

Software Fairness: An Analysis and Survey

EZEKIEL SOREMEKUN, Singapore University of Technology and Design, Singapore

MIKE PAPADAKIS, MAXIME CORDY, and YVES LE TRAON, SnT, University of Luxembourg, Luxembourg

ABSTRACT In the last decade, researchers have studied fairness as a software property. In particular, how to engineer fair software systems. This includes specifying, designing, and validating fairness properties. However, the landscape of works addressing bias as a software engineering concern is unclear, i.e., techniques and studies that analyze the fairness properties of learning-based software. In this work, we provide a clear view of the state-of-the-art in software fairness analysis. To this end, we collect, categorize and conduct in-depth analysis of 164 publications investigating the fairness of learning-based software systems. Specifically, we study the evaluated fairness measure, the studied tasks, the type of fairness analysis, the main idea of the proposed approaches and the access level (e.g., black, white or grey box). Our findings include the following: (1) Fairness concerns (such as fairness specification and requirements engineering) are under-studied; (2) Fairness measures such as conditional, sequential and intersectional fairness are under-explored; (3) Semi-structured datasets (e.g., audio, image, code and text) are barely studied for fairness analysis in the SE community; and (4) Software fairness analysis techniques hardly employ white-box, in-processing machine learning (ML) analysis methods. In summary, we observed several open challenges including the need to study intersectional/sequential bias, policy-based bias handling and human-in-the-loop, socio-technical bias mitigation.

Additional Key Words and Phrases: Software Fairness, Software Analysis, Bias, Discrimination, Artificial intelligence, Machine learning

ACM Reference Format:

Ezekiel Soremekun, Mike Papadakis, Maxime Cordy, and Yves Le Traon. 2025. Software Fairness: An Analysis and Survey. 1, 1 (August 2025), 35 pages. <https://doi.org/10.1145/mnnnnnn.mnnnnnn>

1 INTRODUCTION

Software fairness is a property of learning-based systems which aims to ensure that the software does not exhibit biases [132]. Given a set of inputs, a fair software should not result in discriminatory outputs or behaviors for inputs relating to certain groups or individuals. In essence, the goal is to ensure that software systems exhibit fair behavior for all inputs that are similar for the task-at-hand. For instance, *discriminatory inputs*, inputs that are similar for a task but only differ in *sensitive attributes* (e.g., gender or race), should produce similar outputs or induce similar behaviors [55].

Specifically, the goal of *software fairness analysis* is to ensure a given system produces the same results or exhibit similar behaviors for a number of *discriminatory inputs*. As an example, consider a sentiment analyzer software that determines the emotional state or situation in a text, which outputs either a positive emotion (e.g., text describing excitement) or negative emotion (e.g., text portraying sadness or anger). Figure 1 shows an example of a bias in such a

Authors' addresses: Ezekiel Soremekun, ezekiel_soremekun@sutd.edu.sg, Singapore University of Technology and Design, 8 Somapah Rd, Singapore, Singapore, Singapore, 487372; Mike Papadakis, michail.papadakis@uni.lu; Maxime Cordy, maxime.cordy@uni.lu; Yves Le Traon, Yves.LeTraon@uni.lu, SnT, University of Luxembourg, 6, rue Richard Coudenhove-Kalergi, Luxembourg, Luxembourg, Luxembourg, L-1359.

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.

© 2025 Association for Computing Machinery.

Manuscript submitted to ACM

Manuscript submitted to ACM

53 system, where the output (sentiment) is different when the gender of the noun in the text is a “man” compared to a
54 “woman”. This illustrates (gender) bias found in real-world natural language processing (NLP) systems [9, 122].
55

56 Several researchers have studied software fairness analysis, with the aim to address a fundamental question: How to
57 engineer fair software systems? Such works consider fairness as a non-functional software property and a software
58 engineering (SE) concern. This includes work studying how to specify [118], test [55], and mitigate [31] fairness
59 properties in software systems. However, despite several works on software fairness analysis, it is difficult to understand
60 the state of research practice: Specifically, what fairness concerns have been addressed? How have they been addressed?
61 What are the open problems and challenges?

62 In this paper, we aim to provide a clear view of the state-of-the-art in software fairness analysis, i.e., techniques
63 and studies that analyse the fairness properties of learning-based software. This paper aims to analyze the trends in
64 software fairness analysis, the available techniques, the focus of the research community, the problems that have been
65 addressed and the open research problems. To this end, we perform a systematic analysis of the literature, where we
66 studied 164 papers studying *fairness as a software property*. These papers are mostly published in fairness-related or
67 software engineering (SE) venues, as well as venues focused on machine learning (ML), artificial intelligence (AI),
68 security, computer vision (CV) and natural language processing (NLP). Particularly, we conduct an in-depth study of
69 the set of publications that explore fairness with the lens of software engineering, e.g., in terms of software quality
70 control, requirement engineering, design and development. We then characterize several factors encapsulated in these
71 research papers, including the evaluated fairness measure (e.g., individual, group, causal or conditional fairness), the
72 studied tasks (e.g., credit rating, CV, NLP), the type of fairness analysis (e.g., testing or mitigation), the main idea of the
73 proposed approach and the level of access (e.g., black, white or grey box).

74 Overall, we observe that the research community is facing several open challenges in addressing the following
75 fairness concerns: compounding or intersectional bias, verification of fairness properties, sequential bias, equity-based
76 bias handling and socio-technical solutions to bias. The key findings of this work includes the following:
77

- 78 • Software Fairness analysis is mostly performed to validate or mitigate biases, i.e., the focus of the community has
79 been to test and reduce unfair program behaviors, other fairness concerns (such as requirements engineering
80 and verification) are under-studied;
- 81 • The most studied fairness measure include individual, group and causal fairness, measures such as conditional,
82 sequential and intersectional fairness remain under-explored;
- 83 • The most examined tasks involve tabular datasets (such as the adult income dataset), while semi-structured
84 datasets (e.g., audio, image and text) are barely studied in the SE community;
- 85 • The most employed techniques in the SE community are mutation (e.g., input perturbation) and specification
86 (e.g., using input templates, schemas or grammars) -based techniques, other approaches such as in-processing
87 machine learning (ML) analysis methods are hardly employed for software fairness analysis;
- 88 • Most approaches support the analysis of an atomic attribute (e.g. race, or gender), very few approaches study the
89 combination or sequence of attributes (e.g., race × gender), the compounding effect of multiple instances of an
90 attribute, or complex attributes (such as non-binary gender);
- 91 • We found little or very few works tackling the fairness concerns like fairness test metrics/adequacy, automatic
92 repair of biased classifiers or time-based fairness concerns (e.g., sequential or regression fairness bugs).

93 The rest of this paper is organised as follows: We provide background on software fairness in [section 2](#). We discuss
94 the process of collecting/analyzing publications in [section 3](#). [Section 4.3](#) highlights our research questions. [Section 5](#)

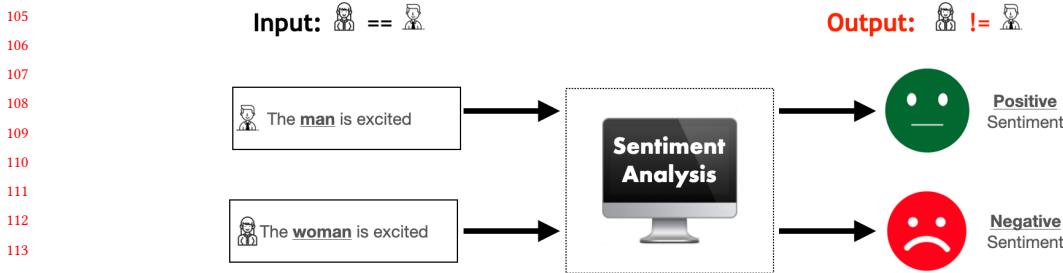


Fig. 1. Example of Individual Fairness Violation using a Sentiment Analysis AI System

provides in-depth analysis of our findings, section 6 discusses the limitations of our study and section 7 presents recent advances. Finally, we discuss open research challenges (section 8) in software fairness analysis and conclude (section 9).

2 BACKGROUND

We provide background on software fairness, definition of terms, illustrative examples and closely-related works.

2.1 Definition of Terms

We describe the main terms used in this paper, and the context in which they apply to our analysis and survey.

Bias: In this paper, we refer to *bias* in terms of algorithmic bias, specifically, bias occurs when a software system *systematically* and *unfairly* discriminate against certain individuals or groups of individuals in favour of others [53]. Algorithmic bias causes discrimination against certain people and can lead to real-world harms in terms of the representation or allocation of resources to such individuals or groups [37].

Software Fairness: There has been a significant work in understanding and defining *fairness* as a software behavior [132]. In this work, we examine *fairness as a software property*. The aim is to focus on works examining *software fairness via the lens of software engineering*, e.g., how to engineer bias-aware software or prevent bias in software systems. To this end, we study papers that study fairness as a SE concern including in terms of a requirements engineering [14], software quality control [55], software design [77], and software verification [3].

2.2 The interplay of ML Fairness and Software Fairness

In this section, we discuss the interplay (similarities, differences and intersections) of ML fairness and software fairness, e.g., in terms of their *design (components and pipelines)*, *terminology*, and *analysis*. We provide a venn diagram describing the taxonomy and interplay of ML and software fairness in our appendix or supplementary material (Figure 1).

Program Fairness: This is the fairness property of a computer program. A program may exhibit unfair behaviors or produce unfair outputs. As an example, consider the rule-based credit rating system in Figure 2. Let us assume that the fairness policy of an institution or a country protects “age” and “marital status”, such that they can not be used to determine the credit approval rating of an individual. Then, we say this program is *unfair* because its credit approval is computed using protected attributes – *age* and *marital status*. Particularly, the conditional checks (“if” condition statements in lines eight (8) and twelve (12) in the program shown in Figure 2) violate the aforementioned fairness policy. Our appendix (Figure 2) further shows an example input where the program violates the defined fairness policy.

```

157 1 public static boolean approve_credit(int age, boolean is_married, int salary,
158 2     int years_of_experience) {
159 3     int salary_threshold = 50000;
160 4     int min_age = 30;
161 5     int max_age = 60;
162 6     int experience_threshold = 5;
163 7     boolean credit_approval = false;
164 8
165 9     if (!is_married || years_of_experience < experience_threshold){
166 10         return credit_approval;
167 11     }
168 12
169 13     if (salary > salary_threshold && age > min_age && age < max_age){
170 14         credit_approval = true;
171 15         return credit_approval;
172 16     }
173 17
174 18     return credit_approval;
175 19 }
176 20
177 21
178 22 Fig. 2. Illustration of an unfair rule-based credit approval program which uses protected attributes (“age”, “marital status”)
179 23
180 24
181 25 ML Fairness: ML fairness refers to the fairness property of a machine learning (ML) model, e.g., an unfair ML model
182 26 may produce unfair predictions/classifications. For instance, an ML classifier can be considered unfair if its prediction
183 27 for a protected group is significantly different from that of other groups. Consider a sentiment analyzer, let us assume
184 28 it is only an ML model (with no additional components), then we say that the ML model is unfair since it produces
185 29 significantly different outcomes for female-biased occupations ( $G_1$ ) versus male-biased occupations ( $G_2$ ). We illustrate
186 30 this example in our appendix (Figure 4).
187 31
188 32 Software Fairness: (Un)fairness may stem from one or more components of a software system. Indeed, unfair behaviors
189 33 may be from the ML model itself, other software components, multiple components or the interaction of the ML model
190 34 with other software components (see Figure 3). For instance, consider the sample ML-based software in Figure 3. On one
191 35 hand, ML fairness in this system may be caused by biases inherent in the ML model, e.g., its training data or algorithms.
192 36 On the other hand, software unfairness may stem from any of the components of the system, including the ML models,
193 37 and other software components – e.g., input/output validators, or servers. Besides, unfair software behaviors may stem
194 38 from the interaction of multiple components, e.g., the response processing from the ML model to the response/output
195 39 servers and validators. For instance, a software whose input validator or tokenizer can only process input up to a limited
196 40 input size (e.g., due to memory constraint) may induce unfair behaviors for very long inputs. As an example, some
197 41 NLP systems (e.g., BERT-based models) only process the first  $N$  (e.g., 512) input tokens [141], and such input/memory
198 42 constraints may cause unfair behaviors despite the model’s capability or fairness properties.
199 43
200 44 Similarly, consider the credit approval software, illustrated in the appendix (Figure 2). Let us assume that marital
201 45 status is a protected attribute. Then, we say that the credit approval software is unfair since it treats two individuals
202 46 differently even though they have similar characteristics except for the protected attribute (marital status). In particular,
203 47 the software approves the credit request of the married individual (input_a) but denies the credit request of the single,
204 48 unmarried individual (input_b). We note that this behavior could be as a result of different sources depending on the
205 49
206 50 Manuscript submitted to ACM
207 51
208 52

```

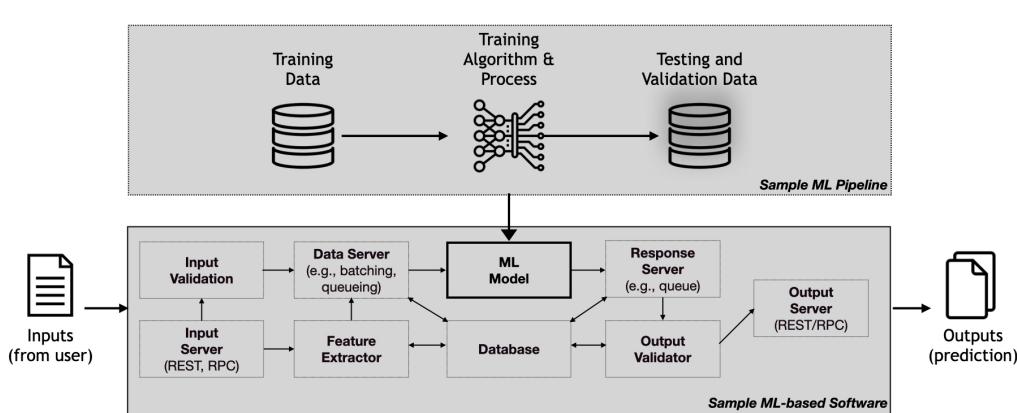


Fig. 3. An illustration of ML-based software systems (adapted from Lewis et al. [83], Muccini and Vaidhyanathan [104])

implementation of the system, i.e., (a) a program implementation (e.g., rule based system in Figure 2), (b) an ML model (e.g., automated credit rating classifier), or (c) a combination of multiple components including programs and ML models (e.g., appendix (Figure 1)). Overall, we consider all aforementioned settings as an unfair software behavior, since it is the behavior of the entire software system regardless of the source. Thus, in our setting, software fairness applies to all of the three aforementioned configurations of the software.

White-box Setting: In a *white-box setting*, where users or developers have access to the software implementation, it is possible to attribute the unfair software behavior to a specific component or analyze a specific sub-component for fairness properties. In this work, if fairness analysis involves examining the internals of a specific sub-component of the software, e.g., *only* the ML model or the rule-based program, then we consider this as *white-box fairness analysis*. This includes the analysis of a specific sub-component, e.g., the ML model (e.g., sentiment analyzer in our appendix (Figure 4)) or a specific program (e.g., the program in Figure 2). For instance, white-box fairness analysis includes examining the training process or architecture of an ML model to improve the fairness properties of the model, as well as the static/dynamic analysis or testing of a program to improve its fairness properties.

Black-box Setting: In a *black-box setting* where a customer or the developer has no access or knowledge of the system that is deployed, it is difficult to determine the component(s) responsible for the unfair behavior. Firstly, in a black-box setting, it is often *unknown* if the system is implemented as a learning-based system (e.g., ML model), a program (e.g., the rule based program in Figure 2), or a combination of both. This is because several ML-based software and AI systems are proprietary, black-box systems (e.g., COMPAS system, Google NLP, Amazon Rekognition, etc). They mostly only provide a web or API access and do not provide information on the implementation of the system, or its architecture. Secondly, in a black-box setting, it is difficult to determine the unfair component in a complex system, even when it is known that the system is composed of an ML model and/or a program. When analyzing an ML-based software system (e.g., Figure 3) in the black-box setting, it is difficult to determine if the ML model or rule-based program is responsible for the unfair behavior, or whether its interaction with other components caused unfairness or even other components are responsible for the unfair behavior. However, in this work we consider the analysis of such systems as software fairness analysis since, despite the lack of knowledge/access, the system is delivered/deployed as a software. For instance, the popular COMPAS system, which has been shown to be biased towards certain race and gender attributes, is a black-box proprietary software system [6]. Our appendix (Figure 3) illustrates the biases

exposed in this system, where *black* and *female* convicts are predicted to have a higher probability of recidivism (risk of re-committing crimes), even when their previous crimes are less serious than their *white* and *male* counterparts [6]. We consider the analysis that reveals biases in COMPAS as a *black-box fairness analysis*, since the researchers had no access to the training process, implementation or software architecture of the system [6].

Grey-box settings refers to methods and tools that employ a combination of the white-box and black-box methods for fairness analysis. Appendix (Table 4) details such grey-box approaches. For instance, Tizpaz-Niari et al. [130] proposed a fairness testing technique that leverages both the input space and ML model for fairness analysis.

White-box versus Black-box in ML Model-only settings: We note that when testing only the ML model, researchers often refer to white-box or black-box analysis as approaches that either inspect the ML model or not, respectively. In such settings, this work considers inspecting model internals (e.g., neurons) as white-box analysis, especially when only the model is tested (see appendix (Table 4)). Conversely, black-box analysis refers to fairness analysis without inspecting the model internals, in ML model-only settings. Subsequently, approaches that combine both approaches (i.e., leveraging analysis/information from both within and outside the ML model) are referred to as grey-box approaches.

2.3 Fairness Metrics

Verma and Rubin [132] provides comprehensive definitions and categories of several fairness metrics employed in the literature. Figure 1 and appendix (Figure 4) exemplify the two most popular metrics, namely *individual fairness* and *group fairness*. In the following, we define these two metrics.

Individual fairness: Dwork et al. [40] provides a comprehensive definition of individual fairness. Individual fairness means that the software should treat similar individuals similarly. Consider the sentiment analysis system in Figure 1, it violates individual fairness since it treats a text with a “man” as a *noun* differently from that with a “woman”. This metric requires that the individuals should be similar for the purposes of the respective task and the outcomes should have similar distributions. Formally, individual fairness is a violation of the following condition:

$$|f(a) - f(a')| \leq \tau \quad (1)$$

Here, a and a' are similar individuals (inputs), f is a software esystem (e.g., automated classifier) and τ is some threshold which is chosen using the inputs and the model as context. In our example (Figure 1), the text containing “man” or “woman” should be treated similarly as individuals, however the output shows that the system provides different outcomes for the similar inputs.

Group fairness: A system satisfies group fairness if subjects in the protected and unprotected groups have equal probability of being assigned a particular outcome [132]. This fairness metric aims to ensure that two or more groups are treated similarly. Let us consider a sentiment analyzer (illustrated in appendix - Figure 4), which violates group fairness because it treats texts containing nouns belonging to *male-biased occupations* (e.g., doctor, farmers etc) differently from similar texts that contain *female-biased occupations* (e.g., nurse, cleaner, etc). Formally, group fairness is maintained if the following condition is true:

$$\Pr(f(a) = + | A = a) = \Pr(f(b) = + | A = b) \quad \forall a, b \in A \quad (2)$$

Given equivalent inputs from different groups a and b , the aforementioned definition checks for the equivalence of the outputs from software f . Here, the choice of a group is determined by random variable A and the positive prediction rate is denoted by $+$. In this example, the groups are *male-biased occupations* and *female-biased occupations*, and the model

313 produces significantly different distribution of predictions for each group. As an example, a group fairness violation is
314 that texts containing male-biased jobs are more (90%) likely to return a positive sentiment than female-biased jobs (30%).
315

316 2.4 Related Work

317

318 Several researchers have surveyed the problem of bias or fairness in learning-based software, but these surveys have
319 been mostly *specialised*, i.e., they studied this problem for a specific domain, bias, or metric. Concretely, previous
320 surveys on fairness analysis target a particular sub-domain (e.g., NLP [20], CV [42], ranking [145] or finance [112]), a
321 specific sensitive attribute or bias (e.g., race [50]), or a specific fairness metric (e.g., causal fairness [99] or intersectional
322 fairness [59]). Thus, these surveys do not account for the advances in other domains, attributes and metrics.
323

324 Other fairness surveys are either more general beyond fairness concerns or are focused on specific goals or goals
325 beyond SE. For instance, some surveys explore general testing of ML systems for several properties [148], other
326 surveys providing a taxonomy or review of bias in ML systems [100] and some explore how to handle bias in special
327 circumstances, e.g., in the absence of demographic information [7]. More importantly, none of these papers have
328 performed a systematic analysis of bias in the context of *software engineering* or *the analysis of software fairness*.
329 Specifically, studies and methods that (*automatically*) evaluate software fairness properties in learning-based systems.
330 Our survey focuses on the literature w.r.t. several aspects of SE, thus involving a significant number of SE-related
331 publications. In the following, we discuss the closely related work, in particular the published surveys in this area.
332

333 There are some *domain-specific* surveys of fairness where researchers have surveyed fairness issues in a specific
334 sub-domain of learning-based software, e.g., recommendation systems [41], ranking-based systems [145], sequential
335 decision systems [153], NLP [20], CV [42], or financial services [112]. Blodgett et al. [20] comprehensively analyzed 146
336 papers addressing fairness of NLP systems, especially quantitative techniques for measuring or mitigating bias. The
337 authors found that almost all surveyed papers are poorly matched to their motivations and do not engage with the
338 relevant literature outside of NLP. The authors recommend that fairness analysis in NLP systems should be based on
339 characterizing system behaviors that are harmful, by centering bias mitigation around people and their experiences.
340 Fabbrizzi et al. [42] provides a survey on bias in visual datasets, i.e., for CV applications. The authors describe different
341 biases in visual datasets, the proposed methods for detecting and measuring these biases and existing bias-aware datasets.
342 The authors concluded that there is no bias-free dataset and detecting biases in visual datasets is an open problem and
343 recommend a checklist to help practitioners spot different biases in their dataset, in order to make bias explicit. Ashurst
344 and Weller [8] surveyed approaches for fairness without demographic data, discussing the benefits and limitations of
345 the approaches. Fabris et al. [45][44] also surveyed the datasets used in algorithmic fairness research. In addition, the
346 authors present a search engine to search amongst the surveyed datasets for algorithmic fairness [43]. Similar to this
347 paper, the authors found that Adult, COMPAS, and German Credit are the three most popular fairness dataset. The
348 paper further highlights the weaknesses of these datasets and propose alternative datasets for fairness practitioners and
349 researchers. Meanwhile, Zehlike et al. [145] studied the application of fairness-enhancing interventions in ranking
350 algorithms, especially examining the literature on incorporating fairness requirements into algorithmic rankers. The
351 authors surveyed papers from several venues, including data management, algorithms, information retrieval, and
352 recommendation systems. The goal of the survey is to provide a new framework that unifies fairness mitigation
353 objectives and ranking requirements, such that it allows to examine the trade-off between both goals. Ekstrand et al.
354 [41] studied how algorithmic fairness issues applies to information retrieval and recommendation systems, especially
355 addressing how to translate algorithmic fairness from classification, scoring, and ranking settings into recommendation
356 and information retrieval settings. In addition, Zhang and Liu [153] provides a literature review of fairness in sequential
357

365 learning-based decision-making systems, i.e., systems where decision-making are not a one time event, but rather occur
366 in a sequential nature, such that decisions made in the past may have an impact on future data. Unlike these papers,
367 this work is not *domain-specific*, instead, we study the problem of fairness analysis across several (sub-)domains.
368

369 Moreover, other surveys on fairness are *bias-specific* or *metric-specific*, they either focus on a specific sensitive
370 attribute (e.g., race) [50], or a specific fairness metric (e.g., causal fairness) [99], respectively. For instance, Field et al. [50]
371 provides a survey of race-related bias in the stages of NLP model development. The authors surveyed 79 papers with
372 the goal of understanding the gaps between *race-related bias* analysis in NLP and other related fields. The authors found
373 that race has been siloed as a niche topic in the NLP community and often ignored in many NLP tasks. The authors
374 also emphasize the need for racial inclusion in NLP research. In addition, Makhlof et al. [99] study the application
375 of causality to address the problem of fairness, by studying papers that examine *causal fairness* properties and their
376 applicability in real-world scenarios. The authors employed *identifiability theory* to determine the criteria for the
377 real-world applicability of *causal fairness*. However, in this work, we examine papers examining several sensitive
378 attributes, biases and fairness metrics. In particular, unlike these metric or bias -specific surveys, we do not focus on
379 papers examining a single type of bias or fairness metric. Instead, we evaluate the literature across different biases and
380 fairness metrics, including race and causal fairness, respectively.
381

382 A few researchers have conducted surveys of the literature on software fairness, albeit mostly targeting fair prediction
383 in machine learning [54] general ML testing [148] or fairness notions across domains [71] or in specific domains, e.g.,
384 concurrent systems [81]. For instance, Gajane and Pechenizkiy [54] studied the *formalization* of fairness in the machine
385 learning literature for prediction tasks, especially examining how this relates to the notions of distributive justice in
386 the social science literature. In their study, the authors proposed that two notions of distributive justice be formalised
387 for fairness in ML, namely *equality of resources* and *equality of capability of functioning*. Likewise, Ntoutsi et al. [106]
388 provides an introductory survey on the technical challenges and available solutions to bias in data-driven AI, with
389 a focus on the legal grounds for bias challenges and the societal implications of these solutions. Boykin et al. [22]
390 surveyed the intersection of ML and psychology, specifically with a focus on psychological mechanisms underlying
391 fairness preferences. Cheng et al. [35] conducted a systematic review of socially responsible AI algorithms with the aim
392 of examining the literature beyond algorithmic fairness and connecting major aspects of AI to societal well-being.
393

394 Mehrabi et al. [100] examined the fairness issues in real-world applications, specifically investigating the different
395 sources of bias in AI systems and providing a taxonomy of fairness definitions in the ML community. Similarly, Caton
396 and Haas [28] provide an overview of fairness mitigation approaches for ML, by categorising mitigation techniques into
397 several stages in the ML model development pipeline. The authors highlight 11 mitigation methods categorised into
398 three areas, namely pre-processing, in-processing, and post-processing methods. In addition, Zhang et al. [148] and Chen
399 et al. [32] surveyed the testing of several properties in machine learning (ML), including fairness properties. In contrast
400 to these works, we focus on fairness properties as it relates to the SE development pipeline. Beyond fairness mitigation
401 and testing, we examine works studying fairness formalization, empirical evaluation, detection and improvement.
402

403 Some researchers, Hutchinson and Mitchell [71] and Kwiatkowska [81], have also studied the notions of fairness
404 properties across domains, and for concurrent systems, respectively. Notably, Hutchinson and Mitchell [71] surveyed
405 the history of the definitions of fairness properties over the last 50 years across multiple disciplines, including education,
406 hiring, and machine learning. The authors compared past and current notions of fairness along several dimensions,
407 including the fairness criteria, the purpose of the criteria (e.g., testing) and how it relates to the mathematical method
408 for measuring fairness (e.g., classification, regression) and people (individuals, groups, and subgroups). Unlike the
409

⁴¹⁷ aforementioned works, our survey of fairness is *more general*, we study software fairness beyond advances in specialised
⁴¹⁸ ML communities. In this work, we additionally examine several papers from security, CHI, PL, CV and NLP venues.
⁴¹⁹

⁴²⁰ Overall, we provide a systematic literature review of the analysis of software fairness properties for learning-based
⁴²¹ systems. To the best of our knowledge, this work is the only systematic literature review concerned with the *analysis of*
⁴²² *fairness a (non-functional) software property of learning-based systems*. We are not aware of any other survey that focuses
⁴²³ on this research area, i.e., software fairness analysis, especially providing a comprehensive survey of its formalization,
⁴²⁴ testing, diagnosis and mitigation across several fairness metrics, biases and domains.
⁴²⁵

⁴²⁶ 3 METHODOLOGY

⁴²⁸ The research methodology employed in this work is based on the methodology detailed in Kitchenham [80]. In the
⁴²⁹ following, we provide the details of our research protocol:
⁴³⁰

- ⁴³¹ (1) **Aim and Scope:** First, we define the goals of this work and the scope of works to be examined. Then, we define
⁴³² the scientific questions, the analysis protocol and the information relevant for our data analysis. To this end, we
⁴³³ define the research questions (see subsection 4.3).
⁴³⁴
- ⁴³⁵ (2) **Detailed Information:** To analyse each paper in-depth, we identify the information necessary to answer all
⁴³⁶ research questions, this includes information that allow us to categorise, understand and describe the problem
⁴³⁷ addressed by each approach, and the technique or results provided by each paper. This informed the publication
⁴³⁸ search (e.g., our focus venues.), as well as the keywords employed in the filtering process used in identifying relevant
⁴³⁹ papers. Among several details, our interests include the studied fairness metrics and biases, the form of fairness
⁴⁴⁰ analysis (e.g., fairness testing) and the level of access (e.g., white, grey or black -box).
⁴⁴¹
- ⁴⁴² (3) **Publication Search:** We curated fairness-related publication via three means, (1) we searched the top venues in SE
⁴⁴³ (such as ICSE, TSE and FSE), Programming Languages (PL) (e.g., PLDI and POPL), Security (e.g., CCS and Euro S &
⁴⁴⁴ P), Artificial Intelligence (AI) (e.g., AAAI), ML (e.g., ICML, NeurIPS), CV (e.g., CVPR) and NLP (e.g., ACL, EMNLP),
⁴⁴⁵ (2) we collected papers from fairness focused conferences and workshops such as FairWare, FAT, FATML, and FaaCT;
⁴⁴⁶ and (3) we conducted a keyword guided publication search of paper repositories (such as ACM Digital library¹,
⁴⁴⁷ IEEE Xplore Digital Library² and Google Scholar³) using popular fairness-related terms such as “bias”, “fairness”,
⁴⁴⁸ “discrimination”, etc. In total we merged all collected papers which amounted to 420 papers from 65 venues. In
⁴⁴⁹ addition, we conducted a focused search of more recent works from SE venues since our initial study. We collected
⁴⁵⁰ an additional 72 publications mostly from the top SE venues from 2022-2023.
⁴⁵¹
- ⁴⁵² (4) **Filtering:** To identify relevant publications we filter out publications that are not relevant to our goals and analysis.
⁴⁵³ Specifically, we filter out papers that (1) do not conduct fairness analysis as a part of the software engineering
⁴⁵⁴ process, i.e., papers that are not relevant to the requirements, design and quality control of AI or ML -based
⁴⁵⁵ software systems; (2) we exclude papers that are not written in English language, duplicate papers, as well as short
⁴⁵⁶ papers or extended abstracts; (3) we also exclude papers that analyze fairness for other systems, i.e., fairness for
⁴⁵⁷ non-learning-based software systems. For instance, we exclude works on fairness of non-software-based systems
⁴⁵⁸ (e.g, transport systems [114], food inspection systems [120]); (4) we also exclude papers that are not focused on
⁴⁵⁹ software fairness properties, e.g., papers studying other properties such as robustness, consistency, security or
⁴⁶⁰ accuracy; (5) finally, we read the abstract of each paper and excluded papers that are not research papers studying
⁴⁶¹

⁴⁶⁵¹<https://dl.acm.org/>

⁴⁶⁶²<https://ieeexplore.ieee.org/Xplore/home.jsp>

⁴⁶⁷³<https://scholar.google.com/>

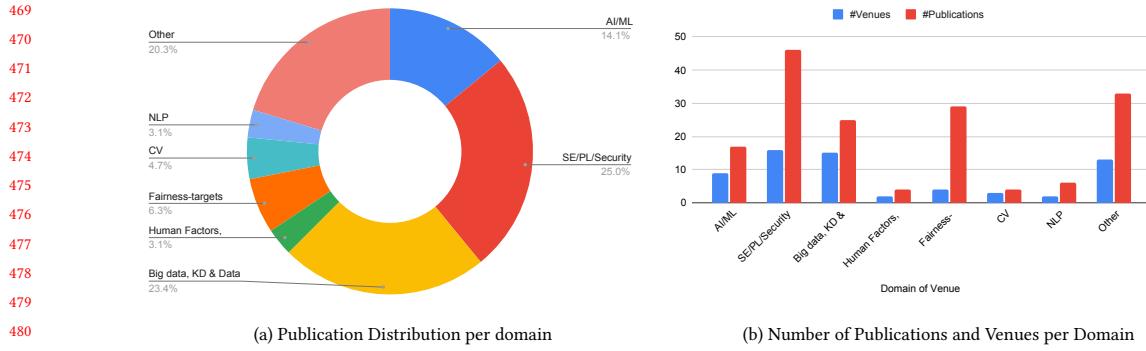


Fig. 4. Details of Publication domains

software fairness, for instance, we excluded surveys, literature reviews and invited lectures. From our initial search, we filtered out 256 papers and were left with 164 papers for our initial analysis. In our follow-up search of the recent SE literature (2022-23), we filtered 31 papers from the 72 collected publications. This resulted in the analysis of 41 additional papers. In the course of both publication searches, i.e., the initial search (2010-2021) and SE-focused search (2022-2023), we found a total of 492 papers (420 in first round and 72 in the second round), and filtered out a total of 287 publications (256 in first round and 31 in the second round). This resulted in a total of 205 analysed papers, i.e., 164 papers in the first round and 41 papers in the second round.

Given the large number of collected papers, we designed a research protocol to analyse each paper and extract certain information from the papers. For each paper, we extract both *the metadata* of the paper, as well as the *detailed research information*. In terms of metadata, we extracted the author details, the paper title, the year and venue of publication and the domain of the publication (e.g., SE, security or PL). For detailed research-relevant information, we studied the following details about the techniques and evaluation of the paper: the employed datasets, the studied fairness metrics, the studied biases (i.e., sensitive attributes, e.g., race), the form of analysis (e.g., fairness testing), the access level, the specific problem addressed, the main idea of the paper, the proposed solution, the resulting findings, as well as the strengths and weaknesses of proposed analysis/solution.

Categorization and Coding Protocol: In our analysis, one researcher collects, analyses and filters the publications following the aforementioned research protocol. For each paper, the researcher collects the metadata and documents all collected papers and their analysis then forwards all details to one other researcher for validation and inspection. For internal validation, at least one other researcher inspects the categorisation of the papers, takes note of conflicts or missing information and informs the researcher. Next, conflicts are resolved by organising a meeting to discuss the differences in the categories, naming of categories, collected data and ascribed descriptions. Finally, to ensure publications are correctly classified, we externally validate our findings and descriptions by sending a draft of our paper to all (cited) authors to provide feedback on mis-catergorization or wrong characterization of their work.

4 EVALUATION SETUP

4.1 Data Collection and Analysis

Our initial analysis involved 164 publications (from 2010 till 2021). Most of the research findings in section 5 (RQ1 to RQ5) are based on the analysis of the initial set of papers. In our follow-up analysis, we examined an additional 41 publications resulting from filtering out 31 papers out of 72 collected papers. These publications were collected from a Manuscript submitted to ACM

521 more focused search of the recent SE literature (2022 till 2023.) We examine these works and shed light on the progress
522 made since our initial analysis. We further discuss these newly published works and mention the domains of these new
523 works and how they relate to our initial findings in our result discussions and tables.

524 In the initial collection, we examine publications (164) that analyse software fairness in learning-based systems. **Table 1**
525 provides an excerpt of some of the collected publications.⁴ We analyse publications from different domains and venues,
526 including SE, PL, AI/ML, CV and NLP (see **Figure 4**). We then perform an *additional paper collection* from 2022 to
527 2023 that is focused on the recent works published in the top SE venues including 2022 and 2023 proceedings of ICSE,
528 ESEC/FSE, TOSEM, ISSTA and EMSE. The goal of this collection is to examine the recent research progress.
529

530 We analyze the volume of publications. We first examine the metadata of our publication corpus to determine the
531 volume of publications in different domains and venues. For this analysis, we collected a corpus of 164 papers from 64
532 different venues and eight (8) different academic domains/fields. **Table 1** and **Figure 6** show the details of publication
533 venues and the distribution of the publications in our corpus by venue. We are also interested in analysing the details of
534 the publications in terms of the type of venues (e.g., conference, journals and workshops) and the domain (or field) of
535 publication venues (e.g., SE or AI/ML). **Figure 4** and **Figure 5** show the details of our publications in terms of domains
536 and venue type, respectively. In **section 4** and **section 5**, we further discuss recent publications (2022-23) in contrast
537 to our initial study (these are also highlighted in brackets in **Table 2**, appendix (**Table 2**) and **Table 5**). We focus on
538 highlighting how such works buttress or address some of our initial findings. Finally, in **section 7**, we reflect on very
539 recent SE-focused works (2024) and how they also relate to our findings.
540

541 4.2 Initial Publication Details

542 **Volume and Domain of Publications:** Our analysis showed that even though algorithmic fairness is studied across
543 different research communities, *the study of fairness as a software property is mostly dominant in the SE and data-centric*
544 *venues*. Notably, *about half (about 48.3%) of all publications on software fairness are in the SE (e.g., FSE) and data-centric*
545 *(i.e., Big data, knowledge discovery and data mining) venues (e.g., KDD)*. Followed by typical top-tier AI/ML domains (e.g.,
546 AAAI and NeurIPS) which account for about 14.1% of all papers on software fairness. **Figure 4(a)** further illustrates that
547 the top software engineering venues (e.g., ICSE, FSE, ISSTA, ASE, TSE, OOPSLA) account for about one in four (about
548 25.0% of) publications, with this community accounting for most of the venues and publications (see **Figure 4(b)**). In
549 addition, we observed that most papers in our corpus are published in three main venues, namely ICSE (6.1%), FSE
550 (7.9%) or TSE (9.8%). Other popular venues include Big Data, Knowledge Discovery and Data Mining domains (23.4%), as
551 well as AI/ML domains (14.1%) (see **Figure 4**). Likewise, fairness-focused venues (such as FaccT, FATE and AAAI AIES)
552 account for about 6.3% of venues and about 17.7% of all publications. Meanwhile, more specialised venues (e.g., ECCV
553 and CAV) had the least number of works on software fairness. This shows that even though research works on software
554 fairness analysis are published across different domains, most publications are in SE venues. This demonstrates the
555 growing interest in software fairness within the SE community. In particular, the need to apply SE processes, methods
556 and tools to study the implications of fairness measures in practice has become paramount in the SE community. In
557 summary, almost half (48%) of software fairness analysis works are published in typical SE venues and Data-centric
558 venues, this is followed by top-tier AI/ML venues which account for about 14.1% of all publications.
559

560 **Publication Trends over time:** Examining the volume of publications over the years, we observed that *the number of*
561 *publications in software fairness analysis has been steadily increasing over the last decade*. **Figure 6(b)** highlights the trend
562

563⁴Our appendix (**Table 1**) also provides a more detailed version of this table.
564

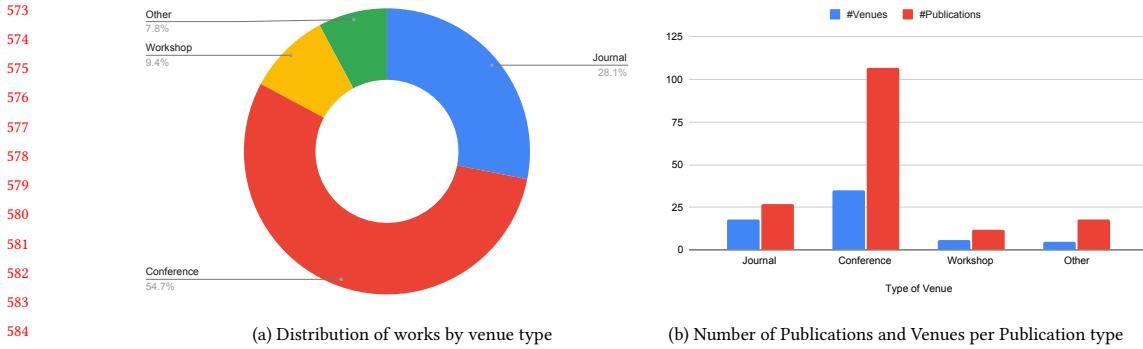


Fig. 5. Details of the type of Publication Venue (e.g., Conference or Journal)

over the last 13 years, it shows that the number of publications in software analysis has been steadily increasing over time, particularly in the last half decade. This trend signifies the growing research interests in fairness analysis in the SE research community. Particularly, a major surge in publications can be observed starting from 2017 till 2021. This is following the publication of Galhotra et al. [55] which *first formalizes causal fairness as a non-functional property of (learning-based) software systems*. Indeed, almost all software fairness papers analysed in this paper (152 papers, 92.68%) were published starting from 2017. This demonstrates the growing interests and increasing number of research output in this area over the years.

Type of Publication Venues: Figure 5 illustrates the distribution of the type of publication venues in our corpus. We observed that the majority (54.7%) of the publications on fairness analysis are in conference proceedings (see Figure 5 (a)). Figure 5(b) also shows that conferences are the most popular venues for software fairness publications, especially top-tier venues like ICSE and FSE, as well as popular conferences focused on Fairness (e.g., FaccT). Journal publications account for about 28.1% of all publications, with TSE, EMSE and JSS being the most popular journal venues. The most popular workshop venue for fairness analysis is FairWare. These findings show that software fairness has (more recently) become an important research area for top-tier SE conferences and journals.

Advances: We analyse the advances made in the analysis of software fairness by focusing on six major areas of fairness concern and their advances over time. Specifically, we analyse the trends in engineering concerns in software fairness such as validation, verification, design, empirical studies, tooling and datasets. Our appendix further provides additional figures illustrating the trends and advances: Appendix (Figure 5) provides details of the trend and advances in publications for all of these six concerns. Additionally, the appendix (Figure 7 and Figure 8) provide detailed distribution of publications in each area over the years.

Generally, the number of publications has increased over time for all six SE areas. However, there has been major increase in publications in the area of fairness validation (see appendix (Figure 7(a))), bias-aware datasets and empirical studies (see appendix (Figure 8 (b) and (c))). This is evident by the number of publications (about eight to 12 publications yearly) in these areas in the last half decade (since 2018). In addition, analysing the span of publications also show that fairness-related validation and datasets have been studied for about seven consecutive years, while empirical studies in fairness have only recently become prominent (i.e., in the last 5 years).

Meanwhile, other SE concerns such as the verification, tooling and design of fair software systems have received little attention. These areas have seen fewer (four to six maximum) publications in the last years. However, the span

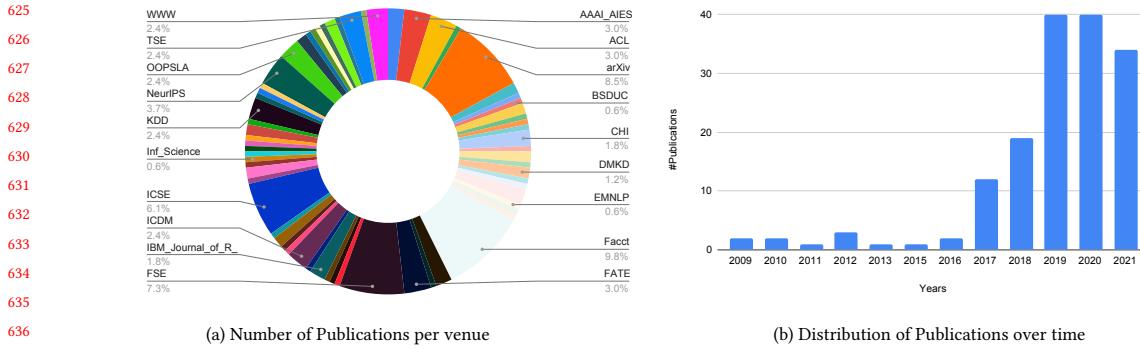


Fig. 6. Distribution of Publication Venues and Year of Publications. (Table 1 provides an overview of venues, note that some sectors are unlabeled in Figure 6(a) due to space constraints.)

Table 1. Excerpt of Publication Details (“#” means “number of”). More detailed table is provided in the appendix (Table 1).

Domain (#Pubs.)	Type	#Venues	Venue Sample Venue(s)	Publications		
				#Pubs	Example	Years
Software Engineering (SE), Prog. Lang. (PL) & Security	Conference	9	ASE, CAV, EuroS&P, FSE, GECCO, ICSE, OOPSLA, TrustCom, ISSTA	28	[36]	2017-21
Natural language processing (NLP)	Journal	4	EMSE, JSS, RE, TSE	7	[14]	2009-21
Artificial Intelligence (AI) & Machine Learning (ML)	Conference	8	AAAI, AISTATS, ICML, NeurIPS, PMLR	17	[146]	2013-21
Computer Vision (CV)	Conference	2	ICCV, CVPR	3	[134]	2019-20
Fairness-targets	Conference	2	AAAI-AIES, Facct	21	[15]	2017-21
Big Data, Data Mining (DM), & Knowledge Discovery (KD)	Conference	8	DMKD, ECML-PKDD, EDBT, KDD, ICDM, ICEDT, ICMD, LAK	18	[78]	2010-21
	Journal	5	Big Data, Inf. Science, JDIQ, KAIS, SIGMOD-Record	5	[76]	2012-19
Human Factors & Usability	Conference	1	CHI	3	[82]	2019-21
	Journal	1	IWC	1	[25]	2016

of publications in the design of fair software is larger, with at least a 10 year spread, showing that there has been a steady interest in this area over the years (see appendix (Figure 7(c))). This implies the need to investigate and examine approaches and methods to address the under-studied areas, especially the verification and tooling of fairness-aware software systems. There has been advances in several areas of software fairness over the years, but most consecutive publications has been in fair learning and validation, with up to 10 yearly publications in the last five years.

4.3 Research Questions

Firstly, we investigate the purpose of these publications including the main idea of the proposed techniques as well as how researchers study software fairness as a software engineering task (RQ1). For instance, we examine if the aim of the proposed method or analysis is to formalize, test, mitigate or diagnose fairness issues in learning-based software systems. Secondly, we analyse the fairness measure studied in the papers, such as individual, group, causal or intersectional fairness (RQ2). Next, we study the bias, i.e., the sensitive or protected attributes (e.g., gender or race) studied in the literature (RQ3). We also examine the tasks and datasets employed in the reviewed publications (RQ4). Finally, we investigate the tools available for software fairness analysis in the literature (RQ5).

Specifically, we aim to address the following research questions (RQs).

Table 2. Purpose of Fairness Analysis (recent works (2022 -2023) are in bracket). More details are provided in the appendix (Table 2).

Categories	Sub-category	Description	#Pubs	Sample Works
Validation	Testing	generating discriminatory test inputs to expose fairness violations	20 (13)	[151] ([137])
	Mitigation	mitigating bias in software systems, e.g., via repair and prevention	14 (13)	[4] ([56])
	Debugging	diagnosis and explanation of fairness violations	8 (3)	[87] ([102])
Verification	Auditing	analysing and measuring bias in software systems	2 (1)	[26] ([150])
	Verifiers	verifying that a system fulfills a fairness metric or goal	12 (1)	[57] ([18])
	Certification	certifying that a system fulfills a fairness goal	4 (1)	[117] ([18])
Design	Proof or Guarantees	providing a formal proof that a system achieves a fairness goal	2 (2)	[101] ([63])
	Requirements	requirement engineering and formalization of fairness properties	4 (2)	[51] ([11])
	Bias-aware Design	designing fair systems and bias-aware software	15 (5)	[27] ([72])
Empirical Evaluation	Analysis	empirical studies about fairness concerns	22 (14)	[17] ([136])
	Benchmarking	providing fair benchmarks or benchmarks for fairness evaluations	4 (3)	[21] ([58])

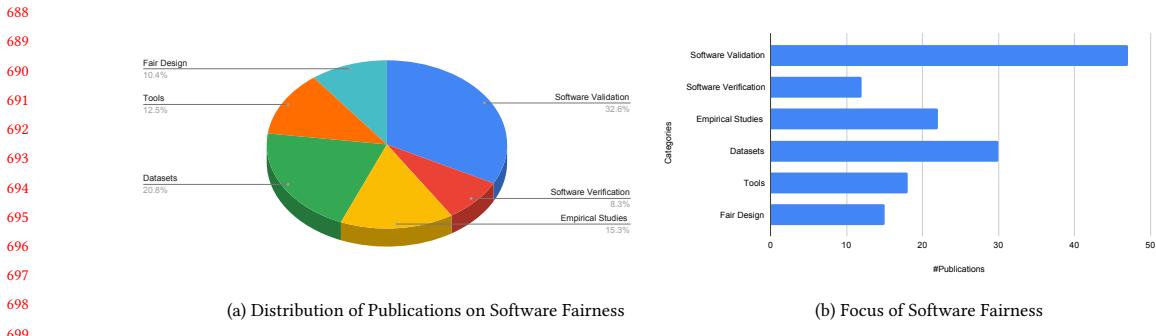


Fig. 7. General Focus of Fairness Analysis

- **RQ1 Purpose of Fairness Analysis.** What is the purpose of fairness analysis in the literature? What fairness problem is addressed or studied in the literature? What is the target or focus of the community in terms of fairness? What areas have been well-examined or un(der-)investigated in the research community?
- **RQ2 Fairness measure.** What are the fairness metrics analysed by the research community (e.g., individual, group or causal fairness)? What metrics are well-studied, under-studied or not investigated?
- **RQ3 Bias and Sensitive Attributes.** Which (societal) biases (e.g., age, race gender or religion) are investigated in the analysis of software fairness, especially w.r.t. to protected or sensitive attributes?
- **RQ4 Datasets and Tasks.** What datasets and tasks are employed for software fairness analysis? What tasks/-datasets have been well or under-studied? What is the distribution of datasets per tasks?
- **RQ5 Tooling.** What are the available fairness analysis tools and frameworks? What problems do they address, what are the analysis goals supported by available tools? Which analysis approaches are employed in the tools? What stage of model processing do these tools support, and what level of model access do they require?

5 RESEARCH FINDINGS

In this study, we investigate the purpose of each paper collected in this survey, particularly, we examine and categorise the fairness problem studied or addressed by each paper. Table 2 provides details of the purpose of software fairness analysis performed in the literature. Overall, we identified six (6) categories for all collected papers. In the following, we discuss the problem addressed in each category, and the notable works that address such issues. In addition, we highlight the gaps in each category and across all identified categories. Figure 7 shows the research focus and purpose of the fairness analysis conducted in the literature. Appendix (Figure 6) also provides more details. We also examine how the research community analyzes software fairness, especially the SE, PL and Security venues. We are interested in

the focus of the community, the areas that are well investigated by the community, and the areas that are not explored. Appendix (Table 4) also provides high level details of publications.

Generally, we observed that *the focus of the research community has been on the validation, design and empirical studies of software fairness*. Particularly, one in three (about 33% of) the collected papers study fairness validation, i.e., the testing, debugging and mitigation of fairness errors (see Figure 7(a)). Analogously, more than one-third of collected papers (together over 36%) either study biases in datasets (about 21%) or conduct empirical evaluation of software fairness (about 15.3%). Appendix (Figure 6) also shows the focus of the papers published in SE venues. An inspection of these works further show that *dataset analysis, empirical studies, testing and tooling of software fairness is popular* in the community (see appendix (Figure 6 (a))). In particular, fairness *debugging, requirements analysis, benchmarking, and auditing are not popularly studied in the community* (see appendix (Figure 6(b))). On one hand, this suggests that the focus of the majority (70%) of the (SE) research work is focused on the validation, empirical evaluation and data(sets) analysis of fairness properties. On the other hand, concerns such as the auditing, debugging and requirements analysis are under-studied.

*Most publications (70%) study the validation, empirical evaluation and data(set) analysis of software fairness:
The design, verification and tooling of fairness-aware software have been under-studied (about 30%).*

Bias Testing, Debugging and Mitigation. The research papers in this category aim to detect, expose, diagnose or mitigate fairness issues in learning-based software systems. Appendix (Figure 6) provides the distribution for each category. The bulk of this work are general-purpose approaches proposed in the SE, PL and security venues, while the rest of the proposed approaches are more *specialised* approaches proposed for analysing applications in specific domains, e.g., CV or NLP venues. Overall, we found that most approaches are focused on fairness testing, some methods are focused on the mitigation of unfairness, while very few techniques are focused on debugging (diagnosing and understanding the root causes of) unfairness. Indeed, we note that the root cause of unfair program behaviors stem from multiple sources, including societal or historical sources. In this work, we refer to debugging in terms of identifying the software components or configuration that induce the fairness behavior in the software. This notion is inline with previous works [30, 130]. In the following, we shed more light on the main idea of available approaches and the gaps in techniques.

Fairness testing techniques employ a plethora of techniques for test generation, including ML, search and program analysis techniques. Appendix (Table 4) provides further details on the fairness testing methods proposed in the literature. These test generation methods include several black box testing approaches (especially, input-based approaches), and few white-box techniques. Notably, *there are even fewer grey-box fairness testing approach*. Thus, most proposed approaches drive the generation of discriminatory inputs either by *analysing the model under test (MUT)* or the *input space*. However, there are very few techniques that leverage both the input space and the model analysis, besides, there are few studies investigating the relationships between both dimensions for testing purposes (e.g., Tizpaz-Niari et al. [130]). Notably, white box approaches (e.g., ADF [151] and EIDIG [149]) mostly employ ML techniques (e.g., gradient computation, and clustering) to drive the generation of discriminatory test cases. Meanwhile, black-box approaches focus on leveraging the knowledge of the input space, program analysis and/or search algorithms to generate discriminatory inputs. They mostly employ templates, schemas, grammar, mutation or search algorithms to drive fairness test generation [122, 131]. Other approaches employ program analysis, e.g., symbolic execution [2] and combinatorial testing [103] to drive the generation of discriminatory test inputs. Notably, *we found few grey-box testing techniques that leverage both the input space and internal model attributes/properties to drive the generation of*

781 782 783 784 785 786 787 788 789 790 791 792 793 794 795 796 797 798 799 800 801 802 803 804 805 806 807 808 809 810 811 812 813 814 815 816 817 818 819 820 821 822 823 824 825 826 827 828 829 830 831 832

discriminatory inputs. Moreover, there is little work that studies the link between the properties of the input space or discriminatory inputs to other internal model attributes/properties for fairness testing (e.g., Tizpaz-Niari et al. [130]).

Fairness testing approaches are mostly black or white box test generation methods: There are few grey-box approaches that leverage (or study) the relationship between the input space and internal model properties.

Fairness Verification, Certification and Proof Guarantees: We examine the literature on the verification, certification and proving of fairness properties in learning-based software systems. Specifically, we examine the categories of verifiers, the main ideas of proposed verification approaches and the gaps in this research area. To this end, we identify three major kinds of fairness verification approaches, namely *distributional verifiers* (e.g., FairSquare [3] and VeriFair [12]), *specialised verifiers* designed for a particular domain, metric or task (e.g., [29]) and sample-based verifiers (e.g., Themis [55]). In our study, most approaches are *distributional-based* verifiers or *specialised* verification techniques, and there are fewer *sample-based verifiers*. Distributional verifiers typically encode fairness metrics as probabilistic properties then verify such properties with respect to the underlying data distribution of the learning-based system. On the other hand, sample-based approaches (e.g., Themis [55]) allows to detect and verify fairness metrics based on a fixed dataset, otherwise they generate counter-examples to refute the satisfaction of the fairness property. Other verification approaches target specific domains (e.g. NLP [95]), specific fairness metrics (e.g., individual fairness [73] and disparate impact [46]), or tasks (fair training [29] and data debiasing [46]).

Some verification approaches encode fairness metrics as probabilistic properties, then provide guarantees over the system's data distribution. These approaches take as input the probability distribution of the attributes in the dataset and the MUT, then verify the fairness of the system with respect to the distribution and MUT. For instance, FairSquare [3] presents a technique for verifying and certifying fairness properties by encoding fairness definitions as probabilistic properties. Moreover, Albarghouthi et al. [4] also proposed an approach that applies *distribution-guided inductive synthesis* to verify (and repair) *unfair* ML classifiers. Likewise, Bastani et al. [12] developed an algorithm for verifying fairness specifications (called VeriFair), the algorithm provides probabilistic guarantees for fairness properties and allows users to verify that the probability of fairness errors is small. Ghosh et al. [57] presents a stochastic satisfiability (SSAT) framework (called Justicia) to formally verify fairness measures of supervised learning algorithms with respect to the underlying data distribution. Justicia is applicable to different fairness metrics including disparate impact, statistical parity, and equalized odds. Compared to previous distribution-based verification approaches (FairSquare [3] and VeriFair [12]), Justicia supports non-Boolean and compound sensitive attribute. It also provides theoretical bound for the finite-sample error of the verified fairness measure, and it is more robust than sample-based verifiers.

There are several verification approaches that are focused on a specific domain, metric, or task (e.g., debiasing datasets [46] or fair training [29]). For domain- or task-specific approaches, Ma et al. [95] provides a black-box technique to enforce fairness guarantee for NLP systems by leveraging advances in certified robustness of machine learning. Their approach employs a neutral phase to piggyback the NLP model to smooth its outputs such that they are certified to preserve individual fairness. Similarly, Liu et al. [91] presents a (semi-)automated verification framework (called FairCon) to ascertain the fairness of smart contracts, their approach can refute false claims with concrete examples, or certify that contract implementation fulfil desired fairness properties.

For metric-specific approaches, John et al. [73] proposes sound (but incomplete) verifiers for proving individual fairness of models by employing appropriate relaxations of the problem, specifically for linear classifiers and kernelized polynomial/radial basis function classifiers. Likewise, Feldman et al. [46] presents a verification technique for certifying the (im)possibility of disparate impact on a data set by employing a regression algorithm that minimizes the balanced

833 error rate (BER) of the dataset. The goal is to verify that a protected or sensitive attribute can not be predicted from the
 834 other attributes in the dataset by ascertaining if there is sufficient information about the dataset to detect sensitive
 835 attributes from the data. In terms of verifying fair training, Celis et al. [29] propose a technique to train fair classifiers
 836 with theoretical guarantees, using a meta-algorithm for classification that can take as input a general class of fairness
 837 constraints with respect to multiple non-disjoint and multi-valued sensitive attributes. Notably, this approach can
 838 handle non-convex fairness constraints such as predictive parity.
 839

840 A different line of fairness verification techniques are sample-based verifiers such as AIF360 [13] and Themis [55].
 841 Typically, these approaches leverage software testing techniques to verify fairness properties on fixed data sample.
 842 AIF360 [13] is an extensible open source toolkit for detecting and verifying fairness properties, particularly for a fixed
 843 data sample. It provides an array of methods to detect and report several fairness performance metrics. This includes 71
 844 bias detection metrics, and nine bias mitigation algorithms. Meanwhile, Themis [55] allows developers to verify that a
 845 fixed sample do not discriminate against specific sensitive attributes (e.g., race and gender) by automatically generating
 846 discriminatory that verifies that changing the instance of the attribute does not cause a change in the output of the
 847 learning-based system. Overall, these approaches leverage advances in software testing to measure and verify that a
 848 fixed sample data fulfills specific fairness properties.
 849

850 There are other general verification approaches that verify multiple (user-defined) fairness constraints or other
 851 properties beyond, but including, fairness properties. As an example, Metevier et al. [101] addressed the problem of
 852 verifying multiple fairness definitions as well as user-defined fairness metrics for learning-based systems. The authors
 853 proposed a verification approach (called RobinHood) which employs an offline contextual bandit algorithm determine
 854 the satisfiability of several fairness constraints. Morevoer, RobinHood provides a probabilisitic guarantee of fairness by
 855 ensuring that it does not return a solution with a probability greater than a user-defined threshold. Besides, Sharma
 856 et al. [118] is a more general verification approach which is also applicable to fairness properties. The authors propose
 857 a (white-box) verification approach that employs the knowledge of the internal structure of the model to verify that ML
 858 models fulfil several properties (including fairness), in particular by training a shadow model that approximates the
 859 MUT by using the prediction of the original model as training data. It employs a property specification language to test
 860 and verify model properties of learning-based software and provides counter-examples (i.e., test cases violating the
 861 property) if the property is not fulfilled.
 862

863 *Most fairness verifiers are distribution-based, sample-based or specialised for a specific fairness measure, domain or task.*

864 *There are very few verifiers that support multiple or user-defined fairness constraints.*

865 **Others:** Our investigation showed that the SE research community has *mostly* focused on analysing software fairness
 866 as a *fairness validation* (i.e., *testing and debugging*) problem [131] and as a *fair system design* problem, especially to
 867 mitigate biases [31]. However, other aspects of SE concerns are under-studied, such as the formalization of fairness as a
 868 *software requirement* [24, 51], the *verification* of fairness properties [3], and *empirical evaluation* of software fairness
 869 properties [67]. Furthermore, some aspects of SE concerns are hardly studied by the community, namely, empirical
 870 evaluation of fairness properties (especially human factors in software fairness [49, 60]) and the maintenance of fairness
 871 properties as the software evolves (e.g., because of model re-training, model compression [65] or software regression).
 872

873 In 2017, Brun and Meliou [24] formalised software fairness as a software engineering problem that needs to be
 874 tackled by all aspects of software engineering. These aspects include steps in the software development life cycle
 875 such as requirements engineering, design, testing, verification and maintenance. Since their publication, the bulk of
 876 the published papers have *focused on the design, testing and mitigation* of software fairness, these areas have been
 877

well explored and investigated by these communities. For instance, several papers have explored the problem of *fairness testing* by employing random test generation (Themis) [5], local search algorithms (Aequitas) [131], gradient computation [152], mutation testing (TransRepair) [125], grammar-based testing (Astraea) [122], symbolic execution [2], property-driven testing [118] and schema or template based testing (BiasRV) [142]. In addition, the problem of fair system design to mitigate biases has been studied by a few researchers via several techniques including behavior mutation [67], and feature or dataset manipulation [147]. Researchers have also proposed search-based optimization methods for the fairness-aware design, testing and mitigation. For instance, Perera et al. [109] proposed a search-based method for regression fairness testing and Hort et al. [68] presents a search-based method for the repair of fairness and accuracy properties of ML-based software. Hort et al. [66] have proposed a multi-objective search method to determine a balance between gender-fairness and semantic correctness of word embeddings using Word2Vec.

Very few researchers have conducted empirical evaluation of software fairness properties in real-world applications. Notably, Biswas and Rajan [16] conducted an empirical study to study several fairness mitigation techniques including their impact on performance. Likewise, Hort et al. [67] conducted a large scale empirical study to test the effectiveness of 12 widely-studied bias mitigation methods. Meanwhile, Zhang and Harman [147] empirically evaluated how feature set and training data affect fairness. The focus of most studies has been on fairness mitigation, except Zhang and Harman [147] that empirically studied the impact of features and datasets on fairness properties. Besides these two concerns, we have found few studies in these fields empirically evaluating other SE concerns (such as human factors [64]) or concerns relevant to steps of the model/software development pipeline.

However, some aspects of software fairness analysis remain under-investigated. Firstly, there is little work addressing concerns about the maintenance of software fairness, for instance, as the software or model evolves over time (i.e., software regression analysis) or is optimized (e.g., via model compression for edge devices). For instance, there is a recent paper investigating the impact of model compression on software fairness [65]. We also found few works supporting the software engineering activities around such changes. Notably, fairness testing for changes in ML-based software, i.e., in a regression scenario has been examined by Perera et al. [109]. The authors proposed a search-based approach to address regression in ML-based systems. Likewise, there are very few empirical studies on software fairness properties, particularly, there has been few empirical evaluation in this area involving humans, which is vital to determine the harm caused by unfair software behavior [64] and practitioners (developers) perception of software fairness properties [49]. Secondly, the formalization and definition of fairness as a software requirement, metric or measure has only been performed by a few researchers. We found few papers in this area including Verma and Rubin [132] and Finkelstein et al. [51].

The SE research community mostly study software fairness properties as mitigation, design and testing problems. Very few works have studied fairness as a requirements engineering or verification problem, and fewer works empirically studying human factors of fairness properties and how to maintain fairness as software changes/evolves.

5.1 RQ2 Fairness measure.

In this research question (RQ2), we examine the fairness metrics analyzed by the studies and methods proposed for software fairness analysis. We investigate the number of studies examining the different classes of fairness metrics (e.g., statistical measures, similarity-based measures and time-based measures), as well as the distribution of specific metrics, such as individual fairness, group fairness and causal fairness. Figure 8 highlights our analysis of the distribution of these metrics across proposed methods and conducted studies.

Generally, we observed that most studies examine statistical or similarity based measures (such as individual, group or causal fairness), while time-based metrics (e.g., sequential or long-term fairness) or measures based on causal reasoning (e.g., fair inference) are not (yet) studied in the SE community as software fairness metrics. Statistical and similarity based fairness metrics (such as group, individual, and intersectional fairness) are the most examined metrics in software engineering (see Figure 8(a)). For instance, individual fairness and group fairness account for more than three in four (76.5%) studies and methods in software fairness analysis. Causal reasoning based fairness metrics (such as causal fairness) are also commonly studied. However, time-based metrics were not found in the SE literature [153].

Statistical and similarity -based fairness metrics are popularly studied (86.5%) in fairness analysis, but metrics hinged on causal reasoning are under-studied (about 10%) and we could not find a single work studying time-based metrics.

In addition, we observed that fairness metrics such as *individual and group fairness metrics are well studied in the SE community*. Concretely, about 42.5% and 35% of publications in software engineering venues study individual fairness and group fairness, respectively. The focus of most of these studies is in terms of testing, validating and mitigating software to fulfill these metrics.

Several of these techniques are tailored towards testing, discovering and mitigating one or both of these metrics. Concretely, 27% of examined studies study only individual discrimination [131, 149, 151]. For instance, Zhang et al. [149] proposed EIDIG (Efficient Individual Discriminatory Instances Generator), a scalable and efficient approach to systematically generate test cases that violate the individual fairness for DNN models. Likewise, Zhang et al. [151] proposed approaches to search for individual discriminatory instances of DNN, using lightweight procedures like gradient computation and clustering. Overall, we found that 23% of examined papers in the community study both individual and group fairness, 17% examine only group fairness and 27% study only individual fairness.

In our analysis, causal fairness is also a commonly studied fairness measure in the community, it accounts for one in ten publications. Notably, Galhotra et al. [55] and Biswas and Rajan [17] have studied the testing and debugging of causal fairness. Meanwhile, other metrics such as intersectional bias, local and group fairness (especially in ML or data pipelines) are the least studied in the community, with each accounting for less 2.5 to 7.5 percent of all examined works. As an example, Cabrera et al. [26] studied the analysis of intersectional bias in SE venues, specifically, in terms of testing, auditing and visual analysis of intersectional bias, respectively.

A recent survey by Gohar and Cheng [59] examines the state-of-the-art in intersectional fairness and presents a taxonomy for intersectional notions of fairness and mitigation. Overall, there are still several open issues in this area. Particularly, how to effectively test, mitigate and debug intersectional software bias and how intersectional bias interacts with other fairness metrics (e.g., group and individual fairness).

We also examine the composition, sequential and long-term aspects of fairness properties. Notably, Gohar et al. [58] have recently studied fairness composition in ensemble models. The authors show that ensembles can be designed to be fairer without using mitigation techniques and provide a benchmark for studying fair ensembles. Analogously, Chen et al. [33] proposed an ensemble method (called MAAT) for improving fairness-performance trade-off for ML software. MAAT combines models optimized for the two objectives – fairness and ML performance. Finally, on the extreme end, we found almost no work in the SE community studying time-based fairness metrics such as sequential and long-term fairness [116, 153]. These metrics are important because of maintaining fairness properties as the software evolves. As an example, consider an automated classifier that is re-trained periodically due to new data requirements, how do we ensure that such data and the resulting classifier maintains the fairness property defined in the previous system?

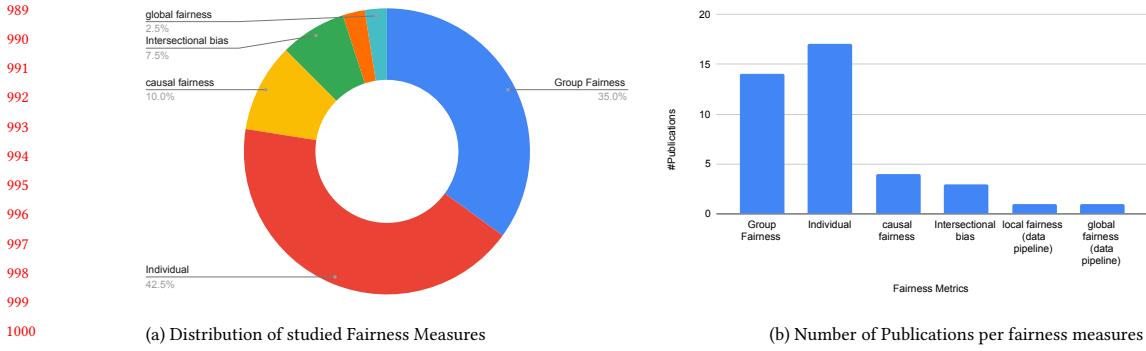


Fig. 8. Details of Fairness Measures

Overall, this shows that the focus of the community is on a set of statistical and similarity-based fairness metrics (such as individual, group, causal), hence, ignoring other metrics, e.g., time related bias concerns.

Most SE studies (86.5%) investigate individual fairness, group fairness and causal fairness metrics, metrics such as intersectional fairness and time-based metrics (sequential and long-term fairness) are under-studied in the SE community.

5.2 RQ3 Bias and Sensitive Attributes.

Let us investigate the societal biases studied in software fairness. In particular, we analyze the sensitive or protected attributes (e.g., age, race, gender etc) that researchers study or develop techniques for. Figure 8 and appendix (Figure 9) illustrate the distribution of sensitive attributes in the literature. In Table 3, we exemplify and illustrate some of these sensitive attributes. We note that the examples in Table 3 are not real violations. However, they are adapted from fairness violations discovered by state-of-the-art NLP bias testing tools [95, 122].

We observed that *about four in every five studies examine four main protected attributes, namely age, race, gender and country*. These four attributes are studied in about 79.5% of works, with age, race and gender accounting for the majority (about 72.5% of all works). We believe this is due to the availability of these attributes in several datasets. For instance, attributes such as age, gender and race are popular features in tabular datasets (e.g., German Credit, Adult Census Income and Bank Marketing). These attributes are popularly studied due to the simplicity of manipulating them. They typically have a tractable input space, bounded constraints or limited range of values. As an example, it is easy to manipulate gender features in tabular data, since they have a small number of possible values (e.g., male, female or non-binary). Meanwhile, rare protected attributes (e.g., language features [39, 92]) either possess a complex input space, are uncommon in popular tabular datasets, or are more relevant to semi-structured datasets (e.g., text, audio and images). Such attributes are more difficult to study or manipulate: For instance, consider an image dataset where the gender-related sensitive attribute are related to pixel values [86]. Despite the complexity of rare attributes, some researchers have explored their fairness concerns. Notably, Black et al. [19] and Lipton et al. [90] both employed hair length, work experience and academic performance in their works for fairness analysis. More recently, researchers have also explored fairness concerns in speech tone [143]. Overall, this analysis suggests that there is a gap in fairness analysis of rare sensitivie attributes. Indeed, the focus of the community has been on the sensitive attributes that are available in less complex tasks or tabular datasets.

Table 3. Illustration of studied biases, i.e., sensitive/protective attributes. These examples are adapted fairness violations discovered by state-of-the-art bias testing tools [9, 95, 113, 122]. Sentiment analysis (SA) systems detect the emotional situation or state in a sentence, \odot indicates a positive sentiment (e.g., happy), \ominus indicates a negative sentiment (e.g., sad), bias-inducing inputs are underlined and unfair outputs are in red (e.g. $\textcolor{red}{\odot}$).

Sensitive Attributes	#Pubs	Illustrative Example for text-based SA system	Sample Works
<i>Gender</i>	24	“The {man/woman} is happy.” = \odot/\odot	[30, 31, 131]
<i>Race</i>	21	“The {white man/ black man} is happy.” = \odot/\ominus	[55, 151, 152]
<i>Age</i>	15	“The {young man/ old man} is happy.” = \odot/\odot	[16, 147, 149]
<i>Country</i>	6	“The {american man/ chinese man} is happy.” = \odot/\ominus	[111, 152]
<i>Language</i>	2	“The {english orator/ spanish orator} is happy.” = \odot/\odot	[125, 126]
<i>Occupation</i>	2	“The {manager/ <u>cleaner</u> } is happy.” = \odot/\odot	[9, 122]
<i>Religion</i>	2	“The {nun/ atheist} is happy.” = \odot/\odot	[122, 152]
<i>Ethnicity/Accents</i>	2	“The {european man/ asian man} is happy.” = \odot/\odot	[111, 152]
<i>Class label (e.g., poverty level)</i>	1	“The {rich man/ poor man} is happy.” = \odot/\odot	[30, 129]
<i>Narcotics Arrests /Gang Affiliation</i>	1	“The {convict/ <u>innocent man</u> } is happy.” = \odot/\odot	[19]
<i>Work experience</i>	1	“The {experienced farmer/ novice farmer} is happy.” = \odot/\odot	[19]
<i>Academics (LSAT /GPA)</i>	1	“The {honors student/ failing student} is happy.” = \odot/\odot	[19]
<i>Marital-status/Relationship</i>	1	“The {married man/ single man} is happy.” = \odot/\odot	[26]

Table 4. Excerpt of Tasks and Datasets employed in Software Fairness Analysis. More detailed table is in appendix (Table 3).

Type	Task Category	#Data.	#Pubs	#Tasks	Tasks (Example Pubs)	Datasets (#Pubs)
Tabular	Education	2	3	1	academic performance [19]	Law School (1)
	Finance	6	54	4	income [2] credit default [17] potential buyers [17] fraud [2]	Adult Census (19) German Credit (17), Default Credit (3), Home Credit (3) Bank Marketing (11) Fraud Detection (1)
	Legal	3	10	3	recidivism [31] arrests [19] US executions [2]	COMPAS (8) Chicago Strategic Subject List (SSL) (1) US Executions (1)
	CV (image)	7	10	2	face detection image recognition	ClbA-IN (1), PPB (1), LFW (2) COCO (2), imSitu (1), CIFAR (2), ImageNet (1)
	NLP (text, speech)	12	10	6	SA [9], CoRef [122], MLM [122] toxicity machine translation	Twitter (1), IMDB (1), EEC Dataset (1), Labor statistics (1) Wiki Comment (1), Jigsaw Comments (1) News Commentary (2)
				ASR [111]	Speech Accent Archive (1), RAVDESS (1), Multi speaker Corpora of the English Accents in the British Isles (1), Nigerian English speech dataset (1)	
				SE (code)	Programs [36]	Bug2Commit (1), Diff Review (1), Code AutoComplete (1), Oncall Recommendation (1)

About four in five (79.5%) works study age, gender or race as sensitive attributes, while model class-label and specific attributes such as relationship (status), class (academics, work) and religion are understudied (less than 2.5% each).

5.3 RQ4 Datasets and Tasks.

For this research question, we examine the datasets examined or employed in the collected publications. Table 4 provides details of the tasks and datasets employed in (the evaluation of) software fairness analysis. Furthermore, Figure 10 and appendix (Table 3, Figure 10 and Figure 11) highlight the type of the dataset (e.g., tabular, or semi-structured), the distribution of examined datasets, and the task associated with each dataset, respectively. We also examine the task category associated with the examined datasets (e.g., NLP, CV, etc.), the volume of publications associated with each task category, and dataset (see appendix (Figure 12, Figure 13 and Figure 14, respectively)).

Volume of Publications per dataset: We found that *more than half of the examined studies employ four (4) major datasets, most of which are tabular datasets for finance-related tasks.* Appendix (Figure 10) also illustrates the distribution of publications using each dataset. Notably, the most common datasets are the Adult Census Income, German Credit, Bank Marketing and COMPAS dataset, they are employed in over half (52%) of all (SE) papers. These datasets account for most of the software fairness works we examined (see appendix (Figure 10 (b))). Most of these datasets are tabular datasets for finance-related tasks such as income prediction, except COMPAS – a popular legal dataset for recidivism. However, semi-structured datasets such as image (COCO, CIFAR) and text/NLP (News Commentary) are among the least studied datasets. This suggests that most research work employ similar datasets, which are mostly tabular or finance related. Hence, implying there is the need for a more diverse evaluation of fairness analysis techniques that are general and cuts across several datasets and tasks.

*Tabular datasets (e.g., Adult Census Income) are the most employed datasets in software fairness analysis.
Semi-structured datasets (e.g., image, text, code and speech) are less popularly studied.*

Type of examined datasets: We observed that *the majority (64%) of the research works study software fairness for tabular datasets.* Appendix (Figure 10(a)) shows that about two-third of the examined papers study fairness properties relating to tabular datasets (e.g., Adult Income Census). *Research works that generalise to both tabular and semi-structured datasets are few (about 14%). Fewer studies (about 21% of papers) solely analyse fairness properties for semi-structured dataset (e.g., text, image, speech or code).* Appendix (Figure 10(b)) provides the distribution of works studying semi-structured datasets. Inspecting works studying fairness properties in semi-structured datasets, we observed that most (about 72%) are focused on text (i.e., NLP-related) datasets or image (i.e., CV) datasets. *Datasets relating to speech (i.e., audio) and code (i.e., programs) are the least studied in software fairness,* accounting for about 18% of the examined papers. Generally, the focus of the community has been on tabular datasets. Overall, this suggests that the SE research community has been focused on studying fairness properties as related to tabular (relatively less complex) datasets, while more semi-structured datasets are under-explored. This suggests the need to develop fairness analysis techniques that are general and agnostic, i.e., applicable to both tabular and semi-structured datasets.

Researchers have recently taken interest in analyzing fairness properties in program-based, and SE-task centric activities. For example, researchers have examined fairness properties for SE-centric tasks such as issue assignment [105], buggy commit classification [36], code review recommendation [97], crowd worker testing recommendation [133], automatic code generation [70] and software development pipelines [61, 136]. Similarly, there has been an increasing attention to bias analysis of LLMs [1, 70, 115], especially due to their importance and recent popularity.

*Fairness concerns relating to semi-structured datasets such as code and speech (audio) are rarely studied,
they account for only about 18% of all inspected papers.*

Task Categories: Figure 10 and appendix (Figure 12) highlight the details of the task categories and the datasets employed for each task (category). We found that *tasks involving finance, computer vision (CV) and natural language process (NLP) are the most studied in software fairness publications, both in terms of the number of tasks and the number of publications.* However, software fairness concerns for tasks involving *education and code analysis* are less frequently studied. Appendix (Figure 12) further illustrates that all of the well studied task categories (CV, NLP and finance) have up to two to four tasks each, while education had one task. As an example, finance-related publications contribute over 50 papers (see appendix (Figure 12(a))), with about four distinct tasks examined (appendix (Figure 12(b))). These tasks include predicting income, credit default, fraud and potential buyers (Figure 10). In addition, despite fewer publications,

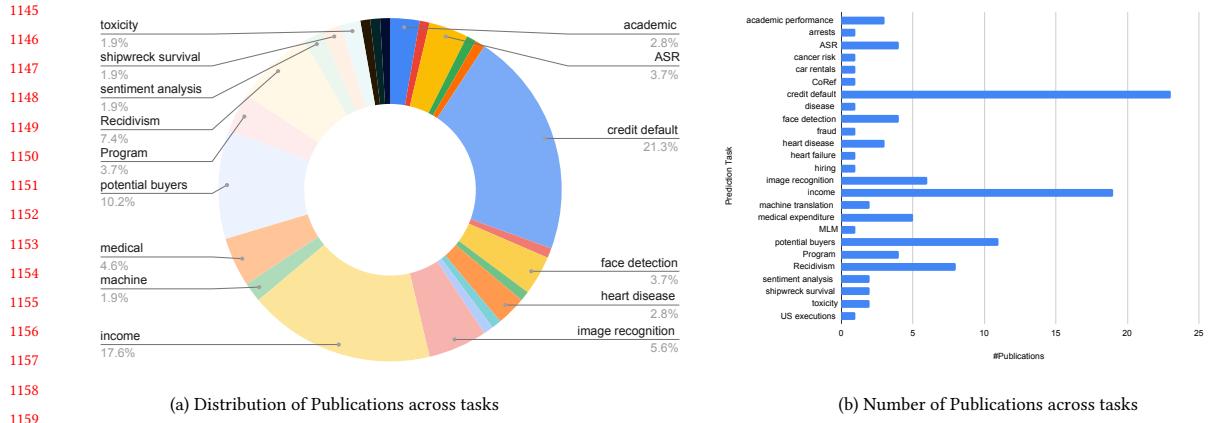


Fig. 9. Details of Publication across tasks

we observed that some task categories have more examined tasks. For instance, consider NLP task category, over six tasks have been studied in the literature, e.g., sentiment analysis, CoRef, MLM, ASR etc. This is despite very few publications for fairness analysis of NLP systems.

Our analysis of the tasks investigated in the collected papers confirm that *financial, CV and NLP tasks are well studied for software fairness analysis*. Appendix (Figure 13 and Figure 14) further shows the distribution of publications based on the number of tasks, and datasets. Similarly, Figure 9 shows that finance related tasks are the most studied, especially tasks such as credit prediction, income prediction and fraud detection. Figure 10 provides more fine grained analysis of the datasets relating to each task. Evidently, tasks involving finance, CV and NLP account for most examined datasets, about seven (7) to nine (9) percent each. For instance, tasks involving face detection, credit default, image recognition, speech recognition (ASR) and programs contribute about four (4) datasets each (see Figure 10(b)). These findings suggests that the community has been focused on investigating bias in specific sectors (CV, NLP and finance) while ignoring other categories (e.g., education and code).

The (SE) research community has been focused on studying software fairness concerns for (three) specific tasks (finance, CV and NLP tasks), but fairness concerns in areas like education and code analysis have been largely ignored.

5.4 RQ5 Tooling

In this evaluation, we inspect the papers that propose a tool, framework or library to enable software fairness analysis. Appendix (Table 5) provides details of some of the tools found in the literature. We categorize the goal of the analysis performed by each tool, the addressed problem, the main approach employed and their processing stage (pre, in and post -processing), as well as the software access required by the tool (black, grey or white box access).

Our analysis showed that *post-deployment validation of AI software (e.g.. testing, auditing, analysis and mitigation) is the most prominent tool support available for software fairness analysis*. Appendix (Table 5) shows that most tools are black-box post-validation tools. This is followed by support for fair learning (i.e., designing fair software systems). These includes measuring and analyzing the trade-off between fairness and accuracy metrics of the software, e.g., AIF360 [13] and POF [15]. Overall, *there is very little tool support for the specification, formalization and verification of*

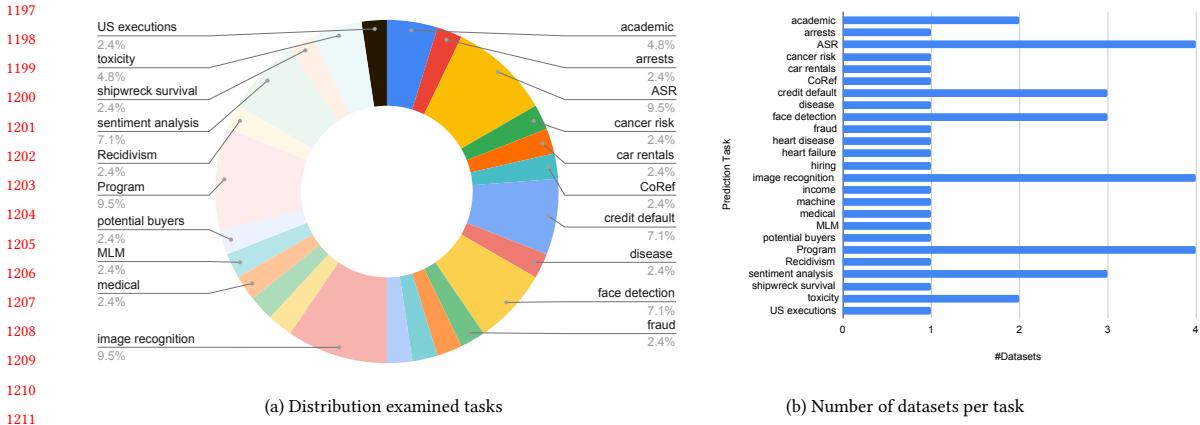


Fig. 10. Details of Datasets examined for each Task

fairness metrics. Some tools support the verification of fairness properties in the post-processing stage, i.e., verifying trained model (e.g., FairSquare [3]). Meanwhile tools like VeriFair [12] verify fairness properties in the pre-processing stage, i.e., verifying if datasets fulfill a fairness property. There is also *little support for in-processing and white-box analysis of software fairness*. While there are some works that support all processing stages (i.e., pre, in and post) of software fairness [10, 13, 74, 75, 121], we found low support for tools specifically focused on the in-processing stage.

Generally, fairness analysis tools are mostly automated, they provide a software architecture, API or framework which implements several (mitigation) algorithms that enable developers to conduct bias analysis. For instance, AIF360 [13] provides an extensible architecture, FAT-Forensics [121] offers a python framework providing several mitigation algorithms, and Themis-ML [10] provides an API for analysis. Interestingly, some approaches allow to visualise fairness metrics while auditing or analyzing fairness properties (e.g., Fairkit-learn [74, 75]), and others enable fairness test experimentation. Notably, 2AFC [89] tests for (un)fairness via psycho-physical experimentation.

Fairness Analysis tools are mostly focused on the post validation of AI-based software, with little (grey) or no access (black) to the AI model. There is a need for tools that support white-box, in-processing stages of fairness analysis, especially for specification, formalization and verification.

6 LIMITATIONS AND THREATS TO VALIDITY

Collection, Filtering and Analysis of Publications: In this work, we have focused on in-depth analysis of publications exploring fairness as a software property or conducting fairness analysis via the lens of software engineering (SE). This scope means that we may have missed or filtered out papers that study fairness in other aspects, e.g., as a legal, ethical or transparency concern. Hence, this work is limited to the analysis of fairness property as an SE concern.

Manual Publication Analysis/Interpretation: The analysis of the publications studied in this paper are potentially open to human bias since they were manually coded and analyzed. However, to mitigate this threat we conduct both internal and external validation of the work. First, we ensure that the in-depth analysis and categorization of each paper is validated by at least one other researcher. In addition, we provide a copy of the paper and data to all (cited) authors for inspection and feedback. Feedback from cited authors improved our publication search and analysis. For instance,

1249 Table 5. Open Research Problems and Opportunities (Recent works (2022-2023) are in bracket). Full table is in the appendix (Table 6).

1250 Open Problems	1251 Problem Description	1252 Potential Solutions	1253 Sample Related Work
Fairness Test Metrics and Adequacy	Measuring when fairness testing is sufficient/enough	Design of fairness test metrics and adequacy criteria	[94] ([156])
Automatic Repair of Biased Classifiers	How to automatically repair biased classifiers to be (less or) un-biased?	Automatic Program Repair for fairness property	[4] ([68])
Tooling for Fairness Property Specification	Specifying and engineering fairness properties for learning-based systems	Requirement Engineering tool support for Fairness properties	[118] ([48])
Unexplored or Poorly Understood Biases	Analyzing rare biases (e.g., age), complex or intersectional biases (e.g., age \times gender)?	Fairness Analysis Support for rare, complex or intersectional Biases	[23] ([34])
Sequential and Long-term Fairness concerns	How to analyse/maintain fairness as the AI system evolves over time?	Techniques to support analysis of sequential and long-term fairness	[153] ([58])
Human factors in fairness analysis	E.g., evaluating the harm induced by fairness violations to humans/society	Empirical studies of Human Factors in Fairness Analysis	[82] ([49])
Non-Specific/Holistic mitigation approaches	Designing bias mitigation methods that are agnostic of tasks, domains or datasets	General (i.e., task, domain and dataset -agnostic) bias analysis techniques	[152]
Fair Policy, Legalisation, and Compliance	How to design fairness analysis tools for policy makers and compliance officers?	Fairness Analysis Tool Support for Policy and Compliance Analysis	[89]

1264
 1265 authors pointed out arxiv papers that have now been published (e.g., FairKit [74]), new publications missed by our
 1266 search, due to string search (e.g., Thomas et al. [128]), and comprehensive or newer versions of cited papers (e.g., Fabris
 1267 et al. [44]). Finally, we provide both the paper and data online to support scrutiny and reuse.

1270 7 REFLECTION ON RECENT ADVANCES

1271 In this section, we reflect on the recent advances in software fairness, since our study. We discuss the advancements
 1272 made by recent works, how they relate to our previous findings and address previously open research challenges. We
 1273 reflect on our initial findings with respect to recent publications in top-tier SE-specific literature, since our analysis. To
 1274 this end, we highlight the following (13) recent (2024) publications, including papers from ICSE 2024 [34, 48, 84, 135],
 1275 FSE 2024 [138], ISSTA 2024 [38, 140], TSE (2024) [144, 155], JSS 2024 [110] and TOSEM (2024) [32, 124, 124].

1276 **Intersectional Bias Testing:** Our initial findings (RQ2) highlight the gap in fairness analysis of multiple attributes
 1277 (aka intersectional bias) and unexplored complex biases (Table 5). In recent works, researchers have begun to tackle
 1278 this problem. Notably, Chen et al. [34] presented an empirical study on fairness improvement for multiple protected
 1279 attributes which shows that improving fairness for a single protected attribute can decrease fairness for other un-
 1280 considered protected attributes by up to 88.3%. This finding emphasizes the importance of studying intersectional
 1281 fairness properties. Chen et al. [34] further conducted a large-scale benchmarking study to evaluate the effectiveness
 1282 of state-of-the-art bias mitigation methods in improving intersectional fairness and assessing the trade-off between
 1283 machine learning performance and intersectional fairness.

1284 **Fairness Requirements Engineering:** Our work emphasized the low number of support for fairness requirements
 1285 engineering (RQ1) and the lack of tooling for fairness specification (Table 5). However, some recent works have made
 1286 significant effort to improve the state-of-the-art in this area. In 2024, Ferrara et al. [48] proposed a context-aware
 1287 requirements engineering framework (called ReFair) that classifies sensitive features from user stories. ReFair employs
 1288 NLP and word embedding analysis to recommend protected features to be considered during ML implementation.

1289 **Fairness-related Empirical Analysis:** We had emphasized the low number of empirical analysis on software fairness
 1290 metrics (RQ1), in particular, human studies (Table 5). We note that, in recent work, researchers have conducted empirical
 1291 studies and human studies examining fairness concerns in bias mitigation and software engineering publication. Yang
 1292 et al. [140] presents a recent empirical study comparing the performance of state-of-the-art fairness mitigation techniques

for image classification. Using three datasets, five performance metrics and 13 mitigation methods, the authors found that pre-processing methods and in-processing methods outperform post-processing methods, with pre-processing methods exhibiting the best performance. Liang et al. [88] investigates the effect of demographic information (gender and age) on the evaluation of technical articles on software engineering and potential behavioral differences among participants. The authors conducted a survey and human study, involving 540 participants, to investigate developers' evaluation of technical articles for software engineering. Results show that participants provide more positive content depth evaluations for younger male authors when compared to older male authors, male participants evaluate faster than female participants, and there is no significant difference in the genders of authors on the evaluation outcome of technical articles in SE.

Fairness Mitigation: In line with our previous observation (**RQ5**), bias mitigation methods remain a top focus in recent SE literature [38, 85, 138, 144]. Recently (2024), Xiao et al. [138] proposed MirrorFair an approach that employs a combination of model ensembling and counterfactual analysis to mitigate unfairness. MirrorFair constructs a counterfactual dataset from the original data and trains two models, one model on the original dataset and a second model on the counterfactual dataset. It then employs an ensemble of both model predictions to generate fairer decisions. Results show that MirrorFair outperforms baselines in fairness improvement, performance preservation, and trade-off metrics. Dasu et al. [38] also presented a fairness mitigation method (NeuFair) that employ a set of randomized algorithms that uses neuron dropout as a post-processing bias mitigation method. Results show that NeuFair is efficient and effective in improving fairness (up to 69%) with minimal performance degradation. Similarly, Li et al. [85] proposed model repair techniques that employ neural condition synthesis to repair biased systems. Yu et al. [144] proposed FairBalance, a pre-processing bias mitigation algorithm which balances the class distribution in each demographic group by assigning calculated weights to the training data. Fairbalance aims to identify the root cause of equalized odd violations and mitigate it without modifying the normal training process. Specifically, the authors observed that equalizing the class distribution in each demographic group with sample weights is a necessary condition for achieving equalized odds. Results show that FairBalance outperforms state-of-the-art approaches in terms of equalized odds.

Tabular and Semi-structured Datasets: Li et al. [84] addresses one of the open problems highlighted in our work (**Table 5**) – dataset-agnostic fairness analysis that cater for both tabular and semi-structured datasets. To achieve this, the authors propose a white-box fairness analysis approach called Responsible UNfair NEuron Repair (RUNNER). RUNNER improves the efficiency, effectiveness and generalization of existing works via Importance-based Neuron Diagnosis and Neuron Stabilizing Retraining. More importantly, RUNNER generalizes to both structured tabular data and large-scale unstructured image data, a combination which we observed as uncommon in previous works (**RQ4** and **Table 5**).

Fairness Testing: In our analysis, we emphasized previous works on fairness testing (**RQ1** and appendix (**Table 5**)) and fairness metrics (e.g. fairness-performance trade-off (**RQ2**)). There are some more recent works in both areas. In particular, Wang et al. [135] present a novel black-box individual fairness testing method called Model-Agnostic Fairness Testing (MAFT) which employs lightweight methods like gradient estimation and attribute perturbation for fairness testing. The authors show that MAFT achieves the same effectiveness as state-of-the-art white-box methods whilst improving the applicability to large-scale networks. Researchers have also proposed DistroFair to test class-level fairness in the presence of distributional shifts (i.e., out-of-distribution datasets) [110]. The authors concretize their approach for image datasets and demonstrate that it exposes class-level fairness in image datasets. Sun et al. [124] recently proposed FairMT an automated fairness testing approach for machine translation systems. FairMT aims to ensure that translations of semantically similar sentences containing protected attributes from distinct demographic

1353 groups maintain comparable meanings. The authors found that fair translations tend to exhibit superior translation
1354 performance, which challenges the existing literature on the existence of a tradeoff between fairness and performance.
1355

1356 **Fairness Test Adequacy:** We highlight the need for fairness test adequacy criteria (Table 5). A recent work studied
1357 this concern: Zheng et al. [155] conducted an empirical analysis to investigate the correlation between fairness property
1358 and test coverage criteria. Their study employed seven (7) coverage criteria, six (6) fairness metrics, three (3) testing
1359 techniques, five (5) bias mitigation methods, five (5) DNN models and nine (9) fairness datasets. Similar to our work, the
1360 authors found several open challenges in this area: They observed that (a) there is limited correlation between coverage
1361 criteria and fairness, (b) coverage and fairness metrics change as the test suite increases, and (c) models debiased by
1362 bias mitigation methods have a lower correlation, between coverage and fairness, compared to the original models.
1363
1364

1365 **Recent Closely-related Survey:** In section 2, we compare our work to fairness-related surveys. However, more
1366 recently (2024), Chen et al. [32] presented a closely-related survey on fairness testing. This work is related to our
1367 work since fairness testing is one of the software development steps, and testing is inter-related to other software
1368 development steps, e.g., exposing fairness bugs is a check that fairness requirements are met. Unlike our work, Chen
1369 et al. [32] examine the literature with a focus on fairness testing workflow (how to test) and testing components (what
1370 to test). Similar to our work, Chen et al. [32] investigates the fairness-related literature, datasets, open source tools,
1371 research distributions, and research opportunities, but as it relates to fairness testing, instead of the entire SE pipeline
1372 (as done in this work). Notably, their work [32] presents several results that corroborate our findings. For instance,
1373 they also highlight the absence of works on intersectional bias (multiple sensitive attributes), the lack of test adequacy
1374 criteria for fairness testing and the lack of socio-technical and human-in-the-loop (stakeholder) support for fairness
1375 testing. Some of their findings are also complementary to ours. As an example, similar to our finding on unexplored
1376 or poorly understood biases (Table 5), Chen et al. [32] also highlight the difficulty of testing rare attributes and the
1377 challenge of fairness testing when sensitive attributes are absent (e.g., due to GDPR regulation). Unlike Chen et al. [32],
1378 we study software fairness beyond testing, we focus on its impact on the entire software development pipeline.
1379
1380

1381 8 FUTURE OUTLOOK

1382 This paper presents a comprehensive analysis and survey of publications on software fairness. In particular, we study
1383 publications that study *fairness as a software property*, works that examine fairness with a software engineering lens, or
1384 study how to engineer fair learning-based software systems. We identified several open problems and gaps. Table 5
1385 highlights the details of the open problems identified in this work. Specifically, we highlight open problems in the
1386 engineering of fair software systems, including concerns in the areas of fairness testing, verification and empirical
1387 evaluation of fairness properties. In the following we discuss the gaps and open problems we elicited from our analysis.
1388

1389 Even though there are several papers exploring fairness testing concerns [2, 131, 149] as well as works investigating
1390 test metrics for learning-based systems [69, 93, 123, 139], there is still the *problem of measuring test adequacy for*
1391 *fairness testing*. In our study, we found a few recent works [98, 156] *exploring the test adequacy of fairness testing*
1392 *approaches*: Zheng et al. [156] empirically studied the relationship between DNN fairness and neuron coverage test
1393 adequacy criteria. Their experiments showed that there is a limited statistical correlation between neuron coverage
1394 criteria and DNN fairness, and these coverage criteria may be invalid for DNN fairness. Similarly, Majumder et al. [98]
1395 explored how to satisfy diverse fairness metrics and proposed a reduced set of fairness metrics. In particular, they
1396 cluster 30 metrics into seven clusters of classification metrics and three clusters of dataset metrics.
1397
1398

1405 However, several questions concerning the adequacy of fairness testing remain unanswered and unexplored: When
 1406 is fairness testing (in)sufficient? How can fairness testing be guided to reduce redundant test cases? Beyond fairness
 1407 violations, are there other test metrics that indicate (in)sufficient testing has been performed? These are open problems
 1408 in the area of fairness testing.
 1409

1410 Despite the advances in test metrics for traditional and learning-based software, determining the appropriate test
 1411 metrics or adequacy criteria for fairness evaluation remains an open problem. Researchers have proposed several test
 1412 adequacy criteria for learning-based software, such as *neuron coverage* [108] and *surprise adequacy* [79]. *Neuron coverage*
 1413 (from DeepXplore [108]) measures the parts of a learning-based systems that is exercised by a test input, it employs
 1414 the ratio of neurons whose activation values were above a predefined threshold to measure the diversity of neuron
 1415 behaviour and guide test generation. Likewise, *surprise adequacy* evaluates the behaviour of learning-based system
 1416 with respect to their training data by measuring the surprise of an input as the difference in the system's behaviour
 1417 between the input and the training data, such that a good test input should be sufficiently but not overtly surprising
 1418 compared to the training data. Other test metrics for learning-based systems include DeepGini [47], Multiple-Boundary
 1419 Clustering and Prioritization (MCP) [119] and Maximum Probability (MaxP) [96]. Despite the availability of several test
 1420 metrics and adequacy criteria for functional testing of learning-based systems, there is no *indication or evaluation that*
 1421 *demonstrates these test metrics are applicable for the fairness testing of learning-based systems*. The test metric typically
 1422 employed by the fairness testing literature remains unfair behavior/outputs characterizing fairness violations. We
 1423 encourage researchers to further investigate this vital challenge of fairness testing.
 1424

1425 Although there are several fairness mitigation approaches, *repairing unfair learning-based systems to fulfill software*
 1426 *fairness properties (i.e., to be fair or unbiased) has been hardly explored*, except for few works (e.g., Albarghouthi et al.
 1427 [4], Li et al. [85], Tao et al. [127]). Albarghouthi et al. [4] proposed an approach that applies *distribution-guided inductive*
 1428 *synthesis* to repair *unfair* ML classifiers with the goal of making them fair, it also verifies that the repaired classifiers is
 1429 semantically close to the original (unfair) classifier. Their approach formulates the problem as a probability distribution
 1430 problem, such that the repaired classifier needs to satisfy a probabilistic Boolean expression. Most importantly, this
 1431 work is one of the few approaches we have found that aims to directly repair ML classifiers to fulfill fairness properties,
 1432 without model re-training. Additionally, Tao et al. [127] have recently proposed model repair techniques that employ
 1433 adversarial training to repair biased systems. Similarly, there is *low support for verification of fairness properties* (see
 1434 RQ5). We had observed that verifying, certifying and providing guarantees for software fairness in learning-based
 1435 software is under-explored.
 1436

1437 Additional open problems include the *low number of empirical evaluation of fairness properties, especially as they relate*
 1438 *to human factors and societal policies*. There are limited empirical studies studying or measuring the harm caused by
 1439 unfair software behavior (to humans) and the impact of fairness mitigation on marginalized individuals and communities.
 1440 Blodgett et al. [20] emphasizes the importance of evaluating the harms induced by unfair (NLP) systems on humans.
 1441

1442 Furthermore, there is *a need to develop approaches and tool support that are task, and dataset agnostic*. For instance,
 1443 we observed that there is a lack of general, non-specific fairness mitigation approaches, most works target a specific
 1444 domain/task (e.g., NLP, CV) with little work that is generally or demonstrably applicable across domains, tasks or
 1445 datasets (except for few works like ADF [152]).
 1446

1447 Besides, *there is a need to provide tools to support the automatic specification and requirements engineering of fairness*
 1448 *properties*. Closely related works in this area includes MLCheck [118] and VeriFair [12] which allows to specify and
 1449 check fairness specifications. However, these tools do not provide support for non-technical experts, e.g., compliance
 1450 officers and policy makers. Indeed, such tools should support the definition and formalization of specific fairness
 1451 Manuscript submitted to ACM
 1452
 1453
 1454
 1455
 1456

1457 properties as a socio-technical issue, not just a technical issue. These tools should allow not only allow to test fairness,
1458 but enable compliance officers and policy makers to audit and derive bias-preserving policies: For instance, equity-based
1459 policies have been proven to be useful in other domains, e.g., education (e.g., affirmative action) [107], toxic language
1460 detection [62] and transport systems [114].
1461

1462 In RQ3, we demonstrate that there are some *poorly explored and understood biases, especially in terms of protected*
1463 *attributes*. Generally, we observed that certain tasks, datasets and biases are poorly studied. For instance, protected
1464 attributes relating to marital status, sexuality, non-binary gender, and education are poorly explored. Relating to this
1465 issue is the fact that the interactions of protected attributes is under-explored, especially in terms of the compounding
1466 or intersectional effect of multiple protected attributes, e.g., biases triggered by a combination of attributes (e.g., age
1467 × gender). While works like FairVis [26] allow to visualize and analyze intersectional bias, there is little support for
1468 specifying, testing and mitigating such biases.
1469

1470 Sequential and long-term fairness concerns also remain an open problem [153]. For instance, how do we mitigate
1471 fairness properties as the software system evolves? Do the proposed mitigation approaches for one-time fairness
1472 analysis scale to sequential or long-term concerns.
1473

1474 In addition, there is the socio-technical and ethical concern of fairness analysis: How do we ensure that our fairness
1475 mitigation approaches do not induce new and unintended biases? How do our proposed mitigation and analysis
1476 approaches translate to real-world intervention for fairness, e.g., in terms of equity-based mitigation approaches
1477 typically employed in public policy (e.g., affirmative action [107])? In summary, we posit that there is a need for
1478 socio-technical, human-in-the-loop bias analysis approaches that translate to the mitigation of real-world harms to
1479 humans and society [52]. There is a need to provide bias analysis methods to support developers and policy makers.
1480

1481 9 CONCLUSION

1482 This paper presents a survey of the literature on software fairness analysis. It is particularly focused on examining the
1483 landscape of the research works that explore fairness with the lens of software engineering. We initially collected 164
1484 papers in our initial analysis (2010-2021) and examine several aspects of the literature including the goals of the available
1485 work. This includes the focus of the literature in terms of the domains, measures, attributes and datasets that are
1486 explored. Additionally, we examine the under-examined or unexplored areas and discuss open research opportunities
1487 in software fairness. We also collected an additional 41 papers (2022-2023), from the dominant (SE) venues. We reflect
1488 on the most recent advancements in software fairness analysis since our initial study. We provide further details of
1489 collected papers and our analysis in our artefact: <https://github.com/ezekiel-soremekun/Software-Fairness-Analysis>
1490

1491 1492 ACKNOWLEDGMENTS

1493 We would like to thank the reviewers for their valuable feedback, which has improved our analysis and paper. In
1494 addition, we would like to thank all (cited) authors who provided feedback on previous drafts of this paper.
1495

1496 1497 REFERENCES

- 1500 [1] Abubakar Abid, Maheen Farooqi, and James Zou. 2021. Persistent anti-muslim bias in large language models. In *Proceedings of the 2021 AAAI/ACM*
1501 *Conference on AI, Ethics, and Society*. 298–306.
1502 [2] Aniya Aggarwal, Pranay Lohia, Seema Nagar, Kuntal Dey, and Diptikalyan Saha. 2019. Black box fairness testing of machine learning models.
1503 In *Proceedings of the 2019 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software*
1504 *Engineering*. 625–635.

- [3] Aws Albarghouthi, Loris D'Antoni, Samuel Drews, and Aditya V Nori. 2017. Fairsquare: probabilistic verification of program fairness. *Proceedings of the ACM on Programming Languages* 1, OOPSLA (2017), 1–30.
- [4] Aws Albarghouthi, Loris D'Antoni, and Samuel Drews. 2017. Repairing decision-making programs under uncertainty. In *International Conference on Computer Aided Verification*. Springer, 181–200.
- [5] Rico Angell, Brittany Johnson, Yuriy Brun, and Alexandra Meliou. 2018. Themis: Automatically testing software for discrimination. In *Proceedings of the 2018 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 871–875.
- [6] Julia Angwin, Jeff Larson, Surya Mattu, and Lauren Kirchner. 2022. Machine bias. In *Ethics of data and analytics*. Auerbach Publications, 254–264.
- [7] Carolyn Ashurst and Adrian Weller. 2023. Fairness Without Demographic Data: A Survey of Approaches. In *Proceedings of the 3rd ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization*. 1–12.
- [8] Carolyn Ashurst and Adrian Weller. 2023. Fairness Without Demographic Data: A Survey of Approaches. In *Proceedings of the 3rd ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization (EAAMO '23)*. Association for Computing Machinery, New York, NY, USA, Article 14, 12 pages. <https://doi.org/10.1145/3617694.3623234>
- [9] Muhammad Hilmi Asyrofi, Zhou Yang, Imam Nur Bani Yusuf, Hong Jin Kang, Ferdinand Thung, and David Lo. 2021. Biasfinder: Metamorphic test generation to uncover bias for sentiment analysis systems. *IEEE Transactions on Software Engineering* (2021).
- [10] Niels Bantilan. 2018. Themis-ml: A fairness-aware machine learning interface for end-to-end discrimination discovery and mitigation. *Journal of Technology in Human Services* 36, 1 (2018), 15–30.
- [11] Luciano Baresi, Chiara Criscuolo, and Carlo Ghezzi. 2023. Understanding fairness requirements for ml-based software. In *2023 IEEE 31st International Requirements Engineering Conference (RE)*. IEEE, 341–346.
- [12] Osbert Bastani, Xin Zhang, and Armando Solar-Lezama. 2019. Probabilistic verification of fairness properties via concentration. *Proceedings of the ACM on Programming Languages* 3, OOPSLA (2019), 1–27.
- [13] Rachel KE Bellamy, Kuntal Dey, Michael Hind, Samuel C Hoffman, Stephanie Houde, Kalapriya Kannan, Pranay Lohia, Jacquelyn Martino, Sameep Mehta, Aleksandra Mojsilović, et al. 2019. AI Fairness 360: An extensible toolkit for detecting and mitigating algorithmic bias. *IBM Journal of Research and Development* 63, 4/5 (2019), 4–1.
- [14] Nelly Bencomo, Jin LC Guo, Rachel Harrison, Hans-Martin Heyn, and Tim Menzies. 2021. The Secret to Better AI and Better Software (Is Requirements Engineering). *IEEE Software* 39, 1 (2021), 105–110.
- [15] Richard Berk, Hoda Heidari, Shahin Jabbari, Matthew Joseph, Michael Kearns, Jamie Morgenstern, Seth Neel, and Aaron Roth. 2017. A Convex Framework for Fair Regression. *Fairness, Accountability, and Transparency in Machine Learning* (2017).
- [16] Sumon Biswas and Hridesh Rajan. 2020. Do the machine learning models on a crowd sourced platform exhibit bias? an empirical study on model fairness. In *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 642–653.
- [17] Sumon Biswas and Hridesh Rajan. 2021. Fair Preprocessing: Towards Understanding Compositional Fairness of Data Transformers in Machine Learning Pipeline. *arXiv preprint arXiv:2106.06054* (2021).
- [18] Sumon Biswas and Hridesh Rajan. 2023. Fairify: Fairness verification of neural networks. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*. IEEE, 1546–1558.
- [19] Emily Black, Samuel Yeom, and Matt Fredrikson. 2020. Fliptest: fairness testing via optimal transport. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*. 111–121.
- [20] Su Lin Blodgett, Solon Barocas, Hal Daumé III, and Hanna Wallach. 2020. Language (Technology) is Power: A Critical Survey of “Bias” in NLP. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 5454–5476.
- [21] Su Lin Blodgett, Gilsinia Lopez, Alexandra Olteanu, Robert Sim, and Hanna Wallach. 2021. Stereotyping Norwegian salmon: an inventory of pitfalls in fairness benchmark datasets. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. 1004–1015.
- [22] C. Malik Boykin, Sophia T. Dasch, Vincent Rice Jr., Venkat R. Lakshminarayanan, Taiwo A. Togun, and Sarah M. Brown. 2021. Opportunities for a More Interdisciplinary Approach to Measuring Perceptions of Fairness in Machine Learning. In *Proceedings of the 1st ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization* (–, NY, USA) (EAAMO '21). Association for Computing Machinery, New York, NY, USA, Article 1, 9 pages. <https://doi.org/10.1145/3465416.3483302>
- [23] Martin Brandao. 2019. Age and gender bias in pedestrian detection algorithms. In *Proceedings of the Workshop on Fairness Accountability Transparency and Ethics in Computer Vision at IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*.
- [24] Yuriy Brun and Alexandra Meliou. 2018. Software fairness. In *Proceedings of the 2018 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 754–759.
- [25] Margaret Burnett, Simone Stumpf, Jamie Macbeth, Stephan Makri, Laura Beckwith, Irwin Kwan, Anicia Peters, and William Jernigan. 2016. GenderMag: A method for evaluating software’s gender inclusiveness. *Interacting with Computers* 28, 6 (2016), 760–787.
- [26] Ángel Alexander Cabrera, Will Epperson, Fred Hohman, Minsuk Kahng, Jamie Morgenstern, and Duen Horng Chau. 2019. FairVis: Visual analytics for discovering intersectional bias in machine learning. In *2019 IEEE Conference on Visual Analytics Science and Technology (VAST)*. IEEE, 46–56.
- [27] Toon Calders and Sicco Verwer. 2010. Three naive Bayes approaches for discrimination-free classification. *Data mining and knowledge discovery* 21, 2 (2010), 277–292.
- [28] Simon Caton and Christian Haas. 2020. Fairness in machine learning: A survey. *arXiv preprint arXiv:2010.04053* (2020).

- [29] L Elisa Celis, Lingxiao Huang, Vijay Keswani, and Nisheeth K Vishnoi. 2019. Classification with fairness constraints: A meta-algorithm with provable guarantees. In *Proceedings of the conference on fairness, accountability, and transparency*. 319–328.
- [30] Joymallya Chakraborty, Suvodeep Majumder, and Tim Menzies. 2021. Bias in machine learning software: why? how? what to do?. In *Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 429–440.
- [31] Joymallya Chakraborty, Suvodeep Majumder, Zhe Yu, and Tim Menzies. 2020. Fairway: A way to build fair ml software. In *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 654–665.
- [32] Zhenpeng Chen, Jie M Zhang, Max Hort, Mark Harman, and Federica Sarro. 2024. Fairness testing: A comprehensive survey and analysis of trends. *ACM Transactions on Software Engineering and Methodology* 33, 5 (2024), 1–59.
- [33] Zhenpeng Chen, Jie M Zhang, Federica Sarro, and Mark Harman. 2022. MAAT: a novel ensemble approach to addressing fairness and performance bugs for machine learning software. In *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 1122–1134.
- [34] Zhenpeng Chen, Jie M Zhang, Federica Sarro, and Mark Harman. 2024. Fairness improvement with multiple protected attributes: How far are we? (2024), 1–13.
- [35] Lu Cheng, Kush R Varshney, and Huan Liu. 2021. Socially responsible ai algorithms: Issues, purposes, and challenges. *Journal of Artificial Intelligence Research* 71 (2021), 1137–1181.
- [36] Jürgen Cito, Isil Dillig, Seohyun Kim, Vijayaraghavan Murali, and Satish Chandra. 2021. Explaining mispredictions of machine learning models using rule induction. In *Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 716–727.
- [37] Kate Crawford. 2017. The trouble with bias. In *Conference on Neural Information Processing Systems, Invited Speaker*. https://www.youtube.com/watch?v=fMym_BKWQzk
- [38] Vishnu Asutosh Dasu, Ashish Kumar, Saeid Tizpaz-Niari, and Gang Tan. 2024. NeuFair: Neural Network Fairness Repair with Dropout. In *Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis*. 1541–1553.
- [39] Isin Demirsahin, Oddur Kjartansson, Alexander Gutkin, and Clara Rivera. 2020. Open-source multi-speaker corpora of the english accents in the british isles. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*. 6532–6541.
- [40] Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. 2012. Fairness through awareness. In *Proceedings of the 3rd innovations in theoretical computer science conference*. 214–226.
- [41] Michael D Ekstrand, Robin Burke, and Fernando Diaz. 2019. Fairness and discrimination in recommendation and retrieval. In *Proceedings of the 13th ACM Conference on Recommender Systems*. 576–577.
- [42] Simone Fabbrizzi, Symeon Papadopoulos, Eirini Ntoutsi, and Ioannis Kompatsiaris. 2021. A Survey on Bias in Visual Datasets. *arXiv preprint arXiv:2107.07919* (2021).
- [43] Alessandro Fabris, Fabio Giachelle, Alberto Piva, Gianmaria Silvello, and Gian Antonio Susto. 2023. A Search Engine for Algorithmic Fairness Datasets. In *Proceedings of the 2nd European Workshop on Algorithmic Fairness, Winterthur, Switzerland, June 7th to 9th, 2023 (CEUR Workshop Proceedings, Vol. 3442)*, Jose M. Alvarez, Alessandro Fabris, Christoph Heitz, Corinna Hertweck, Michele Loi, and Meike Zehlike (Eds.). CEUR-WS.org. <https://ceur-ws.org/Vol-3442/paper-08.pdf>
- [44] Alessandro Fabris, Stefano Messina, Gianmaria Silvello, and Gian Antonio Susto. 2022. Algorithmic fairness datasets: the story so far. *Data Min. Knowl. Discov.* 36, 6 (2022), 2074–2152. <https://doi.org/10.1007/S10618-022-00854-Z>
- [45] Alessandro Fabris, Stefano Messina, Gianmaria Silvello, and Gian Antonio Susto. 2022. Tackling Documentation Debt: A Survey on Algorithmic Fairness Datasets. In *Proceedings of the 2nd ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization (EAAMO '22)*. Association for Computing Machinery, New York, NY, USA, Article 2, 13 pages. <https://doi.org/10.1145/3551624.3555286>
- [46] Michael Feldman, Sorelle A Friedler, John Moeller, Carlos Scheidegger, and Suresh Venkatasubramanian. 2015. Certifying and removing disparate impact. In *proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*. 259–268.
- [47] Yang Feng, Qingkai Shi, Xinyu Gao, Jun Wan, Chunrong Fang, and Zhenyu Chen. 2020. DeepGini: Prioritizing Massive Tests to Enhance the Robustness of Deep Neural Networks. In *Proceedings of the 29th ACM SIGSOFT International Symposium on Software Testing and Analysis (Virtual Event, USA) (ISSTA 2020)*. Association for Computing Machinery, New York, NY, USA, 177–188. <https://doi.org/10.1145/3395363.3397357>
- [48] Carmine Ferrara, Francesco Casillo, Carmine Gravino, Andrea De Lucia, and Fabio Palomba. 2024. Refair: Toward a context-aware recommender for fairness requirements engineering. (2024), 1–12.
- [49] Carmine Ferrara, Giulia Sellitto, Filomena Ferrucci, Fabio Palomba, and Andrea De Lucia. 2023. Fairness-aware machine learning engineering: how far are we? *Empirical Software Engineering* 29, 1 (2023), 9.
- [50] Anjalie Field, Su Lin Blodgett, Zeerak Waseem, and Yulia Tsvetkov. 2021. A Survey of Race, Racism, and Anti-Racism in NLP. *arXiv preprint arXiv:2106.11410* (2021).
- [51] Anthony Finkelstein, Mark Harman, S Afshin Mansouri, Jian Ren, and Yuanyuan Zhang. 2009. A search based approach to fairness analysis in requirement assignments to aid negotiation, mediation and decision making. *Requirements engineering* 14, 4 (2009), 231–245.
- [52] Claudia Flores-Saviaga, Christopher Curtis, and Saiph Savage. 2023. Inclusive Portraits: Race-Aware Human-in-the-Loop Technology. In *Proceedings of the 3rd ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization (EAAMO '23)*. Association for Computing Machinery, New York, NY, USA, Article 15, 11 pages. <https://doi.org/10.1145/3617694.3623235>
- [53] Batya Friedman and Helen Nissenbaum. 1996. Bias in computer systems. *ACM Transactions on Information Systems (TOIS)* 14, 3 (1996), 330–347.

- [54] Pratik Gajane and Mykola Pechenizkiy. 2017. On formalizing fairness in prediction with machine learning. *arXiv preprint arXiv:1710.03184* (2017).
- [55] Sainyam Galhotra, Yuriy Brun, and Alexandra Meliou. 2017. Fairness testing: testing software for discrimination. In *Proceedings of the 2017 11th Joint Meeting on Foundations of Software Engineering*. 498–510.
- [56] Xuanqi Gao, Juan Zhai, Shiqing Ma, Chao Shen, Yufei Chen, and Qian Wang. 2022. FairNeuron: Improving Deep Neural Network Fairness with Adversary Games on Selective Neurons. In *2022 IEEE/ACM 44th International Conference on Software Engineering (ICSE)*. IEEE.
- [57] Bishwamitra Ghosh, Debabrata Basu, and Kuldeep S Meel. 2021. Justicia: A Stochastic SAT Approach to Formally Verify Fairness. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 7554–7563.
- [58] Usman Gohar, Sumon Biswas, and Hridesh Rajan. 2023. Towards understanding fairness and its composition in ensemble machine learning. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*. IEEE, 1533–1545.
- [59] Usman Gohar and Lu Cheng. 2023. A Survey on Intersectional Fairness in Machine Learning: Notions, Mitigation, and Challenges. *arXiv preprint arXiv:2305.06969* (2023).
- [60] Nima Grgic-Hlaca, Elissa M Redmiles, Krishna P Gummadi, and Adrian Weller. 2018. Human perceptions of fairness in algorithmic decision making: A case study of criminal risk prediction. In *Proceedings of the 2018 world wide web conference*. 903–912.
- [61] Emitzá Guzmán, Ricardo Anna-Lena Fischer, and Janev Kok. 2023. Mind the gap: gender, micro-inequities and barriers in software development. *Empirical Software Engineering* 29, 1 (2023), 17.
- [62] Matan Halevy, Camille Harris, Amy Bruckman, Diyi Yang, and Ayanna Howard. 2021. Mitigating Racial Biases in Toxic Language Detection with an Equity-Based Ensemble Framework. In *Proceedings of the 1st ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization* (–, NY, USA) (EAAMO ’21). Association for Computing Machinery, New York, NY, USA, Article 7, 11 pages. <https://doi.org/10.1145/3465416.3483299>
- [63] Austin Hoag, James E. Kostas, Bruno Castro da Silva, Philip S. Thomas, and Yuriy Brun. 2023. Seldonian Toolkit: Building Software with Safe and Fair Machine Learning. In *2023 IEEE/ACM 45th International Conference on Software Engineering: Companion Proceedings (ICSE-Companion)*. 107–111. <https://doi.org/10.1109/ICSE-Companion58688.2023.00035>
- [64] Kenneth Holstein, Jennifer Wortman Vaughan, Hal Daumé III, Miro Dudík, and Hanna Wallach. 2019. Improving fairness in machine learning systems: What do industry practitioners need?. In *Proceedings of the 2019 CHI conference on human factors in computing systems*. 1–16.
- [65] Sara Hooker, Nyalleng Moorosi, Gregory Clark, Samy Bengio, and Emily Denton. 2020. Characterising bias in compressed models. *arXiv preprint arXiv:2010.03058* (2020).
- [66] Max Hort, Rebecca Moussa, and Federica Sarro. 2023. Multi-objective search for gender-fair and semantically correct word embeddings. *Applied Soft Computing* 133 (2023), 109916.
- [67] Max Hort, Jie M Zhang, Federica Sarro, and Mark Harman. 2021. Fairea: a model behaviour mutation approach to benchmarking bias mitigation methods. In *Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 994–1006.
- [68] Max Hort, Jie M. Zhang, Federica Sarro, and Mark Harman. 2024. Search-based Automatic Repair for Fairness and Accuracy in Decision-making Software. *Empirical Software Engineering* 29, 1 (2024), 36. <https://doi.org/10.1007/s10664-023-10419-3>
- [69] Chao Huang, Junbo Zhang, Yu Zheng, and Nitesh V. Chawla. 2018. DeepCrime: Attentive Hierarchical Recurrent Networks for Crime Prediction. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management* (Torino, Italy) (CIKM ’18). Association for Computing Machinery, New York, NY, USA, 1423–1432. <https://doi.org/10.1145/3269206.3271793>
- [70] Dong Huang, Qingwen Bu, Jie Zhang, Xiaofei Xie, Junjie Chen, and Heming Cui. 2023. Bias assessment and mitigation in llm-based code generation. *arXiv preprint arXiv:2309.14345* (2023).
- [71] Ben Hutchinson and Margaret Mitchell. 2019. 50 years of test (un) fairness: Lessons for machine learning. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*. 49–58.
- [72] Wiebke (Toussaint) Hutiri, Aaron Yi Ding, Fahim Kawsar, and Akhil Mathur. 2023. Tiny, Always-on, and Fragile: Bias Propagation through Design Choices in On-Device Machine Learning Workflows. *ACM Trans. Softw. Eng. Methodol.* 32, 6, Article 155 (sep 2023), 37 pages. <https://doi.org/10.1145/3591867>
- [73] Philips George John, Deepak Vijaykeerthy, and Diptikalyan Saha. 2020. Verifying individual fairness in machine learning models. In *Conference on Uncertainty in Artificial Intelligence*. PMLR, 749–758.
- [74] Brittany Johnson, Jesse Bartola, Rico Angell, Sam Witty, Stephen Giguere, and Yuriy Brun. 2023. Fairkit, fairkit, on the wall, who's the fairest of them all? Supporting fairness-related decision-making. *EURO Journal on Decision Processes* 11 (2023), 100031. <https://doi.org/10.1016/j.ejdp.2023.100031>
- [75] Brittany Johnson and Yuriy Brun. 2022. Fairkit-Learn: A Fairness Evaluation and Comparison Toolkit. In *Proceedings of the ACM/IEEE 44th International Conference on Software Engineering: Companion Proceedings* (Pittsburgh, Pennsylvania) (ICSE ’22). Association for Computing Machinery, New York, NY, USA, 70–74. <https://doi.org/10.1145/3510454.3516830>
- [76] Faisal Kamiran and Toon Calders. 2012. Data preprocessing techniques for classification without discrimination. *Knowledge and information systems* 33, 1 (2012), 1–33.
- [77] Faisal Kamiran, Toon Calders, and Mykola Pechenizkiy. 2010. Discrimination aware decision tree learning. In *2010 IEEE International Conference on Data Mining*. IEEE, 869–874.
- [78] Faisal Kamiran, Asim Karim, and Xiangliang Zhang. 2012. Decision theory for discrimination-aware classification. In *2012 IEEE 12th International Conference on Data Mining*. IEEE, 924–929.

- 1665 [79] Jinhan Kim, Robert Feldt, and Shin Yoo. 2019. Guiding deep learning system testing using surprise adequacy. In *2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE)*. IEEE, 1039–1049.
- 1666 [80] Barbara Kitchenham. 2004. Procedures for performing systematic reviews. *Keele, UK, Keele University* 33, 2004 (2004), 1–26.
- 1668 [81] Marta Z Kwiatkowska. 1989. Survey of fairness notions. *Information and Software Technology* 31, 7 (1989), 371–386.
- 1669 [82] Michelle Seng Ah Lee and Jat Singh. 2021. The landscape and gaps in open source fairness toolkits. In *Proceedings of the 2021 CHI conference on human factors in computing systems*. 1–13.
- 1670 [83] Grace A Lewis, Ipek Ozkaya, and Xiwei Xu. 2021. Software architecture challenges for ml systems. In *2021 IEEE International Conference on Software Maintenance and Evolution (ICSME)*. IEEE, 634–638.
- 1672 [84] Tianlin Li, Yue Cao, Jian Zhang, Shiqian Zhao, Yihao Huang, Aishan Liu, Qing Guo, and Yang Liu. 2024. Runner: Responsible unfair neuron repair for enhancing deep neural network fairness. In *Proceedings of the 46th IEEE/ACM International Conference on Software Engineering*. 1–13.
- 1674 [85] Tianlin Li, Xiaofei Xie, Jian Wang, Qing Guo, Aishan Liu, Lei Ma, and Yang Liu. 2023. Faire: Repairing Fairness of Neural Networks via Neuron Condition Synthesis. *ACM Trans. Softw. Eng. Methodol.* 33, 1, Article 21 (nov 2023), 24 pages. <https://doi.org/10.1145/3617168>
- 1676 [86] Xinyue Li, Zhenpeng Chen, Jie M Zhang, Federica Sarro, Ying Zhang, and Xuanzhe Liu. 2023. Dark-skin individuals are at more risk on the street: Unmasking fairness issues of autonomous driving systems. *arXiv preprint arXiv:2308.02935* (2023).
- 1678 [87] Yanhui Li, Linghan Meng, Lin Chen, Li Yu, Di Wu, Yuming Zhou, and Baowen Xu. 2022. Training Data Debugging for the Fairness of Machine Learning Software. In *2022 IEEE/ACM 44th International Conference on Software Engineering (ICSE)*. IEEE.
- 1680 [88] Anda Liang, Emerson Murphy-Hill, Westley Weimer, and Yu Huang. 2024. A Controlled Experiment in Age and Gender Bias When Reading Technical Articles in Software Engineering. *IEEE Trans. Softw. Eng.* 50, 10 (Oct. 2024), 2498–2511. <https://doi.org/10.1109/TSE.2024.3437355>
- 1682 [89] Lizhen Liang and Daniel E Acuna. 2020. Artificial mental phenomena: Psychophysics as a framework to detect perception biases in AI models. In *Proceedings of the 2020 conference on fairness, accountability, and transparency*. 403–412.
- 1684 [90] Zachary Lipton, Julian McAuley, and Alexandra Chouldechova. 2018. Does mitigating ML’s impact disparity require treatment disparity? *Advances in neural information processing systems* 31 (2018).
- 1686 [91] Ye Liu, Yi Li, Shang-Wei Lin, and Rong Zhao. 2020. Towards automated verification of smart contract fairness. In *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 666–677.
- 1688 [92] Steven R Livingstone and Frank A Russo. 2018. The Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS): A dynamic, multimodal set of facial and vocal expressions in North American English. *PloS one* 13, 5 (2018), e0196391.
- 1690 [93] Lei Ma, Felix Juefei-Xu, Fuyuan Zhang, Jiyuan Sun, Minhui Xue, Bo Li, Chunyang Chen, Ting Su, Li Li, Yang Liu, Jianjun Zhao, and Yadong Wang. 2018. DeepGauge: Multi-Granularity Testing Criteria for Deep Learning Systems. In *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering* (Montpellier, France) (ASE 2018). Association for Computing Machinery, New York, NY, USA, 120–131. <https://doi.org/10.1145/3238147.3238202>
- 1692 [94] Lei Ma, Fuyuan Zhang, Jiyuan Sun, Minhui Xue, Bo Li, Felix Juefei-Xu, Chao Xie, Li Li, Yang Liu, Jianjun Zhao, et al. 2018. Deepmutation: Mutation testing of deep learning systems. In *2018 IEEE 29th International Symposium on Software Reliability Engineering (ISSRE)*. IEEE, 100–111.
- 1694 [95] Pingchuan Ma, Shuai Wang, and Jin Liu. 2020. Metamorphic Testing and Certified Mitigation of Fairness Violations in NLP Models. In *IJCAI*. 458–465.
- 1696 [96] Wei Ma, Mike Papadakis, Anestis Tsakmalis, Maxime Cordy, and Yves Le Traon. 2021. Test Selection for Deep Learning Systems. *ACM Trans. Softw. Eng. Methodol.* 30, 2, Article 13 (Jan. 2021), 22 pages. <https://doi.org/10.1145/3417330>
- 1698 [97] Mohammad Mahdi Mohajer, Alvine Boaye Belle, Junjie Wang, Hadi Hemmati, Song Wang, Zhen Ming, et al. 2023. A First Look at Fairness of Machine Learning Based Code Reviewer Recommendation. *arXiv e-prints* (2023), arXiv–2307.
- 1700 [98] Suvodeep Majumder, Joymallya Chakraborty, Gina R. Bai, Kathryn T. Stolee, and Tim Menzies. 2023. Fair Enough: Searching for Sufficient Measures of Fairness. *ACM Trans. Softw. Eng. Methodol.* 32, 6, Article 134 (sep 2023), 22 pages. <https://doi.org/10.1145/3585006>
- 1702 [99] Karima Makhlof, Sami Zhioua, and Catuscia Palamidessi. 2020. Survey on Causal-based Machine Learning Fairness Notions. *arXiv preprint arXiv:2010.09553* (2020).
- 1704 [100] Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. 2021. A survey on bias and fairness in machine learning. *ACM Computing Surveys (CSUR)* 54, 6 (2021), 1–35.
- 1706 [101] Blossom Metevier, Stephen Giguerre, Sarah Brockman, Ari Kobren, Yuriy Brun, Emma Brunskill, and Philip Thomas. 2019. Offline contextual bandits with high probability fairness guarantees. *Advances in neural information processing systems* 32 (2019).
- 1708 [102] Verya Monjezi, Ashutosh Trivedi, Gang Tan, and Saeid Tizpaz-Niari. 2023. Information-Theoretic Testing and Debugging of Fairness Defects in Deep Neural Networks. *arXiv preprint arXiv:2304.04199* (2023).
- 1710 [103] Daniel Perez Morales, Takashi Kitamura, and Shingo Takada. 2021. Coverage-Guided Fairness Testing. In *International Conference on Intelligence Science*. Springer, 183–199.
- 1712 [104] Henry Muccini and Karthik Vaidhyanathan. 2021. Software architecture for ML-based systems: What exists and what lies ahead. In *2021 IEEE/ACM 1st Workshop on AI Engineering-Software Engineering for AI (WAIN)*. IEEE, 121–128.
- 1714 [105] Aniruddhan Murali, Gaurav Sahu, Kishanthan Thangarajah, Brian Zimmerman, Gema Rodriguez-Pérez, and Meiyappan Nagappan. 2024. Diversity in issue assignment: humans vs bots. *Empirical Software Engineering* 29, 2 (2024), 37.
- 1716 [106] Eirini Ntoutsi, Pavlos Falafios, Ujwal Gadiraju, Vasileios Iosifidis, Wolfgang Nejdl, Maria-Ester Vidal, Salvatore Ruggieri, Franco Turini, Symeon Papadopoulos, Emmanouil Krasanakis, et al. 2020. Bias in data-driven artificial intelligence systems—An introductory survey. *Wiley Interdisciplinary*

- 1717 *Reviews: Data Mining and Knowledge Discovery* 10, 3 (2020), e1356.
- 1718 [107] U.S. Department of Labor. Accessed: 25.04.2022. Affirmative Action. <https://www.dol.gov/general/topic/hiring/affirmativeact>
- 1719 [108] Kexin Pei, Yinzheng Cao, Junfeng Yang, and Suman Jana. 2017. Deepxplore: Automated whitebox testing of deep learning systems. In *proceedings of the 26th Symposium on Operating Systems Principles*. 1–18.
- 1720 [109] Anjana Perera, Aldeida Aleti, Chakkrit Tantithamthavorn, Jirayus Jiarpakdee, Burak Turhan, Lisa Kuhn, and Katie Walker. 2022. Search-based fairness testing for regression-based machine learning systems. *Empirical Software Engineering* 27, 3 (2022), 79.
- 1721 [110] Sai Sathiesh Rajan, Ezekiel Soremekun, Yves Le Traon, and Sudipta Chattopadhyay. 2024. Distribution-aware fairness test generation. *Journal of Systems and Software* 215 (2024), 112090. <https://doi.org/10.1016/j.jss.2024.112090>
- 1722 [111] Sai Sathiesh Rajan, Sakshi Udeshi, and Sudipta Chattopadhyay. 2022. AequeVox: Automated Fairness Testing of Speech Recognition Systems. (2022).
- 1723 [112] Stephen Rea. 2020. A Survey of Fair and Responsible Machine Learning and Artificial Intelligence: Implications of Consumer Financial Services. Available at SSRN 3527034 (2020).
- 1724 [113] Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. 2020. Beyond Accuracy: Behavioral Testing of NLP Models with CheckList. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 4902–4912.
- 1725 [114] Adam Rumpf and Hemanshu Kaul. 2021. A Public Transit Network Optimization Model for Equitable Access to Social Services. In *Proceedings of the 1st ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization* (–, NY, USA) (EAAMO ’21). Association for Computing Machinery, New York, NY, USA, Article 16, 17 pages. <https://doi.org/10.1145/3465416.3483288>
- 1726 [115] Abel Salinas, Parth Shah, Yuzhong Huang, Robert McCormack, and Fred Morstatter. 2023. The Unequal Opportunities of Large Language Models: Examining Demographic Biases in Job Recommendations by ChatGPT and LLaMA. In *Proceedings of the 3rd ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization* (EAAMO ’23). Association for Computing Machinery, New York, NY, USA, Article 34, 15 pages. <https://doi.org/10.1145/3617694.3623257>
- 1727 [116] Meirav Segal, Anne-Marie George, and Christos Dimitrakakis. 2023. Policy Fairness and Unknown Bias Dynamics in Sequential Allocations. In *Proceedings of the 3rd ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization* (EAAMO ’23). Association for Computing Machinery, New York, NY, USA, Article 38, 10 pages. <https://doi.org/10.1145/3617694.3623262>
- 1728 [117] Shahar Segal, Yossi Adi, Benny Pinkas, Carsten Baum, Chaya Ganesh, and Joseph Keshet. 2021. Fairness in the eyes of the data: Certifying machine-learning models. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*. 926–935.
- 1729 [118] Arnab Sharma, Caglar Demir, Axel-Cyrille Ngonga Ngomo, and Heike Wehrheim. 2021. MLCheck-Property-Driven Testing of Machine Learning Models. *arXiv preprint arXiv:2105.00741* (2021).
- 1730 [119] Weijun Shen, Yanhui Li, Lin Chen, Yuanlei Han, Yuming Zhou, and Baowen Xu. 2020. Multiple-Boundary Clustering and Prioritization to Promote Neural Network Retraining. In *2020 35th IEEE/ACM International Conference on Automated Software Engineering* (ASE). 410–422.
- 1731 [120] Shubham Singh, Bhuvni Shah, Chris Kanich, and Ian A. Kash. 2022. Fair Decision-Making for Food Inspections. In *Proceedings of the 2nd ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization* (EAAMO ’22). Association for Computing Machinery, New York, NY, USA, Article 5, 11 pages. <https://doi.org/10.1145/3551624.3555289>
- 1732 [121] Kacper Sokol, Alexander Hepburn, Rafael Poyiadzi, Matthew Clifford, Raul Santos-Rodriguez, and Peter Flach. 2020. FAT forensics: a python toolbox for implementing and deploying fairness, accountability and transparency algorithms in predictive systems. *Journal of Open Source Software* 5, 49 (2020), 1904.
- 1733 [122] Ezekiel Soremekun, Sakshi Udeshi, and Sudipta Chattopadhyay. 2022. Astraea: Grammar-based fairness testing. *IEEE Transactions on Software Engineering* (2022).
- 1734 [123] Youcheng Sun, Xiaowei Huang, Daniel Kroening, James Sharp, Matthew Hill, and Rob Ashmore. 2019. DeepConcolic: Testing and Debugging Deep Neural Networks. In *2019 IEEE/ACM 41st International Conference on Software Engineering: Companion Proceedings (ICSE-Companion)*. 111–114. <https://doi.org/10.1109/ICSE-Companion.2019.00051>
- 1735 [124] Zeyu Sun, Zhenpeng Chen, Jie Zhang, and Dan Hao. 2024. Fairness Testing of Machine Translation Systems. *ACM Trans. Softw. Eng. Methodol.* 33, 6, Article 156 (June 2024), 27 pages. <https://doi.org/10.1145/3664608>
- 1736 [125] Zeyu Sun, Jie M Zhang, Mark Harman, Mike Papadakis, and Lu Zhang. 2020. Automatic testing and improvement of machine translation. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 974–985.
- 1737 [126] Zeyu Sun, Jie M Zhang, Yingfei Xiong, Mark Harman, Mike Papadakis, and Lu Zhang. 2022. Improving Machine Translation Systems via Isotopic Replacement. In *2022 IEEE/ACM 44th International Conference on Software Engineering (ICSE)*. IEEE.
- 1738 [127] Guanhong Tao, Weisong Sun, Tingxu Han, Chunrong Fang, and Xiangyu Zhang. 2022. RULER: discriminative and iterative adversarial training for deep neural network fairness. In *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 1173–1184.
- 1739 [128] Philip S. Thomas, Bruno Castro da Silva, Andrew G. Barto, Stephen Giguere, Yuriy Brun, and Emma Brunskill. 2019. Preventing undesirable behavior of intelligent machines. *Science* 366, 6468 (2019), 999–1004. <https://doi.org/10.1126/science.aag3311> arXiv:<https://www.science.org/doi/pdf/10.1126/science.aag3311>
- 1740 [129] Yuchi Tian, Ziyuan Zhong, Vicente Ordonez, Gail Kaiser, and Baishakhi Ray. 2020. Testing DNN image classifiers for confusion & bias errors. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 1122–1134.
- 1741 [129] Yuchi Tian, Ziyuan Zhong, Vicente Ordonez, Gail Kaiser, and Baishakhi Ray. 2020. Testing DNN image classifiers for confusion & bias errors. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 1122–1134.
- 1742 [129] Yuchi Tian, Ziyuan Zhong, Vicente Ordonez, Gail Kaiser, and Baishakhi Ray. 2020. Testing DNN image classifiers for confusion & bias errors. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 1122–1134.
- 1743 [129] Yuchi Tian, Ziyuan Zhong, Vicente Ordonez, Gail Kaiser, and Baishakhi Ray. 2020. Testing DNN image classifiers for confusion & bias errors. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 1122–1134.
- 1744 [129] Yuchi Tian, Ziyuan Zhong, Vicente Ordonez, Gail Kaiser, and Baishakhi Ray. 2020. Testing DNN image classifiers for confusion & bias errors. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 1122–1134.
- 1745 [129] Yuchi Tian, Ziyuan Zhong, Vicente Ordonez, Gail Kaiser, and Baishakhi Ray. 2020. Testing DNN image classifiers for confusion & bias errors. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 1122–1134.
- 1746 [129] Yuchi Tian, Ziyuan Zhong, Vicente Ordonez, Gail Kaiser, and Baishakhi Ray. 2020. Testing DNN image classifiers for confusion & bias errors. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 1122–1134.
- 1747 [129] Yuchi Tian, Ziyuan Zhong, Vicente Ordonez, Gail Kaiser, and Baishakhi Ray. 2020. Testing DNN image classifiers for confusion & bias errors. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 1122–1134.
- 1748 [129] Yuchi Tian, Ziyuan Zhong, Vicente Ordonez, Gail Kaiser, and Baishakhi Ray. 2020. Testing DNN image classifiers for confusion & bias errors. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 1122–1134.
- 1749 [129] Yuchi Tian, Ziyuan Zhong, Vicente Ordonez, Gail Kaiser, and Baishakhi Ray. 2020. Testing DNN image classifiers for confusion & bias errors. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 1122–1134.
- 1750 [129] Yuchi Tian, Ziyuan Zhong, Vicente Ordonez, Gail Kaiser, and Baishakhi Ray. 2020. Testing DNN image classifiers for confusion & bias errors. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 1122–1134.
- 1751 [129] Yuchi Tian, Ziyuan Zhong, Vicente Ordonez, Gail Kaiser, and Baishakhi Ray. 2020. Testing DNN image classifiers for confusion & bias errors. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 1122–1134.
- 1752 [129] Yuchi Tian, Ziyuan Zhong, Vicente Ordonez, Gail Kaiser, and Baishakhi Ray. 2020. Testing DNN image classifiers for confusion & bias errors. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 1122–1134.
- 1753 [129] Yuchi Tian, Ziyuan Zhong, Vicente Ordonez, Gail Kaiser, and Baishakhi Ray. 2020. Testing DNN image classifiers for confusion & bias errors. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 1122–1134.
- 1754 [129] Yuchi Tian, Ziyuan Zhong, Vicente Ordonez, Gail Kaiser, and Baishakhi Ray. 2020. Testing DNN image classifiers for confusion & bias errors. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 1122–1134.
- 1755 [129] Yuchi Tian, Ziyuan Zhong, Vicente Ordonez, Gail Kaiser, and Baishakhi Ray. 2020. Testing DNN image classifiers for confusion & bias errors. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 1122–1134.
- 1756 [129] Yuchi Tian, Ziyuan Zhong, Vicente Ordonez, Gail Kaiser, and Baishakhi Ray. 2020. Testing DNN image classifiers for confusion & bias errors. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 1122–1134.
- 1757 [129] Yuchi Tian, Ziyuan Zhong, Vicente Ordonez, Gail Kaiser, and Baishakhi Ray. 2020. Testing DNN image classifiers for confusion & bias errors. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 1122–1134.
- 1758 [129] Yuchi Tian, Ziyuan Zhong, Vicente Ordonez, Gail Kaiser, and Baishakhi Ray. 2020. Testing DNN image classifiers for confusion & bias errors. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 1122–1134.
- 1759 [129] Yuchi Tian, Ziyuan Zhong, Vicente Ordonez, Gail Kaiser, and Baishakhi Ray. 2020. Testing DNN image classifiers for confusion & bias errors. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 1122–1134.
- 1760 [129] Yuchi Tian, Ziyuan Zhong, Vicente Ordonez, Gail Kaiser, and Baishakhi Ray. 2020. Testing DNN image classifiers for confusion & bias errors. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 1122–1134.
- 1761 [129] Yuchi Tian, Ziyuan Zhong, Vicente Ordonez, Gail Kaiser, and Baishakhi Ray. 2020. Testing DNN image classifiers for confusion & bias errors. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 1122–1134.
- 1762 [129] Yuchi Tian, Ziyuan Zhong, Vicente Ordonez, Gail Kaiser, and Baishakhi Ray. 2020. Testing DNN image classifiers for confusion & bias errors. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 1122–1134.
- 1763 [129] Yuchi Tian, Ziyuan Zhong, Vicente Ordonez, Gail Kaiser, and Baishakhi Ray. 2020. Testing DNN image classifiers for confusion & bias errors. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 1122–1134.
- 1764 [129] Yuchi Tian, Ziyuan Zhong, Vicente Ordonez, Gail Kaiser, and Baishakhi Ray. 2020. Testing DNN image classifiers for confusion & bias errors. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 1122–1134.
- 1765 [129] Yuchi Tian, Ziyuan Zhong, Vicente Ordonez, Gail Kaiser, and Baishakhi Ray. 2020. Testing DNN image classifiers for confusion & bias errors. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 1122–1134.
- 1766 [129] Yuchi Tian, Ziyuan Zhong, Vicente Ordonez, Gail Kaiser, and Baishakhi Ray. 2020. Testing DNN image classifiers for confusion & bias errors. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 1122–1134.
- 1767 [129] Yuchi Tian, Ziyuan Zhong, Vicente Ordonez, Gail Kaiser, and Baishakhi Ray. 2020. Testing DNN image classifiers for confusion & bias errors. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 1122–1134.
- 1768 [129] Yuchi Tian, Ziyuan Zhong, Vicente Ordonez, Gail Kaiser, and Baishakhi Ray. 2020. Testing DNN image classifiers for confusion & bias errors. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 1122–1134.

- [1769] [130] Saeid Tizpaz-Niari, Ashish Kumar, Gang Tan, and Ashutosh Trivedi. 2022. Fairness-aware Configuration of Machine Learning Libraries. In *2022 IEEE/ACM 44th International Conference on Software Engineering (ICSE)*. IEEE.
- [1770] [131] Sakshi Udeshi, Pryanshu Arora, and Sudipta Chattopadhyay. 2018. Automated directed fairness testing. In *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering*. 98–108.
- [1771] [132] Sahil Verma and Julia Rubin. 2018. Fairness definitions explained. In *2018 ieee/acm international workshop on software fairness (fairware)*. IEEE, 1–7.
- [1772] [133] Junjie Wang, Ye Yang, Song Wang, Jun Hu, and Qing Wang. 2022. Context- and Fairness-Aware In-Process Crowdworker Recommendation. *ACM Trans. Softw. Eng. Methodol.* 31, 3, Article 35 (mar 2022), 31 pages. <https://doi.org/10.1145/3487571>
- [1773] [134] Tianlu Wang, Jieyu Zhao, Mark Yatskar, Kai-Wei Chang, and Vicente Ordonez. 2019. Balanced datasets are not enough: Estimating and mitigating gender bias in deep image representations. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5310–5319.
- [1774] [135] Zhaohui Wang, Min Zhang, Jingran Yang, Bojie Shao, and Min Zhang. 2024. MAFT: Efficient Model-Agnostic Fairness Testing for Deep Neural Networks via Zero-Order Gradient Search. In *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering*. 1–12.
- [1775] [136] Nimmri Rashiniika Weeraddana, Xiaoyan Xu, Mahmoud Alfadel, Shane McIntosh, and Meiyappan Nagappan. 2023. An empirical comparison of ethnic and gender diversity of DevOps and non-DevOps contributions to open-source projects. *Empirical Software Engineering* 28, 6 (2023), 150.
- [1776] [137] Yisong Xiao, Aishan Liu, Tianlin Li, and Xianglong Liu. 2023. Latent Imitator: Generating Natural Individual Discriminatory Instances for Black-Box Fairness Testing. *arXiv preprint arXiv:2305.11602* (2023).
- [1777] [138] Ying Xiao, Jie M Zhang, Yeqang Liu, Mohammad Reza Mousavi, Sicen Liu, and Dingyuan Xue. 2024. MirrorFair: Fixing fairness bugs in machine learning software via counterfactual predictions. *Proceedings of the ACM on Software Engineering* 1, FSE (2024), 2121–2143.
- [1778] [139] Xiaofei Xie, Lei Ma, Felix Juefei-Xu, Minhui Xue, Hongxu Chen, Yang Liu, Jianjun Zhao, Bo Li, Jianxiong Yin, and Simon See. 2019. DeepHunter: A Coverage-Guided Fuzz Testing Framework for Deep Neural Networks. In *Proceedings of the 28th ACM SIGSOFT International Symposium on Software Testing and Analysis* (Beijing, China) (ISSTA 2019). Association for Computing Machinery, New York, NY, USA, 146–157. <https://doi.org/10.1145/3293882.3330579>
- [1779] [140] Junjie Yang, Jiajun Jiang, Zeyu Sun, and Junjie Chen. 2024. A large-scale empirical study on improving the fairness of image classification models. In *Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis*. 210–222.
- [1780] [141] Liu Yang, Mingyang Zhang, Cheng Li, Michael Bendersky, and Marc Najork. 2020. Beyond 512 tokens: Siamese multi-depth transformer-based hierarchical encoder for long-form document matching. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*. 1725–1734.
- [1781] [142] Zhou Yang, Muhammad Hilmi Asyrofi, and David Lo. 2021. BiasRV: Uncovering biased sentiment predictions at runtime. In *Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 1540–1544.
- [1782] [143] Christina Yeung, Umar Iqbal, Tadayoshi Kohno, and Franziska Roesner. 2023. Gender Biases in Tone Analysis: A Case Study of a Commercial Wearable. In *Proceedings of the 3rd ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization* (EAAMO '23). Association for Computing Machinery, New York, NY, USA, Article 21, 12 pages. <https://doi.org/10.1145/3617694.3623241>
- [1783] [144] Zhe Yu, Joymallya Chakraborty, and Tim Menzies. 2024. FairBalance: How to Achieve Equalized Odds With Data Pre-Processing. *IEEE Trans. Softw. Eng.* 50, 9 (Sept. 2024), 2294–2312. <https://doi.org/10.1109/TSE.2024.3431445>
- [1784] [145] Meike Zehlike, Ke Yang, and Julia Stoyanovich. 2021. Fairness in Ranking: A Survey. *arXiv preprint arXiv:2103.14000* (2021).
- [1785] [146] Rich Zemel, Yu Wu, Kevin Swersky, Toni Pitassi, and Cynthia Dwork. 2013. Learning fair representations. In *International conference on machine learning*. PMLR, 325–333.
- [1786] [147] Jie M Zhang and Mark Harman. 2021. “Ignorance and Prejudice” in Software Fairness. In *2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE)*. IEEE, 1436–1447.
- [1787] [148] Jie M Zhang, Mark Harman, Lei Ma, and Yang Liu. 2020. Machine learning testing: Survey, landscapes and horizons. *IEEE Transactions on Software Engineering* (2020).
- [1788] [149] Lingfeng Zhang, Yueling Zhang, and Min Zhang. 2021. Efficient white-box fairness testing through gradient search. In *Proceedings of the 30th ACM SIGSOFT International Symposium on Software Testing and Analysis*. 103–114.
- [1789] [150] Mengdi Zhang, Jun Sun, Jingyi Wang, and Bing Sun. 2023. TestSGD: Interpretable Testing of Neural Networks against Subtle Group Discrimination. *ACM Trans. Softw. Eng. Methodol.* 32, 6, Article 137 (sep 2023), 24 pages. <https://doi.org/10.1145/3591869>
- [1790] [151] Peixin Zhang, Jingyi Wang, Jun Sun, Guoliang Dong, Xinyu Wang, Xingen Wang, Jin Song Dong, and Ting Dai. 2020. White-box fairness testing through adversarial sampling. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 949–960.
- [1791] [152] Peixin Zhang, Jingyi Wang, Jun Sun, Xinyu Wang, Guoliang Dong, Xingen Wang, Ting Dai, and Jin Song Dong. 2021. Automatic Fairness Testing of Neural Classifiers through Adversarial Sampling. *IEEE Transactions on Software Engineering* (2021).
- [1792] [153] Xueru Zhang and Mingyan Liu. 2021. Fairness in learning-based sequential decision algorithms: A survey. In *Handbook of Reinforcement Learning and Control*. Springer, 525–555.
- [1793] [154] Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2017. Men also like shopping: Reducing gender bias amplification using corpus-level constraints. *arXiv preprint arXiv:1707.09457* (2017).
- [1794] [155] Wei Zheng, Lidan Lin, Xiaoxue Wu, and Xiang Chen. 2024. An Empirical Study on Correlations Between Deep Neural Network Fairness and Neuron Coverage Criteria. *IEEE Trans. Softw. Eng.* 50, 3 (March 2024), 391–412. <https://doi.org/10.1109/TSE.2023.3349001>
- [1795] [156] W. Zheng, L. Lin, X. Wu, and X. Chen. 5555. An Empirical Study on Correlations between Deep Neural Network Fairness and Neuron Coverage Criteria. *IEEE Transactions on Software Engineering* 01 (jan 5555), 1–22. <https://doi.org/10.1109/TSE.2023.3349001>

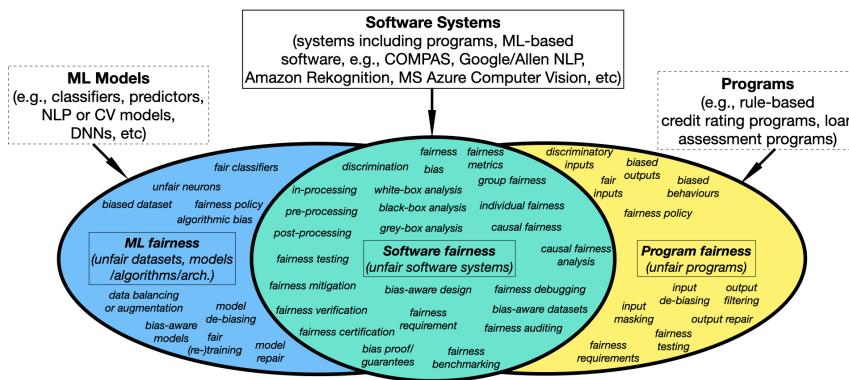
1 **Appendix: Supplementary Material for Paper titled “Software Fairness: An**
2 **Analysis and Survey”**

5 EZEKIEL SOREMEKUN, Singapore University of Technology and Design, Singapore
6

7 MIKE PAPADAKIS, MAXIME CORDY, and YVES LE TRAON, SnT, University of Luxembourg, Luxembourg
8

9 **Summary**

10 This document is the appendix (or supplementary material) for the paper titled “Software Fairness: An Analysis and Survey”. It
11 provides additional or more detailed tables and figures to support the reported findings in the original paper.
12



31 Fig. 1. Taxonomy and interplay of ML fairness and Software Fairness
32

49 Authors' addresses: Ezekiel Soremekun, ezekiel_soremekun@sutd.edu.sg, Singapore University of Technology and Design, 8 Somapah Rd, Singapore,
50 Singapore, Singapore, 487372; Mike Papadakis, michail.papadakis@uni.lu; Maxime Cordy, maxime.cordy@uni.lu; Yves Le Traon, Yves.LeTraon@uni.lu,
51 SnT, University of Luxembourg, 6, rue Richard Coudenhove-Kalergi, Luxembourg, Luxembourg, Luxembourg, L-1359.

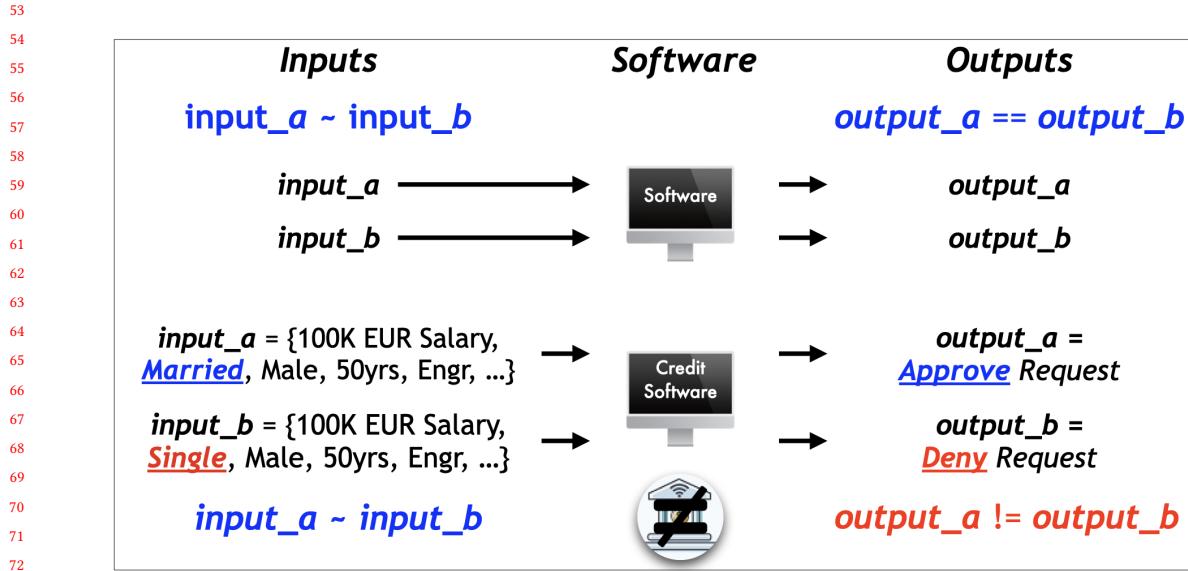


Fig. 2. Program fairness property showing fair program behavior ($\text{output_a} == \text{output_b}$) and unfair program behavior ($\text{output_a} != \text{output_b}$) using a (credit approval) software

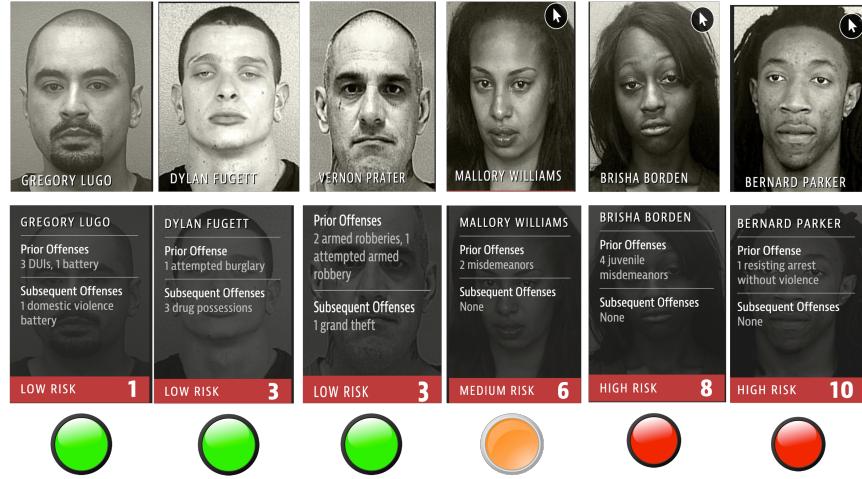


Fig. 3. An illustrative example of real-world discrimination (unfair software behavior) in the COMPAS software, a program that predicts recidivism – the likelihood of committing a future crime (adapted from Angwin et al. [12, 13])

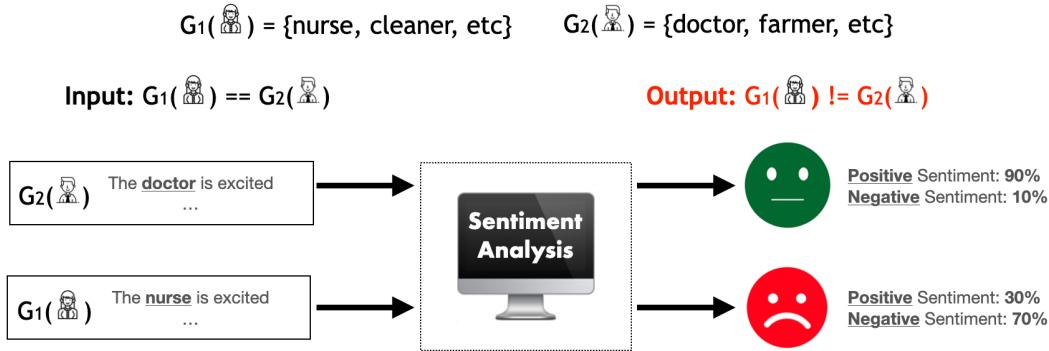


Fig. 4. Example of Group Fairness Violation using a Sentiment Analysis AI System

Table 1. Excerpt of Publication Details (“#” means “number of”)

Domain (#Pubs.)	Type	#Venues	Venue Sample Venue(s)	#Pubs	Publications Example Publications	Years
Software Engineering (SE), Programming Languages (PL) & Security	Conference	9	ASE, CAV, EuroS&P, FSE, GECCO, ICSE, OOPSLA, TrustCom, ISSTA	28	[24, 42, 50, 61, 62, 78, 101, 141]	2017-21
	Journal	4	EMSE, JSS, RE, TSE	7	[14, 20, 21, 57, 136, 137, 164]	2009-21
	Other	1	ICSE-C	1	[6]	2021
Natural language processing (NLP)	Conference	2	ACL, EMNLP	6	[28, 29, 56, 123, 129, 166]	2017-21
Artificial Intelligence (AI) & Machine Learning (ML)	Conference	8	AAAI, AISTATS, ICML, NeurIPS, PMLR	17	[4, 96, 109, 158] [3, 35, 83, 92, 121, 154]	2013-21
	Other	1	HRLC	1	[165]	2021
Computer Vision (CV)	Conference	2	ICCV, CVPR	3	[94, 150, 151]	2019-20
	Workshop	1	ECCV	1	[157]	2020
Fairness-targets	Conference	2	AAAI-AIES, Facet	21	[9, 22, 27, 36, 58, 73, 93, 125, 134]	2017-21
	Workshop	2	FairWare, FATE	8	[15, 30, 46, 47, 80, 148, 153]	2018-19
Big Data, Data Mining (DM), & Knowledge Discovery (KD)	Conference	8	DMKD, ECML-PKDD, EDBT, KDD, ICDM, ICEDT, ICMD, LAK	18	[34, 89, 99, 169] [43, 52, 88, 91, 139]	2010-21
	Journal	5	Big Data, Inf. Science, JDIQ, KAIS, SIGMOD-Record	5	[1, 48, 87] [90, 128]	2012-19
	Workshop	2	BSDUC, KDD-XAI	2	[69, 124]	2018-19
Human Factors & Usability	Conference	1	CHI	3	[44, 75, 98]	2019-21
	Journal	1	IWC	1	[32]	2016
Others	Conference	3	CCCT, VAST, WWW	6	[33, 51, 70, 86, 97]	2009-20
	Journal	6	CACM, DGRP, Scientific-data, SSRN, IBM Journal of R & D	10	[19, 64, 82, 130, 131]	2018-21
	Workshop	1	CEUR Workshop	1	[138]	2019
	Other	3	arXiv, HRDAG, MS Tech. Report	15	[23, 45, 76, 84, 104, 127, 133]	2019-21

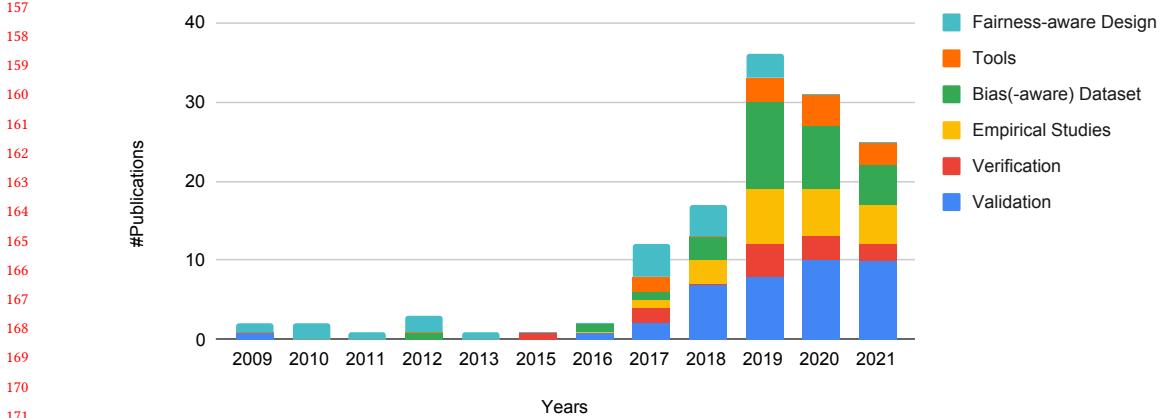


Fig. 5. Detailed Publication Trend by Year

Table 2. Purpose of Software Fairness Analysis (recent (2022 and 2023) publications are in bracket).

Categories	Sub-category	Description	#Pubs	Sample Works
Validation	Testing	generating discriminatory test inputs to expose fairness violations	20 (13)	[160, 163] ([10, 39, 100, 155])
	Mitigation	mitigating bias in software systems, e.g., via repair and prevention	14 (13)	[8] ([40, 62, 113, 118, 142, 161])
	Debugging	diagnosis and explanation of fairness violations	8 (3)	[42, 101, 137] ([110, 116])
	Auditing	analysing and measuring bias in software systems	2 (1)	[33, 96] ([162])
Verification	Verifiers	verifying that a system fulfills a fairness metric or goal	12 (1)	[7, 18, 65, 83] ([26])
	Certification	certifying that a system fulfills a fairness goal	4 (1)	[52, 132] ([26])
Design	Proof or Guarantees	providing a formal proof that a system achieves a fairness goal	2 (2)	[36, 109] ([66, 74])
	Requirements	requirement engineering and formalization of fairness properties	4 (2)	[21, 57, 103] ([17, 120])
Empirical Evaluation	Bias-aware Design	designing fair systems and bias-aware software	15 (5)	[34, 80, 88] ([63, 81, 143])
	Analysis	empirical studies about fairness concerns	22 (14)	[25, 30, 47, 166] ([72, 112, 152])
Datasets	Benchmarking	providing fair benchmarks or benchmarks for fairness evaluations	4 (3)	[24, 29, 78, 151] ([67, 77, 149])
	Bias in Datasets	studying biases in training and evaluation datasets	30	[35, 64, 150, 157]
Tooling	Bias-aware Datasets	developing unbiased or bias-aware datasets for better evaluation	5 (1)	[30, 124, 129, 130] ([149])
	Automatic	providing fully automatic tools for fairness analysis	18	[7, 19, 23, 65]
	Semi-automatic	building tools that require human interaction for fairness analysis	2	[33, 102]

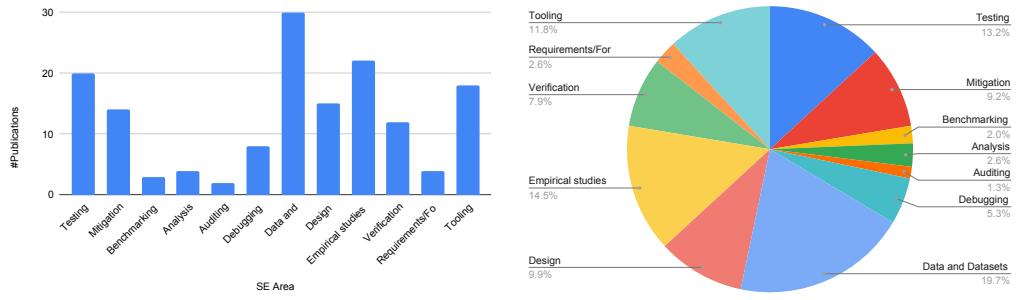


Fig. 6. Purpose of Fairness Analysis in SE community

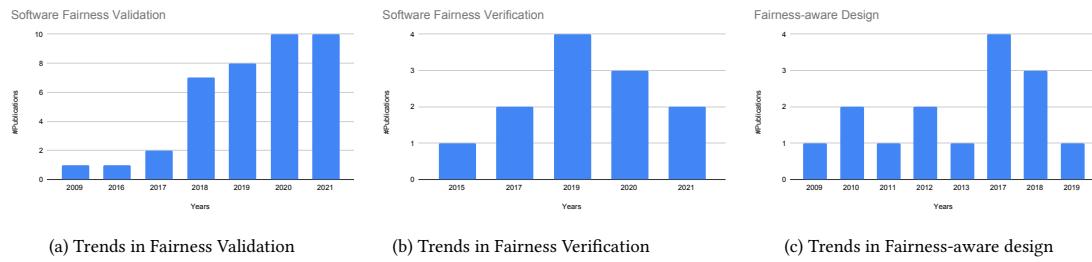


Fig. 7. Details of trends in Fairness Verification, Validation and Design

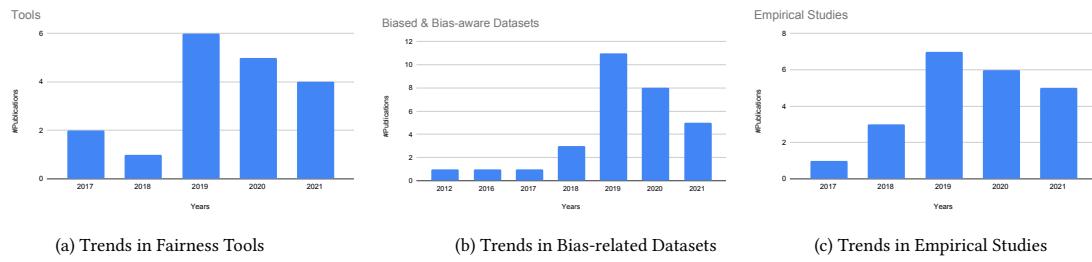


Fig. 8. Details of trends in Fairness Tooling, Datasets and Empirical Studies

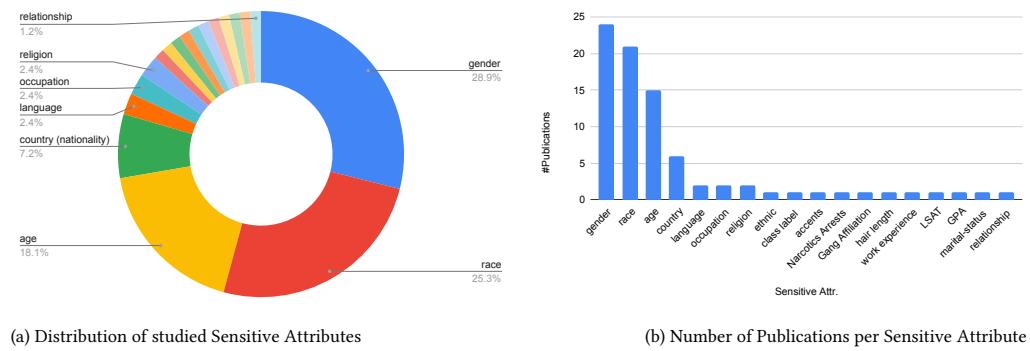


Fig. 9. Details of Biases, i.e., Sensitive/Protected attributes

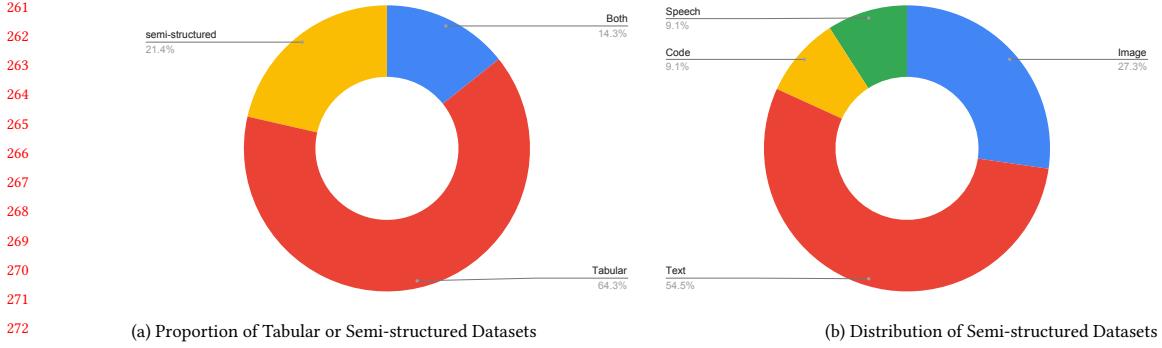


Fig. 10. Type of Studied Datasets

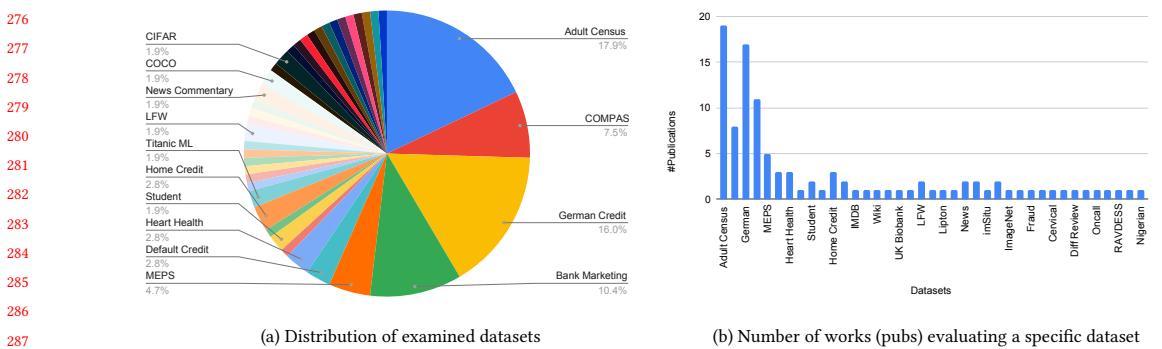


Fig. 11. Details of examined Datasets

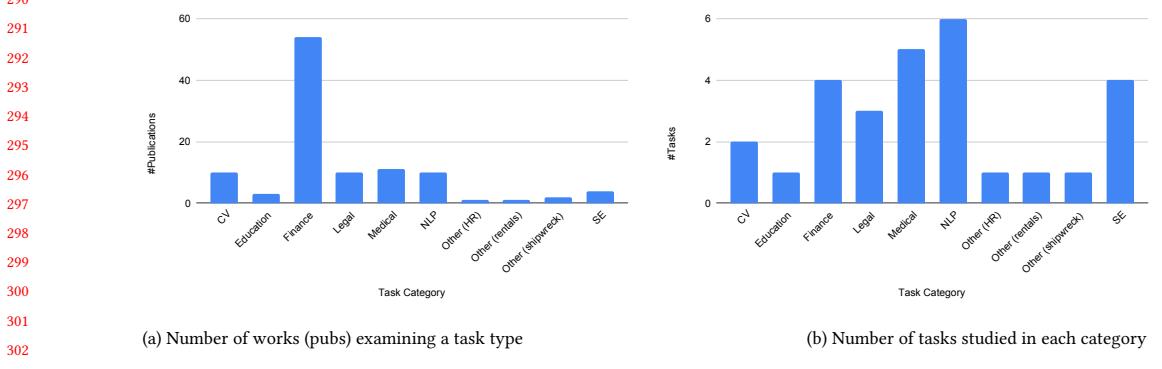


Fig. 12. Details of Task Categories

Table 3. Details of Tasks and Datasets employed in Software Fairness Analysis

Type	Task Category	#Data	#Pubs	#Tasks	Tasks (Example Pubs)	Datasets (#Pubs)
Tabular	Education	2	3	1	academic performance [27]	Law School (1)
	Finance	6	54	4	income [5, 159, 163]	Adult Census (19)
					credit default [11, 25, 38]	German Credit (17), Default Credit (3), Home Credit (3)
					potential buyers [25, 37, 101]	Bank Marketing (11)
	Legal	3	10	3	fraud [5]	Fraud Detection (1)
					recidivism [33, 38, 78]	COMPAS (8)
					arrests [27]	Chicago Strategic Subject List (SSL) (1)
	Medical	5	11	5	US executions [5]	US Executions (1)
					medical expenditure [37, 101, 167]	MEPS (5)
					heart disease [37, 38]	Heart Health (3)
					heart failure [42]	Heart Failure (1)
	Other (HR)	1	1	1	cancer risk [42]	Cervical Cancer (1)
					hiring [27]	Lipton (1)
					car rentals [5]	Raw Car Rentals (1)
	Other (shipwreck)	1	2	1	shipwreck survival	Titanic ML (2)
Semi-structured	CV (image)	7	10	2	face detection	ClbA-IN (1), PPB (1), LFW (2)
					image recognition	COCO (2), imSitu (1), CIFAR (2), ImageNet (1)
	NLP (text, speech)	12	10	6	SA [14], CoRef [137], MLM [137]	Twitter (1), IMDB (1), EEC Dataset (1), Labor statistics (1)
					toxicity	Wiki Comment (1), Jigsaw Comments (1)
					machine translation	News Commentary (2)
					ASR [122]	Speech Accent Archive (1), RAVDESS (1), Multi speaker Corpora of the English Accents in the British Isles (1), Nigerian English speech dataset (1)
	SE (code)	4	4	4	Programs [42]	Bug2Commit (1), Diff Review (1), Code AutoComplete (1), Oncall Recommendation (1)

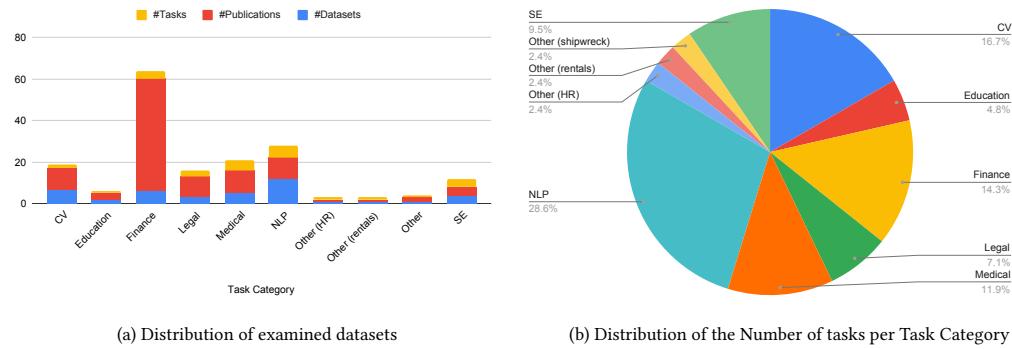


Fig. 13. Details of Publications and Datasets for each Task category

Table 4. Excerpt of works performing Software Fairness Analysis (“Acc.” means “level of software access”)

Acc.	Approach	Goal	Problem	Main Idea	Core Technique
Black-box	LTDD [101]	Debugging	identifying biased features in training data	debugging biased features to build fair ML software.	data debugging, linear-regression
	Multiacc. Boost [96]	Auditing, Analysis, Mitigation	audit/mitigate multiaccuracy, i.e., group fairness for all subgroups	perform multiaccuracy audit and post-process models to achieve it	multiaccuracy auditing, post-processing
	Cito et al. [42]	Debugging, Mitigation	debug and isolate the cause of mispredictions in ML models	characterize the data on which the model performs poorly	rule induction
	ASTR-AEA [137]	Debugging, Testing, Mitigation	performing fairness testing without existing datasets	discover and diagnose fairness violations in NLP software	grammar-based testing
	Fair-Way [38]	Debugging, Testing, Mitigation	detect and explain how ML model acquires bias from training data	identify how ground truth bias affects ML fairness	multi-objective optimization, pre/in -processing
	FairSM-OTE [37]	Debugging, Testing, Mitigation	finding biased labels in training data generation	remove biased labels and balance data using sensitive attribute	situation testing, data balancing
	Fair-Vis [33]	Auditing, Analysis	auditing and analysing group fairness in ML model	visual analytics for the discovery and audit of (sub)group fairness	visual analytics, domain knowledge
	Flip-test [27]	Testing, Mitigation	testing individual fairness – similar treatment of protected statuses	discover individual (un)fairness and its associated features	optimal transport, flipset, dist. sampling
	Aequitas [147]	Testing, Mitigation	validation of fairness for arbitrary ML models?	generating discriminatory inputs to uncover fairness violations	directed testing, probabilistic search
	Themis [11, 31, 61]	Formaliz., Testing	formalize software fairness testing for causal discovering of discrimination	measure causal discrimination in software to direct fairness testing	input schema, causal relationships
	Aequ Vox [122]	Debugging, Testing	testing group fairness for Automatic Speech Recognition (ASR) systems	group fairness testing by simulating different environments	ML robustness, test simulation, fault localization
	ExpGA [50]	Testing	current individual fairness testing methods suffer poor efficiency, effectiveness, and model specificity	fairness testing by modifying feature values using explanation results and genetic algorithm (GA)	genetic algorithm, feature mutation, search based testing
	CGFT [111]	Testing	Uneven distribution of fairness tests and variations in execution results	leverage combinatorial testing to generate evenly-distributed test suites	combinatorial testing, input coverage
	SG [5]	Testing	detecting the presence of individual discrimination in ML models.	auto-generation of test inputs for detecting individual discrimination.	symbolic execution, local explainability
	Bias-Finder [14]	Testing	Bias testers for SA systems rely on small, short, predefined templates	discover biased predictions in SA systems via metamorphic testing.	template curation, NLP techniques, metamorphic testing
	Biswas and Rajan [24]	Mitigation	understanding fairness characteristics in ML models from practice	empirical evaluation of fairness and mitigations on real-world ML models	empirical study
	Fairea [78]	Mitigation	what is the SE trade-off between accuracy and fairness?	benchmarking and quantifying the fairness-accuracy trade-off achieved by bias mitigation methods	model behaviour mutation
White-box	ADF [163, 164]	Testing	searching individual discriminatory instances	generating discriminatory inputs violating individual fairness via ML	gradient computation and clustering
	Deep-Inspect [144]	Testing	detecting confusion and bias errors at class-level	expose confusion and bias errors in image classifiers	class property violations, robustness
	EIDIG [160]	Testing, Mitigation	how to detect and improve individual fairness of a model	generating test cases that violate individual fairness	gradient descent, global/local search
	Neuron-Fair [167]	Testing, Mitigation, Analysis	interpretability, performance, and generalizability in bias testing	identifying biased neurons, i.e., neurons that cause discrimination	neuron activation, adversarial attacks
	Fair-Neuron [62]	Mitigation, Analysis	balancing accuracy-fairness trade-off without additional model(s)	detect neurons with contradictory optimization directions, and achieve trade-off via selective dropout	joint-optimization, adversarial game
Grey	Tizpaz-Niari [145]	Debugging, Testing, Mitigation	explaining fairness impact of hyper-parameters	identify the effect of parameters on software fairness	search based testing, statistical debugging
	CAT/TransRepair [140, 141]	Testing, Mitigation	detecting inconsistency in machine translation (MT)	detect inconsistency bugs without access to human oracles	mutation testing, metamorphic testing, language model (BERT)

Table 5. Excerpt of Fairness Analysis Tools

Tool (Paper)	Goal	Addressed Problem	Process. Stage	Approach	Access
Fairkit-learn [84, 85]	Fair learning, Analysis	how to reason about and determine the trade-off between model quality (accuracy) and fairness	pre, & post	model search, visualisation	Grey
AIF360 [19]	Fair learning, Analysis	understanding how, when and why to use different bias handling algorithms in the model life-cycle	pre, in, & post	extensible architecture for analysing fairness metrics	Grey
POF [22]	Fair learning, Analysis	how to compute the “Pareto curve” of the trade-off between accuracy and fairness in the <i>regression</i> settings (<i>continuous</i> prediction/targeted values)	pre	fairness regularizers, Price of Fairness (PoF) metric	Grey
AITEST [6]	Testing	how to detect the presence of individual discrimination in ML models	post	symbolic execution and local explainability	Black
2AFC [102]	Testing	how to relate unobservable phenomena deep inside models with observable, outside quantities that we can measure from inputs and outputs	post	Test Experiments, Experimental Psychology, Psychophysics, two-alternative forced choice (2AFC)	Black
Pc-fairness [154]	Formalization, Analysis	how to bound path-specific counterfactual fairness, address their <i>identifiability</i> , i.e., whether they can be uniquely measured from observational data	post	parameterized causal modelling, linear programming, response-function variables, constraints	Black
BiasRV [156]	Testing	how to monitor and uncover biased predictions at runtime	post	automatic template generation, mutation, metamorphic relations	Black
Themis [11]	Formalization, Testing	how to formally define and test software fairness using a causality-based measure of discrimination	post	causal inference, schema-based test generation	Black
Fair-Square [7]	Formalization, Verification, Certification	how to verify or certify that a program meets a given fairness property	post	probabilistic reasoning, SMT solving, symbolic weighted -volume-computation algorithm	Black
FAT Forensics [136]	Analysis, Auditing, certifying	how to inspect datasets (features), models and their prediction for fairness metrics	pre, in, & post	an inter-operable Python framework for fairness (FAT) algorithms	White, Black, & Grey
VeriFair [18]	Verification, Specification	how to verify fairness specifications, i.e., fairness properties of ML programs	post	adaptive concentration inequalities	Black
Justicia [65]	Verification	how to formally verify the fairness metrics are satisfied by different algorithms on different datasets	pre	stochastic satisfiability (SSAT)	Black
Checklist [123]	Testing	how to test (fairness) behaviors of NLP systems	post	behavioral testing, template-based test generation	Black
FairML [2]	Auditing, Analysis	how to determine the significance of inputs in assessing the fairness of black-box models	post	model compression, input ranking algorithms	Black
MT-NLP [106]	Testing, Mitigation	how to determine if NLP models are free of unfair bias toward certain sub-populations/groups	post	metamorphic testing	Black
ASTRAEA [137]	Debugging, Testing, Mitigation	how to perform fairness testing without an existing dataset, i.e., no training data access	post	grammar-based testing, metamorphic relations	Black
Aequitas [126]	Auditing, Analysis	how to audit for bias and fairness when developing and deploying algorithmic decision making systems	post	bias audit toolkit to support many bias metrics	Black
FairTest [146]	Testing, Debugging	how to detect unwarranted associations (UA) (disparate impact, offensive labels, and uneven error rates) between model outcomes and data attributes	post	unwarranted associations (UA) framework to determine UA between outcomes and attributes	Black
Themis-ml [16]	Fair Learning, Mitigation, Analysis, Auditing	how to measure, understand, and mitigate the implicit historical biases in socially sensitive data	pre, in & post	API for Fair ML for simple binary classifier	White, Black, & Grey
Fairify [26]	Verification Certification	how to verify individual fairness property in neural network (NN) models	in	SMT, formal analysis, pruning, input partitioning, interval arithmetic and activation heuristic	White

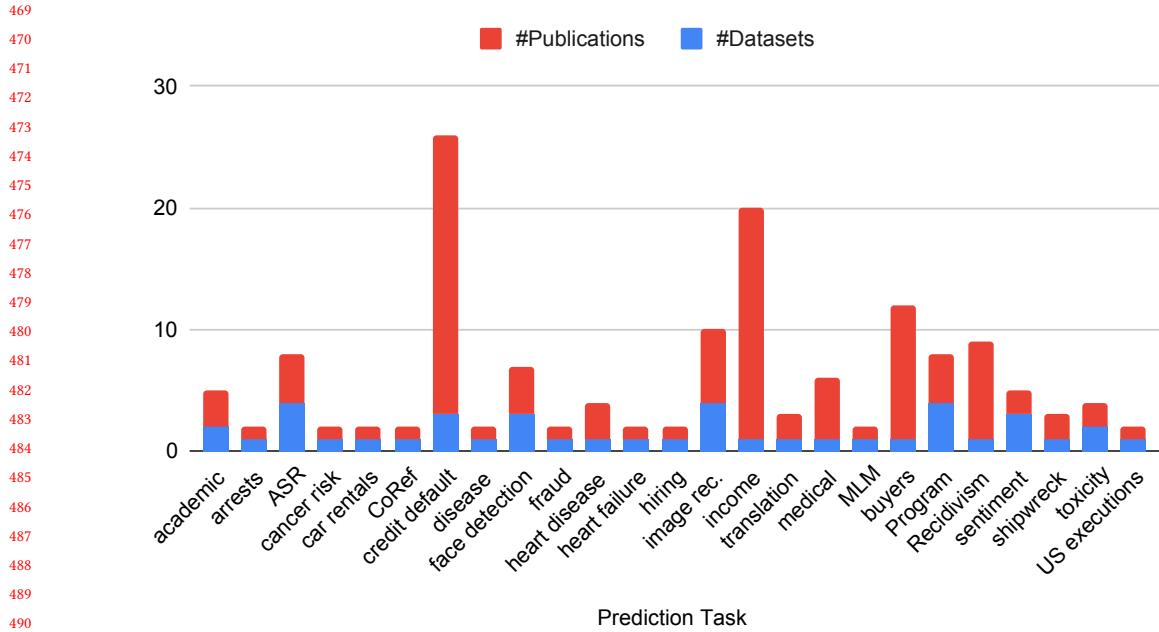


Fig. 14. Details of Publications and Datasets per Prediction Task

Table 6. Details of Open Problems and Future Research Opportunities (more recent (2022 and 2023) solutions are in bracket)

Open Problems	Problem Description	Potential Solutions	Sample Related Work
Fairness Test Metrics and Adequacy	Measuring when fairness testing is sufficient/enough	Design of fairness test metrics and adequacy criteria	[49, 53, 95, 105, 117, 135] [107] ([108, 168])
Automatic Repair of Biased Classifiers	How to automatically repair biased classifiers to be (less or) un-biased?	Automatic Program Repair for fairness property	[8, 128, 140] ([79, 142])
Tooling for Fairness Property Specification	Specifying and engineering fairness properties for learning-based systems	Requirement Engineering tool support for Fairness properties	[18, 133] ([17, 54])
Unexplored or Poorly Understood Biases	Analyzing rare biases (e.g., age), complex or intersectional biases (e.g., age × gender)?	Fairness Analysis Support for rare, complex or intersectional Biases	[30, 33] ([41, 68])
Sequential and Long-term Fairness concerns	How to analyse/maintain fairness as the AI system evolves over time?	Techniques to support analysis of sequential and long-term fairness	[165] ([67, 119])
Human factors in fairness analysis	E.g., evaluating the harm induced by fairness violations to humans/society	Empirical studies of Human Factors in Fairness Analysis	[44, 75, 98] ([55, 59, 60, 71, 114])
Non-Specific/Holistic mitigation approaches	Designing bias mitigation methods that are agnostic of tasks, domains or datasets	General (i.e., task, domain and dataset -agnostic) bias analysis techniques	[164]
Fair Policy, Legalisation, and Compliance	How to design fairness analysis tools for policy makers and compliance officers?	Fairness Analysis Tool Support for Policy and Compliance Analysis	[102, 115]

521 REFERENCES

- 522 [1] Serge Abiteboul and Julia Stoyanovich. 2019. Transparency, fairness, data protection, neutrality: Data management challenges in the face of new
523 regulation. *Journal of Data and Information Quality (JDIQ)* 11, 3 (2019), 1–9.
- 524 [2] Julius A Adebayo et al. 2016. *FairML: ToolBox for diagnosing bias in predictive modeling*. Ph.D. Dissertation. Massachusetts Institute of Technology.
- 525 [3] Alekh Agarwal, Alina Beygelzimer, Miroslav Dudík, John Langford, and Hanna Wallach. 2018. A reductions approach to fair classification. In
526 *International Conference on Machine Learning*. PMLR, 60–69.
- 527 [4] Chirag Agarwal, Himabindu Lakkaraju, and Marinka Zitnik. 2021. Towards a unified framework for fair and stable graph representation learning.
528 In *Uncertainty in Artificial Intelligence*. PMLR, 2114–2124.
- 529 [5] Aniya Aggarwal, Pranay Lohia, Seema Nagar, Kuntal Dey, and Diptikalyan Saha. 2019. Black box fairness testing of machine learning models.
530 In *Proceedings of the 2019 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software
Engineering*. 625–635.
- 531 [6] Aniya Aggarwal, Samiulla Shaikh, Sandeep Hans, Swastik Halder, Rema Ananthanarayanan, and Diptikalyan Saha. 2021. Testing framework for
532 black-box AI models. In *2021 IEEE/ACM 43rd International Conference on Software Engineering: Companion Proceedings (ICSE-Companion)*. IEEE,
533 81–84.
- 534 [7] Aws Albarghouthi, Loris D’Antoni, Samuel Drews, and Aditya V Nori. 2017. Fairsquare: probabilistic verification of program fairness. *Proceedings
535 of the ACM on Programming Languages* 1, OOPSLA (2017), 1–30.
- 536 [8] Aws Albarghouthi, Loris D’Antoni, and Samuel Drews. 2017. Repairing decision-making programs under uncertainty. In *International Conference
537 on Computer Aided Verification*. Springer, 181–200.
- 538 [9] Aws Albarghouthi and Samuel Vinitsky. 2019. Fairness-aware programming. In *Proceedings of the Conference on Fairness, Accountability, and
539 Transparency*. 211–219.
- 540 [10] Jose Manuel Alvarez and Salvatore Ruggieri. 2023. Counterfactual Situation Testing: Uncovering Discrimination under Fairness given the Difference.
541 In *Proceedings of the 3rd ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization (EAAMO ’23)*. Association for Computing
542 Machinery, New York, NY, USA, Article 2, 11 pages. <https://doi.org/10.1145/3617694.3623222>
- 543 [11] Rico Angell, Brittany Johnson, Yuriy Brun, and Alexandra Meliou. 2018. Themis: Automatically testing software for discrimination. In *Proceedings of
544 the 2018 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 871–875.
- 545 [12] Julia Angwin, Jeff Larson, Surya Mattu, and Lauren Kirchner. 2022. Machine bias. In *Ethics of data and analytics*. Auerbach Publications, 254–264.
- 546 [13] Julia Angwin, Jeff Larson, Surya Mattu, and Lauren Kirchner. Accessed:09.11.2023. Machine Bias: There’s software used across the country to
547 predict future criminals. And it’s biased against blacks. <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>
- 548 [14] Muhammad Hilmi Asyrofi, Zhou Yang, Imam Nur Bani Yusuf, Hong Jin Kang, Ferdinand Thung, and David Lo. 2021. Biasfinder: Metamorphic test
549 generation to uncover bias for sentiment analysis systems. *IEEE Transactions on Software Engineering* (2021).
- 550 [15] Fatma Basak Aydemir and Fabiano Dalpiaz. 2018. A roadmap for ethics-aware software engineering. In *2018 IEEE/ACM International Workshop on
551 Software Fairness (FairWare)*. IEEE, 15–21.
- 552 [16] Niels Bantilan. 2018. Themis-ml: A fairness-aware machine learning interface for end-to-end discrimination discovery and mitigation. *Journal of
553 Technology in Human Services* 36, 1 (2018), 15–30.
- 554 [17] Luciano Baresi, Chiara Criscuolo, and Carlo Ghezzi. 2023. Understanding fairness requirements for ml-based software. In *2023 IEEE 31st International
555 Requirements Engineering Conference (RE)*. IEEE, 341–346.
- 556 [18] Osbert Bastani, Xin Zhang, and Armando Solar-Lezama. 2019. Probabilistic verification of fairness properties via concentration. *Proceedings of the
557 ACM on Programming Languages* 3, OOPSLA (2019), 1–27.
- 558 [19] Rachel KE Bellamy, Kuntal Dey, Michael Hind, Samuel C Hoffman, Stephanie Houde, Kalapriya Kannan, Pranay Lohia, Jacquelyn Martino, Sameep
559 Mehta, Aleksandra Mojsilović, et al. 2019. AI Fairness 360: An extensible toolkit for detecting and mitigating algorithmic bias. *IBM Journal of
560 Research and Development* 63, 4/5 (2019), 4–1.
- 561 [20] Rachel KE Bellamy, Kuntal Dey, Michael Hind, Samuel C Hoffman, Stephanie Houde, Kalapriya Kannan, Pranay Lohia, Sameep Mehta, Aleksandra
562 Mojsilovic, Seema Nagar, et al. 2019. Think your artificial intelligence software is fair? Think again. *IEEE Software* 36, 4 (2019), 76–80.
- 563 [21] Nelly Bencomo, Jin LC Guo, Rachel Harrison, Hans-Martin Heyn, and Tim Menzies. 2021. The Secret to Better AI and Better Software (Is
564 Requirements Engineering). *IEEE Software* 39, 1 (2021), 105–110.
- 565 [22] Richard Berk, Hoda Heidari, Shahin Jabbari, Matthew Joseph, Michael Kearns, Jamie Morgenstern, Seth Neel, and Aaron Roth. 2017. A Convex
566 Framework for Fair Regression. *Fairness, Accountability, and Transparency in Machine Learning* (2017).
- 567 [23] Sarah Bird, Miro Dudík, Richard Edgar, Brandon Horn, Roman Lutz, Vanessa Milan, Mehrnoosh Sameki, Hanna Wallach, and Kathleen Walker.
568 2020. Fairlearn: A toolkit for assessing and improving fairness in AI. *Microsoft, Tech. Rep. MSR-TR-2020-32* (2020).
- 569 [24] Sumon Biswas and Hridesh Rajan. 2020. Do the machine learning models on a crowd sourced platform exhibit bias? an empirical study on model
570 fairness. In *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software
Engineering*. 642–653.
- 571 [25] Sumon Biswas and Hridesh Rajan. 2021. Fair Preprocessing: Towards Understanding Compositional Fairness of Data Transformers in Machine
572 Learning Pipeline. *arXiv preprint arXiv:2106.06054* (2021).

- [573] [26] Sumon Biswas and Hridesh Rajan. 2023. Fairify: Fairness verification of neural networks. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*. IEEE, 1546–1558.
- [574] [27] Emily Black, Samuel Yeom, and Matt Fredrikson. 2020. Fliptest: fairness testing via optimal transport. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*. 111–121.
- [575] [28] Su Lin Blodgett, Solon Barocas, Hal Daumé III, and Hanna Wallach. 2020. Language (Technology) is Power: A Critical Survey of “Bias” in NLP. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 5454–5476.
- [576] [29] Su Lin Blodgett, Gilsinia Lopez, Alexandra Olteanu, Robert Sim, and Hanna Wallach. 2021. Stereotyping Norwegian salmon: an inventory of pitfalls in fairness benchmark datasets. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. 1004–1015.
- [577] [30] Martin Brandao. 2019. Age and gender bias in pedestrian detection algorithms. In *Proceedings of the Workshop on Fairness Accountability Transparency and Ethics in Computer Vision at IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*.
- [578] [31] Yuri Brun and Alexandra Meliou. 2018. Software fairness. In *Proceedings of the 2018 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 754–759.
- [579] [32] Margaret Burnett, Simone Stumpf, Jamie Macbeth, Stephan Makri, Laura Beckwith, Irwin Kwan, Anicia Peters, and William Jernigan. 2016. GenderMag: A method for evaluating software’s gender inclusiveness. *Interacting with Computers* 28, 6 (2016), 760–787.
- [580] [33] Ángel Alexander Cabrera, Will Epperson, Fred Hohman, Minsuk Kahng, Jamie Morgenstern, and Duen Horng Chau. 2019. FairVis: Visual analytics for discovering intersectional bias in machine learning. In *2019 IEEE Conference on Visual Analytics Science and Technology (VAST)*. IEEE, 46–56.
- [581] [34] Toon Calders and Sicco Verwer. 2010. Three naive Bayes approaches for discrimination-free classification. *Data mining and knowledge discovery* 21, 2 (2010), 277–292.
- [582] [35] Flavio Calmon, Dennis Wei, Bhanukiran Vinzamuri, Karthikeyan Natesan Ramamurthy, and Kush R Varshney. 2017. Optimized pre-processing for discrimination prevention. *Advances in neural information processing systems* 30 (2017).
- [583] [36] L Elisa Celis, Lingxiao Huang, Vijay Keswani, and Nisheeth K Vishnoi. 2019. Classification with fairness constraints: A meta-algorithm with provable guarantees. In *Proceedings of the conference on fairness, accountability, and transparency*. 319–328.
- [584] [37] Joymallya Chakraborty, Suvodeep Majumder, and Tim Menzies. 2021. Bias in machine learning software: why? how? what to do?. In *Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 429–440.
- [585] [38] Joymallya Chakraborty, Suvodeep Majumder, Zhe Yu, and Tim Menzies. 2020. Fairway: A way to build fair ml software. In *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 654–665.
- [586] [39] Jialuo Chen, Jingyi Wang, Xingjun Ma, Youcheng Sun, Jun Sun, Peixin Zhang, and Peng Cheng. 2023. QuoTe: Quality-Oriented Testing for Deep Learning Systems. *ACM Trans. Softw. Eng. Methodol.* 32, 5, Article 125 (jul 2023), 33 pages. <https://doi.org/10.1145/3582573>
- [587] [40] Zhenpeng Chen, Jie M. Zhang, Federica Sarro, and Mark Harman. 2023. A Comprehensive Empirical Study of Bias Mitigation Methods for Machine Learning Classifiers. *ACM Trans. Softw. Eng. Methodol.* 32, 4, Article 106 (may 2023), 30 pages. <https://doi.org/10.1145/3583561>
- [588] [41] Zhenpeng Chen, Jie M Zhang, Federica Sarro, and Mark Harman. 2024. Fairness improvement with multiple protected attributes: How far are we? (2024), 1–13.
- [589] [42] Jürgen Cito, Isil Dillig, Seohyun Kim, Vijayaraghavan Murali, and Satish Chandra. 2021. Explaining mispredictions of machine learning models using rule induction. In *Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 716–727.
- [590] [43] Sam Corbett-Davies, Emma Pierson, Avi Feller, Sharad Goel, and Aziz Huq. 2017. Algorithmic decision making and the cost of fairness. In *Proceedings of the 23rd acm sigkdd international conference on knowledge discovery and data mining*. 797–806.
- [591] [44] Jenna Cryan, Shiliang Tang, Xinyi Zhang, Miriam Metzger, Haitao Zheng, and Ben Y Zhao. 2020. Detecting gender stereotypes: lexicon vs. supervised learning methods. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*. 1–11.
- [592] [45] Anubrata Das and Matthew Lease. 2019. A Conceptual Framework for Evaluating Fairness in Search. *arXiv preprint arXiv:1907.09328* (2019).
- [593] [46] Terrance De Vries, Ishan Misra, Changhan Wang, and Laurens Van der Maaten. 2019. Does object recognition work for everyone?. In *Proceedings of the Workshop on Fairness Accountability Transparency and Ethics in Computer Vision at IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*. 52–59.
- [594] [47] Emily Denton, Ben Hutchinson, Margaret Mitchell, and Timnit Gebru. 2019. Detecting bias with generative counterfactual face attribute augmentation. (2019).
- [595] [48] Marina Drosou, HV Jagadish, Evangelia Pitoura, and Julia Stoyanovich. 2017. Diversity in big data: A review. *Big data* 5, 2 (2017), 73–84.
- [596] [49] Xiaoning Du, Xiaofei Xie, Yi Li, Lei Ma, Yang Liu, and Jianjun Zhao. 2019. DeepStellar: Model-Based Quantitative Analysis of Stateful Deep Learning Systems. In *Proceedings of the 2019 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering* (Tallinn, Estonia) (ESEC/FSE 2019). Association for Computing Machinery, New York, NY, USA, 477–487. <https://doi.org/10.1145/3338906.3338954>
- [597] [50] Ming Fan, Wenyng Wei, Wuxia Jin, Zijiang Yang, and Ting Liu. 2022. Explanation-Guided Fairness Testing through Genetic Algorithm. In *2022 IEEE/ACM 44th International Conference on Software Engineering (ICSE)*. IEEE.
- [598] [51] Golnoosh Farnad, Behrouz Babaki, and Michel Gendreau. 2020. A unifying framework for fairness-aware influence maximization. In *Companion Proceedings of the Web Conference 2020*. 714–722.

- [52] Michael Feldman, Sorelle A Friedler, John Moeller, Carlos Scheidegger, and Suresh Venkatasubramanian. 2015. Certifying and removing disparate impact. In *proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*. 259–268.
- [53] Yang Feng, Qingkai Shi, Xinyu Gao, Jun Wan, Chunrong Fang, and Zhenyu Chen. 2020. DeepGini: Prioritizing Massive Tests to Enhance the Robustness of Deep Neural Networks. In *Proceedings of the 29th ACM SIGSOFT International Symposium on Software Testing and Analysis (Virtual Event, USA) (ISSTA 2020)*. Association for Computing Machinery, New York, NY, USA, 177–188. <https://doi.org/10.1145/3395363.3397357>
- [54] Carmine Ferrara, Francesco Casillo, Carmine Gravino, Andrea De Lucia, and Fabio Palomba. 2024. Refair: Toward a context-aware recommender for fairness requirements engineering. (2024), 1–12.
- [55] Carmine Ferrara, Giulia Sellitto, Filomena Ferrucci, Fabio Palomba, and Andrea De Lucia. 2023. Fairness-aware machine learning engineering: how far are we? *Empirical Software Engineering* 29, 1 (2023), 9.
- [56] Anjalie Field, Su Lin Blodgett, Zeerak Waseem, and Yulia Tsvetkov. 2021. A Survey of Race, Racism, and Anti-Racism in NLP. *arXiv preprint arXiv:2106.11410* (2021).
- [57] Anthony Finkelstein, Mark Harman, S Afshin Mansouri, Jian Ren, and Yuanyuan Zhang. 2009. A search based approach to fairness analysis in requirement assignments to aid negotiation, mediation and decision making. *Requirements engineering* 14, 4 (2009), 231–245.
- [58] Jessie Finocchiaro, Roland Maio, Faidra Monachou, Gourab K Patro, Manish Raghavan, Ana-Andreea Stoica, and Stratis Tsirtsis. 2021. Bridging machine learning and mechanism design towards algorithmic fairness. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*. 489–503.
- [59] Claudia Flores-Saviaga, Christopher Curtis, and Saiph Savage. 2023. Inclusive Portraits: Race-Aware Human-in-the-Loop Technology. In *Proceedings of the 3rd ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization (EAAMO '23)*. Association for Computing Machinery, New York, NY, USA, Article 15, 11 pages. <https://doi.org/10.1145/3617694.3623235>
- [60] Aimen Gaba, Zhanna Kaufman, Jason Cheung, Marie Shvakel, Kyle Wm Hall, Yuriy Brun, and Cindy Xiong Bearfield. 2023. My Model is Unfair, Do People Even Care? Visual Design Affects Trust and Perceived Bias in Machine Learning. *IEEE Transactions on Visualization and Computer Graphics* (2023).
- [61] Sainyam Galhotra, Yuriy Brun, and Alexandra Meliou. 2017. Fairness testing: testing software for discrimination. In *Proceedings of the 2017 11th Joint Meeting on Foundations of Software Engineering*. 498–510.
- [62] Xuanqi Gao, Juan Zhai, Shiqing Ma, Chao Shen, Yufei Chen, and Qian Wang. 2022. FairNeuron: Improving Deep Neural Network Fairness with Adversary Games on Selective Neurons. In *2022 IEEE/ACM 44th International Conference on Software Engineering (ICSE)*. IEEE.
- [63] Xuanqi Gao, Juan Zhai, Shiqing Ma, Chao Shen, Yufei Chen, and Shiwei Wang. 2023. CILIATE: Towards Fairer Class-based Incremental Learning by Dataset and Training Refinement. *arXiv preprint arXiv:2304.04222* (2023).
- [64] Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé III, and Kate Crawford. 2021. Datasheets for datasets. *Commun. ACM* 64, 12 (2021), 86–92.
- [65] Bishwamitra Ghosh, Debabrota Basu, and Kuldeep S Meel. 2021. Justicia: A Stochastic SAT Approach to Formally Verify Fairness. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 7554–7563.
- [66] Stephen Giguerre, Blossom Metevier, Yuriy Brun, Bruno Castro Da Silva, Philip S Thomas, and Scott Niekum. 2022. Fairness guarantees under demographic shift. In *Proceedings of the 10th International Conference on Learning Representations (ICLR)*.
- [67] Usman Gohar, Sumon Biswas, and Hridesh Rajan. 2023. Towards understanding fairness and its composition in ensemble machine learning. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*. IEEE, 1533–1545.
- [68] Usman Gohar and Lu Cheng. 2023. A Survey on Intersectional Fairness in Machine Learning: Notions, Mitigation, and Challenges. *arXiv preprint arXiv:2305.06969* (2023).
- [69] CV González Zelaya, P Missier, and D Prangle. 2019. Parametrised data sampling for fairness optimisation. In *2019 XAI Workshop at SIGKDD, Anchorage, AK, USA*.
- [70] Nina Grgić-Hlaca, Elissa M Redmiles, Krishna P Gummadi, and Adrian Weller. 2018. Human perceptions of fairness in algorithmic decision making: A case study of criminal risk prediction. In *Proceedings of the 2018 world wide web conference*. 903–912.
- [71] Nina Grgić-Hlaca, Gabriel Lima, Adrian Weller, and Elissa M. Redmiles. 2022. Dimensions of Diversity in Human Perceptions of Algorithmic Fairness. In *Proceedings of the 2nd ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization (EAAMO '22)*. Association for Computing Machinery, New York, NY, USA, Article 21, 12 pages. <https://doi.org/10.1145/3551624.3555306>
- [72] Emītā Guzmnān, Ricarda Anna-Lena Fischer, and Janey Kok. 2023. Mind the gap: gender, micro-inequities and barriers in software development. *Empirical Software Engineering* 29, 1 (2023), 17.
- [73] 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*. 392–402.
- [74] Austin Hoag, James E. Kostas, Bruno Castro da Silva, Philip S. Thomas, and Yuriy Brun. 2023. Seldonian Toolkit: Building Software with Safe and Fair Machine Learning. In *2023 IEEE/ACM 45th International Conference on Software Engineering: Companion Proceedings (ICSE-Companion)*. 107–111. <https://doi.org/10.1109/ICSE-Companion58688.2023.00035>
- [75] Kenneth Holstein, Jennifer Wortman Vaughan, Hal Daumé III, Miro Dudik, and Hanna Wallach. 2019. Improving fairness in machine learning systems: What do industry practitioners need?. In *Proceedings of the 2019 CHI conference on human factors in computing systems*. 1–16.
- [76] Sara Hooker, Nyalleng Moorosi, Gregory Clark, Samy Bengio, and Emily Denton. 2020. Characterising bias in compressed models. *arXiv preprint arXiv:2010.03058* (2020).

- [677] Max Hort, Rebecca Moussa, and Federica Sarro. 2023. Multi-objective search for gender-fair and semantically correct word embeddings. *Applied Soft Computing* 133 (2023), 109916.
- [678] Max Hort, Jie M Zhang, Federica Sarro, and Mark Harman. 2021. Fairea: a model behaviour mutation approach to benchmarking bias mitigation methods. In *Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 994–1006.
- [679] Max Hort, Jie M. Zhang, Federica Sarro, and Mark Harman. 2024. Search-based Automatic Repair for Fairness and Accuracy in Decision-making Software. *Empirical Software Engineering* 29, 1 (2024), 36. <https://doi.org/10.1007/s10664-023-10419-3>
- [680] Waqar Hussain, Davoud Mougouei, and Jon Whittle. 2018. Integrating social values into software design patterns. In *2018 IEEE/ACM International Workshop on Software Fairness (FairWare)*. IEEE, 8–14.
- [681] Wiebke (Toussaint) Hutiri, Aaron Yi Ding, Fahim Kawsar, and Akhil Mathur. 2023. Tiny, Always-on, and Fragile: Bias Propagation through Design Choices in On-Device Machine Learning Workflows. *ACM Trans. Softw. Eng. Methodol.* 32, 6, Article 155 (sep 2023), 37 pages. <https://doi.org/10.1145/3591867>
- [682] HV Jagadish, Julia Stoyanovich, and Bill Howe. 2021. Covid-19 brings data equity challenges to the fore. *Digital Government: Research and Practice* 2, 2 (2021), 1–7.
- [683] Philips George John, Deepak Vijayakeerthy, and Diptikalyan Saha. 2020. Verifying individual fairness in machine learning models. In *Conference on Uncertainty in Artificial Intelligence*. PMLR, 749–758.
- [684] Brittany Johnson, Jesse Bartola, Rico Angell, Sam Witty, Stephen Giguere, and Yuriy Brun. 2023. Fairkit, fairkit, on the wall, who's the fairest of them all? Supporting fairness-related decision-making. *EURO Journal on Decision Processes* 11 (2023), 100031. <https://doi.org/10.1016/j.ejdp.2023.100031>
- [685] Brittany Johnson and Yuriy Brun. 2022. Fairkit-Learn: A Fairness Evaluation and Comparison Toolkit. In *Proceedings of the ACM/IEEE 44th International Conference on Software Engineering: Companion Proceedings* (Pittsburgh, Pennsylvania) (ICSE '22). Association for Computing Machinery, New York, NY, USA, 70–74. <https://doi.org/10.1145/3510454.3516830>
- [686] Faisal Kamiran and Toon Calders. 2009. Classifying without discriminating. In *2009 2nd international conference on computer, control and communication*. IEEE, 1–6.
- [687] Faisal Kamiran and Toon Calders. 2012. Data preprocessing techniques for classification without discrimination. *Knowledge and information systems* 33, 1 (2012), 1–33.
- [688] Faisal Kamiran, Toon Calders, and Mykola Pechenizkiy. 2010. Discrimination aware decision tree learning. In *2010 IEEE International Conference on Data Mining*. IEEE, 869–874.
- [689] Faisal Kamiran, Asim Karim, and Xiangliang Zhang. 2012. Decision theory for discrimination-aware classification. In *2012 IEEE 12th International Conference on Data Mining*. IEEE, 924–929.
- [690] Faisal Kamiran, Sameen Mansha, Asim Karim, and Xiangliang Zhang. 2018. Exploiting reject option in classification for social discrimination control. *Information Sciences* 425 (2018), 18–33.
- [691] Toshihiro Kamishima, Shotaro Akaho, Hideki Asoh, and Jun Sakuma. 2012. Fairness-aware classifier with prejudice remover regularizer. In *Joint European conference on machine learning and knowledge discovery in databases*. Springer, 35–50.
- [692] Michael Kearns, Seth Neel, Aaron Roth, and Zhiwei Steven Wu. 2018. Preventing fairness gerrymandering: Auditing and learning for subgroup fairness. In *International Conference on Machine Learning*. PMLR, 2564–2572.
- [693] Michael Kearns, Seth Neel, Aaron Roth, and Zhiwei Steven Wu. 2019. An empirical study of rich subgroup fairness for machine learning. In *Proceedings of the conference on fairness, accountability, and transparency*. 100–109.
- [694] Byungju Kim, Hyunwoo Kim, Kyungsu Kim, Sungjin Kim, and Junmo Kim. 2019. Learning not to learn: Training deep neural networks with biased data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 9012–9020.
- [695] Jinhan Kim, Robert Feldt, and Shin Yoo. 2019. Guiding deep learning system testing using surprise adequacy. In *2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE)*. IEEE, 1039–1049.
- [696] Michael P Kim, Amirata Ghorbani, and James Zou. 2019. Multiaccuracy: Black-box post-processing for fairness in classification. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*. 247–254.
- [697] Caitlin Kuhlman, MaryAnn VanValkenburg, and Elke Rundensteiner. 2019. Fare: Diagnostics for fair ranking using pairwise error metrics. In *The World Wide Web Conference*. 2936–2942.
- [698] Michelle Seng Ah Lee and Jat Singh. 2021. The landscape and gaps in open source fairness toolkits. In *Proceedings of the 2021 CHI conference on human factors in computing systems*. 1–13.
- [699] Chenglu Li, Wanli Xing, and Walter Leite. 2021. Yet Another Predictive Model? Fair Predictions of Students' Learning Outcomes in an Online Math Learning Platform. In *LAK21: 11th International Learning Analytics and Knowledge Conference*. 572–578.
- [700] Xinyue Li, Zhenpeng Chen, Jie M Zhang, Federica Sarro, Ying Zhang, and Xuanzhe Liu. 2023. Dark-skin individuals are at more risk on the street: Unmasking fairness issues of autonomous driving systems. *arXiv preprint arXiv:2308.02935* (2023).
- [701] Yanhui Li, Linghan Meng, Lin Chen, Li Yu, Di Wu, Yuming Zhou, and Baowen Xu. 2022. Training Data Debugging for the Fairness of Machine Learning Software. In *2022 IEEE/ACM 44th International Conference on Software Engineering (ICSE)*. IEEE.
- [702] Lizhen Liang and Daniel E Acuna. 2020. Artificial mental phenomena: Psychophysics as a framework to detect perception biases in AI models. In *Proceedings of the 2020 conference on fairness, accountability, and transparency*. 403–412.

- [103] Zachary Lipton, Julian McAuley, and Alexandra Chouldechova. 2018. Does mitigating ML’s impact disparity require treatment disparity? *Advances in neural information processing systems* 31 (2018).
- [104] Kristian Lum and Tarak Shah. 2019. Measures of fairness for New York City’s Supervised Release Risk Assessment Tool. *Human Rights Data Analytics Group* (2019), 21.
- [105] Lei Ma, Fuyuan Zhang, Jiyuan Sun, Minhui Xue, Bo Li, Felix Juefei-Xu, Chao Xie, Li Li, Yang Liu, Jianjun Zhao, et al. 2018. Deepmutation: Mutation testing of deep learning systems. In *2018 IEEE 29th International Symposium on Software Reliability Engineering (ISSRE)*. IEEE, 100–111.
- [106] Pingchuan Ma, Shuai Wang, and Jin Liu. 2020. Metamorphic Testing and Certified Mitigation of Fairness Violations in NLP Models.. In *IJCAI*. 458–465.
- [107] Wei Ma, Mike Papadakis, Anestis Tsakmalis, Maxime Cordy, and Yves Le Traon. 2021. Test Selection for Deep Learning Systems. *ACM Trans. Softw. Eng. Methodol.* 30, 2, Article 13 (Jan. 2021). 22 pages. <https://doi.org/10.1145/3417330>
- [108] Suvodeep Majumder, Joymallya Chakraborty, Gina R. Bai, Kathryn T. Stolee, and Tim Menzies. 2023. Fair Enough: Searching for Sufficient Measures of Fairness. *ACM Trans. Softw. Eng. Methodol.* 32, 6, Article 134 (sep 2023). 22 pages. <https://doi.org/10.1145/3585006>
- [109] Blossom Metevier, Stephen Giguere, Sarah Brockman, Ari Kobren, Yuriy Brun, Emma Brunskill, and Philip Thomas. 2019. Offline contextual bandits with high probability fairness guarantees. *Advances in neural information processing systems* 32 (2019).
- [110] Verya Monjezi, Ashutosh Trivedi, Gang Tan, and Saeid Tizpaz-Niari. 2023. Information-Theoretic Testing and Debugging of Fairness Defects in Deep Neural Networks. *arXiv preprint arXiv:2304.04199* (2023).
- [111] Daniel Perez Morales, Takashi Kitamura, and Shingo Takada. 2021. Coverage-Guided Fairness Testing. In *International Conference on Intelligence Science*. Springer, 183–199.
- [112] Aniruddhan Murali, Gaurav Sahu, Kishanthan Thangarajah, Brian Zimmerman, Gema Rodriguez-Pérez, and Meiyappan Nagappan. 2024. Diversity in issue assignment: humans vs bots. *Empirical Software Engineering* 29, 2 (2024), 37.
- [113] Giang Nguyen, Sumon Biswas, and Hradesh Rajan. 2023. Fix Fairness, Don’t Ruin Accuracy: Performance Aware Fairness Repair using AutoML. *arXiv preprint arXiv:2306.09297* (2023).
- [114] Julian Nyarko, Sharad Goel, and Roseanna Sommers. 2021. Breaking Taboos in Fair Machine Learning: An Experimental Study. In *Proceedings of the 1st ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization* (–, NY, USA) (EAAMO ’21). Association for Computing Machinery, New York, NY, USA, Article 14, 11 pages. <https://doi.org/10.1145/3465416.3483291>
- [115] U.S. Department of Labor. Accessed: 25.04.2022. Affirmative Action. <https://www.dol.gov/general/topic/hiring/affirmativeact>
- [116] Moses Openja, Gabriel Laberge, and Foutse Khomh. 2023. Detection and evaluation of bias-inducing features in machine learning. *Empirical Software Engineering* 29, 1 (2023), 22.
- [117] Kexin Pei, Yinzhai Cao, Junfeng Yang, and Suman Jana. 2017. Deepxplore: Automated whitebox testing of deep learning systems. In *proceedings of the 26th Symposium on Operating Systems Principles*. 1–18.
- [118] Kewen Peng, Joymallya Chakraborty, and Tim Menzies. 2022. FairMask: Better fairness via model-based rebalancing of protected attributes. *IEEE Transactions on Software Engineering* 49, 4 (2022), 2426–2439.
- [119] Anjana Perera, Aldeida Aleti, Chakkrit Tantithamthavorn, Jirayus Jiarpakdee, Burak Turhan, Lisa Kuhn, and Katie Walker. 2022. Search-based fairness testing for regression-based machine learning systems. *Empirical Software Engineering* 27, 3 (2022), 79.
- [120] Nga Pham, Hung Pham-Ngoc, and Anh Nguyen-Duc. 2023. Fairness Requirement in AI Engineering—A Review on Current Research and Future Directions. In *International Conference on Sustainability in Software Engineering & Business Information Management*. Springer, 3–13.
- [121] Geoff Pleiss, Manish Raghavan, Felix Wu, Jon Kleinberg, and Kilian Q Weinberger. 2017. On fairness and calibration. *Advances in neural information processing systems* 30 (2017).
- [122] Sai Sathiesh Rajan, Sakshi Udeshi, and Sudipta Chattopadhyay. 2022. AequiVox: Automated Fairness Testing of Speech Recognition Systems. (2022).
- [123] Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. 2020. Beyond Accuracy: Behavioral Testing of NLP Models with CheckList. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 4902–4912.
- [124] Luke Rodriguez, Babak Salimi, Haoyue Ping, Julia Stoyanovich, and Bill Howe. 2018. MobilityMirror: Bias-adjusted transportation datasets. In *Workshop on Big Social Data and Urban Computing*. Springer, 18–39.
- [125] Debjani Saha, Candice Schumann, Duncan C McElfresh, John P Dickerson, Michelle L Mazurek, and Michael Carl Tschantz. 2020. Human comprehension of fairness in machine learning. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*. 152–152.
- [126] Pedro Saleiro, Benedict Kuester, Loren Hinkson, Jesse London, Abby Stevens, Ari Anisfeld, Kit T Rodolfa, and Rayid Ghani. 2018. Aequitas: A bias and fairness audit toolkit. *arXiv preprint arXiv:1811.05577* (2018).
- [127] Babak Salimi, Bill Howe, and Dan Suciu. 2019. Data management for causal algorithmic fairness. *arXiv preprint arXiv:1908.07924* (2019).
- [128] Babak Salimi, Bill Howe, and Dan Suciu. 2020. Database repair meets algorithmic fairness. *ACM SIGMOD Record* 49, 1 (2020), 34–41.
- [129] Maarten Sap, Saadia Gabriel, Lianhui Qin, Dan Jurafsky, Noah A Smith, and Yejin Choi. 2020. Social Bias Frames: Reasoning about Social and Power Implications of Language. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 5477–5490.
- [130] Prasanna Sattigeri, Samuel C Hoffman, Vijil Chenthamarakshan, and Kush R Varshney. 2019. Fairness GAN: Generating datasets with fairness properties using a generative adversarial network. *IBM Journal of Research and Development* 63, 4/5 (2019), 3–1.
- [131] Daniel Schvarcz. 2021. Health-Based Proxy Discrimination, Artificial Intelligence, and Big Data. *Artificial Intelligence, and Big Data (March 3, 2021)*. *Houston Journal of Health Law and Policy* (2021).

- [781] Shahar Segal, Yossi Adi, Benny Pinkas, Carsten Baum, Chaya Ganesh, and Joseph Keshet. 2021. Fairness in the eyes of the data: Certifying machine-learning models. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*. 926–935.
- [782] Arnab Sharma, Caglar Demir, Axel-Cyrille Ngonga Ngomo, and Heike Wehrheim. 2021. MLCheck-Property-Driven Testing of Machine Learning Models. *arXiv preprint arXiv:2105.00741* (2021).
- [783] Shubham Sharma, Yunfeng Zhang, Jesús M Ríos Aliaga, Djallel Bouneffouf, Vinod Muthusamy, and Kush R Varshney. 2020. Data augmentation for discrimination prevention and bias disambiguation. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*. 358–364.
- [784] Weijun Shen, Yanhui Li, Lin Chen, Yuanlei Han, Yuming Zhou, and Baowen Xu. 2020. Multiple-Boundary Clustering and Prioritization to Promote Neural Network Retraining. In *2020 35th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. 410–422.
- [785] Kacper Sokol, Alexander Hepburn, Rafael Poyiadzi, Matthew Clifford, Raul Santos-Rodriguez, and Peter Flach. 2020. FAT forensics: a python toolbox for implementing and deploying fairness, accountability and transparency algorithms in predictive systems. *Journal of Open Source Software* 5, 49 (2020), 1904.
- [786] Ezekiel Soremekun, Sakshi Udeshi, and Sudipta Chattopadhyay. 2022. Astraea: Grammar-based fairness testing. *IEEE Transactions on Software Engineering* (2022).
- [787] Julia Stoyanovich. 2019. TransFAT: Translating fairness, accountability and transparency into data science practice. In *CEUR Workshop Proceedings*, Vol. 2417. CEUR-WS.
- [788] Julia Stoyanovich, Serge Abiteboul, and Gerome Miklau. 2016. Data, responsibly: Fairness, neutrality and transparency in data analysis. In *International Conference on Extending Database Technology*.
- [789] Zeyu Sun, Jie M Zhang, Mark Harman, Mike Papadakis, and Lu Zhang. 2020. Automatic testing and improvement of machine translation. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 974–985.
- [790] Zeyu Sun, Jie M Zhang, Yingfei Xiong, Mark Harman, Mike Papadakis, and Lu Zhang. 2022. Improving Machine Translation Systems via Isotopic Replacement. In *2022 IEEE/ACM 44th International Conference on Software Engineering (ICSE)*. IEEE.
- [791] Guanhong Tao, Weisong Sun, Tingxu Han, Chunrong Fang, and Xiangyu Zhang. 2022. RULER: discriminative and iterative adversarial training for deep neural network fairness. In *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 1173–1184.
- [792] Philip S. Thomas, Bruno Castro da Silva, Andrew G. Barto, Stephen Giguere, Yuriy Brun, and Emma Brunskill. 2019. Preventing undesirable behavior of intelligent machines. *Science* 366, 6468 (2019), 999–1004. <https://doi.org/10.1126/science.aag3311> arXiv:<https://www.science.org/doi/pdf/10.1126/science.aag3311>
- [793] Yuchi Tian, Ziyuan Zhong, Vicente Ordonez, Gail Kaiser, and Baishakhi Ray. 2020. Testing DNN image classifiers for confusion & bias errors. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 1122–1134.
- [794] Saeid Tizpaz-Niari, Ashish Kumar, Gang Tan, and Ashutosh Trivedi. 2022. Fairness-aware Configuration of Machine Learning Libraries. In *2022 IEEE/ACM 44th International Conference on Software Engineering (ICSE)*. IEEE.
- [795] Florian Tramer, Vaggelis Atlidakis, Roxana Geambasu, Daniel Hsu, Jean-Pierre Hubaux, Mathias Humbert, Ari Juels, and Huang Lin. 2017. Fairtest: Discovering unwarranted associations in data-driven applications. In *2017 IEEE European Symposium on Security and Privacy (EuroS&P)*. IEEE, 401–416.
- [796] Sakshi Udeshi, Pryanshu Arora, and Sudipta Chattopadhyay. 2018. Automated directed fairness testing. In *Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering*. 98–108.
- [797] Sahil Verma and Julia Rubin. 2018. Fairness definitions explained. In *2018 ieee/acm international workshop on software fairness (fairware)*. IEEE, 1–7.
- [798] Yuxuan Wan, Wenxuan Wang, Pinjia He, Jiazheng Gu, Haonan Bai, and Michael R Lyu. 2023. Biasasker: Measuring the bias in conversational ai system. In *Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 515–527.
- [799] Tianlu Wang, Jieyu Zhao, Mark Yatskar, Kai-Wei Chang, and Vicente Ordonez. 2019. Balanced datasets are not enough: Estimating and mitigating gender bias in deep image representations. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5310–5319.
- [800] Zeyu Wang, Clint Qinami, Ioannis Christos Karakozis, Kyle Genova, Prem Nair, Kenji Hata, and Olga Russakovsky. 2020. Towards fairness in visual recognition: Effective strategies for bias mitigation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 8919–8928.
- [801] Nimmi Rashinika Weeraddana, Xiaoyan Xu, Mahmoud Alfadel, Shane McIntosh, and Meyappan Nagappan. 2023. An empirical comparison of ethnic and gender diversity of DevOps and non-DevOps contributions to open-source projects. *Empirical Software Engineering* 28, 6 (2023), 150.
- [802] Benjamin Wilson, Judy Hoffman, and Jamie Morgenstern. 2019. Predictive inequity in object detection. (2019).
- [803] Yongkai Wu, Lu Zhang, Xintao Wu, and Hanghang Tong. 2019. P_c-fairness: A unified framework for measuring causality-based fairness. *Advances in Neural Information Processing Systems* 32 (2019).
- [804] Yisong Xiao, Aishan Liu, Tianlin Li, and Xianglong Liu. 2023. Latent Imitator: Generating Natural Individual Discriminatory Instances for Black-Box Fairness Testing. *arXiv preprint arXiv:2305.11602* (2023).
- [805] Zhou Yang, Muhammad Hilmi Asyrofi, and David Lo. 2021. BiasRV: Uncovering biased sentiment predictions at runtime. In *Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 1540–1544.
- [806] Jun Yu, Xinlong Hao, Haonian Xie, and Yu Yu. 2020. Fair face recognition using data balancing, enhancement and fusion. In *European Conference on Computer Vision*. Springer, 492–505.

- 833 [158] Rich Zemel, Yu Wu, Kevin Swersky, Toni Pitassi, and Cynthia Dwork. 2013. Learning fair representations. In *International conference on machine*
834 *learning*. PMLR, 325–333.
- 835 [159] Jie M Zhang and Mark Harman. 2021. “Ignorance and Prejudice” in Software Fairness. In *2021 IEEE/ACM 43rd International Conference on Software*
836 *Engineering (ICSE)*. IEEE, 1436–1447.
- 837 [160] Lingfeng Zhang, Yueling Zhang, and Min Zhang. 2021. Efficient white-box fairness testing through gradient search. In *Proceedings of the 30th ACM*
838 *SIGSOFT International Symposium on Software Testing and Analysis*. 103–114.
- 839 [161] Mengdi Zhang and Jun Sun. 2022. Adaptive fairness improvement based on causality analysis. In *Proceedings of the 30th ACM Joint European*
840 *Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 6–17.
- 841 [162] Mengdi Zhang, Jun Sun, Jingyi Wang, and Bing Sun. 2023. TestSGD: Interpretable Testing of Neural Networks against Subtle Group Discrimination.
842 *ACM Trans. Softw. Eng. Methodol.* 32, 6, Article 137 (sep 2023), 24 pages. <https://doi.org/10.1145/3591869>
- 843 [163] Peixin Zhang, Jingyi Wang, Jun Sun, Guoliang Dong, Xinyu Wang, Xingen Wang, Jin Song Dong, and Ting Dai. 2020. White-box fairness testing
844 through adversarial sampling. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 949–960.
- 845 [164] Peixin Zhang, Jingyi Wang, Jun Sun, Xinyu Wang, Guoliang Dong, Xingen Wang, Ting Dai, and Jin Song Dong. 2021. Automatic Fairness Testing
846 of Neural Classifiers through Adversarial Sampling. *IEEE Transactions on Software Engineering* (2021).
- 847 [165] Xueru Zhang and Mingyan Liu. 2021. Fairness in learning-based sequential decision algorithms: A survey. In *Handbook of Reinforcement Learning*
848 and Control. Springer, 525–555.
- 849 [166] Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2017. Men also like shopping: Reducing gender bias amplification
850 using corpus-level constraints. *arXiv preprint arXiv:1707.09457* (2017).
- 851 [167] Haibin Zheng, Zhiqing Chen, Tianyu Du, Xuhong Zhang, Yao Cheng, Shouling Ji, Jingyi Wang, Yue Yu, and Jinyin Chen. 2022. NeuronFair:
852 Interpretable White-Box Fairness Testing through Biased Neuron Identification. In *2022 IEEE/ACM 44th International Conference on Software*
853 *Engineering (ICSE)*. IEEE.
- 854 [168] W. Zheng, L. Lin, X. Wu, and X. Chen. 5555. An Empirical Study on Correlations between Deep Neural Network Fairness and Neuron Coverage
855 Criteria. *IEEE Transactions on Software Engineering* 01 (jan 5555), 1–22. <https://doi.org/10.1109/TSE.2023.3349001>
- 856 [169] Indre Žliobaite, Faisal Kamiran, and Toon Calders. 2011. Handling conditional discrimination. In *2011 IEEE 11th International Conference on Data*
857 *Mining*. IEEE, 992–1001.
- 858
- 859
- 860
- 861
- 862
- 863
- 864
- 865
- 866
- 867
- 868
- 869
- 870
- 871
- 872
- 873
- 874
- 875
- 876
- 877
- 878
- 879
- 880
- 881
- 882
- 883
- 884