134 [176] Ibo van de Poel. 2013. Translating Values into Design Requirements. 253–266. https://doi.org/10.1007/978-94-007-7762-0_{ }20
- 1135 [177] Sahil Verma and Julia Rubin. 2018. Fairness definitions explained. In *Proceedings of the International Workshop on Software Fairness*. ACM, New
1136 York, NY, USA, 1–7. <https://doi.org/10.1145/3194770.3194776>
- 1137 [178] Sandra Wachter and Brent Mittelstadt. 2019. Right to Reasonable Inferences: Re-Thinking Data Protection Law in the Age of Big Data and AI.
1138 *Columbia Business Law Review* 2 (2019), 494–620.
- 1139 [179] Zezhong Wang, Jacob Ritchie, Jingtao Zhou, Fanny Chevalier, and Benjamin Bach. 2021. Data Comics for Reporting Controlled User Studies in Human-
1140 Computer Interaction. *IEEE Transactions on Visualization and Computer Graphics* 27, 2 (2 2021), 967–977. <https://doi.org/10.1109/TVCG.2020.3030433>
- 1141 [180] James Wexler, Mahima Pushkarna, Tolga Bolukbasi, Martin Wattenberg, Fernanda Viegas, and Jimbo Wilson. 2019. The What-If Tool: Interactive
1142 Probing of Machine Learning Models. (7 2019). <https://doi.org/10.1109/TVCG.2019.2934619>
- 1143 [181] Christo Wilson, Avijit Ghosh, Shan Jiang, Alan Mislove, Lewis Baker, Janelle Szary, Kelly Trindel, and Frida Polli. 2021. Building and Auditing
1144 Fair Algorithms: A Case Study in Candidate Screening. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*
(FAcCT '21). Association for Computing Machinery, New York, NY, USA, 666–677. <https://doi.org/10.1145/3442188.3445928>
- [182] Langdon Winner. 1980. Do Artifacts Have Politics? *Daedalus* 109, 1 (1980), 121–136. <http://www.jstor.org/stable/20024652>
- [183] Pulei Xiong, Scott Buffett, Shahrear Iqbal, Philippe Lamontagne, Mohammad Mamun, and Heather Molyneaux. 2021. Towards a Robust and
Trustworthy Machine Learning System Development. (1 2021).

- 1145 [184] Depeng Xu, Shuhan Yuan, Lu Zhang, and Xintao Wu. 2018. FairGAN: Fairness-aware Generative Adversarial Networks. In *2018 IEEE International*
1146 *Conference on Big Data (Big Data)*. 570–575. <https://doi.org/10.1109/BigData.2018.8622525>
- 1147 [185] An Yan and Bill Howe. 2020. Fairness-Aware Demand Prediction for New Mobility. *Proceedings of the AAAI Conference on Artificial Intelligence*
1148 34, 01 (4 2020), 1079–1087. <https://doi.org/10.1609/aaai.v34i01.5458>
- 1149 [186] Vahid Yazdanpanah, Enrico Gerding, Sebastian Stein, Mehdi Dastani, Catholijn M Jonker, and Timothy Norman. 2021. Responsibility Research
1150 for Trustworthy Autonomous Systems. In *20th International Conference on Autonomous Agents and Multiagent Systems (03/05/21 - 07/05/21)*. 57–62.
1151 <https://eprints.soton.ac.uk/447511/>
- 1152 [187] Jiayuan Ye, Aadyaa Maddi, Sasi Kumar Murakonda, and Reza Shokri. 2021. Enhanced Membership Inference Attacks against Machine Learning
1153 Models. (11 2021).
- 1154 [188] Brian Hu Zhang, Blake Lemoine, and Margaret Mitchell. 2018. Mitigating Unwanted Biases with Adversarial Learning. In *Proceedings of*
1155 *the 2018 AAAI/ACM Conference on AI, Ethics, and Society (AIES '18)*. Association for Computing Machinery, New York, NY, USA, 335–340.
<https://doi.org/10.1145/3278721.3278779>
- 1156 [189] Angela Zhou, David Madras, Inioluwa Raji Raji, Bogdan Kulynych, Smitha Mili, and Richard Zemel. [n. d.]. Call for participation: Participatory
1157 Approaches to Machine Learning. <https://participatoryml.github.io/>
- 1158 [190] Liming Zhu, Xiwei Xu, Qinghua Lu, Guido Governatori, and Jon Whittle. 2021. AI and Ethics – Operationalising Responsible AI. (5 2021).
- 1159 [191] Christian Zimmermann, Rafael Accorsi, and Gunter Muller. 2014. Privacy Dashboards: Reconciling Data-Driven Business Models and Privacy.
1160 In *2014 Ninth International Conference on Availability, Reliability and Security*. IEEE, 152–157. <https://doi.org/10.1109/ARES.2014.27>
- 1161 [192] Arkaitz Zubiaga, Bo Wang, Maria Liakata, and Rob Procter. 2019. Political Homophily in Independence Movements: Analyzing and Classifying
1162 Social Media Users by National Identity. *IEEE Intelligent Systems* 34, 6 (2019), 34–42. <https://doi.org/10.1109/MIS.2019.2958393>
- 1163
- 1164
- 1165
- 1166
- 1167
- 1168
- 1169
- 1170
- 1171
- 1172
- 1173
- 1174
- 1175
- 1176
- 1177
- 1178
- 1179
- 1180
- 1181
- 1182
- 1183
- 1184
- 1185
- 1186
- 1187
- 1188
- 1189
- 1190
- 1191
- 1192
- 1193
- 1194
- 1195
- 1196

A FROM VALUES TO SPECIFIC CRITERIA MANIFESTATIONS

Value	Criteria	Manifestations
Privacy	(1) Consent for data usage [3, 56, 60]	• Written declaration of consent [56]
	(2) Data protection [3, 60, 61]	• Description of what data is collected [125]
	(3) Control over data / ability to restrict processing [56, 60]	• Description of how data is handled [125]
	(4) Right to rectification [3, 56, 60]	• Purpose statement of data collection [125]
	(5) Right to erase the data [3, 56, 60]	• Statement of how long the data is kept [125]
	(6) Right of access by data subject, data agency [56, 167]	• Form and submission mechanisms to object data collection and to make complaints [27]
		• Obfuscation of data [3]
		AGAINST INTEGRITY THREATS [183]:
		• Training time [183] Ex.:
		• Data sanitization ¹² [23, 40]
		• Robust learning ¹³ [23, 73]
		• Prediction time [183]
		• Model enhancement [23, 74, 122, 141] Ex.:
		• Adversarial Learning ¹⁴
		• Gradient masking ¹⁵
		• Defensive Distillation ¹⁶
Security	(1) Resilience to attacks : protection of privacy [86, 127, 178], vulnerabilities, fallback plans [3, 60, 75, 133]	
	(2) Predictability [3, 57, 60]	AGAINST PRIVACY THREATS [183]:
	(3) Robustness / reliability : prevent manipulation [3]	• Mitigation techniques [136]:
		• Restrict prediction vector to top k classes ¹⁷ [161]
		• Coarsen the precision of the prediction vector ¹⁸ [161]
		• Increase entropy of the prediction vector ¹⁹ [161]
		• Use regularization ²⁰ [101, 161]
		• Differential privacy mechanisms [136]:
		• Differential privacy ²¹ [53, 187]. Ex.:
		• Adversarial regularization ²² [136]
		• MemGuard ²³ [97]

¹²It ensures data soundness by identifying abnormal input samples and by removing them [183].

¹³It ensures that algorithms are trained on statistically robust datasets, with little sensitivity to outliers [183].

¹⁴Adversarial samples are introduced to the training set [183].

¹⁵Input gradients are modified to enhance model robustness [183].

¹⁶The dimensionality of the network is reduced [183].

¹⁷Applicable when the number of classes is very large. Even if the model only outputs the most likely k classes, it will still be useful [161].

¹⁸It consists in rounding the classification probabilities down [161].

¹⁹Modification of the softmax layer (in neural networks) to increase its normalizing temperature [161].

²⁰Technique to avoid overfitting in ML that penalizes large parameters by adding a regularization factor λ to the loss function [161].

²¹It prevents any adversary from distinguishing the predictions of a model when its training dataset is used compared to when other dataset is used [187]

²²Membership privacy is modeled as a min-max optimization problem, where a model is trained to achieve minimum loss of accuracy and maximum robustness against the strongest inference attack [136].

²³Noise is added to the confidence vector of the attacker so as to mislead the attacker's classifier [97]

	Value	Criteria	Manifestations
1249	Conservation	Performance	<ul style="list-style-type: none"> • Accuracy (for classification, sum of true positive and true negative rates) [130, 180] • False Positive and False Negative rates [130, 180] • False Discovery and Omission Rate [130] • Mean and median error [180] • R2 score [25] • Precision and recall rates [180] • Area under ROC curve (AUC) [25] • Estimation of energy consumption through [68]: <ul style="list-style-type: none"> • performance counters • simulation • instruction- or architecture-level estimations • real-time estimation • Estimation of GPU memory consumption [67, 123] • Wall-clock training time [14, 41]
1250			
1251			
1252			
1253			
1254	Respect for public interest	(1) Correctness of predictions [26, 57, 60] (2) Memory efficiency [3, 26] (3) Training efficiency [26] (4) Energy efficiency [3, 26] (5) Data efficiency [26]	<ul style="list-style-type: none"> • Diverse and inclusive forum for discussion [60, 129] • Measure of social and environmental impact [21, 133, 147]
1255			
1256			
1257	Universalism	Fairness	<ul style="list-style-type: none"> • Accuracy across groups (for classification, sum of true positive and true negative rates) [37, 80, 105, 133] • False positive and negative rates across groups [37, 105, 126, 151, 179] • False discovery and omission rates across groups [130, 151] • Pinned AUC [48, 130] • Debiasing algorithms [19] • Election of protected classes based on user considerations [77]
1258			
1259			
1260			
1261			
1262			
1263			
1264			
1265	Non-discrimination	(1) Desirability of technology [1, 34, 104] (2) Benefit to society [60–62, 133] (3) Environmental impact [3, 21]	<ul style="list-style-type: none"> • Inclusive data generation process [3, 34, 70, 133] • Analysis of data for potential biases, data quality assessment [3, 60, 69, 86, 126] • Diversity of participant in development process [3, 60, 114, 189] • Access to code and technology to all [3, 26, 60, 133]
1266			
1267			
1268	Individual fairness ²⁴ Demographic parity ²⁵ Conditional Statistical parity ²⁶ Equality of opportunity ²⁷ Equalized odds ²⁸ Treatment equality ²⁹ Test fairness ³⁰ Procedural fairness ³¹	(1) Individual fairness [18, 52, 110, 126] (2) Demographic parity [18, 52, 80, 86, 102, 110, 126, 163, 177] (3) Conditional Statistical parity [126, 177] (4) Equality of opportunity [79, 126, 175] (5) Equalized odds [126] (6) Treatment equality [22, 126] (7) Test fairness [37, 126, 177] (8) Procedural fairness [77, 110, 126]	<ul style="list-style-type: none"> • Accuracy across groups (for classification, sum of true positive and true negative rates) [37, 80, 105, 133] • False positive and negative rates across groups [37, 105, 126, 151, 179] • False discovery and omission rates across groups [130, 151] • Pinned AUC [48, 130] • Debiasing algorithms [19] • Election of protected classes based on user considerations [77]
1269			
1270			
1271			
1272			
1273			
1274			
1275			
1276			
1277			
1278			
1279			
1280			
1281			
1282			
1283			
1284			
1285			
1286			
1287			
1288			
1289			

²⁴ Similar individuals should be treated in a similar way. Diverging definitions state that: two individuals that are similar with respect to a common metric should receive the same outcome (*fairness through awareness*); or any protected attribute should not be used when making a decision (*fairness through unawareness*); or the outcome obtained by an individual should be the same if this individual belonged to a counterfactual world or group (*counterfactual fairness*) [126].

²⁵ The probability of getting a positive outcome should be the same whether the individual belongs to a protected group or not [126].

²⁶ Given a set of factors L, individuals belonging to the protected or unprotected group should have the same probability of getting a positive outcome [126].

²⁷ The probability for a person from class A (positive class) of getting a positive outcome, which should be the same regardless of the group (protected group or not) that the individual belongs to [126].

²⁸ The probability for a person from class A (positive class) of getting a positive outcome and the probability for a person from class B (negative class) of getting a negative outcome should be the same [126].

²⁹ The ratio of false positives and negatives has to be the same for both groups [126].

³⁰ For any probability score S, the probability of correctly belonging to the positive class should be the same for both the protected and unprotected group [126].

³¹ It deals with the fairness of the decision-making process that leads to the outcome in question [77].

	Value	Criteria	Manifestations	
1301 1302 1303 1304 1305 1306 1307 1308 1309	Openness	Transparency	(1) Interpretability of data and models [26, 168]	• Description of data generation process [3, 20, 34, 69, 70, 133]
			(2) Enabling human oversight of operations [60, 133]	• Disclosure of origin and properties of models and data [3, 130, 168]
			(3) Accessibility of data and algorithm [3, 60, 168]	• Open access to data and algorithm [3, 26, 60, 168]
			(4) Traceability [133]	• Notification of usage/interaction [60]
			(5) Reproducibility [26]	• Regular reporting [60]
1310 1311 1312 1313 1314	Openness	Explainability	(1) Ability to understand AI systems and the decision reached [26, 57, 61, 62, 139, 168]	• Interpretability by design [18]
			(2) Traceability [133]	• Post-hoc explanations [18]
			(3) Enable evaluation [60, 133]	
1315 1316 1317 1318 1319 1320 1321 1322 1323	Individual empowerment	Contestability	(1) Enable argumentation / negotiation against a decision [6, 16, 57, 60, 100, 113, 121, 168]	• Information of who determines and what constitutes a contestable decision and who is accountable [121]
			(2) Citizen empowerment [16, 57, 100]	• Determination of who can contest the decision (subject or representative) [121]
				• Indication of type of review in place [121]
				• Information regarding the contestability workflow [121]
				• Mechanisms for users to ask questions and record disagreements with system behavior [87, 131]
1324 1325 1326 1327 1328	Individual empowerment	Human Control	(1) User/collective influence [26, 113]	• Continuous monitoring of system to intervene [57, 60, 166]
			(2) Human review of automated decision [60]	• Establishment levels of human discretion during the use of the system [57, 127]
			(3) Choice of how and whether to delegate [60]	• Ability to override the decision made by a system [57]
1329 1330 1331 1332 1333 1334	Individual empowerment	Human agency	(1) Respect for human autonomy [57, 60, 133]	• Give knowledge and tools to comprehend and interact with AI system [57]
			(2) Power to decide. Ability to make informed autonomous decision [26, 57]	• Opportunity to self-assess the system [57]
			(3) Ability to opt out of an automated decision [57, 60]	

Table 3. Summary of the specific criteria that relate to each value considered in our ML assessment framework. These criteria are then translated into specific manifestations in the form of signifiers (orange), process-oriented practices (olive) or quantifiable indicators (magenta).

1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352

B MAPPING STAKEHOLDERS

Stakeholder	Mapping [164]	Nature of knowledge	Purpose of insight
Development team	ML, Formal + Instrumental + Personal	<ul style="list-style-type: none"> • “Knowledge of the math behind the architecture” [164] • “Stakeholder involved in an ex-ante impact assessment of the automatic decision system” [84] 	<ul style="list-style-type: none"> • Ensure/improve product efficiency and debug [18] • Research new functionalities [18]
Auditing team	Milieu, Formal + Instrumental	<ul style="list-style-type: none"> • “Familiarity with broader ML-enabled systems” [164] • “Experts who intervene wither upstream or downstream” [84] 	<ul style="list-style-type: none"> • Verify model compliance with legislation [18]
Data domain experts	Data domain, Formal + Instrumental	<ul style="list-style-type: none"> • “Theories relevant to the data domain” [164] • “Professional involved in the operational phase of the automatic decision system” [84] 	<ul style="list-style-type: none"> • Gain scientific or domain-specific knowledge [18, 164] • Trust the model [18, 164] • Act based on the output [164]
Decision subjects	Data domain + Milieu, Personal	<ul style="list-style-type: none"> • “Lived experience and cultural knowledge” [164] • “Layperson affected by the outcomes of the automatic decision system” [84] 	<ul style="list-style-type: none"> • Understand their situation [18] • Verify fair decision [18] • Contest decision [164] • Understand how one’s data is being used [164]

Table 4. Description of potential stakeholders that can be brought together as part of our value-based framework. These stakeholders have been mapped following the two dimensional criteria (type of knowledge —formal, instrumental or personal— and contexts in which this knowledge manifests —ML, data domain, milieu—) outlined by Suresh et al. [164]. The nature of their knowledge and the purpose of gaining insight for each of them have also been defined.

C TAILORED COMMUNICATION OF SYSTEM-RELATED INFORMATION

		Development team	Auditing team	Data Domain experts	Decision subjects
Conservation	Privacy	[K]	[K]		[A] [B]
	Security	[K] [W] [AB]	[K] [W]		
	Performance	[F] [G] [H] [Y] [Z] [AE]	[G] [H] [Y] [Z] [AE]	[I] [J]	[J]
Universalism	Respect for public interest	[E] [AE]	[E] [AE]	[E]	[C] [D]
	Fairness	[G] [H] [K] [X] [Y] [Z] [AD]	[G] [H] [K] [X] [Y] [Z] [AD]	[I] [J]	[J]
	Non- discrimination	[H] [K] [X] [Y] [AD]	[H] [K] [X] [Y] [AD]	[L]	[J] [L]
Openness	Transparency	[H] [K] [M]	[H] [K] [M] [X] [Y]	[I] [L] [M]	[B] [J] [L] [M]
	Explainability	[M] [N] [O] [Q] [AC] [AD] [P]	[M] [N] [O] [Q] [AC] [AD] [P]	[M] [N] [O] [Q] [P]	[J] [M] [N] [O] [Q] [R] [S] [P]
Individual empowerment	Contestability	[U]	[U]	[T] [U]	[T] [AF]
	Human Control	[V]	[V]	[T] [V]	[C] [T] [V]
	Human Agency				[B] [AA]

Table 5. Mapping of available means for transmitting value-specific manifestations to different stakeholders based on the purpose of their insight and the nature of their knowledge —more details about stakeholders can be found on appendix B, table 4—. These means have been classified into three main categories: descriptive documents specifying whether/how a value manifestation is fulfilled (red), strategies for fulfilling value manifestations (blue), and complete tools for enabling the fulfillment of value manifestations (green). This table aims at facilitating the navigation of table 6, where each means is documented.

1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497

Means	Value	Manifestation(s)	Stakeholder				Application (model)	Approach	Visual elements	Additional details
			DT	AT	DE	DS				
Iconsets for data privacy [A]	Privacy	<ul style="list-style-type: none"> Description of what data is collected Description of how data is handled Purpose statement of data collection Statement of how long the data is kept 				✓	Agnostic	Iconsets		
		<ul style="list-style-type: none"> Description of what data is collected Description of how data is handled Purpose statement of data collection Statement of how long the data is kept 								
Privacy dashboards [B]	Privacy	<ul style="list-style-type: none"> Description of what data is collected Description of how data is handled Purpose statement of data collection Opportunity to self-assess the system 				✓	Agnostic	<ul style="list-style-type: none"> Timelines Bar charts Maps Network graphs 		
		<ul style="list-style-type: none"> Opportunity to self-assess the system Disclosure of origin and properties of data 								
Risk matrix [C]	Respect for public interest	<ul style="list-style-type: none"> Measure of social impact Ability to override the decision made by a system 				✓	Agnostic	<ul style="list-style-type: none"> Two-dimensional space (vulnerability vs dependence of the decision) 		
		<ul style="list-style-type: none"> Measure of social impact Ability to override the decision made by a system 								
Moral space [D]	Respect for public interest	<ul style="list-style-type: none"> Measure of social impact 				✓	Agnostic	Based on human judgement	<ul style="list-style-type: none"> Three-dimensional moral space. Wrongness as a function of intention and harm 	
		<ul style="list-style-type: none"> Measure of social impact 								
Social impact assessment [E]	Respect for public interest	<ul style="list-style-type: none"> Measure of social impact 				✓	Agnostic	Anticipate scenarios		
		<ul style="list-style-type: none"> Measure of social impact 								

1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538

Means	Value	Manifestation(s)	Stakeholder			Application (model)	Approach	Visual elements	Additional details
			DT	AT	DE				
Model Tracker in- teractive [F] visual- ization [9]	Perfor- mance	<ul style="list-style-type: none"> • Accuracy • False Positive and Negative rates 			✓	Classification tasks		<ul style="list-style-type: none"> • Summary statistics • Confusion matrices • Labels chart • Precision-recall curves • Connector lines to identify similar examples in feature space • Highlighted boxes for correlations between features and target classes 	
Model cards for models [G] [130]	Perfor- mance Fairness	<ul style="list-style-type: none"> • Accuracy • False Positive and Negative rates • False Discovery and omission rates • Accuracy across groups • False Positive and Negative rates across groups • False Discovery and omission rates across groups 			✓	Agnostic		<ul style="list-style-type: none"> • Confidence bars • Bar charts 	

1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565
1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579

Means	Value	Manifestation(s)	Stakeholder				Application (model)	Approach	Visual elements	Additional details
			DT	AT	DE	DS				
What-if [H] tool ³² [180]	Performance	<ul style="list-style-type: none"> Accuracy False Positive and Negative Rates False Discovery and omission rates 						<ul style="list-style-type: none"> Confusion matrices (Two-dimensional) Histograms Scatterplots Summary statistics of datasets Partial dependence plots 	<ul style="list-style-type: none"> Interactive modules include: list of feature values, inference values, and counterfactual controls 	
	Fairness	<ul style="list-style-type: none"> Accuracy across groups False Positive and Negative Rates across groups False Discovery and omission rates across groups 			✓		Classification tasks, Regression tasks			
	Transparency	<ul style="list-style-type: none"> Disclosure of origin and properties of data 								
Interactive transfer learning tools [128]	Performance	<ul style="list-style-type: none"> Accuracy False Positive and Negative Rates 						<ul style="list-style-type: none"> Confusion matrices Z-scored of each filter Bar charts Activation heatmaps t-SNE clusters 		
	Fairness	<ul style="list-style-type: none"> Accuracy across groups False Positive and Negative Rates across groups 			✓		Convolutional Neural Networks			
	Transparency	<ul style="list-style-type: none"> Disclosure of properties of data 								
Question- Driven XAIDe- sign [118]	Performance	<ul style="list-style-type: none"> Accuracy 						<ul style="list-style-type: none"> Summary statistics (percentage scores) for data explanations and performance metrics Feature importance Contrastive explanations 	<ul style="list-style-type: none"> End users were more interested in the limitation of the model: uncertainty 	
	Fairness	<ul style="list-style-type: none"> Accuracy across groups 								
	Transparency	<ul style="list-style-type: none"> Disclosure of origin and properties of data Analysis of data for potential biases, data quality assessment 			✓	✓	Agnostic			

³²<https://github.com/pair-code/what-if-tool>

1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620

Means	Value	Manifestation(s)	Stakeholder				Application (model)	Approach	Visual elements	Additional details
			DT	AT	DE	DS				
Explainability	Transparency	<ul style="list-style-type: none"> • Post-hoc explanations 								
		<ul style="list-style-type: none"> • Description of data generation process • Disclosure of origin properties of models and data 								
Datasheets for [K] datasets [69]	Non-discrimination	<ul style="list-style-type: none"> • Analysis of data for potential biases, data quality assessment 								
		<ul style="list-style-type: none"> • Written declaration of consent • Description of what data is collected • Description of how data is handled 	✓				Agnostic	<ul style="list-style-type: none"> • Summary statistics • Visual examples of datasets (if images, for instance) 		
Data centric explanations [L] [12]	Transparency	<ul style="list-style-type: none"> • Purpose statement of data collection • Statement of how long the data is kept 								
		<ul style="list-style-type: none"> • Election of protected classes • Membership inference 								
Non-discrimination	Transparency	<ul style="list-style-type: none"> • Description of data generation process • Disclosure of origin and properties of the models and data 								
		<ul style="list-style-type: none"> • Analysis of data for potential biases, data quality assessment 								

1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661

Means	Value	Manifestation(s)	Stakeholder				Application (model)	Approach	Visual elements	Additional details
			DT	AT	DE	DS				
[M] Example-based explanations [18, 24, 32, 49, 98, 117, 118]	Transparency	<ul style="list-style-type: none"> • Disclosure of properties of data 					<ul style="list-style-type: none"> • Similar example • Typical example • Counter-factual example 	<ul style="list-style-type: none"> • Example images from dataset if in the visual domain 	<ul style="list-style-type: none"> • Normative vs comparative explanations [52] 	
	Explainability	<ul style="list-style-type: none"> • Post-hoc explanations • Post-hoc explanations 	✓	✓	✓	✓	Agnostic			
[N] Explanation by simplification [18, 98]	Explainability	<ul style="list-style-type: none"> • Post-hoc explanations 	✓	✓	✓	✓	<ul style="list-style-type: none"> • Decision rule • Decision tree 			
	Feature relevance						<ul style="list-style-type: none"> • Feature attribute • Feature shape • Feature interaction • Sensitivity / perturbation-based • Saliency maps (visual domain) 	<ul style="list-style-type: none"> • Bar charts • Visualization of element importance, saliency (visual domain) 	<ul style="list-style-type: none"> • Usability of saliency maps for non-experts [7]. They should be accompanied by global descriptors 	
[O] Feature relevance explanation [7, 18, 24, 49, 98, 118]	Explainability	<ul style="list-style-type: none"> • Post-hoc explanations 	✓	✓	✓	✓	<ul style="list-style-type: none"> • Agnostic 			
	Contrastive explanations [47, 118, 134]						<ul style="list-style-type: none"> • Example of minimum change that leads to different outcomes 			
[P] Text-based explanation [18, 175]	Explainability	<ul style="list-style-type: none"> • Post-hoc explanations 	✓	✓	✓	✓	<ul style="list-style-type: none"> • Agnostic 			
	Text-based explanation [18, 175]	<ul style="list-style-type: none"> • Post-hoc explanations 	✓	✓	✓	✓	<ul style="list-style-type: none"> • With or without outcome comparison 			

1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702

Means	Value	Manifestation(s)	Stakeholder				Application (model)	Approach	Visual elements	Additional details
			DT	AT	DE	DS				
Interactive demonstrations [R] [120]	Explainability	<ul style="list-style-type: none"> Post-hoc explanations 				✓	Agnostic			
Experiential AI [82]	Explainability	<ul style="list-style-type: none"> Post-hoc explanations 				✓	Agnostic	<ul style="list-style-type: none"> Art mediated between computer code and human comprehension 		
Interactive contestations [T] [84, 106]	Contestability	<ul style="list-style-type: none"> Mechanisms for users to ask questions and record disagreements with system behavior 						<ul style="list-style-type: none"> Statements restricted to natural language 		
Challenge justifications provided by operator using the same means [84]	Human Control	<ul style="list-style-type: none"> Ability to override the decision made by the system 				✓	Agnostic			
	Human agency	<ul style="list-style-type: none"> Opportunity to self-assess the system 								
Mapping of actors and tasks depending on automation level [33]	Contestability	<ul style="list-style-type: none"> Mechanisms for users to ask questions and record disagreements with system behavior 				✓	Agnostic	<ul style="list-style-type: none"> Further testing Verification 		
	Human Control	<ul style="list-style-type: none"> Establishment of levels of human discretion during the use of the system 				✓	Agnostic		<ul style="list-style-type: none"> Relationship diagrams 	

1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743

Means	Value	Manifestation(s)	Stakeholder				Application (model)	Approach	Visual elements	Additional details
			DT	AT	DE	DS				
Failure Modes and Effects Analysis [W] [147]	Security	<ul style="list-style-type: none"> Threats against integrity (adversarial learning) and mitigation techniques 								
	Fairness	<ul style="list-style-type: none"> Accuracy across groups False positives and negatives across groups 	✓				Agnostic			
Aequitas [X] 33 [151]	Fairness	<ul style="list-style-type: none"> Accuracy across groups False Positive and Negative rates across groups False Discovery and Omission rates across groups Counterfactual examples 								
	Non-discrimination	<ul style="list-style-type: none"> Analysis of data for potential biases, data quality assessment 								
AI Fairness [Y] 360 34 [19]	Performance	<ul style="list-style-type: none"> False Positive and Negative rates 								
	Fairness	<ul style="list-style-type: none"> False positive and negative rates across groups Debiasing algorithms 							<ul style="list-style-type: none"> Bar charts Confidence bars 	
	Non-discrimination	<ul style="list-style-type: none"> Analysis of data for potential biases, data quality assessment 							Classifiers: logistic regression, random forest classifier and neural networks	

³³<https://github.com/dsrg/aequitas>

³⁴<https://github.com/Trusted-AI/AIF360>

1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781
1782
1783
1784

Means	Value	Manifestation(s)	Stakeholder			Application (model)	Approach	Visual elements	Additional details
			DT	AT	DE				
Fairlearn [Z] ³⁵ [25]	Performance	<ul style="list-style-type: none"> • Accuracy • False Positive and False Negative rates • Precision and recall rates 	✓	✓		Agnostic		<ul style="list-style-type: none"> • Bar charts • Pie charts 	
	Fairness	<ul style="list-style-type: none"> • Accuracy across groups • False negative and false positive rates across groups • Debiasing algorithms 							
Playbook [AA] AI ³⁶ [90]	Human agency	<ul style="list-style-type: none"> • Give knowledge and tools to comprehend and interact with AI systems • Opportunity to self-assess the system 			✓	NLP	Early AI prototyping	<ul style="list-style-type: none"> • Interactive survey 	
Counterfit [AB] ³⁷	Security	<ul style="list-style-type: none"> • Defence against integrity threats • Defence against privacy threats 			✓			Agnostic	
InterpretML [AC] ^{38 39} [134, 137]	Explainability	<ul style="list-style-type: none"> • Interpretability by design • Post-hoc explanations 			✓			Both white-box and blackbox models	<ul style="list-style-type: none"> • Bar charts • Line charts • Decision trees

³⁵<https://github.com/fairlearn/fairlearn>
³⁶<https://github.com/microsoft/HAXIPlaybook>
³⁷<https://github.com/Azure/counterfit>
³⁸<https://github.com/interpretml/interpret>
³⁹<https://github.com/interpretml/DICE>

1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825

Means	Value	Manifestation(s)	Stakeholder			Application (model)	Approach	Visual elements	Additional details
			DT	AT	DE				
Error analysis dashboard ⁴⁰	Non-discrimination	• Analysis of data for potential biases, data quality assessment							
	Explainability	• Post-hoc explanations	✓	✓		Agnostic		• Decision tree • Error heatmap	
Impact tracker ⁴¹ [83]	Fairness	• Accuracy across groups							
	Performance	• Estimation of energy consumption • Estimation of GPU memory consumption						• Dot plots • Bar charts	
Representative contestations [174]	Respect for public interest	• Measure of environmental impact							
	Contestability	• Mechanisms for users to ask questions and record disagreement with system behaviour				✓	Agnostic		

Table 6. Mapping of available means for transmitting value-specific manifestations to different stakeholders based on the purpose of their insight and the nature of their knowledge (DT = Development Team; AT = Auditing Team; DE = Data Domain Experts; DS = Decision Subjects). The identification and color code correspond to those on table 5. Each means is linked to the value and criteria manifestations that they communicate, the stakeholders that the original papers address, model specificity, deployed approach, visual elements and any additional details.

⁴⁰<https://github.com/microsoft/responsible-ai-toolbox/blob/main/docs/erroranalysis-dashbord-README.md>

⁴¹<https://github.com/Breakend/experiment-impact-tracker>