

Mapping the Potential of Explainable AI for Fairness Along the AI Lifecycle

Luca Deck^{1,2}, Astrid Schomäcker¹, Timo Speith¹, Jakob Schöffner³, Lena Kästner¹ and Niklas Kühl^{1,2}

¹University of Bayreuth, Germany

²Fraunhofer FIT, Germany

³University of Texas at Austin, USA

Abstract

The widespread use of artificial intelligence (AI) systems across various domains is increasingly surfacing issues related to algorithmic fairness, especially in high-stakes scenarios. Thus, critical considerations of how fairness in AI systems might be improved—and what measures are available to aid this process—are overdue. Many researchers and policymakers see explainable AI (XAI) as a promising way to increase fairness in AI systems. However, there is a wide variety of XAI methods and fairness conceptions expressing different desiderata, and the precise connections between XAI and fairness remain largely nebulous. Besides, different measures to increase algorithmic fairness might be applicable at different points throughout an AI system's lifecycle. Yet, there currently is no coherent mapping of fairness desiderata along the AI lifecycle. In this paper, we distill eight fairness desiderata, map them along the AI lifecycle, and discuss how XAI could help address each of them. We hope to provide orientation for practical applications and to inspire XAI research specifically focused on these fairness desiderata.

Keywords

Explainable AI, Algorithmic Fairness, Fairness Desiderata, AI Lifecycle

1. Introduction

The emergence and widespread use of artificial intelligence (AI) systems across various sectors and domains is increasingly shifting attention from considerations of mere performance to considerations about algorithmic fairness. This is particularly relevant for systems employed in high-stakes scenarios and especially pressing in contexts prone to harmful societal biases [1, 2]. This has sparked a growing demand for approaches to scrutinize and improve fairness in AI systems. In the literature, there is a common recognition that fairness in AI systems demands various perspectives and measures (e.g., [3, 4, 5]). In this paper, we set out to integrate two strategies that have been suggested: First, a growing community of researchers proposes explainable AI (XAI) as a versatile and powerful tool to combat unfairness [6]. Second, others have focused on the AI lifecycle and tried to determine where fairness issues originate (e.g.,

EWF'24: European Workshop on Algorithmic Fairness, July 01–03, 2024, Mainz, Germany

✉ luca.deck@uni-bayreuth.de (L. Deck); astrid.schomaecker@uni-bayreuth.de (A. Schomäcker); timo.speith@uni-bayreuth.de (T. Speith); schoeffner@utexas.edu (J. Schöffner); lena.kaestner@uni-bayreuth.de (L. Kästner); kuehl@uni-bayreuth.de (N. Kühl)



© 2024 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).



CEUR Workshop Proceedings (CEUR-WS.org)

[7, 8]). They hope that once identified, fairness issues can be mitigated by taking appropriate steps in the relevant phase within the lifecycle.

While both these strategies seem intuitively promising, neither currently presents a satisfactory and comprehensive picture. A primary reason why neither approach currently fulfills their potential is, we take it, that there are various types and kinds of fairness discussed in the literature, yielding different *fairness desiderata*. Differentiating between these desiderata is crucial to gain an overall picture of which measures are most promising to address fairness in which contexts. Without such differentiation, the utility of XAI is not as straightforward as commonly claimed. Accordingly, there is a need for more clarity about how, exactly, XAI contributes to which fairness desideratum [6]. Existing discussions on this matter suffer, by and large, from both a too narrow conception of XAI and a mostly under-specified notion of fairness. Similarly, existing attempts to map measures for achieving fairness onto the AI lifecycle remain crucially incomplete, for they have limited their attention to only a subset of the relevant fairness desiderata. To overcome these limitations, we distill eight fairness desiderata, map them along the AI lifecycle, and discuss how XAI can help address each of them.

For the purposes of this paper, we will focus specifically on how XAI can be utilized to improve AI systems' fairness throughout their lifecycle. Concerning this question, we aim to develop a holistic account based on the fairness desiderata we distill. While our proposal is not intended as a guideline, we believe it will stimulate scientific discussions about and practical applications of fairness measures and, as such, will be a valuable starting point to explore fairness opportunities for researchers, developers, and regulators alike.

We begin by introducing preliminaries and diagnosing the problem in Section 2. Importantly, we do not limit our understanding of algorithmic fairness to computational perspectives. In Section 3 we propose eight fairness desiderata from interdisciplinary literature: fairness understanding, data fairness, formal fairness, perceived fairness, fairness with human oversight, empowering fairness, long-term fairness, and informational fairness (contribution #1). We map our fairness desiderata along the different stages within the AI lifecycle and suggest that each fairness desideratum affords a different entry point for taking measures to improve fairness. This closes the first gap in the current literature as it highlights where in the lifecycle different fairness desiderata become especially relevant (contribution #2). Utilizing this mapping, we discuss how XAI can be leveraged to address the different fairness desiderata at the respective points throughout the AI lifecycle. This closes the second gap in the current literature as it allows us to systematically examine the potential of XAI to address algorithmic fairness in different circumstances (contribution #3). To illustrate the utility of our approach, we discuss its application to the COMPAS case (see [9]) in Section 4. Before closing, we point to some avenues for future research in Section 5.

2. Background

In this article, we contribute to closing two research gaps: 1) considerations of fairness along the AI lifecycle are incomplete, and 2) it is unclear how exactly XAI can help foster fairness. Before closing these gaps, we retrace contributions and debates of prior works.

2.1. Algorithmic Fairness

The debate on algorithmic fairness draws inspiration from various disciplines. Several scandals surrounding AI systems that disadvantage marginalized groups [9, 10, 11] have shattered early hopes put into the “neutrality” of AI [1]. In reaction, scholars from computer science, philosophy, social science, law, and psychology have been engaging in debates on what it means for algorithmic decision-making to be fair (see [12] for interdisciplinary perspectives).

The technical debate on algorithmic fairness primarily focuses on formal measures of fairness [13, 14, 15]. After it has proven unhelpful to remove information about membership in marginalized groups from training data (“fairness through unawareness” [16]), most approaches from the field of computer science rely on comparing the outcomes for people from different groups. An important finding is that the different formal measures for fairness can, in most cases, not be fulfilled simultaneously [17, 18]. Thus, (formally) fair AI systems require crucial design choices about which formal measures to apply in which context.

To aid such decisions, philosophers have connected formal measures to different philosophical theories [19, 20, 21]. In general, “fairness” is a normative concept—and debates about fairness can be understood as a way to discuss what is morally right or wrong in a given case (see [12]). More specifically, fairness can be understood as an issue of justice, non-discrimination, or equality [19], which connects it to numerous strands of philosophical discussion.

Beyond that, social scientists have pointed out shortcomings of existing formal methods, e.g., problems in detecting issues of intersectionality [22] and their reliance on constructed categories like race and gender [23]. Legal scholars are primarily interested in how legislation on discrimination can be applied to algorithmic decision-making and whether additional regulation is required [24, 1]. Psychological research is generally interested in whether algorithmic decision-making is perceived as fair [4, 25]. For example, Colquitt and Rodell [26] distinguish four different psychological dimensions of fairness: Whether the outcome is fair (viz., distributive justice), whether the process that leads to the outcome is fair (viz., procedural justice), whether the information about the decision is communicated truthfully and thoroughly (viz., informational justice), and whether individuals are treated respectfully (viz., interpersonal fairness).

Overall, algorithmic fairness can be addressed from several different angles—some of which are complementary, some of which are contradictory [27, 18]. This stresses the need to be clear about the specific desideratum behind any attempt to improve fairness.

2.2. Considerations of Fairness Along the AI Lifecycle Remain Incomplete

A common strategy to determine potential entry points to address fairness issues involves examining each step of the AI lifecycle. There have been several proposals of prototypical AI lifecycles (e.g., [28, 29, 30, 31, 32]), each with slightly different stages or adapted to different application areas. The AI lifecycle used in this article is a combination of the proposals by Quemy [31] and Wang et al. [30] involving the following stages: 1) problem formulation, 2) data collection, 3) data analysis, 4) feature selection, 5) model construction, 6) model evaluation, 7) deployment, and 8) inference and usage (see Figure 1).

First, *problem formulation* involves abstraction and formulation of the problem, which is to be solved using AI. Afterward, *data collection* aims to establish a representative data set suitable

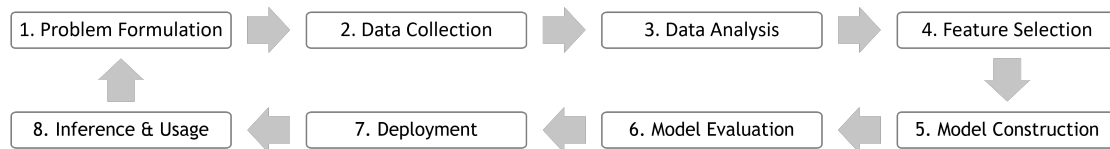


Figure 1: The AI Lifecycle we use, combined from Quemy [31] and Wang et al. [30].

to develop an AI system able to solve the defined problem. *Data analysis* includes descriptive statistics to understand the characteristics of the data at hand and data pre-processing to prepare the data for further operations. Subsequently, in the *feature selection* stage, features are excluded, transformed, or aggregated for effective and efficient handling in the training process. *Model construction*, then, includes the selection of the training algorithm and the training itself. Iteratively, performance objectives are tested during *model evaluation* and optimized through algorithm adjustments, e.g., parameter tuning. *Deployment* refers to the integration of the AI system into a productive environment once the model achieves sufficient performance. Finally, *inference & usage* describes concrete output generated by the AI system and its impact on the surroundings, such as the business context, society, and environment. Note that these stages are often not entered sequentially, leaving room for iterations and loops between stages [33].

Much research has focused on what types of bias can emerge in or affect different steps of the AI lifecycle (e.g., [3, 8, 34, 35]). For example, Singh et al. [35] take the CRISP-DM model [28] and review the types of bias that can occur at different stages of this model. However, these considerations remain incomplete, as most of these works have a strong technical focus and deal only with formal notions of fairness, e.g., by providing guidance for choosing fairness metrics or debating the role of potential fairness-utility tradeoffs (e.g., Castelnovo et al. [36], De-Arteaga et al. [3]). Against this background, it is important to develop a view that also maps other conceptions of fairness into the lifecycle.

2.3. Considerations of XAI for Fairness Remain Incomplete

To amend the lack of understanding of AI-based systems, their reasoning processes, and their outputs, the research field of XAI emerged in recent years [37, 38, 39, 40, 41]. Generally, XAI aims to provide a stakeholder with information about some aspect of an AI system to facilitate their understanding of this aspect [39]. Understanding (some aspect of) a system is often just an intermediary step to other goals, such as fairness or appropriate reliance [42, 43]. For example, by understanding how a particular system’s output came to be, prior work has argued that a person should be enabled to assess whether this output was based on valid criteria or not [38, 44]. If an unfavorable decision was based on (a proxy for) the skin color of a person, this person should supposedly be able to recognize that they were treated unfairly [45, 46].

However, the goal of XAI is often narrowly construed as the development of methods to create interpretable surrogate models of black box models [47, 48]. Such a narrow view of XAI artificially restricts the space of options that can be chosen to facilitate stakeholders’ understanding of AI systems. Take, e.g., model cards for model reporting [49] or datasheets for data sets [50]. Although these approaches provide important information about certain

aspects of an AI model¹ (e.g., which training data was used) and thus improve stakeholders’ understanding of these aspects, they would rarely be counted as XAI. Further, XAI methods alone do not provide other kinds of highly relevant information, such as the normative motivations guiding the development of an AI model [52], the analysis of the social context in which it is deployed [53], or the descriptive outcome statistics of the deployed model [54].

In a similar vein, the high hopes placed in XAI to mitigate issues of fairness often remain vague or unfulfilled [6]. Prior works have criticized the explanatory value [47, 55], susceptibility to manipulations [56, 57], and unsatisfactory interpretations [58, 6] of XAI methods. For example, popular feature-based explanations like LIME [59] or SHAP [60] highlighting the use of sensitive features (e.g., gender and race) provide little information about fairness. This is because these features are often correlated with proxy variables and embedded in use case-specific normative contexts [16, 61, 62].

Instead, broadening the conception of XAI enables a more holistic and meaningful mapping of how information about an AI system can contribute to various fairness desiderata. Against this background, we propose to distinguish the *narrow* conception of XAI, solely pertaining to the inner workings of an AI system, from a *broad* conception of XAI that provides information beyond these inner workings and includes explanations of the broader socio-technical system [63]. Specifically, we suggest that the broad conception of XAI includes all types of information that increases stakeholders’ understanding of (aspects of) an AI system. For the remainder of this paper, if not indicated otherwise, we refer to the broad conception when we mention XAI.

3. How XAI Can Be Leveraged for Fairness Along the AI Lifecycle

Grounded in a broad interdisciplinary literature review (summarized in Table 1), we distill eight categories of what previous work on fairness has called for. We call these categories *fairness desiderata*, and each desideratum can be instantiated with specific *fairness objectives*. For the purposes of this paper, we focus specifically on fairness desiderata connected to XAI. While various stakeholders are involved in utilizing XAI for fairness, an extensive discussion of their various roles is beyond the scope of this paper. When we mention stakeholders, we rely on the taxonomy provided by Langer et al. [40].

In what follows, we introduce our eight fairness desiderata, describe the underlying core ideas for each, discuss how they relate to similar concepts, map them onto the AI lifecycle, and elaborate on how we think XAI could contribute to their satisfaction (see Figure 2). Note that our proposal does not necessarily present a comprehensive list (see Section 5).

3.1. Fairness Understanding

“Accounting for bias not only requires an understanding of the different sources, that is, data, knowledge bases, and algorithms, but more importantly, it demands the interpretation and description of the meaning, potential side effects, provenance, and context of bias.” [2, p. 8]

¹Throughout this paper, we use the term “AI models” synonymously for “machine learning model” or “ML model”. While we acknowledge the technical distinction between AI and Machine Learning (ML) as discussed in [51], we adopt the use of the broader “AI” term as the prevalent terminology established in the scientific community. This choice reflects the contemporary linguistic trend rather than a lack of distinction between the two phenomena.

Table 1

Fairness desiderata and their related concepts in existing interdisciplinary literature.

Fairness desiderata	Related concepts
Fairness understanding Gaining higher-level insights on fairness and the socio-technical challenges surrounding the development and deployment of an AI application to specify concrete fairness objectives.	Understanding bias [2] Interdisciplinary fairness conceptualization [12] Lessons from political philosophy [19] Socio-technical perspective [64]
Data fairness Identifying and addressing flaws in the data set that might be unfair themselves or potentially lead to downstream violations of fairness objectives.	Sampling bias [65, 3] Data errors and bias [66] Data-centric factors in algorithmic fairness [67]
Formal fairness Identifying and addressing model properties leading to violations of formal fairness objectives.	Algorithmic bias [3] Formal fairness definitions [68] Fairness metrics [15] Statistical fairness criteria [13] Disparate impact & disparate treatment [1]
Perceived fairness Providing affected parties with explanations and justifications to improve or “calibrate” fairness perceptions.	Fairness perceptions [4, 69] Perceptions of justice [70, 71] Fairness judgment [72] Society-in-the loop [73]
Fairness with human oversight Supporting human decision-makers interacting with an AI system to effectively align human discretion with fairness objectives.	Human-ML augmentation for fairness [74] Appropriate reliance [75] Human-in-the-loop [72] Domain expertise [76]
Empowering fairness Providing affected parties with practical information to foster contestability and recourse.	Self-informed advocacy [77] Procedural fairness [78] Counterfactual explanations [79] Fair and adequate explanations [80] The explanation game [81]
Long-term fairness Monitoring and analyzing the socio-technical long-term impacts of algorithmic decision-making to adjust unfair repercussions over time.	Long-term effects of algorithmic fairness [3] Fairness drift [82] Fairness through time [83] Fairness monitoring [84, 85]
Informational fairness Providing truthful, understandable, and relevant information about all fairness desiderata across the AI lifecycle.	Informational fairness [70] Design publicity [52] Model cards [49] Outward transparency [86]

Before mitigating algorithmic unfairness, developers should reflect upon the multidimensional and conflicting nature of fairness, define concrete fairness objectives, and understand how these might be achieved. In light of this, *fairness understanding* relates to the specification of concrete fairness objectives and knowledge about the socio-technical challenges surrounding the development and deployment of an AI application.

Among others, this may include a sound understanding of fairness schools of thought [12], existing social inequalities [53], stakeholder-centered fairness requirements [40], or relevant legal frameworks [87]. As the following fairness desiderata will show, there is no one-size-fits-all way in which a system can be fair, both across (e.g., a system can be perceived as fair despite not being formally fair) as well as within fairness desiderata (e.g., formal fairness criteria pose inherent tradeoffs [27]).

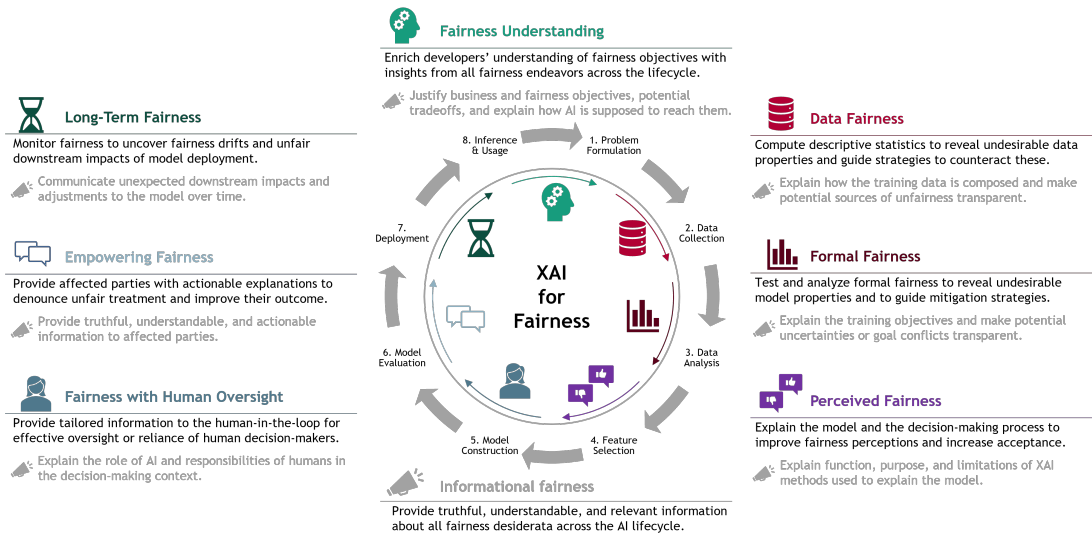


Figure 2: Fairness desiderata along the AI lifecycle depicting how XAI may directly contribute to these desiderata and how it may contribute to informational fairness across all desiderata (symbolized by speakerphones).

Stage of the Lifecycle. In the AI lifecycle, we map fairness understanding to the initial problem formulation step because all subsequent steps are guided by the fairness objectives discussed and defined in the early stages. We note, however, that fairness understanding is developed iteratively, based on trials and errors, interdisciplinary exchange, and stakeholder feedback from various development stages—similarly to how a traditional business problem can be better understood and adjusted over the entire course of an AI project.

The Role of XAI. XAI can contribute to a better understanding of fairness across all fairness desiderata which, e.g., may lead to a re-evaluation of earlier fairness objectives. Data-centric explanations [88] might reveal pre-existing group disparities (e.g., in the form of differing base rates such as the gender pay gap), which are crucial factors in determining the fairness objectives and how to achieve them. Further, simulating and testing various formal fairness metrics promotes an understanding of conflicting objectives that may lead to adjustments or re-weightings of fairness objectives [3]. Similarly, XAI may change stakeholders' view of proactively using sensitive features in a specific decision-making context [62]. Concerning “substantive” fairness [53, 21], XAI may also be applied to normatively gauge the legitimacy of certain features depending on the societal context.

3.2. Data Fairness

“Explanations [...] are crucial for helping system developers and ML practitioners to debug ML algorithms by identifying data errors and bias in training data, such as measurement errors and misclassifications, data imbalance, missing data and selection bias, covariate shift, technical biases introduced during data preparation, and poisonous data points injected through adversarial attacks.” [66, p. 248]

Undesirable model behavior often stems from flawed data that contains misrepresentations of the world (e.g., erroneous, mislabeled, or imbalanced data) or accurate representations that are societally undesirable (e.g., historical inequalities). The goal of *data fairness* is to identify and address such flaws in the data set used to train an AI model.

Prior studies have pointed out data issues as drivers of unfairness [89, 65, 67] or highlighted data as a starting point for unfairness mitigation [90, 91, 92]. Mehrabi et al. [65] provide a comprehensive list of data biases that may introduce unfairness early on in AI development. Awareness and understanding of these biases are crucial to derive strategies to handle them.

Stage of the Lifecycle. Data fairness relates to data collection (e.g., in the form of labeling errors, sampling bias, imbalances, etc.) and data analysis, which aims to identify and potentially mitigate unfair data characteristics early on. However, data fairness also reaches into feature selection which is usually part of an iterative loop together with model construction. For example, features might be dropped if they are not justifiably task-relevant. Importantly, in the context of sensitive attributes, developers should always be aware of the flaws of the idea of “fairness through unawareness” [16, 62].

The Role of XAI. As data is a main source of unfairness [65], XAI is suited to identify potential disparities, imbalances, or abnormalities manifested in the available data early on. Descriptive statistics are a natural first step to explore pre-existing disparities [67]. Anik and Bunt [88] and Mitchell et al. [49] provide two examples of how to present simple descriptions and visualizations about the collection, feature distributions, and patterns of a dataset. XAI techniques have further been claimed to reveal instances and features in the data that have undesirable effects on the model output [93, 94, 66] which, however, are often subject to questionable causal and normative assumptions [95, 96]. These approaches indicate a strong connection between data fairness and *formal fairness*.

3.3. Formal Fairness

“To counteract biases, it is, therefore, crucial to enable their detection. Explainability approaches may aid in this regard by providing means to track down factors that may have contributed to unfair and unethical decision-making processes and either to eliminate such factors, to mitigate them, or at least to be aware of them.” [40, p. 6]

The most common fairness desideratum is concerned with formal model properties. By *formal fairness*, we refer to the vast array of formal criteria that have been proposed as mathematical and statistical measures of fairness [13, 15].

Consistent with Verma and Rubin [14], this includes all fairness definitions based on statistical measures (e.g., demographic parity), similarity measures (e.g., fairness through awareness), or based on causal reasoning (e.g., counterfactual fairness). Formal fairness notions are often distinguished into group and individual fairness. Group fairness criteria typically require a form of parity between demographic groups, e.g., along sensitive attributes like gender or race [54]. Individual fairness criteria typically demand to treat similar people alike [16].

Stage of the Lifecycle. Formal fairness is particularly relevant for the iterative loop of model construction and model evaluation. As soon as the first model prototype is ready, fairness metrics can be evaluated. Based on the evaluation, unfairness mitigation techniques can be implemented and validated iteratively [97].

The Role of XAI. Regarding formal fairness, XAI could be of exploratory value, offering a plethora of tools and a novel perspective from which to explore formal fairness from multiple angles. Again, descriptive statistics present a natural first step to identify potential disparities, imbalances, or abnormalities outputted by the model [54, 98]. This is especially useful if formal fairness objectives are already specified, e.g., when testing whether formal group fairness metrics are sufficiently satisfied [97]. XAI might also point towards specific features driving the violation of group fairness metrics [99] or shed light on more subtle forms of formal fairness such as fairness of recourse [100]. XAI has further been claimed to elucidate the complex interplay between “task-relevant” and correlated “protected” features (e.g., [101]) which, again, relies on causal and normative assumptions [95, 96]. Technical possibilities of XAI to explore formal fairness are manifold but require utmost caution when applied for unfairness mitigation in specific application contexts and should not be misinterpreted as fairness “proofs” or “guarantees” [6]. Particularly, when interpreting the legitimacy of using sensitive information, it is paramount to account for the “fairness through unawareness” fallacy [16, 61] and for differential subgroup validity [102, 103].

3.4. Perceived Fairness

“We argue that only fair systems that are also perceived to be fair by their users should be accepted and employed in practice.” [104, p. 5]

Another common fairness desideratum refers to positive perceptions of stakeholders, which is often related to trust and acceptance [104, 105]. In this sense, *perceived fairness* captures how stakeholders, particularly those affected by a decision, perceive the fairness of the AI system.

Measures of perceived fairness can be derived from the justice constructs of Colquitt [70], which decompose fairness perceptions into a procedural, distributive, and informational dimension (see also [71, 72, 69]). Accounting for fairness perceptions promotes a valuable means to design AI systems based on human needs and ideals while giving voice to societal values, which Rahwan [73] coined as “society-in-the-loop”. Notably, the desirable aspect of perceived fairness does not necessarily or solely lie in positive fairness perceptions but in “appropriate” fairness perceptions [106].

Stage of the Lifecycle. Although fairness perceptions can be relevant at any stage of the lifecycle (e.g., to evaluate data fairness [88]), we map fairness perceptions to the model evaluation, deployment, and inference & usage stage. Most prior works have measured fairness perceptions after a (mock-up) model has been developed to ask stakeholders about certain outcomes or the AI system itself [4]. Beyond that, there are approaches trying to embed stakeholders’ values throughout the conception and development of AI systems (e.g., [107, 108, 109, 110]).

The Role of XAI. Explanations may serve stakeholders as a cue to confirm whether their ideas of fairness are implemented in a system. This idea has also been commonly expressed in prior literature [111, 104, 112, 105]. For example, Lee et al. [113] indicate that explanations can decrease fairness perceptions when they reveal information that stands in conflict with peoples’ fairness beliefs. However, the effect of XAI on fairness perceptions is highly context-dependent and moderated by human factors like political ideology and self-interest [114]. This indicates that stakeholders require tailored information addressing their case-specific concerns. Affected

parties, e.g., seem to appreciate information beyond a model’s inner workings, such as system context, usage, and data [115]. We note that optimizing XAI to stimulate positive fairness perceptions isolated from complementary desiderata (e.g., trustworthiness or formal fairness) can lead to undesirable effects such as placebic explanations [116], fairwashing [56] or deception [117, 118].

3.5. Fairness With Human Oversight

“Research suggests that neither humans nor ML models are likely to achieve fairness working alone. Instead, human–ML augmentation, where humans and technology work together to perform organizational tasks jointly, is the most promising path to achieving fairness.” [74, p. 2]

Beyond the (formal) fairness of a model itself, fairness can also refer to the decision-making process in which the model is embedded. *Fairness with human oversight* aims at installing and supporting a human decision-maker to realize case and context-sensitive fairness objectives through human oversight or human discretion.

The human-AI setting can take various forms but usually involves overseeing and overruling unfair outputs or fostering effective reliance behavior [74]. Although legal and ethical guidelines often demand human oversight [78, 119], the effect of human oversight on fairness is not necessarily beneficial and not well understood yet [86, 120, 74, 43]. Accordingly, this desideratum does not necessarily capture fairness *through* human oversight, but potentially also fairness *despite* human oversight.

Stage of the Lifecycle. Fairness with human oversight becomes relevant after the model has been deployed, i.e., during inference & usage.

The Role of XAI. Where formal implementation of fairness during model development is difficult or human oversight is required, XAI is crucial to inform human discretion so that fairness objectives can still be realized. In this sense, XAI is commonly proposed to support human decision-makers and domain experts in fostering fairer decisions (e.g., [109, 121, 122]). Beyond simple recommendations, such information may include a comparison to similar instances [76], disclosure of uncertainty [123], or conditional heatmaps in computer vision tasks [124]. To tackle automation bias, Miller [125]’s concept of “evaluative AI” also highlights providing not only explanations for a certain recommendation but rather balanced evidence for multiple possible outcomes. However, both designing XAI and training human decision-makers to interact with the explanations is challenging [126]. As feature importance explanations may even hinder fairness of human decisions [43], we are in need of more conceptual and empirical research on how to design XAI towards fairness objectives in human-in-the-loop settings (e.g., how to effectively override certain types of AI recommendations that violate fairness objectives).

3.6. Empowering Fairness

“From the perspective of individuals affected by automated decision-making, we propose three aims for explanations: (1) to inform and help the individual understand why a particular decision was reached, (2) to provide grounds to contest the decision if the outcome is undesired, and (3) to understand what could be changed to receive a desired result in the future, based on the current decision-making model.” [79, p. 2]

Fairness considerations do not end after an AI model has been designed and decisions have been made. *Empowering fairness* refers to the ability of affected parties to take effective actions regarding the outcome of a particular decision, e.g., by contesting decisions or seeking recourse.

In her article on the right to explanation, Vredenburg [77] proposes two types of *self-informed advocacy*: retrospectively, affected parties should be able to identify the accountable entity to demand remedy for unfair treatment (viz., responsibility); prospectively, contesting and recourse options should enable affected parties to actively improve upon their possibly unfair outcome (viz., agency). Our conception of empowering fairness strongly relates to Vredenburg’s forward-looking self-informed advocacy. It more generally relates to attempts to conceptualize what makes a *fair explanation* in the context of the right to explanation [87].

Stage of the Lifecycle. Recourse and contesting is only possible after a decision has been made. Accordingly, we map empowering fairness to inference & usage.

The Role of XAI. XAI could be useful for empowering fairness because it may help acquire the information to be communicated. Contrastive explanations (answering the question “Why P rather than Q?”) are commonly proposed to provide intuitive entry points to engage with affected parties [79, 127]. For instance, counterfactual explanations have been claimed to provide a promising solution to inform and empower affected parties by clarifying what combination of feature values would lead to a different outcome. [79, 80]. Therefore, they can provide valuable information required for algorithmic recourse, e.g., by providing actionable recommendations to a loan applicant on what to do to be granted a loan in the future [128, 129]. While research on counterfactual explanations is growing rapidly, they do not come without limitations regarding their actionability [128], validity [130], underlying assumptions [131], or susceptibility to manipulations [132]. Extending the scope beyond model-centered explanations, XAI might also involve practical information about responsible contact persons, guidance on how to seek redress, or collaborative platforms for affected parties to share experienced outcomes [117].

3.7. Long-Term Fairness

“Accuracy, discrimination, and security characteristics of a system can change over time as well. Simply testing for these problems at training time [...] is not adequate for high-stakes, human-centered, or regulated ML systems. Accuracy, discrimination, and security should be monitored in real-time and over time, as long as a model is deployed.” [133, p. 18]

Fairness remains relevant over the entire AI lifecycle, even after deployment. Hence, *long-term fairness* captures the dynamic interplay of an AI system with the socio-technical system it is deployed in over time.

The long-term impact of an AI model can be analyzed from several perspectives. Arif Khan et al. [134] contrast “formal” and “substantive” equality of opportunity where a forward-facing view of fair life chances also accounts for affected parties’ future prospects of success (as opposed to ensuring fair contests at a discrete point in time). Since such conceptions make assumptions about structural disadvantages and future prospects of affected parties, they require much broader and longitudinal evaluation. Further, due to concept drift [135], a model’s formal fairness properties can change over time and should be monitored.

Stage of the Lifecycle. We map long-term fairness primarily onto inference & usage, noting that it may encompass all future iterations of the AI lifecycle until the AI system is shut down.

The Role of XAI. Monitoring tools may help to track changes in fairness metrics and identify situations where interventions are necessary [83, 85, 82]. Another long-term fairness impact is strategic gaming behavior that arises from transparent models [136]. For example, loan applicants who receive counterfactual explanations exactly describing how to reach the decision threshold are prone to create a game-theoretic situation where information itself can be unfairly distributed among clients. We suspect several other forms of long-term fairness issues will emerge that are currently under-explored in the literature. For example, Liu et al. [137] model the impact of formal fairness on the underlying population over time, which Hardt et al. [138] coined as “performative power”. Novel forms of XAI may help to anticipate and evaluate such dynamics.

3.8. Informational Fairness

“Transparency [...] is valuable because and in so far as it enables the individuals, who are subjected to algorithmic decision-making, to assess whether these decisions are morally and politically justifiable.” [52, p. 254]

All of the listed fairness desiderata can be augmented with a meta desideratum targeted at transparency *about* fairness, which we label as *informational fairness*.

Informational fairness was originally introduced as a psychological construct in the context of organizational justice [70] to test whether the communication accompanying a decision is candid, truthful, reasonable, timely, and specific. Our conception of informational fairness is inspired by this construct and corresponds to a great extent to Loi et al. [52]’s concept of design publicity. In line with our idea of broad XAI, Loi et al. [52] demand a form of transparency that explains not only a model and its underlying functioning but also the goals and values that went into its design and how these are embedded in the model. A similar distinction is made by Walmsley [86], who differentiates between functional transparency concerned with the inner workings of an AI model and outward transparency related to communication with stakeholders.

Stage of the Lifecycle. We conceive informational fairness as a meta desideratum that applies to all other fairness desiderata (see Figure 2). Accordingly, informational fairness can be considered across all stages of the AI lifecycle from problem formulation to inference & usage.

The Role of XAI. In the case of informational fairness, XAI is not an aid to but the desideratum itself. Following this idea, we provide examples of what could be made transparent and communicated to affected parties for each fairness desideratum. Regarding fairness understanding, developers and deployers could justify their fairness objectives, delineate potential tradeoffs, and elucidate the mechanisms through which AI aims to achieve them [52]. For data fairness, they could explain how the training data is composed, highlighting potential sources of bias [88]. Explaining the training objectives and outlining potential uncertainties or goal conflicts might be appropriate to address formal fairness [49]. Regarding perceived fairness, explaining the functions and limitations of XAI-generated information could enhance understanding to approach “appropriate fairness perceptions” [106]. Transparency about the role of AI and the responsibilities of a human decision-maker within the decision-making context might be desirable for fairness with human oversight [139]. The idea of empowering fairness relies on truthful, understandable, and actionable information [81]. Lastly, regarding long-term fairness,

unexpected downstream impacts and the factors driving potentially unfair dynamics can be communicated to stakeholders [121].

4. The COMPAS Case

To illustrate the utility of our mapping, we describe how XAI could help address fairness throughout the lifecycle of a high-stakes AI system. As an example, we consider the recidivism prediction software COMPAS developed by Northpointe (today rebranded to *equivant Supervision*). COMPAS has been installed in many US-American jurisdictions in order to predict whether a defendant will likely commit another crime in the near future. Judges use this information, e.g., for decisions about who they release on bail. However, COMPAS has been criticized by investigative journalists at ProPublica for disadvantaging Black people [9]. This is especially problematic given the history of systematic discrimination and marginalization of Black people in the US [140]. In this section, we illustrate both how Northpointe could have addressed different fairness desiderata during system development and how they could still, given COMPAS' continued use, address some of them. Notably, our illustration also presupposes an inherent motivation to actually address fairness desiderata, and XAI is not the only means to this end.

First, Northpointe could have used XAI to ensure that the development is based on an appropriate *fairness understanding*. ProPublica's analysis of COMPAS shows that although it was tested for some fairness objective (viz., predictive parity), it falls short on other fairness objectives (viz., equalized odds) in problematic ways [54]. Thus, predictive parity may have been an inadequate fairness objective. Northpointe could also have based its development on insights about the predictive value and correlations of certain features. It has been shown that defendants' age and number of previous crimes are most predictive for recidivism [141, 47] but also highly correlated to race.

To address *data fairness*, Northpointe would have needed to ensure that the data set adequately represents demographic groups targeted by the system and to be aware of existing structural relationships such as statistically higher crime rates in predominantly Black neighborhoods [19]. Descriptive statistics could have already pointed them towards biases introduced in the data collection process—a reason to reiterate the data collection stage to apply new data collection strategies [67]. Statistical analysis could also have helped unveil traces of systematic discrimination in the data (e.g. that Black people are more likely to be arrested for minor offenses due to increased police presence in Black neighborhoods). Northpointe could have used such insights either to change their data collection strategies (e.g., collecting features that are less correlated to race), or they could have used them during the feature selection phase (e.g., to mitigate systematic discrimination at the data-level [66]).

Regarding *formal fairness*, Northpointe tested COMPAS for equality of error rates (viz., predictive parity). By contrast, journalists at ProPublica tested for the balance of true and false positives (viz., equalized odds) [54]. Notably, many formal fairness metrics cannot be satisfied simultaneously [27], and selecting appropriate formal fairness metrics is an intricate endeavor [3]. However, had Northpointe tested COMPAS for equalized odds, too, they might have anticipated the social backlash and prevented some of the harmful impact. Today, developers

can draw on an extensive suite of fairness testing tools for a comprehensive formal fairness assessment [142, 97, 6].

Northpointe could have also detected flaws by engaging with stakeholders and gathering feedback on *perceived fairness*. Outcome statistics, fairness metrics, or model explanations (e.g., in the form of counterfactuals) could have been presented to a focus group of diverse stakeholders to assess and discuss fairness perceptions (see [109]). Potential concerns could be handled by reconciling the different perspectives, and even today (after deployment) Northpointe could continue to monitor fairness perceptions by establishing feedback channels.

Our discussion thus far has primarily focused on measures that COMPAS' developers could have taken early in the AI lifecycle. But even now, during COMPAS' continued usage, Northpointe could still address several fairness issues via XAI. First, they could seek *fairness with human oversight*. This is a particularly relevant desideratum because COMPAS does not make autonomous decisions but informs the decisions of judges. So, Northpointe could (if not already happening) familiarize judges with the basic functioning, strengths, limitations, and uncertainties of the system as well as the factors driving the risk score to ensure judges will use it appropriately.

Northpointe could also foster *empowering fairness* by providing defendants with two kinds of explanations: First, global explanations [72] about the general functioning of COMPAS could be communicated. Where relevant factors are under a person's control (as, e.g., with defendants' history of misdemeanor or substance abuse), insights about their impact can allow people to make life choices decreasing their risk scores. Second, local explanations [72] could provide information about an individual's risk score granting the opportunity to (justifiably) contest the score, e.g., by denouncing discriminatory treatment [118], or to (effectively) seek recourse, e.g., by signaling plans to go into rehab [79].

Moreover, Northpointe could monitor *long-term fairness* by tracking the specified formal fairness objective(s) over time and accordingly adjust the model over several iterations of the AI lifecycle [83]. Further, changes in defendants' behavior due to specified fairness objectives or increased transparency could be examined. For example, precise global or counterfactual explanations might allow defendants to game the system in unexpected ways strategically [136], which might—similar to unfair recourse [100]—be easier for some than for others.

Finally, Northpointe could have aimed and still could aim, to ensure *informational fairness*. To this end, they would have needed to document and communicate fairness-related information accordingly. For example, Northpointe could have used model (fairness) cards similar to [49], which would have clarified many of the issues that were only revealed due to the investigative work of Angwin et al. [9]. Most importantly, perhaps, Northpointe could have been candid about the underlying business and fairness objectives from the start, ideally making transparent decisions about tradeoffs (e.g., when looking at predictive parity rather than equalized odds) and justifying why an AI system was an appropriate tool for criminal risk prediction in the first place. In communication with affected parties, COMPAS could be complemented with truthful, understandable, and relevant information [80] explaining the underlying logic. Overall, the strategic use of XAI along the AI lifecycle could benefit all fairness desiderata regarding a model like COMPAS—given sincere intentions to actually pursue these desiderata.

5. Conclusion and Outlook

We distilled eight fairness desiderata from interdisciplinary literature, mapped them onto the AI lifecycle, and discussed how XAI might be leveraged to contribute to fulfilling them. Finally, we illustrated the utility of our approach by applying it to the COMPAS case. The overall picture we paint (see Figure 2) highlights that to design and utilize XAI for fairness effectively, it is paramount to reflect on which fairness desideratum one seeks to fulfill. Before closing, we shall briefly comment on some limitations of our proposal and point to avenues for future work.

Conceptualization and Validation. We do not claim our proposed fairness desiderata to be mutually exclusive or collectively exhaustive. As our discussion in Section 3 has shown, different terms elicit different associations in different communities and conceptual overlap between different fairness desiderata seems difficult to avoid altogether (e.g., long-term fairness to some extent includes formal fairness, etc.). Yet, we are confident that our application to the COMPAS case highlights the general usefulness of our work. Thus, we firmly believe that our eight fairness desiderata provide a valuable starting point for refinement in follow-up work. We expect that future interdisciplinary research can—through both conceptual work and validation by application to real-world cases—provide more sharply distinguished categories, capture more fine-grained distinctions, and incorporate an even broader spectrum of perspectives on fairness. Specific open issues include, e.g., where to locate procedural fairness in our account. For now, we deliberately excluded that notion as it has inherently differing meanings across (and sometimes even within) disciplines. Another open question concerns the status of privacy. Legal scholars might consider data privacy a form of fairness; conceived this way, privacy might qualify as a form data fairness.

Generalizability. Our discussion has focused primarily on high-risk applications and more traditional decision-support systems. To what extent does our proposal generalize to low-risk applications? And what about generative AI? Consider Google’s multi-modal AI-chatbot Gemini, which has recently received bad publicity for its inaccurate historical depictions—which were driven by misdirected fairness interventions [143]. We suspect that the usage of (seemingly low-risk) AI systems at scale will accumulate to significant societal impacts bearing more subtle threats to fairness such as, e.g., representational harms [89] in the Gemini case. Thus, although contemporary regulation (like the European AI Act) focuses predominantly on high-risk systems, fairness is an important concern for low-risk applications, too. Beyond that, the Gemini case illustrates that fairness considerations affect generative AI just as much as more traditional AI systems for decision-making. Given the rapid advancement and broad adoption of generative AI, we believe that fairness considerations will be indispensable in this context, though they only start to gain traction [144, 145].

Actionability. We acknowledge that our proposal does not present an actionable process model for fair software engineering; nor does it provide regulatory guidelines to be enforced by legal institutions. Still, it may serve as a useful roadmap for researchers, developers, and regulators. For researchers, it provides a starting point for a truly interdisciplinary discourse going beyond discussions of algorithmic fairness focusing on technical aspects of AI systems. For developers, our mapping provides urgently needed guidance as to what fairness challenges should be addressed when developing specific guidelines and requirements (see Hacker [146]). More specifically, it may help determine under which circumstances human oversight is benefi-

cial [74, 43] or how responsibilities can be attributed along the AI lifecycle [147]. Naturally, the implications of our approach will vary with the precise settings; so fine-tuning may be needed for different business, societal, and legal contexts.

References

- [1] S. Barocas, A. D. Selbst, Big Data's Disparate Impact, SSRN Electronic Journal (2016).
- [2] E. Ntoutsi, P. Fafalios, U. Gadiraju, V. Iosifidis, W. Nejdl, M.-E. Vidal, S. Ruggieri, F. Turini, S. Papadopoulos, E. Krasanakis, I. Kompatsiaris, K. Kinder-Kurlanda, C. Wagner, F. Karimi, M. Fernandez, H. Alani, B. Berendt, T. Kruegel, C. Heinze, K. Broelemann, G. Kasneci, T. Tiropanis, S. Staab, Bias in data-driven artificial intelligence systems—An introductory survey, *WIREs Data Mining and Knowledge Discovery* 10 (2020).
- [3] M. De-Arteaga, S. Feuerriegel, M. Saar-Tsechansky, Algorithmic fairness in business analytics: Directions for research and practice, *Production and Operations Management* 31 (2022) 3749–3770.
- [4] C. Starke, J. Baleis, B. Keller, F. Marcinkowski, Fairness perceptions of algorithmic decision-making: A systematic review of the empirical literature, *Big Data & Society* 9 (2022) 1–16.
- [5] A. Aler Tubella, D. Coelho Mollo, A. Dahlgren Lindström, H. Devinney, V. Dignum, P. Ericson, A. Jonsson, T. Kampik, T. Lenaerts, J. A. Mendez, J. C. Nieves, ACROCPoLis: A Descriptive Framework for Making Sense of Fairness, in: *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*, ACM Digital Library, Association for Computing Machinery, 2023, pp. 1014–1025.
- [6] L. Deck, J. Schoeffler, M. De-Arteaga, N. Köhl, A Critical Survey on Fairness Benefits of XAI, in: *ACM Conference on Fairness, Accountability, and Transparency (ACM FAccT '24)*, 2024.
- [7] S. Fazelpour, D. Danks, Algorithmic bias: Senses, sources, solutions, *Philosophy Compass* 16 (2021) e12760.
- [8] A. Agarwal, H. Agarwal, A seven-layer model with checklists for standardising fairness assessment throughout the ai lifecycle, *AI and Ethics* (2023) 1–16.
- [9] J. Angwin, J. Larson, S. Mattu, L. Kirchner, Machine Bias: There's software used across the country to predict future criminals. And it's biased against blacks., in: K. Martin (Ed.), *Ethics of data and analytics*, An Auerbach Book, CRC Press Taylor & Francis Group, Boca Raton and London and New York, 2022, pp. 254–264.
- [10] J. Dastin, Amazon Scraps Secret AI Recruiting Tool that Showed Bias against Women *, in: K. Martin (Ed.), *Ethics of data and analytics*, An Auerbach Book, CRC Press Taylor & Francis Group, Boca Raton and London and New York, 2022, pp. 296–299.
- [11] A. Fuster, P. Goldsmith-Pinkham, T. Ramadoral, A. Walther, Predictably unequal? the effects of machine learning on credit markets, *The Journal of Finance* 77 (2022) 5–47.
- [12] D. K. Mulligan, J. A. Kroll, N. Kohli, R. Y. Wong, This Thing Called Fairness, *Proceedings of the ACM on Human-Computer Interaction* 3 (2019) 1–36.
- [13] S. Barocas, M. Hardt, A. Narayanan, *Fairness and Machine Learning: Limitations and Opportunities*, fairmlbook.org, 2019.

- [14] S. Verma, J. Rubin, Fairness definitions explained, in: Proceedings of the International Workshop on Software Fairness, ACM Conferences, ACM, 2018, pp. 1–7.
- [15] A. Castelnovo, R. Crupi, G. Greco, D. Regoli, I. G. Penco, A. C. Cosentini, A clarification of the nuances in the fairness metrics landscape, *Scientific Reports* 12 (2022) 4209.
- [16] C. Dwork, M. Hardt, T. Pitassi, O. Reingold, R. Zemel, Fairness through awareness, in: Proceedings of the 3rd Innovations in Theoretical Computer Science Conference on - ITCS '12, ACM Press, 2012.
- [17] J. Kleinberg, S. Mullainathan, M. Raghavan, Inherent Trade-Offs in the Fair Determination of Risk Scores: 67(43): 1–23, *Leibniz International Proceedings in Informatics* 67 (2017) 1–23.
- [18] M. DeFrance, T. de Bie, Maximal fairness, in: Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency, ACM Digital Library, Association for Computing Machinery, 2023, pp. 851–880.
- [19] R. Binns, Fairness in Machine Learning: Lessons from Political Philosophy, *Conference on Fairness, Accountability and Transparency* 81 (2018) 149–159.
- [20] M. S. A. Lee, L. Floridi, J. Singh, Formalising trade-offs beyond algorithmic fairness: lessons from ethical philosophy and welfare economics, *AI and Ethics* 1 (2021) 529–544. URL: <https://link.springer.com/10.1007/s43681-021-00067-y>. doi:10.1007/s43681-021-00067-y.
- [21] A. A. Khan, S. Badshah, P. Liang, M. Waseem, B. Khan, A. Ahmad, M. Fahmideh, M. Niazi, M. A. Akbar, Ethics of AI: A Systematic Literature Review of Principles and Challenges, in: The International Conference on Evaluation and Assessment in Software Engineering 2022, ACM Digital Library, Association for Computing Machinery, 2022, pp. 383–392.
- [22] K. Crenshaw, Mapping the Margins: Intersectionality, Identity Politics, and Violence against Women of Color, *Stanford Law Review* 43 (1991) 1241–1299. URL: <https://www.jstor.org/stable/1229039>. doi:10.2307/1229039, publisher: Stanford Law Review.
- [23] A. Hanna, E. Denton, A. Smart, J. Smith-Loud, Towards a critical race methodology in algorithmic fairness, in: Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency, FAT* '20, Association for Computing Machinery, New York, NY, USA, 2020, pp. 501–512. URL: <https://dl.acm.org/doi/10.1145/3351095.3372826>. doi:10.1145/3351095.3372826.
- [24] S. Wachter, B. Mittelstadt, C. Russell, Why fairness cannot be automated: Bridging the gap between eu non-discrimination law and aiq1bcomputer law & security review; h-index: 41 sjr: Q1 core: Na abdc: B ft50: Na c2computer law & security review; h-index: 41 vhb: Na fnege: Na conrs: Na hcere: Na ccf: C bfi: 2 ajg: Na +, *Computer Law & Security Review* 41 (2021).
- [25] M. K. Lee, Understanding perception of algorithmic decisions: Fairness, trust, and emotion in response to algorithmic management, *Big Data & Society* 5 (2018) 205395171875668.
- [26] J. A. Colquitt, J. B. Rodell, Measuring Justice and Fairness, in: R. Cropanzano, R. S. Cropanzano, M. L. Ambrose (Eds.), *The Oxford handbook of justice in the workplace*, Oxford library of psychology, Oxford University Press, Oxford, 2015.
- [27] S. A. Friedler, C. Scheidegger, S. Venkatasubramanian, The (Im)possibility of fairness, *Communications of the ACM* 64 (2021) 136–143.
- [28] R. Wirth, J. Hipp, Crisp-dm: towards a standard process modell for data mining, 2000. URL: <https://api.semanticscholar.org/CorpusID:1211505>.

- [29] U. Fayyad, G. Piatetsky-Shapiro, P. Smyth, From data mining to knowledge discovery in databases, *AI magazine* 17 (1996) 37–37.
- [30] M. Wang, Y. Cui, X. Wang, S. Xiao, J. Jiang, Machine learning for networking: Workflow, advances and opportunities, *Ieee Network* 32 (2017) 92–99.
- [31] A. Quemy, Two-stage optimization for machine learning workflow, *Information Systems* 92 (2020) 101483.
- [32] N. Kühn, R. Hirt, L. Baier, B. Schmitz, G. Satzger, How to conduct rigorous supervised machine learning in information systems research: the supervised machine learning report card, *Communications of the Association for Information Systems* 48 (2021) 46.
- [33] D. Kreuzberger, N. Kühn, S. Hirschl, Machine learning operations (mlops): Overview, definition, and architecture, *IEEE Access* 11 (2023) 31866–31879. doi:10.1109/ACCESS.2023.3262138.
- [34] M. S. A. Lee, J. Singh, Risk identification questionnaire for detecting unintended bias in the machine learning development lifecycle, in: *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, 2021, pp. 704–714.
- [35] V. Singh, A. Singh, K. Joshi, Fair crisp-dm: Embedding fairness in machine learning (ml) development life cycle (2022).
- [36] A. Castelnovo, R. Crupi, G. D. Gamba, G. Greco, A. Naseer, D. Regoli, B. S. Miguel Gonzalez, BeFair: Addressing Fairness in the Banking Sector, in: *2020 IEEE International Conference on Big Data*, IEEE, 2020, pp. 3652–3661.
- [37] A. Barredo Arrieta, N. Díaz-Rodríguez, J. Del Ser, A. Bennetot, S. Tabik, A. Barbado, S. Garcia, S. Gil-Lopez, D. Molina, R. Benjamins, R. Chatila, F. Herrera, Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI, *Information Fusion* 58 (2020) 1–72.
- [38] A. Páez, The pragmatic turn in explainable artificial intelligence (XAI), *Minds and Machines* 29 (2019) 441–459. doi:10.1007/s11023-019-09502-w.
- [39] L. Chazette, W. Brunotte, T. Speith, Exploring explainability: A definition, a model, and a knowledge catalogue, in: *Proceedings of the 29th IEEE International Requirements Engineering Conference*, IEEE, Piscataway, NJ, USA, 2021, pp. 197–208. doi:10.1109/RE51729.2021.00025.
- [40] M. Langer, D. Oster, T. Speith, H. Hermanns, L. Kästner, E. Schmidt, A. Sesing, K. Baum, What do we want from Explainable Artificial Intelligence (XAI)? – A stakeholder perspective on XAI and a conceptual model guiding interdisciplinary XAI research, *Artificial Intelligence* 296 (2021) 1–24.
- [41] M. A. Köhl, K. Baum, M. Langer, D. Oster, T. Speith, D. Bohlender, Explainability as a non-functional requirement, in: D. E. Damian, A. Perini, S. Lee (Eds.), *Proceedings of the 27th IEEE International Requirements Engineering Conference*, IEEE, Piscataway, NJ, USA, 2019, pp. 363–368. doi:10.1109/RE.2019.00046.
- [42] M. Langer, R. N. Landers, The future of artificial intelligence at work: A review on effects of decision automation and augmentation on workers targeted by algorithms and third-party observers, *Computers in Human Behavior* 123 (2021) 106878.
- [43] J. Schoeffler, M. De-Arteaga, N. Kuehl, Explanations, fairness, and appropriate reliance in human-ai decision-making, in: *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, 2024.

- [44] L. H. Gilpin, C. Testart, N. Fruchter, J. Adebayo, Explaining explanations to society, in: C. Bakalar, S. Bird, T. Caetano, E. Felten, D. Garcia-Garcia, I. Kloumann, F. Lattimore, S. Mullainathan, D. Sculley (Eds.), Proceedings of the NeurIPS 2018 Workshop on Ethical, Social and Governance Issues in AI, 2019, pp. 1–6. [arXiv:1901.06560](https://arxiv.org/abs/1901.06560).
- [45] D. V. Carvalho, E. M. Pereira, J. S. Cardoso, Machine learning interpretability: A survey on methods and metrics, *Electronics* 8 (2019). doi:10.3390/electronics8080832.
- [46] M. Langer, K. Baum, K. Hartmann, S. Hessel, T. Speith, J. Wahl, Explainability auditing for intelligent systems: A rationale for multi-disciplinary perspectives, in: T. Yue, M. Mirakhorli (Eds.), 29th IEEE International Requirements Engineering Conference Workshops, REW 2021, IEEE, Piscataway, NJ, USA, 2021, pp. 164–168. doi:10.1109/REW53955.2021.00030.
- [47] C. Rudin, Stop Explaining Black Box Machine Learning Models for High Stakes Decisions and Use Interpretable Models Instead, *Nature Machine Intelligence* 1 (2019) 206–215.
- [48] R. Dwivedi, D. Dave, H. Naik, S. Singhal, R. Omer, P. Patel, B. Qian, Z. Wen, T. Shah, G. Morgan, et al., Explainable ai (xai): Core ideas, techniques, and solutions, *ACM Computing Surveys* 55 (2023) 1–33.
- [49] M. Mitchell, S. Wu, A. Zaldivar, P. Barnes, L. Vasserman, B. Hutchinson, E. Spitzer, I. D. Raji, T. Gebru, Model Cards for Model Reporting, in: Proceedings of the Conference on Fairness, Accountability, and Transparency, ACM Conferences, ACM, 2019, pp. 220–229.
- [50] T. Gebru, J. Morgenstern, B. Vecchione, J. W. Vaughan, H. Wallach, H. D. Iii, K. Crawford, Datasheets for datasets, *Communications of the ACM* 64 (2021) 86–92.
- [51] N. Kühl, M. Schemmer, M. Goutier, G. Satzger, Artificial intelligence and machine learning, *Electronic Markets* 32 (2022) 2235–2244.
- [52] M. Loi, A. Ferrario, E. Viganò, Transparency as design publicity: Explaining and justifying inscrutable algorithms, *Ethics and information technology* 23 (2021) 253–263.
- [53] B. Green, Escaping the Impossibility of Fairness: From Formal to Substantive Algorithmic Fairness, *Philosophy & Technology* 35 (2022).
- [54] A. Chouldechova, Fair Prediction with Disparate Impact: A Study of Bias in Recidivism Prediction Instruments, *Big Data* 5 (2017) 153–163.
- [55] B. Herman, The Promise and Peril of Human Evaluation for Model Interpretability, in: 31st Conference on Neural Information Processing Systems (NIPS 2017), 2017.
- [56] U. Aïvodji, H. Arai, O. Fortineau, S. Gambs, S. Hara, A. Tapp, Fairwashing: The risk of rationalization, in: Proceedings of the 36th International Conference on Machine Learning, 2019, pp. 161–170.
- [57] H. Lakkaraju, O. Bastani, "How do I fool you?" Manipulating User Trust via Misleading Black Box Explanations, in: Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society, 2020.
- [58] E. Balkir, S. Kiritchenko, I. Nejadgholi, K. C. Fraser, Challenges in Applying Explainability Methods to Improve the Fairness of NLP Models, in: Proceedings of the 2nd Workshop on Trustworthy Natural Language Processing (TrustNLP 2022), 2022.
- [59] M. T. Ribeiro, S. Singh, C. Guestrin, "Why should I trust you?" Explaining the predictions of any classifier, in: Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining, 2016.
- [60] S. M. Lundberg, Su-In Lee, A unified approach to interpreting model predictions, in: 31st

- Conference on Neural Information Processing Systems, 2017.
- [61] Z. C. Lipton, A. Chouldechova, J. McAuley, Does mitigating ML’s impact disparity require treatment disparity?, in: Proceedings of the 32nd International Conference on Neural Information Processing Systems, 2018, pp. 8136–8146.
 - [62] J. Nyarko, S. Goel, R. Sommers, Breaking Taboos in Fair Machine Learning: An Experimental Study, in: Equity and Access in Algorithms, Mechanisms, and Optimization, ACM Digital Library, Association for Computing Machinery, 2021, pp. 1–11.
 - [63] S. Dhanorkar, C. T. Wolf, K. Qian, A. Xu, L. Popa, Y. Li, Who needs to know what, when?: Broadening the explainable ai (xai) design space by looking at explanations across the ai lifecycle, in: Proceedings of the 2021 ACM Designing Interactive Systems Conference, 2021, pp. 1591–1602.
 - [64] L. Sartori, A. Theodorou, A sociotechnical perspective for the future of AI: Narratives, inequalities, and human control, *Ethics and Information Technology* 24 (2022) 1–11.
 - [65] N. Mehrabi, F. Morstatter, N. Saxena, K. Lerman, A. Galstyan, A Survey on Bias and Fairness in Machine Learning, *ACM Computer Surveys* 54 (2021).
 - [66] R. Pradhan, J. Zhu, B. Glavic, B. Salimi, Interpretable Data-Based Explanations for Fairness Debugging, in: Proceedings of the 2022 International Conference on Management of Data, ACM Digital Library, Association for Computing Machinery, 2022, pp. 247–261.
 - [67] N. Li, N. Goel, E. Ash, Data-Centric Factors in Algorithmic Fairness, in: Proceedings of the 2022 AAAI/ACM Conference on AI, Ethics, and Society, ACM Digital Library, Association for Computing Machinery, 2022, pp. 396–410.
 - [68] S. Corbett-Davies, S. Goel, The Measure and Mismeasure of Fairness: A Critical Review of Fair Machine Learning, 2018. URL: <http://arxiv.org/pdf/1808.00023v2>.
 - [69] J. Schoeffer, N. Kuehl, Y. Machowski, “There Is Not Enough Information”: On the Effects of Explanations on Perceptions of Informational Fairness and Trustworthiness in Automated Decision-Making, in: 2022 ACM Conference on Fairness, Accountability, and Transparency (FAccT ’22), ACM, 2022.
 - [70] J. A. Colquitt, On the dimensionality of organizational justice: a construct validation of a measure, *The Journal of applied psychology* 86 (2001) 386–400.
 - [71] R. Binns, M. van Kleek, M. Veale, U. Lyngs, J. Zhao, N. Shadbolt, “It’s Reducing a Human Being to a Percentage” Perceptions of Justice in Algorithmic Decisions, in: Proceedings of the 2018 CHI, ACM, 2018, pp. 1–14.
 - [72] J. Dodge, Q. V. Liao, Y. Zhang, R. K. E. Bellamy, C. Dugan, Explaining Models: An Empirical Study of How Explanations Impact Fairness Judgment, in: Proceedings of the 24th international conference on intelligent user interfaces, 2019, pp. 275–285.
 - [73] I. Rahwan, Society-in-the-Loop: Programming the Algorithmic Social Contract, *Ethics and Information Technology* 20 (2018).
 - [74] M. Teodorescu, L. Morse, Y. Awwad, G. Kane, Failures of fairness in automation require a deeper understanding of human-ml augmentation, *Management Information Systems Quarterly* 45 (2021) 1483–1500.
 - [75] M. Schemmer, N. Kuehl, C. Benz, A. Bartos, G. Satzger, Appropriate reliance on ai advice: Conceptualization and the effect of explanations, in: Proceedings of the 28th International Conference on Intelligent User Interfaces, 2023, pp. 410–422.
 - [76] J. Chakraborty, K. Peng, T. Menzies, Making fair ML software using trustworthy expla-

- nation, in: Proceedings of the 35th IEEE/ACM International Conference on Automated Software Engineering, ACM, 2020, pp. 1229–1233.
- [77] K. Vredenburg, The Right to Explanation*, *Journal of Political Philosophy* 30 (2022) 209–229.
- [78] European Commission. Directorate General for Communications Networks, Content and Technology., High Level Expert Group on Artificial Intelligence., Ethics guidelines for trustworthy AI, 2019. URL: <https://digital-strategy.ec.europa.eu/en/library/ethics-guidelines-trustworthy-ai>.
- [79] S. Wachter, B. Mittelstadt, C. Russell, Counterfactual Explanations Without Opening the Black Box: Automated Decisions and the GDPR, *SSRN Electronic Journal* (2017) 1–47.
- [80] N. Asher, L. de Lara, S. Paul, C. Russell, Counterfactual Models for Fair and Adequate Explanations, *Machine Learning and Knowledge Extraction* 4 (2022) 316–349.
- [81] D. S. Watson, L. Floridi, The Explanation Game: A Formal Framework for Interpretable Machine Learning, in: L. Floridi (Ed.), *Ethics, Governance, and Policies in Artificial Intelligence*, volume 144 of *Springer eBook Collection*, Springer International Publishing and Imprint Springer, Cham, 2021.
- [82] A. Ghosh, A. Shanbhag, C. Wilson, FairCanary: Rapid Continuous Explainable Fairness, in: Proceedings of the 2022 AAAI/ACM Conference on AI, Ethics, and Society (AIES’22), 2022.
- [83] A. Castelnovo, L. Malandri, F. Mercurio, M. Mezzanzanica, A. Cosentini, Towards Fairness Through Time, in: Joint European Conference on Machine Learning and Knowledge Discovery in Databases, Springer, Cham, 2021, pp. 647–663.
- [84] M. Hardt, X. Chen, X. Cheng, M. Donini, J. Gelman, S. Gollaprolu, J. He, P. Larroy, X. Liu, N. McCarthy, A. Rathi, S. Rees, A. Siva, E. Tsai, K. Vasist, P. Yilmaz, M. B. Zafar, S. Das, K. Haas, T. Hill, K. Kenthapadi, Amazon SageMaker Clarify: Machine Learning Bias Detection and Explainability in the Cloud, in: Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD 2021), 2021.
- [85] C. Panigutti, A. Perotti, A. Panisson, P. Bajardi, D. Pedreschi, FairLens: Auditing black-box clinical decision support systems, *Information Processing & Management* 58 (2021) 1–17.
- [86] J. Walmsley, Artificial intelligence and the value of transparency, *AI & SOCIETY* 36 (2021) 585–595.
- [87] P. Hacker, J.-H. Passoth, Varieties of AI Explanations Under the Law. From the GDPR to the AIA, and Beyond, in: A. Holzinger, R. Goebel, R. Fong, T. Moon, K.-R. Müller, W. Samek (Eds.), *xxAI – Beyond Explainable AI*, volume 13200 of *Springer eBook Collection*, Springer International Publishing and Imprint Springer, Cham, 2022, pp. 343–373.
- [88] A. I. Anik, A. Bunt, 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, 2021.
- [89] S. Barocas, K. Crawford, A. Shapiro, H. Wallach, The problem with bias: Allocative versus representational harms in machine learning, in: 9th Annual conference of the special interest group for computing, information and society (SIGCIS), 2017.
- [90] Y. Li, N. Vasconcelos, REPAIR: Removing Representation Bias by Dataset Resampling, in: 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), IEEE, 2019.

- [91] W. Cai, J. Gaebler, N. Garg, S. Goel, Fair Allocation through Selective Information Acquisition, in: Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society, ACM Digital Library, Association for Computing Machinery, 2020, pp. 22–28.
- [92] P. Saleiro, K. T. Rodolfa, R. Ghani, Dealing with bias and fairness in data science systems: A practical hands-on tutorial, in: Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, ACM Digital Library, Association for Computing Machinery, 2020, pp. 3513–3514.
- [93] E. Black, S. Yeom, M. Fredrikson, FlipTest: Fairness Testing via Optimal Transport, in: Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency, ACM, 2020, pp. 111–121.
- [94] M. Fan, W. Wei, W. Jin, Z. Yang, T. Liu, Explanation-Guided Fairness Testing through Genetic Algorithm, in: Proceedings of the 44th International Conference on Software Engineering, ACM Digital Library, Association for Computing Machinery, 2022, pp. 871–882.
- [95] S. Chiappa, Path-specific counterfactual fairness, in: Proceedings of the AAAI Conference on Artificial Intelligence, volume 33, 2019, pp. 7801–7808.
- [96] R. Binns, On the apparent conflict between individual and group fairness, in: Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency, ACM, 2020.
- [97] Z. Chen, J. M. Zhang, M. Hort, M. Harman, F. Sarro, Fairness testing: A comprehensive survey and analysis of trends, ACM Transactions on Software Engineering and Methodology (2023).
- [98] Y. Ahn, Y.-R. Lin, FairSight: Visual Analytics for Fairness in Decision Making, IEEE Transactions on Visualization and Computer Graphics 26 (2020) 1086–1095.
- [99] B. Ghosh, D. Basu, K. S. Meel, “How Biased is Your Feature?”: Computing Fairness Influence Functions with Global Sensitivity Analysis, in: Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency, ACM Digital Library, Association for Computing Machinery, 2023.
- [100] V. Gupta, P. Nokhiz, C. D. Roy, S. Venkatasubramanian, Equalizing Recourse across Groups, arxiv preprint, 2019. URL: <https://arxiv.org/abs/1909.03166>.
- [101] P. Grabowicz, N. Perello, A. Mishra, Marrying Fairness and Explainability in Supervised Learning, in: 2022 ACM Conference on Fairness, Accountability, and Transparency (FAccT ’22), ACM, 2022.
- [102] J. E. Hunter, F. L. Schmidt, R. Hunter, Differential validity of employment tests by race: A comprehensive review and analysis, Psychological Bulletin 86 (1979) 721–735.
- [103] S. Gupta, S. Lee, M. De-Arteaga, M. Lease, Same Same, But Different: Conditional Multi-Task Learning for Demographic-Specific Toxicity Detection, in: Proceedings of the ACM Web Conference 2023, Association for Computing Machinery and Saarländische Universitäts- und Landesbibliothek, 2023, pp. 3689–3700.
- [104] A. Shulner-Tal, T. Kuflik, D. Kliger, Enhancing Fairness Perception – Towards Human-Centred AI and Personalized Explanations Understanding the Factors Influencing Laypeople’s Fairness Perceptions of Algorithmic Decisions, International Journal of Human-Computer Interaction (2022) 1–28.
- [105] A. Papanmeier, G. Englebienne, C. Seifert, How model accuracy and explanation fidelity influence user trust in AI, in: Proceedings of the IJCAI 2019, International Joint

- Conferences on Artificial Intelligence, 2019, pp. 94–100.
- [106] J. Schoeffer, N. Kuehl, Appropriate fairness perceptions? On the effectiveness of explanations in enabling people to assess the fairness of automated decision systems., in: Companion Publication of the 2021 Conference on Computer Supported Cooperative Work and Social Computing, ACM, 2021.
 - [107] R. Noothigattu, S. Gaikwad, E. Awad, S. Dsouza, I. Rahwan, P. Ravikumar, A. Procaccia, A Voting-Based System for Ethical Decision Making, Proceedings of the AAAI Conference on Artificial Intelligence 32 (2018).
 - [108] M. K. Lee, D. Kusbit, A. Kahng, J. T. Kim, X. Yuan, A. Chan, D. See, R. Noothigattu, S. Lee, A. Psomas, A. D. Procaccia, WeBuildAI: Participatory Framework for Algorithmic Governance, Proceedings of the ACM on Human-Computer Interaction 3 (2019) 1–35.
 - [109] S. Stumpf, L. Strappelli, S. Ahmed, Y. Nakao, A. Naseer, G. Del Gamba, D. Regoli, Design Methods for Artificial Intelligence Fairness and Transparency, in: Joint Proceedings of the ACM IUI 2021 Workshops, 2021.
 - [110] M. Yaghini, A. Krause, H. Heidari, A human-in-the-loop framework to construct context-aware mathematical notions of outcome fairness, in: Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society, ACM Digital Library, Association for Computing Machinery, 2021, pp. 1023–1033.
 - [111] G. Ras, M. van Gerven, P. Haselager, Explanation Methods in Deep Learning: Users, Values, Concerns and Challenges, in: Explainable and Interpretable Models in Computer Vision and Machine Learning, Springer, Cham, 2018.
 - [112] A. Shulner-Tal, T. Kuflik, D. Kliger, Fairness, explainability and in-between: Understanding the impact of different explanation methods on non-expert users’ perceptions of fairness toward an algorithmic system, Ethics and Information Technology 24 (2022) 1–13.
 - [113] M. K. Lee, A. Jain, H. J. Cha, S. Ojha, D. Kusbit, Procedural Justice in Algorithmic Fairness: Leveraging Transparency and Outcome Control for Fair Algorithmic Mediation, in: Proceedings of the ACM on Human-Computer Interaction, volume 3, 2019, pp. 1–26.
 - [114] G. Starke, B. Schmidt, E. de Clercq, B. S. Elger, Explainability as fig leaf? An exploration of experts’ ethical expectations towards machine learning in psychiatry, AI and Ethics (2022) 1–12.
 - [115] T. Schmude, L. Koesten, T. Möller, S. Tschatschek, On the Impact of Explanations on Understanding of Algorithmic Decision-Making, in: Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency, ACM Digital Library, Association for Computing Machinery, 2023, pp. 959–970.
 - [116] M. Eiband, D. Buschek, A. Kremer, H. Hussmann, The Impact of Placebic Explanations on Trust in Intelligent Systems, in: Extended Abstracts of the 2019 CHI Conference on Human Factors in Computing Systems, ACM Digital Library, Association for Computing Machinery, 2019, pp. 1–6.
 - [117] E. Le Merrer, G. Trédan, Remote explainability faces the bouncer problem, Nature Machine Intelligence 2 (2020) 529–539.
 - [118] J.-M. John-Mathews, Some critical and ethical perspectives on the empirical turn of AI interpretability, Technological Forecasting and Social Change 174 (2022) 1–29.
 - [119] L. Enqvist, ‘human oversight’ in the eu artificial intelligence act: what, when and by whom?, Law, Innovation and Technology 15 (2023) 508–535.

- [120] M. Langer, Fehlgeleitete Hoffnungen? Grenzen menschlicher Aufsicht beim Einsatz algorithmusbasierter Systeme am Beispiel Personalauswahl, *Psychologische Rundschau* 74 (2023) 211–220.
- [121] D. Slack, S. A. Friedler, E. Givental, Fairness warnings and Fair-MAML: Learning Fairly with Minimal Data, in: *Conference on Fairness, Accountability, and Transparency (FAT '20)*, 2020.
- [122] Y. Nakao, S. Stumpf, S. Ahmed, A. Naseer, L. Strappelli, Toward Involving End-users in Interactive Human-in-the-loop AI Fairness, *ACM Transactions on Interactive Intelligent Systems* 12 (2022) 1–30.
- [123] U. Bhatt, J. Antorán, Y. Zhang, Q. V. Liao, P. Sattigeri, R. Fogliato, G. Melançon, R. Krishnan, J. Stanley, O. Tickoo, L. Nachman, R. Chunara, M. Srikumar, A. Weller, A. Xiang, Uncertainty as a Form of Transparency: Measuring, Communicating, and Using Uncertainty, in: *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, ACM Digital Library, Association for Computing Machinery, 2021, pp. 401–413.
- [124] R. Achibat, M. Dreyer, I. Eisenbraun, S. Bosse, T. Wiegand, W. Samek, S. Lopuschkin, From attribution maps to human-understandable explanations through concept relevance propagation, *Nature Machine Intelligence* 5 (2023) 1006–1019.
- [125] T. Miller, Explainable AI is Dead, Long Live Explainable AI!, in: *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*, ACM Digital Library, Association for Computing Machinery, 2023, pp. 333–342.
- [126] S. Lebovitz, H. Lifshitz-Assaf, N. Levina, To Engage or Not to Engage with AI for Critical Judgments: How Professionals Deal with Opacity When Using AI for Medical Diagnosis, *Organization Science* 33 (2022) 126–148.
- [127] T. Miller, Explanation in artificial intelligence: Insights from the social sciences, *Artificial Intelligence* 267 (2019) 1–38.
- [128] A.-H. Karimi, G. Barthe, B. Schölkopf, I. Valera, A survey of algorithmic recourse: Contrastive explanations and consequential recommendations, *ACM Computing Surveys* (2022) 1–26.
- [129] R. Crupi, A. Castelnovo, D. Regoli, B. San Miguel Gonzalez, Counterfactual explanations as interventions in latent space, *Data Mining and Knowledge Discovery* (2022) 1–37.
- [130] R. Guidotti, Counterfactual explanations and how to find them: literature review and benchmarking, *Data Mining and Knowledge Discovery* (2022) 1–55.
- [131] L. Hu, I. Kohler-Hausmann, What's sex got to do with machine learning?, in: *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, ACM, 2020, p. 513.
- [132] D. Slack, A. Hilgard, H. Lakkaraju, S. Singh, Counterfactual Explanations Can Be Manipulated, in: *Advances in Neural Information Processing Systems*, volume 34, Curran Associates, Inc, 2021, pp. 62–75.
- [133] N. Gill, P. Hall, K. Montgomery, N. Schmidt, A Responsible Machine Learning Workflow with Focus on Interpretable Models, Post-hoc Explanation, and Discrimination Testing, *Information* 11 (2020) 1–32.
- [134] F. Arif Khan, E. Manis, J. Stoyanovich, Towards Substantive Conceptions of Algorithmic Fairness: Normative Guidance from Equal Opportunity Doctrines, in: *Equity and Access in Algorithms, Mechanisms, and Optimization*, ACM Digital Library, Association for Computing Machinery, 2022, pp. 1–10.

- [135] L. Baier, N. Kühl, G. Satzger, How to cope with change? preserving validity of predictive services over time, in: Hawaii International Conference on System Sciences (HICSS-52), Grand Wailea, Maui, Hawaii, Januar 8-11, 2019, 2019, p. 1085.
- [136] B. Lepri, N. Oliver, E. Letouzé, A. Pentland, P. Vinck, Fair, Transparent, and Accountable Algorithmic Decision-making Processes, *Philosophy & Technology* 31 (2018) 611–627.
- [137] L. T. Liu, S. Dean, E. Rolf, M. Simchowicz, M. Hardt, Delayed impact of fair machine learning, *International Conference on Machine Learning* (2018) 3150–3158.
- [138] M. Hardt, M. Jagadeesan, C. Mendler-Dünner, Performative power, *Advances in Neural Information Processing Systems* 35 (2022) 22969–22981.
- [139] K. Martin, Ethical implications and accountability of algorithms, *Journal of Business Ethics* 160 (2019) 835–850.
- [140] B. Pettit, C. Gutierrez, Mass Incarceration and Racial Inequality, *American journal of economics and sociology* 77 (2018) 1153–1182. URL: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC9540942/>. doi:10.1111/ajes.12241.
- [141] J. Dressel, H. Farid, The accuracy, fairness, and limits of predicting recidivism, *Science Advances* 4 (2018).
- [142] P. Linardatos, V. Papastefanopoulos, S. Kotsiantis, Explainable AI: A Review of Machine Learning Interpretability Methods, *Entropy* 23 (2020) 1–45.
- [143] Google pauses ai-generated images of people after ethnicity criticism, *The Guardian* (2024). URL: <https://www.theguardian.com/technology/2024/feb/22/google-pauses-ai-generated-images-of-people-after-ethnicity-criticism>.
- [144] F. Friedrich, M. Brack, L. Struppek, D. Hintersdorf, P. Schramowski, S. Luccioni, K. Kersting, Fair diffusion: Instructing text-to-image generation models on fairness, 2023. URL: <http://arxiv.org/pdf/2302.10893>.
- [145] T. Neumann, S. Lee, M. De-Arteaga, S. Fazelpour, M. Lease, Diverse, but divisive: LLMs can exaggerate gender differences in opinion related to harms of misinformation, *arXiv preprint arXiv:2401.16558* (2024).
- [146] P. Hacker, The european ai liability directives – critique of a half-hearted approach and lessons for the future, 2022. URL: https://papers.ssrn.com/sol3/papers.cfm?abstract_id=4279796.
- [147] I. D. Raji, A. Smart, R. N. White, M. Mitchell, T. Gebru, B. Hutchinson, J. Smith-Loud, D. Theron, P. Barnes, Closing the AI Accountability Gap: Defining an End-to-End Framework for Internal Algorithmic Auditing, in: *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, ACM, 2020, pp. 33–44.