

Beyond Accuracy: A Critical Review of Fairness in Machine Learning for Mobile and Wearable Computing

SOFIA YFANTIDOU*, Aristotle University of Thessaloniki, Greece

MARIOS CONSTANTINIDES†, Nokia Bell Labs, United Kingdom

DIMITRIS SPATHIS†, Nokia Bell Labs, United Kingdom

ATHENA VAKALI, Aristotle University of Thessaloniki, Greece

DANIELE QUERCIA, Nokia Bell Labs, United Kingdom

FAHIM KAWSAR, Nokia Bell Labs, United Kingdom

The field of mobile, wearable, and ubiquitous computing (UbiComp) is undergoing a revolutionary integration of machine learning. Devices can now diagnose diseases, predict heart irregularities, and unlock the full potential of human cognition. However, the underlying algorithms are not immune to biases with respect to sensitive attributes (e.g., gender, race), leading to discriminatory outcomes. The research communities of HCI and AI-Ethics have recently started to explore ways of reporting information about datasets to surface and, eventually, counter those biases. The goal of this work is to explore the extent to which the UbiComp community has adopted such ways of reporting and highlight potential shortcomings. Through a systematic review of papers published in the Proceedings of the ACM Interactive, Mobile, Wearable and Ubiquitous Technologies (IMWUT) journal over the past 5 years (2018-2022), we found that progress on algorithmic fairness within the UbiComp community lags behind. Our findings show that only a small portion (5%) of published papers adheres to modern fairness reporting, while the overwhelming majority thereof focuses on accuracy or error metrics. In light of these findings, our work provides practical guidelines for the design and development of ubiquitous technologies that not only strive for accuracy but also for fairness.

CCS Concepts: • **Human-centered computing** → **Ubiquitous and mobile computing**; • **Applied computing** → *Consumer health*; • **Computing methodologies** → *Artificial intelligence*; • **Social and professional topics** → *Codes of ethics*.

Additional Key Words and Phrases: literature review, survey, machine learning, bias, fairness, responsible artificial intelligence, ubiquitous computing, sensing data

ACM Reference Format:

Sofia Yfantidou, Marios Constantinides, Dimitris Spathis, Athena Vakali, Daniele Quercia, and Fahim Kawsar. 2023. *Beyond Accuracy: A Critical Review of Fairness in Machine Learning for Mobile and Wearable Computing*. 1, 1 (March 2023), 31 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

Tasks once thought impossible or reserved exclusively for humans are now within our grasp thanks to the integration of Machine Learning (ML) in ubiquitous computing (UbiComp). Algorithms deployed on ubiquitous

*Work done at Nokia Bell Labs

† Also affiliated with the University of Cambridge, United Kingdom

Authors' addresses: Sofia Yfantidou, syfantid@csd.auth.gr, Aristotle University of Thessaloniki, Thessaloniki, Greece; Marios Constantinides, marios.constantinides@nokia-bell-labs.com, Nokia Bell Labs, Cambridge, United Kingdom; Dimitris Spathis, dimitrios.spathis@nokia-bell-labs.com, Nokia Bell Labs, Cambridge, United Kingdom; Athena Vakali, avakali@csd.auth.gr, Aristotle University of Thessaloniki, Thessaloniki, Greece; Daniele Quercia, daniele.quercia@nokia-bell-labs.com, Nokia Bell Labs, Cambridge, United Kingdom; Fahim Kawsar, fahim.kawsar@nokia-bell-labs.com, Nokia Bell Labs, Cambridge, United Kingdom.

devices were typically used to recognize human activities [45], facilitate indoor localization [137], detect breathing phases [120] and infer sleep quality [63]. Today, we are also witnessing an increasing trend toward high-stakes applications. For example, detecting Atrial Fibrillation (AFib) [82], diagnosing COVID-19 infection [11], predicting fertility windows [84], and even improving cognitive performance [25]. Independence of healthcare access, individualized health-promoting interventions, and easier dissemination of medical information, to name a few, make up the list of benefits that algorithmic decision-making (i.e., the use of algorithms and mathematical models to automate the process of decision-making in an efficient and objective manner) has enabled [86]. However, the data and algorithms powering these advancements are not immune to biases. With great ethical opportunities come ethical risks, and, similarly to humans, ML algorithms are susceptible to biases rendering their decisions “unfair” [4, 10, 13, 105].

The research community of Fairness, Accountability and Transparency in ML (FAccT, formerly FAT/ML) defines fairness as a principle that “ensures that algorithmic decisions do not create discriminatory or unjust impacts when comparing across different demographics (e.g., race, sex)” [6]. Real-world cases of “unfair” ML algorithms abound. For example, Kamulegeya et al. [56] found that neural network algorithms trained to perform skin lesion classification showed approximately half the original diagnostic accuracy on black patients compared to white patients. At the same time, people of color are consistently misclassified by health sensors such as oximeters as they were scientifically tested on predominantly white populations [121].

As we shall see throughout this critical review, fairness in UbiComp remains relatively unexplored due to the primary focus on accuracy, and particularities of the community. But what makes UbiComp unique compared to other communities? UbiComp data and models have certain particularities, oftentimes not shared with the broader scholarly discourse on ML and AI Ethics (Figure 1). For example, UbiComp typically deals with small-scale studies, proof-of-concept datasets often collected by the authors in-the-lab or in-the-wild, while the broader ML community frequently utilizes popular, medium- to large-scale benchmark datasets such as the UCI Adult, the German Credit, the COMPAS, and the Diversity in faces datasets [29, 62, 67, 92]. Such data are collected once and are immutable, opposite to UbiComp data that are mutable and, by definition, continuously collected. Contrary to the tabular format of such datasets, UbiComp data are mostly sequential in nature, with biases being harder to surface. In other words, while it is relatively straightforward to distinguish a person’s skin tone from a picture, it is much harder to do so from oximetry measurements, necessitating the collection of supplementary metadata, such as demographics; as UbiComp strives to blend technologies in the background, biases are blended, too. However, with a conscious approach, it is possible to create ML models that are both accurate and fair. As the field of ML continues to evolve, the UbiComp community needs to stay vigilant, ensuring that UbiComp technologies are designed and deployed in a responsible and ethical manner.

Building on the footsteps of Human-Computer Interaction (HCI) and FAccT communities, we set out to understand how fairness has been discussed in UbiComp and identify pathways for ensuring that UbiComp technologies do not cause any harm or infringe on any individual rights [24]. Such research communities recently started to explore ways of reporting fairness in data and models to surface and, eventually, counter biases. In this work, we intend to spark a discussion on how the UbiComp community defines, measures, and assesses fairness. While the community has perhaps indirectly adopted (and adapted) the meaning of fairness to capture UbiComp’s particularities, the question remains though: ***Where does UbiComp fairness overlap with other communities, and, more importantly, where does it lag behind?***

To answer this question, we performed a literature review spanning five years (2018-2022) and 523 papers published in UbiComp literature. We targeted papers published in the Proceedings of the ACM Interactive, Mobile, Wearable and Ubiquitous Technologies (IMWUT), a high-quality journal series capturing the emerging trends in

the UbiComp community and bearing an h-index of 58, placing it among the top-3 publications in HCI.¹ In so doing, we made three contributions:

- We conducted the first review of fairness in ML for UbiComp, where we screened 523 and critically reviewed 49 IMWUT papers ($N_{included} = 49$) published at IMWUT between 2018 and 2022 (§4).
- We found that: *a*) Only 5% of all IMWUT papers reported fairness assessments (included papers); *b*) from this proportion of papers, 24% implement fairness enhancement mechanisms, making up only a limited fraction of all IMWUT papers (1%); *c*) Yet, we surfaced biases across several sensitive attributes, otherwise scattered in UbiComp literature; *d*) Included papers predominantly used performance evaluation metrics, rather than fairness ones, in their fairness discourse. Yet, we confirmed assessment gaps in regression and multi-class classification cases; *e*) Similar to other communities, defining fairness in UbiComp was far from straightforward. Ethical risks and opportunities laid at the heart of this decision; *f*) Fairness in UbiComp was often viewed through the lens of generalizability, taking the form of ablation studies, in-the-wild deployments, and personalization; *g*) Measurement inaccuracies and concept drift in audio, video, image, and sensor data led to performance differences across demographics; and *h*) IMWUT papers suffered from a lack of diverse datasets, concealing biases in the absence of heterogeneous demographics (§5).
- In light of these findings, we made ten recommendations to the UbiComp community pertaining to the integration of fairness into the entire ML pipeline of UbiComp studies (§6).


The remainder of the paper is organized as follows. Section 2 examines review and position papers related to machine learning fairness and responsible artificial intelligence. Section 3 outlines fairness definitions and metrics. Section 4 describes the methodology used to conduct the literature review, and Section 5 presents the results obtained. Lastly, Section 6 discusses the findings and limitations of the review, and offers recommendations for fair reporting in the UbiComp community.

2 RELATED WORK

Next, we situate our critical review in previous fairness literature covering three broad areas: *a*) far-reaching yet broad surveys; *b*) surveys targeted to specific domains; and *c*) call-for-action works focusing on diverse, representative, and balanced research samples.

Broad Fairness Surveys. Addressing algorithmic bias in ML has been a longstanding issue [16], despite its recent surge. A number of comprehensive surveys shed light on data and model biases across domains and compared potential mitigation solutions. For example, Caton and Haas [16] and Pessach and Shmueli [109] discussed fairness metrics and categorized mitigation approaches into a widely accepted framework of pre-processing, in-processing, and post-processing methods independently of the application domain. Wan et al. [134] focused exclusively on in-processing modeling methods such as adversarial debiasing, disentangled representations, and fairness-aware data augmentation, while Pessach and Shmueli [109]’s work provides an overview of emerging research trends, including fair adversarial learning, fair word embeddings, fair recommender systems, and fair visual description. More recently, Mehrabi et al. [90] explored data-to-algorithm (e.g., representation bias, measurement bias, aggregation bias, etc.), algorithm-to-user (e.g., popularity bias, user interaction bias, evaluation bias, etc.), and user-to-data (e.g., historical bias, temporal bias, content production bias, etc.) biases, and how these biases are generally encountered in ML practice. Along these lines, Le Quy et al. [69] surveyed available datasets for fairness research, including financial, criminological, healthcare, social, and educational datasets. Yet, despite these surveys’ considerable contributions, they tend to be of generic nature and rarely discuss data, models, and applications related to an individual community.

¹https://scholar.google.co.uk/citations?view_op=top_venues&hl=en&vq=eng_humancomputerinteraction



MACHINE LEARNING FAIRNESS	MOBILE & WEARABLE COMPUTING
Benchmark, standardized datasets	In-the-wild, author-collected datasets
Collected once	Collection ongoing
Tabular, low-dimensional data	Sequential, passively-sensed data
Unchangeable, immutable	Changeable, mutable
Medium-, large-scale data	Small-scale, proof-of-concept data
Human-distinguishable biases	Hard to distinguish biases
Classification/Ranking problems	Regression/Classification problems

Fig. 1. **Conceptual differences between ML Fairness and UbiComp datasets.** UbiComp data and models are oftentimes inherently different from those commonly used by the ML fairness community. Figure style inspired by [108].

Targeted Fairness Surveys. Another line of work took a deep dive into well-defined domains (e.g., recommender systems, social networks, healthcare, etc.). A number of works targeted specific ML paradigms, for example, by focusing on fairness for ML for graphs [23], on exploring notions of fairness in clustering [21], and on studying fairness in recommender systems [72]. Another group of works targeted specific unprivileged groups or high-stakes domains. For example, Olteanu et al. [104] reviewed the literature surrounding social data biases, such as biases in user-generated content, expressed or implicit relations between people, and behavioral traces, while in [127], the authors focused specifically on gender bias in Natural Language Processing (NLP). On a different note, Abdul et al. [1] featured emerging trends for explainable, accountable, and intelligible systems within the CHI community, also discussing notions of fairness. Closer to our work, Mhasawade et al. [94] discussed ML fairness in the domain of public and population health, and Xu et al. [144] explored algorithmic fairness in computational medicine, which only covers a subset of the broad, interdisciplinary UbiComp research domains.

WEIRD Research. Last, another strand of fairness work is concerned with what is coined as WEIRD research. WEIRD research refers to a common criticism in the social sciences that much of the research is conducted on a sample of participants that is Western, Educated, Industrialized, Rich, and Democratic (WEIRD). In particular, a comprehensive study conducted by Henrich et al. [50] in 2010 revealed a significant bias in sample populations. The study found that most research samples come from WEIRD populations, which represent only 12% of the global population but account for 96% of research samples. This criticism suggested that using such a narrow and unrepresentative sample of participants can limit the generalizability of the findings to the broader population. Over the past decade, the *CHI* community, which focuses on human-centered design, and the *FAccT* community, which aims to democratize ML and advance the development of responsible artificial intelligence, have become more aware of the potential biases introduced by WEIRD samples. For instance, Linxen et al. [75] conducted a meta-study on CHI findings from 2016 to 2020, reporting that 73% of CHI studies are based on Western populations,

representing less than 12% of the population worldwide, invariably making CHI “WEIRD”, as it is based on the knowledge and ethics of people who are Western, Educated, Industrialized, Rich, and Democratic. Similarly, a recent meta-study on FAccT proceedings from 2018 to 2021 extracted research topics and identified community values, placing fairness and ML, bias in word embeddings, bias in vision, and racial disparities among the ten largest sub-communities within the conference [68]. Yet again, as highlighted in Introduction (Section 1) in line with prior work [68], “off-the-shelf” benchmark datasets are encountered in the majority of published work, while only a ~ 10% of FAccT papers use original, empirical datasets, let alone UbiComp data.

It is evident that research communities other than UbiComp have recently started to explore ways of reporting data and models in a fair manner to surface and, ultimately, address encountered biases. Yet, the state of fairness in the UbiComp community remains unknown, as, at the time of writing, there exists no other survey or position paper in the intersection between UbiComp and fairness. Are our data susceptible to biases? Do our models discriminate against certain demographics, and if so, how do we make them right? These are just a handful of questions we set to provide answers to in this work.

3 BACKGROUND: FAIRNESS DEFINITIONS & MEASUREMENT

Fairness is a social construct that defies a simple definition [101]. In the legal domain, fairness entails the “protection of individuals and groups from discrimination or mistreatment with a focus on prohibiting behaviors, biases and basing decisions on certain protected factors or social group categories” [122]. Social sciences often consider fairness “in light of social relationships, power dynamics, institutions, and markets” [98], while quantitative fields (e.g., computer science, statistics) view fairness as a mathematical problem of “equal or equitable allocation, representation, or error rates, for a particular task or problem” [101].

Viewed through the lens of quantitative science, ML research has broadly grouped fairness into three categories: *group fairness*, *individual fairness*, and *subgroup fairness* [90]. Group fairness ensures some form of statistical parity (e.g., in terms of positive outcomes or errors) for individuals belonging to different protected groups (i.e., groups characterized by a sensitive attribute, such as gender or race) [30, 66]. On the other hand, individual fairness ensures that “similar” individuals receive similar outcomes [30, 66]. Subgroup fairness aims to exploit the best of both worlds by ensuring some form of statistical parity but holding this constraint over a large collection of subgroups to prevent fairness gerrymandering [57, 58]. While proponents of individual fairness have argued that it should be preferred to other categories for determining fairness, individual fairness has also received criticism: (1) counterexamples show that similar treatment is insufficient to guarantee fairness; (2) similarity metrics are susceptible to encoding implicit human biases; (3) similarity definition assumes prior moral judgments; and (4) the incommensurability of relevant moral values makes similarity metrics infeasible for many tasks [36]. Considering these limitations, we refer to both group and subgroup fairness when the term fairness is used.

In quantifying group fairness, there exist two opposing perspectives: “We’re All Equal” (WAE) and “What You See Is What You Get” (WYSIWYG) [38, 147]. The WAE perspective assumes equal ability across groups to perform the task, and thus it is closely linked with treating equals equally, whereas the WYSIWYG viewpoint assumes that the data themselves reflect a group’s ability with respect to the task, and thus, unequals should not be treated equally. Each perspective is quantified by different fairness metrics [41]. The WAE perspective, for example, uses demographic parity metrics, such as disparate impact and statistical parity difference, while the WYSIWYG perspective uses equality of odds metrics, such as average odds and average absolute odds difference. The two schools of thought find some common ground in equality of opportunity metrics, such as false negative rate, false positive rate, and error rate ratios, among others, where the choice of appropriate fairness metric is often guided by the question “What is the consequence of the predictive outcome?”. In quantifying individual fairness, “similar” individuals should be treated similarly. Dwork et al. [30] formalized this intuition by considering ML models as mappings between input and output metric spaces and defining individual fairness as their Lipschitz

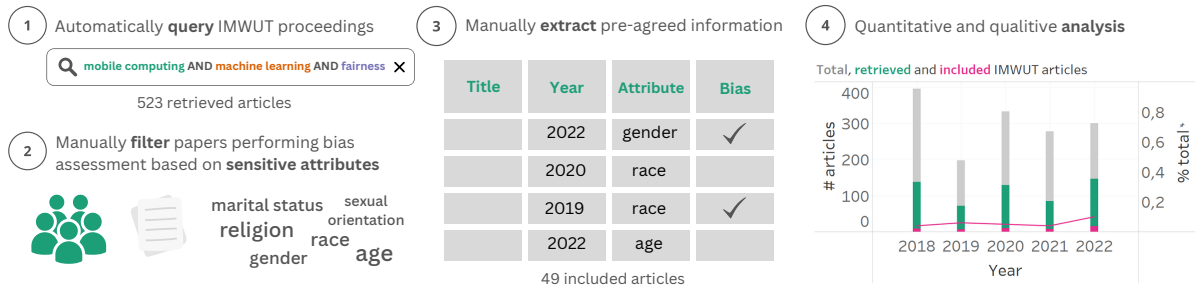


Fig. 2. **Illustration of the literature review methodology.** A high-level overview of the process, including the querying, manual extraction, and filtering, and one of the main results. Note that for readability purposes, we present a simplified version of our query and data extraction. Our query retrieves ~ 55% of all IMWUT publications, while our eligibility assessment filtering ends up with 49 papers (~ 9% of retrieved papers). We notice that only a very small fraction of all IMWUT papers looks at fairness issues, with only a small deviation across years.

continuity. Specifically, a Lipschitz continuity condition requires that for any two individuals u_1 and u_2 , their distance (as defined by a given distance metric) is proportional to the difference between the classifier's output for u_1 and the classifier's output for u_2 . In other words, if x and y are similar according to the distance metric, then their classifier outputs should also be similar, and vice versa. The distance metric on the input space though is crucial to the definition as it encodes individuals' similarity, and the choice was originally deferred to regulatory bodies or civil rights organizations rather than researchers or practitioners [97]. Naturally, different works use different definitions of algorithmic fairness, and although these appear internally consistent, they may also be mutually incompatible, as many quantitative fairness metrics cannot be satisfied simultaneously [38].

Despite the contradictory nature of fairness [38], the common element across traditional, even opposing definitions and metrics is that fairness is defined and assessed with respect to one or more sensitive (a.k.a. protected) attributes, such as gender, race, or age. Specifically, such attributes enable post hoc analyses to evaluate model performance across demographics and, ultimately, model fairness. This is indeed how fairness has been interpreted within the machine learning community [4, 10, 13], but viewed through the lens of the traditional definition, the UbiComp community is significantly lagging behind.

4 METHODOLOGY

Next, we delineate our methodology for conducting this systematic review (§4.1 and §4.2) and provide our positionality statement (§4.3).

4.1 Conducting the Literature Review

We followed an established protocol for conducting systematic reviews introduced by Kitchenham and Charters [61] to ensure the quality of included works and limit the initial retrieved papers. At least three authors were involved at each step to minimize the effects of bias and priming. In accordance with this protocol, we initially identified the need for a systematic review, as discussed in Section 1, namely to explore where UbiComp fairness overlaps with traditional fairness definitions and, more importantly, where it lags behind. Figure 2 provides a high-level overview of the process.

Paper Identification & Screening. UbiComp is a relatively new area in Information and Computer Science that crosses several fields, ranging from HCI, Hardware and Software Systems, and Knowledge Discovery and Data Mining. For the scope of this review, we focused on the Proceedings of the ACM on Interactive, Mobile,

```

(Abstract: ("mobile*" OR "wearable*" OR "smartwatch*" OR "smartphone*" OR "track*") OR
Title: ("mobile*" OR "wearable*" OR "smartwatch*" OR "smartphone*" OR "track*") OR
Keywords: ("mobile*" OR "wearable*" OR "smartwatch*" OR "smartphone*" OR "track*"))
AND
("deep learning" OR "machine learning" OR "artificial intelligence" OR "classification" OR "regression")
AND
("fair*" OR "bias*" OR "parity" OR "ethic*" OR "responsible AI" OR "RAI" OR "discriminat*" OR "non-discriminat*"
OR "equal*" OR "inclusiv*" OR "transparen*")

```

Fig. 3. **The query utilized for recovering relevant papers from the ACM Digital Library.** Terms related to UbiComp are highlighted in green, ML in orange, and fairness in purple.

Wearable and Ubiquitous Technologies (IMWUT), a premier journal in the UbiComp community, also placed among the top-3 publications in HCI. For the search process, we utilized the ACM Digital Library, focusing on papers that were published in the last five years (between 2018 and 2022) to capture emerging trends in fairness and UbiComp research. We also considered a broader search on Google Scholar but opted not to include it due to scale concerns (e.g., the query “machine learning fairness” returns approximately 279,000 results on Google Scholar) and the fact that indexed papers may not have undergone peer review. Apart from year filtering, for the most part, we did not limit our search to meta-data, such as titles, keywords, and abstracts, but rather we expanded it to any searchable field, including full text. That excludes the first part of the query, which tries to match terms such as wearable(s) or mobile(s) only in the papers’ meta-data, as seen in Figure 3.

Query Definition. For the definition of our query, we followed similar terminology with relevant review papers in the fairness literature [16, 69]. Additionally, according to Fjeld et al.’s analysis of prominent AI principles documents, [35], “the fairness and non-discrimination theme is the most highly represented theme in our dataset, with every document referencing at least one of its six principles: “non-discrimination and the prevention of bias”, “representative and high-quality data”, “fairness”, “equality”, “inclusiveness in impact”, and “inclusiveness in design”, mostly included in our query’s coverage. To capture the industrial perspective, we consulted the Responsible Artificial Intelligence (RAI) white papers issued by large tech companies. Specifically, Google’s² and Meta’s³ RAI principles talk about “fairness and inclusion”, Amazon’s⁴ RAI principles promote “diversity, equity, and inclusion” through “detecting bias”. Similarly, Nokia’s⁵ RAI fairness pillar talks about “fairness, non-discrimination, accessibility, and inclusivity”, and Intel’s RAI pillars mention “enabling ethical and equitable AI”. Thus, an iterative refinement process resulted in the query shown in Figure 3.

Eligibility Assessment. To further validate our query, we manually inspected all publications from the latest IMWUT proceedings (Volume 6, Issue 4, published in January 2023) ($N = 56$) to identify eligible papers for inclusion (see inclusion and exclusion criteria below). In total, we identified seven relevant publications, all of which were also returned by our query. This process was irrelevant to our final paper retrieval (Preferred Reporting Items for Systematic Reviews and Meta-Analyses (PRISMA) [96] flow diagram is pictured in Figure 4) and served validation purposes only. To ensure the high quality and relevance of the included papers, we defined appropriate exclusion criteria that helped us determine the included papers:

- (1) Papers that do not provide a quantitative assessment of at least one empirical or artifact contribution in ubiquitous computing (UBI);

²<https://ai.google/responsibilities/responsible-ai-practices/?category=fairness>

³<https://ai.facebook.com/blog/facebooks-five-pillars-of-responsible-ai/>

⁴<https://aws.amazon.com/machine-learning/responsible-machine-learning/>

⁵<https://www.bell-labs.com/institute/blog/introducing-nokias-6-pillars-of-responsible-ai/>

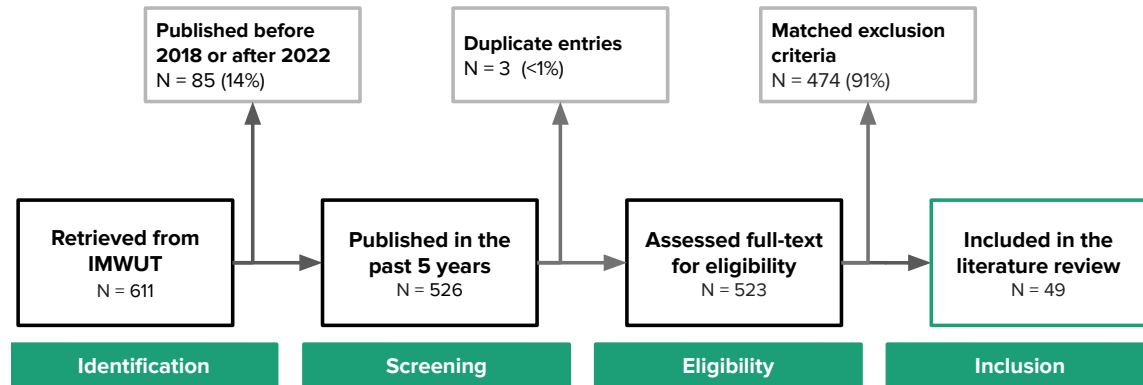


Fig. 4. **PRISMA flow diagram for paper inclusion.** Out of the 611 papers retrieved by our query after the screening, only 9% ($N = 49$) did not check any exclusion criterion and thus were included in the literature review. The majority of the retrieved papers were false positives, highlighting the importance of the eligibility assessment.

- (2) Papers that do not include a quantitative assessment of bias or performance discrepancy in their evaluation with regard to sensitive attributes, such as age, gender, race, disability, religion, and sexual orientation, among others (FAIR);
- (3) Papers that discuss different domains, such as natural language processing or computer vision without incorporating a UbiComp component (DOM);
- (4) Papers that refer to bias in a different context, such as the bias-variance trade-off or the bias parameter in neural networks (CON).

Inclusion & PRISMA Statement. Finally, the sequential execution of the steps above, as depicted in Figure 4, led to our review’s included papers. Overall, we screened 523 papers after date filtering and duplicate elimination. After carefully screening these papers, we excluded 474 based on our exclusion criteria. In detail, 31 papers did not provide any quantitative assessment of a UbiComp contribution (UBI: 6.5%), 394 papers did not provide any fairness assessment (FAIR: 83.1%), 2 papers did not discuss a UbiComp component (DOM: <0.5%), 42 papers referred to bias in a different context (CON: 8.9%), and finally, 5 papers were excluded for other reasons (OTHER: 1.1%). Hence, we included 49 papers in our review⁶.

4.2 Methodological Limitations

While we have made every effort to ensure broad coverage of papers relevant to fairness in UbiComp studies, our search for literature might not be comprehensive and exhaustive. However, covering IMWUT as a prominent academic venue for ubiquitous computing research allowed us to capture emerging trends. Besides, we intended to provide insights and research directions to the IMWUT community; therefore, we narrowed down our research to IMWUT proceedings as opposed to the broader community by including the proceedings of, for example, PerCom, SenSys, MobiSys, or medical journals such as JMIR. We also acknowledge that despite our best efforts to pick literature-driven keywords and manually validate the retrieved results, the output might well have produced both false positives and false negative results.

⁶To foster reproducibility, upon acceptance, we intend to release the review data and codebooks

4.3 Positionality Statement

Understanding researcher positionality is essential to demystifying our lens on data collection and analysis [39, 49]. We situate this review paper in a Western country (REDACTED FOR REVIEW) in the 21st century, writing as authors who primarily work as academic and industry researchers. We identify as two females and four males, and our shared backgrounds include HCI, ML, and ubiquitous computing.

5 RESULTS

In the discussion of our results, we looked into bias in both models and data in the included papers. *Model bias* refers to systematic errors that are introduced into ML algorithms due to the underlying assumptions in the data, algorithms, or the learning process itself, leading to unfair or unequal outcomes for certain sensitive groups. In the following sections, we explore the state of fairness in the UbiComp community (Section 5.1), the ethical risks and opportunities in the domain and how they can inspire the choice of fairness metrics (§5.2), and ultimately, we capture alternative notions of fairness that the UbiComp community has perhaps indirectly adopted (and adapted) to capture the particularities of the domain (§5.3). *Data bias* refers to errors occurring because certain user groups, or generally elements, within a dataset are more heavily weighted and/or represented than others. A biased dataset does not accurately represent a model’s target domain, causing skewed outcomes, low performance, and systematic errors. In the sections below, we also discuss reported or historical biases relevant to the domain’s data modalities (§5.4) and explore diversity in UbiComp datasets and research teams (§5.5).

5.1 What is the State of Fairness in UbiComp?

Before delving into details about fairness assessment, enhancement mechanisms, and metrics documented in the included papers, we summarize the approaches adopted by the UbiComp community in Table 1. The table presents all included papers, categorized per application domain, sensitive attributes, and fairness mechanisms and metrics. In summary, out of the 523 retrieved papers, a small portion of 9% ($N_{included} = 49$) were included in the review, which in turn make up only 5% of all IMWUT publications between 2018 and 2022, highlighting the timeliness and necessity of this work.

Takeaway #1

Out of all papers published at IMWUT between 2018 and 2022, only a small portion of 5% (included papers) adhered to modern fairness reporting.

To identify appropriate application domains, we consulted the past four years (2019-2022) of UbiComp tracks and sessions to identify commonalities in discussed themes. We grouped together tracks’ themes between years based on their similarity and relevance to the included papers. This process led us to the identification of ten domains: *Health; Human-Activity Recognition; Behavioral Sensing & Emotion; Cognition & Attention; Motion, Gaze, Gesture & Touch; Sound, Voice & Hearing; Mobility & Navigation; Privacy & Security; Localization*, and *Miscellaneous*. We encountered all but one theme (localization) in the included papers, which is not unreasonable given localization’s usually low-stakes applications and, thus, less relevance to fairness. Health was the most commonly encountered domain, accounting for more than one in four papers, while Cognition & Attention was the least common, accounting for only ~ 6% of included papers.

Figure 5 provides an overview of the discussed domains, as captured by the papers’ keywords. It also serves as a validation to the domains’ categorization, as many categories (e.g., *Health; Human-Activity Recognition; Motion, Gaze, Gesture & Touch; Sound, Voice & Hearing* and *Privacy & Security*) also appear in the keyword clouds. Deep learning, ML, and human activity recognition are among the most frequently overlapping keywords (colored in dark grey). Over-represented keywords in the retrieved papers (colored in green) include mobile

Table 1. A summary of the included papers categorized by application domain, fairness enhancement mechanism, sensitive attribute, and bias metrics. High-stakes health application papers are the most active in the UbiComp community regarding fairness. While all included papers talk about fairness with respect to one or more sensitive attributes, the vast majority do not offer enhancement mechanisms.

APPLICATION DOMAIN	FAIRNESS MECHANISM	SENSITIVE ATTRIBUTE	FAIRNESS METRICS	PAPERS
Health (26.5%)	Preprocessing	Gender	Accuracy	[150]
		Age	MAE	[74]
	In-processing	Health Condition	Precision, Sensitivity, Specificity	[116]
		Socioeconomic Status	Accuracy, AUC-ROC, F1, Sensitivity, Specificity, MAPE, Error rate	[43, 46, 55, 142, 151]
	None	Gender	Accuracy, AUC-ROC, Sensitivity, Specificity, Pearson's r, Error rate	[7, 55, 138, 142, 151]
		Age	Accuracy, AUC-PRC, AUC-ROC, F1, Precision, Sensitivity, Specificity, RMSE	[46, 47, 151, 152]
		Health Condition	F1, Sensitivity, Specificity	[46]
Privacy & Security (12.2%)	None	Race	AUC-ROC, MAE	[55, 78]
		Age	Error rate	[40]
		Health Condition	Precision, P-value	[3, 132]
		Physiology	Accuracy	[44, 117, 145]
		Religion		[44]
Human-Activity Recognition (10.2%)	Pre-processing	Gender	F1	[124]
		Age		
	In-processing	Physiology		
		Gender	F1	[119]
	None	Gender	Accuracy	[148]
		Physiology	F1	[153]
Behavioral Sensing & Emotion (10.2%)	Pre-processing	Miscellaneous	MSE	[64]
		Gender	Accuracy	[59]
		Age	Accuracy, MAE	[59, 77, 91]
	None	Nationality	Accuracy	[59]
		Gender	Accuracy, Coefficient of determination	[60, 89]
Sound, Voice & Hearing (10.2%)	Pre-processing	Gender	Accuracy	[125]
	None	Gender	Error rate, Mel cepstral distortion	[5, 73, 136]
		Age	Error rate	[136]
		Physiology	Accuracy, MAE	[71]
		Nationality	Error rate	[5]
Motion, Gaze, Gesture & Touch (8.2%)	None	Gender	Accuracy	[76, 81]
		Age		
		Health Condition	Error rate	[53]
		Physiology	Accuracy	[76, 133]
		Gender	Accuracy	[135]
Mobility & Navigation (8.2%)	Pre-processing	Age		
		Physiology		
		Gender	Accuracy, AUC-ROC, F1	[154]
	In-processing	Age		
		Socioeconomic Status		
		Gender	F1	[100]
		Age	F1	[143]
Miscellaneous (8.2%)	Pre-processing	Miscellaneous	Accuracy	[155]
	None	Miscellaneous	Accuracy	[155]
		Gender	Required rate of return	[93]
		Age	Required rate of return	[93]
Cognition & Attention (6.1%)	None	Marital Status	Required rate of return	[93]
		Physiology	P-value	[20]
		Language	Error rate	[111]
		Age	RMSE	[118]
		Physiology	MSE	[139]
		Miscellaneous	Accuracy, AUC-ROC, Precision, Sensitivity, MAE, MSE, Pearson's r	[76, 140]



Fig. 5. **Keyword differences of retrieved (left) and included (right) papers.** Frequent keywords in both retrieved and included papers are colored in dark grey. Over-represented keywords in the retrieved papers are colored in green, while over-represented keywords in the included papers are colored in pink. Even within UbiComp, the privacy, audio, and vision communities are trailblazers in ML fairness.

sensing, wearables, Internet of Things (IoT), Radio-frequency identification (RFID), and self-supervised learning. Over-represented keywords in the included papers (colored in pink) include mobile health, Post-traumatic stress disorder (PTSD), acoustic sensing, computer vision, privacy, and gesture recognition.

5.1.1 Fairness Enhancement Mechanisms. For each application domain, we categorized its papers based on three fairness enhancement mechanisms [16, 109]: a) *pre-processing*; b) *in-processing*; and c) *post-processing* mechanisms. It is often infeasible to eliminate all sources of unfairness and guarantee fairness. Yet, the goal is to surface and mitigate biases as much as possible through fairness enhancement mechanisms.

Pre-processing mechanisms involve altering the training data before feeding it into a ML algorithm. Within UbiComp, preliminary but effective mechanisms include fair data representation. For instance, during data collection, Liaqat et al. [74] equally included both healthy subjects and subjects with Chronic Obstructive Pulmonary Disease (COPD) in their dataset for respiratory rate monitoring using smartwatches, leading to non-significant differences in model performance across health condition. Similarly, Zhou et al. [155] employed a fairness-aware client selection mechanism for federated learning to ensure equal representation for subjects with worse connectivity.⁷ Post data collection, Su et al. [125] performed data balancing, conditioned on the sensitive attribute, managing to narrow the impact of gender voice differences on their speech recognition model. Similarly, a strand of work explored data splitting, conditioned on the sensitive attribute (gender, age, BMI, skin tone, country, and health condition) to enable model personalization [59, 77, 91, 125, 150]. More advanced mechanisms suggest modifying feature representations so that a subsequent classifier will be fairer. For example, Wang et al. [135] improved the performance of their activity detection model by normalizing the window-level features across gender and physiology, yet their model remained dependent on sensitive attributes. Similarly, in line with prior work [79], Su et al. [124] utilized disentangled representations, aiming to isolate relevant activity

⁷While Internet connectivity is not a sensitive attribute per se, it has been linked with socioeconomic status, race, nationality, gender, and age, all of which are sensitive attributes [102].

patterns from redundant noises such as gender, age, and physiological differences, reducing the effect of such covariate factors. However, they did not manage to completely separate the activity signals from the redundancy, attributing it to the diversity limitations of Human-Activity Recognition datasets.

In-processing mechanisms involve modifying the ML algorithms to account for fairness during training. In the included papers, Shahid et al. [116] altered their logistic regression model for Post-traumatic stress disorder (PTSD) screening to include sensitive attributes in its parameters, having observed statistically significant inter-group and intra-group differences based on gender and socioeconomic status. Their alteration led to a statistically significant improvement in performance across groups. On a similar note, Sheng and Huber [119] devised a multi-task loss function consisting of activity, subject, and gender loss. However, they noticed unbalanced performance, with the performance on gender attribute learning being around 11% lower than the performance on the other two tasks. Finally, in quantifying the causal effect of individual mobility on health status, Zhang et al. [154] considered certain sensitive attributes, such as age and socioeconomic status, as confounding variables in their causal model, after establishing that such attributes can affect the correlation between mobility and health. However, they noted that due to dataset privacy constraints in reporting demographic variables, potential unobserved confounding variables, such as occupation, employment, and education, might have been missed, highlighting the conflict between fairness and privacy [18].

Post-processing mechanisms involve altering the output scores of the ML model to make decisions fairer. However, due to the relatively late stage in the learning process in which they are applied, post-processing mechanisms commonly obtain inferior results [141]. They are also considered too invasive or discriminatory since they deliberately damage accuracy for some subjects to compensate others [109]; hence they are less frequently preferred in practice. Perhaps not surprisingly, such mechanisms are not present in the included papers.

Despite the notable efforts of the aforementioned pioneering works in the UbiComp community, 3 out of 4 included papers did not report any fairness enhancement mechanism, regardless of the presence of bias in their models. This is partly due to a lack of consideration for fairness-related harms, but it is also connected with the nature of several UbiComp works: artifact contributions, proof of concept, and early-stage technology development, where performance is prioritized.

— Takeaway #2 —

Papers implementing fairness constitute a small portion of included papers (24%), which, in turn, make up only a limited fraction of all IMWUT publications (1%).

5.1.2 Sensitive Attributes & Biases. The categorization of sensitive attributes is inspired by the EU Charter of Fundamental Rights that prohibits any discrimination based on any ground such as sex, age, race, ethnic or social origin, genetic features, language, and religion or belief, among others [31]. In line with such declarations and prior fairness work [109], UbiComp works investigate a variety of sensitive attributes individually or combined: *gender, age, physiology* (e.g., height, weight), *health, language, nationality, socioeconomic status, religion, race, occupation*, and *marital status*, as seen in Table 1.

Nevertheless, some attributes are better represented than others in the included papers, with gender and age being at the top (mentioned in almost 9 out of 10 included papers), followed by physiology and health condition (mentioned in 4 out of 10 included papers), as seen in Figure 6b. Surprisingly, attributes with long-history of discrimination in ML, such as race, language, and nationality [67, 70, 85], are rarely encountered in the UbiComp literature, with only two papers discussing racial discrepancies in model performance. Yet, UbiComp is far from immune to such discrepancies. Sjoding et al. [121] uncovered racial and ethnic biases in pulse oximetry, while Hutiri and Ding [52] reported language biases in speech recognition.

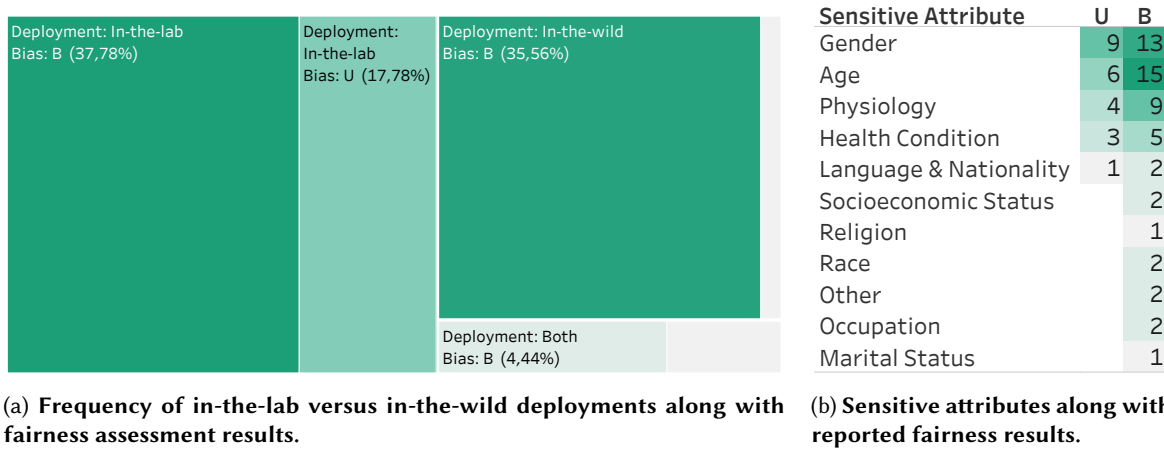


Fig. 6. **Bias results across deployment settings and sensitive attributes.** The figures visualize the fairness assessment results reported in the included papers against deployment setting (left) and sensitive attributes (right). “B” indicates a bias towards one or more sensitive attribute(s), while “U” indicates an unbiased model. In 6a (left), the deployment environment distribution is relatively balanced, but the reported fairness assessments differ with in-the-wild studies reporting significantly less unbiased results. In 6b, gender, age, and physiology are amongst the most frequently assessed attributes, while, surprisingly, race and language are understudied. Note that a single paper might assess more than one sensitive attribute; hence the sum may exceed the number of included papers ($N = 49$).

In line with such findings, our review of UbiComp’s work highlighted biases in ML models across all sensitive attributes and a wide range of UbiComp applications. Gender biases have been reported in monitoring sleep posture with wireless signals [148], opioid usage tracking [46], diaphragmatic breathing monitor based on acoustic signals [43], and speech recognition via accelerometer sensors [125]. Age biases have been reported in medication adherence monitoring through gait assessment [150], fatigue estimation via smartphone tapping frequency [7], mobility purpose and route choice inference [100], and neural activation prediction [55]. Biases based on physiological measurements have been reported by Li et al. [71] in fine-grained activity sensing (e.g., eye blinking, finger tracking) using acoustic signals against people of small stature, by Wang et al. [136] in vital sign monitoring through acoustic sensing against obese or overweight people, and by Griffiths et al. [44] in image processing with binocular thermal cameras against people of non-average height. Similarly, a model for early detection and burden estimation of AFib under-performed for long-term AFib patients compared to their healthy counterparts [151], while a wearable-based clinical opioid use tracker showed bias against chronic opioid users [46]. Regarding less explored sensitive attributes, Griffiths et al. [44] encountered model biases in user authentication via binocular thermal cameras for hijab wearers, a proxy for religion, while Ruan et al. [111] uncovered language biases in speech and keyboard text entry for non-English speakers. It is worth noting that all included papers explore notions of group fairness rather than individual and subgroup fairness (detailed in §3).

Takeaway #3

Across application domains, *gender, age, and physiology* and *health conditions* are commonly taken into account, while *race, nationality, and language* are unfairly overlooked.

5.1.3 Fairness Metrics. Metrics for performance evaluation monopolize UbiComp fairness assessment and reporting as shown in Table 1. Classification metrics, such as Accuracy, Area under the ROC Curve (AUC-ROC), and F1-Score, etc., and regression metrics, such as Root-mean-square Error (RMSE), Mean Absolute Error (MAE), and Mean Absolute Percentage Error (MAPE) highlight the interest of the UbiComp community for such tasks. There exist two challenges, though, with UbiComp’s take on fairness: Firstly, how does one define a threshold in performance evaluation above which a model is considered unfair across sensitive attributes? For instance, if an AFib detection model has 85% accuracy on healthy adults and 80% accuracy on the elderly, should it be considered fair or unfair? This challenge holds for the entirety of ML fairness research, as there are seldom clear-cut answers. Secondly, how does one perform fairness assessment in regression or multi-class classification scenarios, both wildly understudied areas in the fairness domain compared to binary classification [95]? This challenge especially holds for UbiComp. To see how, consider that 1 out of 2 included papers did not discuss binary classification.

To answer the first question, one could employ a statistical hypothesis test, such as the Student’s t-test, for comparing samples’ performance across sensitive attributes, a practice also adopted by a portion of the included papers. Yet the choice of a statistical test is far from straightforward, each incorporating strict assumptions. For example, a key assumption of the paired Student’s t-test is that the observations in each sample are independent, which is not the case in k-fold cross-validation, a common practice in ML model evaluation, leading to an incorrect calculation of the t-statistic and a misleading interpretation of the results and p-value [28]. Better alternatives proposed include McNemar’s Test or 5×2 cross-validation and its refinements [28, 99]. Note that more than 1 in 4 included papers (27.7%) did not use any statistical significance testing in their assessment. Another limitation of performance-based assessment of fairness is the assumption that performance optimization and fairness criteria always overlap. For instance, an AFib detection algorithm might be optimized for accuracy [12], but in fairness assessment, one might also want to ensure equal false negative and false positive rates across groups. An alternative to statistical significance testing on performance metrics, also adopted by the FAccT community, is the usage of fairness metrics; namely several measures, such as demographic parity or equalized odds, that enable the detection of bias in one’s data or model, as briefly discussed in Section 3. Once a fairness metric is obtained, it is common practice to apply the “4/5 rule” or “80% rule”, which states that “the selection rate of any group should be not less than 4/5 than the one of the group with the highest selection rate” [15]. This refers to the guidelines established by the US Equal Employment Opportunity Commission (EEOC) [128], which are frequently cited as one of the few legal frameworks that rely on a specific definition of fairness, particularly the concept of demographic parity. Yet, there is no single fairness definition, metric, or “fair” threshold that will universally apply to different applications. The rule should only be used as a “rule of thumb” and is dependent on the application domain. For example, in high-stakes applications (e.g., health), would we consider “acceptable”, or fair, a model for arterial oxygen saturation estimation with 80% accuracy on Asian, Black, and Hispanic patients and 95% accuracy on White patients—even though it abides by the “4/5 rule”? It is worth noting that we did not identify any fairness metrics in the included papers, indicating the disjointedness between the UbiComp and the FAccT communities.

Regarding the second question, nearly half of the included papers (47%) engage in regression or multi-class classification tasks, such as respiratory rate detection and human-activity recognition, respectively. Yet, the most common ML paradigm explored in fairness research is binary classification [95], with most fairness metrics and enhancement mechanisms specifically targeted to such tasks. In one of the few works about fair regression, Agarwal et al. [2] introduced two definitions of fairness in regression: *statistical parity*, which asks that the prediction be statistically independent of the sensitive attribute, and *bounded group loss*, which asks that the prediction error restricted to any sensitive group remain below some predefined threshold. A popular way to quantify fairness in regression is to compare the outcome distribution across sensitive attributes using the Kullback–Leibler divergence [54], or Kolmogorov–Smirnov test for goodness of fit [88]. If the test fails to reject the null hypothesis that the distributions come from the same population, it is considered fair. Otherwise, it is determined

that there is at least one sensitive group whose distribution does not come from the same population. However, this does not reveal which distributions are different (i.e., which group the model is biased against), and what are the characteristics of such differences (e.g., differences in mean, variance, skewness), requiring a subsequent analysis [2]. Similarly, popular ML fairness libraries, such as AIF360⁸ and FairLearn⁹, at the time of writing, do not include any regression-specific fairness metrics' implementations. In multi-class classification scenarios, computing "standard" fairness metrics such as equalized odds and demographic parity can be challenging due to the lack of a clear definition of "positive" and "negative" classes. To address this issue, one feasible solution is to transform the problem into multiple binary classification problems and then aggregate the results. More recent approaches include the Combined Error Variance (CEV) and Symmetric Distance Error (SDE) metrics [8, 9] to quantitatively evaluate the class-wise bias of multi-class classification models.

Takeaway #4

No UbiComp paper (from the included papers) uses modern fairness metrics, likely due to the lack of widely available regression and multi-class classification metrics.

5.2 Model Consequences: Ethical Risks versus Opportunities

It is no longer a matter of debate that ML has and will have a major impact on UbiComp, and by extension, on society; human authentication [145], early detection and burden estimation of AFib [151], respiratory rate monitoring [151], and opioid use tracking [46] are only a handful of examples that reinforce this argument. The discussion has now shifted to determining the extent and specifics of this impact. In other words, it is given that ubiquitous ML will have an impact on society; what is now being questioned is the specifics of who will feel the effects, how, where, and when they will be felt.

To concretize these questions, the Scientific Committee of the AI4People¹⁰ has categorized the chief ethical opportunities offered by artificial intelligence in "four fundamental points in the understanding of human dignity and flourishing" [37]: "who we want to become" (enabling self-realization), "what we can do" (enhancing human agency), "what we can achieve" (increasing individual and societal capabilities), and "how can we interact with each other and the world" (cultivate societal cohesion). Ethical opportunities within UbiComp span across all four points: UbiComp-based automation of mundane tasks, such as gait-based human authentication [145], speech transcription [111], or gesture recognition [76] may easily mean more time spent more intelligently (self-realization). UbiComp-based augmentation of human intelligence, such as cognitive load measurement [140], or cognitive performance prediction [140] may enable humans to do more, better, and faster (human agency). UbiComp-based innovations in medicine, such as PTSD screening [116], medication adherence monitoring [150], or post-operative complications prediction [152], may reinvent society by radically enhancing what humans are collectively capable of (individual and societal capabilities), while UbiComp-supported cooperative work, such as social context inference [89], may support societal cohesion and collaboration (societal cohesion).

However, we must also consider the ethical risks associated with inadvertent overuse and deliberate or unintended misuse of UbiComp technologies, stemming, for example, from lack of awareness, conflicting interests, greed, or malicious intent. Simply put, ethical risks are the most likely and predictable negative consequence of any action or inaction. And while performance optimization is frequently fueled by the potential of ethical opportunities, fairness assessment is also driven by ethical risks. Oftentimes, fairness considerations in ML systems are influenced by the question: What is the consequence of the predictive outcome? The answer to this

⁸<https://aif360.mybluemix.net/>

⁹<https://fairlearn.org/>

¹⁰An Atomium–European Institute for Science, Media and Democracy initiative.

question can drive the choice of suitable fairness definitions and metrics but is far from straightforward. For instance, in the case of AFib detection, a false negative outcome might prove deadly. Yet, “*deploying a system with a high false alarm rate can add anxiety to people*” [151]. Similarly, in predicting postoperative complications in pancreatic cancer patients, a false negative outcome can deprive a patient of much-needed care. However, “*many false positive errors, [...] means the patients without complications are incorrectly predicted to be at high risk. As a result, clinicians may decide to provide pre-habilitation to reduce their risk of surgical complications or even cancel the surgery. However, pre-habilitation delays surgery and the patients might miss their opportunity for successful recovery and surgery is the only cure for the cancer.*” [152]. On the contrary, in speech-based human identification scenarios, false positive outcomes are critical in preventing unauthorized access, as “*existing voiceprint-based authentication often suffers from various voice spoofing attacks*” [40].

Prioritizing predictive outcomes becomes even more challenging once perceived through their sociotechnical context. Oftentimes UbiComp technologies are built and evaluated as if they were fully autonomous, while in reality, they operate in a complicated sociotechnical system moderated by institutional structures and human stakeholders (the “framing trap” [115]). For instance, in opioid use tracking [46], and drug-seeking behavior sensing [47] applications—both encountered in the included papers—the consequence of a predictive outcome depends on the assumption of punitive or restorative justice [48]. According to traditional punitive justice, punishment serves as a deterrent for wrongdoing, and a means to alter behavior. However, restorative justice takes a different approach, recognizing that punishment alone does not repair the harm caused to the community and relationships. Additionally, restorative believes that relying solely on punishment can result in individuals becoming dependent on external factors rather than internal self-control to modify their behavior [130]. As an example, if a ubiquitous substance abuse detection technology is adopted by a restorative system, a false negative outcome might derive an individual struggling with drug addiction from crucial access to rehabilitation services. For instance, “*if such a device were found to be reliable, it could be used to monitor early treatment response and therefore could allow clinicians to more rapidly optimize patient care*” [47]. On the contrary, if the exact same technology is employed as part of a punitive system, then a false negative outcome might lead to a wrongful accusation or conviction.

Nevertheless, fear, ignorance, and misplaced concerns should not inhibit the UbiComp community from innovation and realizing ethical opportunities for individual and societal good. On the contrary, mindful use of ML is conscious of our commonalities and differences across sensitive groups, as well as the factors within and outside the community’s control, and serves as a framework to act with a sense of responsibility and fostering trust.

Takeaway #5

Given UbiComp’s high-stake applications, fairness reporting and justification help prioritize tradeoffs between ethical risks and opportunities.

5.3 How does UbiComp Capture Alternative Notions of Fairness?

Previously, we have established that only a small fraction of IMWUT works (5%) follow conventional fairness definitions, where fairness is defined with respect to one or more sensitive attributes. Yet, we believe such definitions do not do full justice to the community’s work, which strives for “fairer” models, perhaps not across sensitive attributes, but across differing experimental conditions. In particular, we noticed that, in evaluating new UbiComp systems, artifacts, or applications, the community aims for generalizable and robust models by performing ablations studies, comparing deployment settings, and personalizing models for users and groups. As an indication, in the retrieved papers, almost one out of two papers (44%) reported an ablation study or a deployment setting comparison in their results, while in the included papers, 57% did so.

Ablation Studies. In an ablation study, one or more components of the model are systematically removed or modified, and the performance of the model is evaluated after each change. By comparing the performance of the original model with the performance of the modified models, researchers can determine the importance of each component and gain insights into the functioning of the model across diverse conditions, ensuring its generalizability and robustness. In the included papers, ablation studies take the form of performance evaluation comparisons based on: *user-related components*, *device-related components*, *environmental components*, *experimental components*, and *domain-specific components*. In particular *user-related components* include user motion and orientation during data collection in sleep posture monitoring [148], breathing monitoring [43], gesture recognition [76], user identification [117], and heart activity monitoring [142], as well as aesthetics, such as hair or clothing in fine-grained activity sensing [71], breathing and vital sign monitoring [43, 136], and user identification [44, 145]. *Device-related components* include device type, sampling rate, and operating system, as well as device placement and orientation in activity and gaze tracking [53, 71], vital sign monitoring and physiological sensing [78, 136], speech recognition via built-in sensors and speech synthesis [73, 125], and user behavior sensing [60]. *Environmental components* include ambient noise, light, and temperature that might affect data quality of acoustic [40, 71, 136] or video [44, 78, 133, 139] signals, respectively, or random passers-by that might affect model performance on the individual for human identification [145]. Regarding experimental setup, few included papers studied the effect of equipment placement (i.e., distance, angle) and characteristics (i.e., range) on the model’s robustness in activity sensing [71] and vital sign monitoring using acoustic signals [43, 136]. Apart from such common components, the choice of components to consider in an ablation study is highly domain-dependent. *Domain-specific components* have no limitations and can range from screen size in scrolling interaction experiments [81] to food structure in food-related artifact development [20].

Deployment Setting. Beyond ablation studies, a study’s deployment setting, ranging from in-the-lab to in-the-wild, can significantly impact its outcomes. While laboratory settings can provide controlled environments for experimentation, they may not accurately reflect the complexities of the real world in which the applications are deployed. As a result, in-the-wild (or in-situ) studies have emerged as an alternative approach, focusing on evaluating the situated design experience of UbiComp. Such studies provide insight into how new ubiquitous technologies are adopted in real-world settings [110]. Figure 6a shows the distribution of in-the-lab and in-the-wild studies in the included papers, along with their reported fairness assessment results. We see that perhaps not surprisingly, in-the-lab studies prevail (~ 55%), which can be explained by the nature of numerous IMWUT papers presenting cutting-edge artifacts or early-stage model development work. Nevertheless, ~ 38% of included papers conduct in-the-wild studies, and a small fraction of papers (~ 7%) compare and report results for both deployment settings. An interesting point to be made here is that while 4 out of 10 in-the-lab studies do not identify biases in their models, this number falls to 0.5 in 10 for in-the-wild studies. This confirms our intuition that controlled environments might conceal biases that would emerge once a model is deployed in the real world.

Personalization. In 20% of the papers included, personalization was reported as a commonly used approach for gaining insights on performance differences between individuals. In particular, several works trained separate personalized models for inference on a single subject [5, 117, 150] or a group of subjects sharing a common characteristic. For instance, Liu et al. [77] built personalized models for different age groups “*illustrating differences in communication patterns across age demographics that can impact model performance*”. Similarly, Mendel et al. [91] utilized age-specific models for predicting the right moment for providing mobile safety help, as “*different ages in the sample have a significant influence on supportable moment predictions*”. Liu et al. [78] developed personalized models based on skin tone for camera-based, non-contact photoplethysmography, as “*previous work had already highlighted [skin tone and gender] issues with the Plane-Orthogonal-to-Skin [method]*”. Su et al. [125] developed gender-specific models for speech recognition, as “*women’s voice is generally thinner and higher in pitch*”, while Zhang et al. [153] explored BMI-based models for detecting eating activities via a multi-sensor necklace, due to

“differences in movement patterns while eating, change in the distance of the proximity sensor from the neck, and difference in posture during the eating activity”.

Takeaway #6

Modern fairness reporting aside, the UbiComp community strives for generalizability by conducting and reporting ablation studies, in-the-wild vs. in-the-lab experiments, and personalized models.

5.4 Is UbiComp Susceptible to Data Biases?

UbiComp is inherently multimodal; even the relatively small subset of included papers contained heterogeneous modalities of input data such as audio, video, images, text, and sensor data (e.g., accelerometer, gyroscope, temperature, electrodermal activity (EDA) sensors). These different modalities may be related or complementary and can provide a more complete and nuanced representation of a given phenomenon than any one data modality alone.

Biases in audio (used by 18% of included papers) are well-reported in fairness research communities [52, 87, 129]. UbiComp posed no exception, with such biases surfacing from the included papers. Several works discovered biases in acoustic signals dependent on body size, a potential proxy for gender, physiology, and race. Specifically, Li et al. [71] generically reported “*higher respiration [detection] error [...] due to [...] weaker chest motions and smaller body size*”, while Gong et al. [43] specifically reported bias against women due to “*different physiological structures*”, which led to “*the reflective surface of women [being] smaller than that of men, which encodes less information*”. Similarly, Wang et al. [136] found that acoustic signals for heartbeat monitoring were biased against people with larger BMI, as “larger BMI [would] make the thoracic muscle thicker and block the weak heartbeat signals”.

Biases in video and image (used by 20% of included papers) are also well-studied in fairness literature [27, 103]. Within UbiComp, Kalanadhabhatta et al. [55] attributed performance discrepancies in computer vision models to data-related factors regarding participants’ identity, gender, age, race, ethnicity, and Fitzpatrick skin type [34]. In line with this discourse, Griffiths et al. [44] identified video biases related to physiology, such as height, and appearance features such as hair or head covering, potentially proxies for gender and religion, respectively. Specifically, they encountered certain issues during video capture: “*As [the tallest participants] moved closer to the cameras, both participant’s heads moved above the viewing range of the [...] camera*”; or “*[the algorithm] often measured the participant at their forehead rather than the top of their head [...] due to [the participant’s] long, thick, curly hair*”; and finally, “*when [the participant] turned away from the camera and only her hijab was visible [the algorithm] was unable to detect her presence as different from the background*”.

Some of these biases are easier to distinguish as demographic information is integrated into the data. For example, gender, age, physiology, and race may be inferred from video, while gender, language, and possibly age and race may be inferred from speech signals. However, sensor signals—the most prevalent data modality used by the UbiComp community—are more challenging and equally susceptible to data biases. Prior research has reported racial and ethnic biases in pulse oximeters [121]. Additionally, gender biases may be a concern for electrocardiogram (ECG) quality, given that certain ECG metrics, such as the PR interval, heart rate, QRS duration, and lead voltages, exhibit gender-based differences [146]. Even the most inconspicuous sensors, such as heart rate and acceleration sensors, are shown to be correlated with health, fitness, and demographic characteristics [123].

Biases in sensor signals (used by 51% of included papers) also surfaced during our review. These biases could be attributed either to measurement inaccuracy or concept drift phenomena in signal patterns (i.e., distributions of data change over time, making machine learning models less accurate without updates [80]). For example,

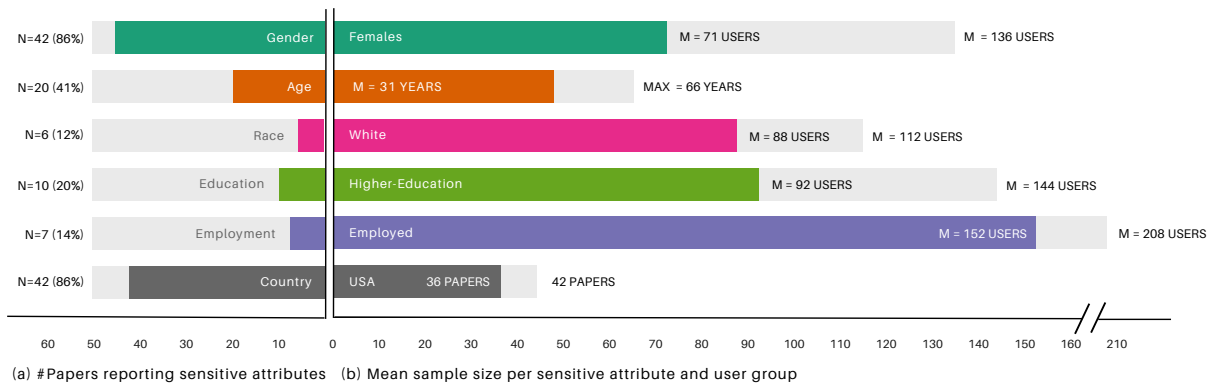


Fig. 7. **Analysis of sensitive attributes and data size.** The bar plots show the percentage of papers reporting certain sensitive attributes (left) and the (mean) sample size in the subset of papers reporting that attribute (right). Note that for age and country, we report the mean sample age and the number of papers originating from the USA, respectively. Sample demographics reporting is not standardized and frequently incomplete, with race, employment status, and education being the least reported sensitive attributes ($\leq 20\%$). While UbiComp samples tend to be gender-balanced, they are still WEIRD, as they consist of predominantly White, highly-educated, and US-based subjects.

Zhang et al. [150] reported decreased accuracy in gait detection via accelerometer and gyroscope measurements for elderly users as “Human gait patterns inevitably change with the increase of age. The so-called aging effect may affect the detection accuracy of [the] system”. Furthermore, Su et al. [125] encountered inaccuracies in accelerometer-based speech recognition, “Since women’s voice is generally thinner and higher in pitch, [and] it may be harder for [the] accelerometer to preserve voice feature”. Similarly, in speech synthesis from accelerometer measurement, Liang et al. [73] also found gender-based differences in the accelerometer’s ability to preserve the voice feature. In particular, “[...] *the frequency bands of male speakers are lower than those of female speakers. The lower pitch of the male speaker makes it practicable to encompass vocal traits with a lower sampling frequency with fewer losses on the high-frequency bands*”. On a different note, Zhang et al. [153] utilized multiple sensors (i.e., inertial measurement unit (IMU), proximity, and ambient light sensors) for eating activity detection. They detected performance discrepancies for participants with large BMI, attributing them to data-related factors, such as “[...] *differences in movement patterns while eating, change in the distance of the proximity sensor from the neck, and difference in posture during the eating activity*”. It is crucial to understand that such biases are not straightforward to distinguish. On the contrary, they remain hidden, blended into time-series signals, or, even worse, propagated to inferred high-level signals such as steps, physical activity, and sleep. Such challenges call for additional care and demographic meta-data for post hoc fairness analysis.

Takeaway #7

Measurement inaccuracies and concept drift phenomena in captured audio, video, image, and sensor signals lead to performance discrepancies across sensitive attributes.

5.5 How WEIRD is UbiComp?

Inspired by previous call-to-action papers appearing in other communities about diversity in datasets and sample demographics [75, 113], we performed an analysis of UbiComp datasets with regard to sensitive attribute

distributions. We chose the WEIRD acronym coined by Henrich et al. [50] as a starting point to inspect how Western, Educated, Industrialized, Rich, or Democratic UbiComp really is.

Figure 5.5 presents the results of our analysis. The bar plot on the left (a) shows the percentage of papers reporting certain sensitive attributes: gender, age, race, education, employment, and sample country. The bar plot on the right (b) gives the mean sample size within each subset of papers reporting that sensitive attribute. Note that the mean was used, despite its sensitivity to outliers, because the median tends to be less accurate and more biased than the mean when sample sizes are small. In particular, gender was the most reported sensitive attribute in IMWUT datasets. 86% of included papers ($N = 42$) disclosed this information. Within these 42 papers, the mean sample size was 136 users, and the mean number of females in the sample was 71. Evidently, the community has made a step in the right direction by engaging in a conscious effort to achieve more balanced, diverse, and representative datasets. Yet, there is plenty of room for improvement. The second most-reported attribute was only reported in 4 out of 10 papers ($N = 20$), and the mean sample age within these papers was 31 years old. For comparison, the maximum mean age reported was 66 years. In a world that is rapidly aging [106], the UbiComp community is predominantly developing and testing on young populations. Such a finding is in line with prior work, reporting that within the fall detection domain, for example, datasets usually comprise imitated falls performed by younger people while they are intended for deployment on older people [126]. Finally, only 12% of included papers ($N = 6$) mentioned the participants' race. Within this subset, the mean sample size was 112 users, 88 of which were White (79%). Even though these numbers should be taken with a grain of salt, due to the small number of papers, it is worth pointing out that not only race is a wildly overlooked sensitive attribute (see results in Section 5.1.2), but non-White populations are significantly underrepresented within UbiComp datasets. There is a risk that models underperform for non-White users, but this fact may go unnoticed as there is no effort to check for it. For instance, Zhang et al. [151] could not assess the impact of skin tone on AFib detection “due to the unbalanced dataset where the majority (88.7%) of participants were White”.

Confirming IMWUT's WEIRDness, out of the 42 papers for which the participants' country (a proxy for Western) is reported or can be inferred, 36 (86%) engaged with US samples. China (26%) and Switzerland (7%) completed the top-3 of country representation. Note that the percentages do not sum up to 100 because of papers with more than one sample country. Concerning education (a proxy for Educated), 20% of included papers ($N = 10$) reported relevant information, with a median sample size of 144 users, 92 of which were college-educated (80%). This is perhaps not surprising, as in the early stages of development in UbiComp, participant recruitment frequently takes place within the universities and from the researchers' close circle. Regarding employment status (a proxy for Industrialized and Rich), it was only reported in 14% of included papers ($N = 7$), with a mean sample size of 208 users, of which 152 were employed.

Digging deeper into IMWUT's WEIRDness, we found that out of the 49 included papers, 32 (65%) included at least one author with a US-based affiliation. To evaluate whether the participants' countries are more diverse than the authors' locations, we analyzed the author affiliations reported in the 42 articles that also contained information about the participants' countries. Out of those, only 3 papers (7%) recruited at least part of their sample from a country different than the authors' location. In the remaining cases, participants were from the same country as at least one of the affiliated institutions. These results demonstrate that the vast majority of IMWUT authors (93%) recruit samples within the country they are located. This proportion is in line—even though larger—with similar analyses in other communities, such as CHI [75]. This is possibly due to the requirements of in-the-lab or artifact-based research common within the UbiComp community.

Takeaway #8

While UbiComp populations are balanced in terms of gender, they are otherwise predominantly young, White, western, highly educated, and employed, calling for more diverse sample recruitment.

6 DISCUSSION

In the following, we discuss our main findings (§6.1) and our work’s implications and recommendations for achieving “Fairness by design” in UbiComp work (§6.2).

6.1 Main Findings


By screening 523 papers published at IMWUT between 2018 and 2022, we found that only a small portion of 5% adhered to fairness reporting, while the overwhelming majority thereof focused on accuracy or error metrics. By delving into the smaller number of 49 papers, we surfaced biases in machine learning data and models across several sensitive attributes and application domains that would otherwise remain scattered in the UbiComp literature. Yet, the identified lack of diverse datasets in IMWUT publications could result in biases remaining undetected in the absence of heterogeneous demographics. To quantify such biases, included papers primarily employed performance evaluation instead of fairness metrics, while challenges in fairness assessment were found in regression and multi-class classification scenarios. Similar to other communities, defining fairness in UbiComp was not a simple task and involved considering its sociotechnical context, its ethical risks, and opportunities. Nevertheless, in an effort to employ fairness in practice —sensitive attributes aside— we found that the community has been striving for generalizability through ablation studies, real-world deployments, and personalization.

6.2 Implications and Recommendations

Drawing from these findings and borrowing from the “Privacy by design” literature [17, 35], we propose a “Fairness by design” equivalent, requiring AI developers and researchers to consider data and model fairness concerns from the very beginning of any AI project or system design. It is, thus, a proactive and preventative approach that prioritizes fairness as a core value in the development and implementation of UbiComp technologies, products, and services. To facilitate the community achieving “Fairness by design”, we next discuss recommendations for integrating fairness into the entire ML pipeline of UbiComp studies. These recommendations span two fronts, one concerning the data and the other the model.

Data Collection. Prior to the problem definition, researchers should identify the types of fairness-related harms relevant to their work (e.g., quality-of-service, allocation, stereotyping, and erasure harms [26]). For example, in an AFib detection application, quality-of-service harms could occur if the model had a substantially different performance for different ages, while allocation harms could occur if such difference led to one group unfairly receiving better care than another. Additionally, it is important to consider the demographic groups —including historically marginalized groups (e.g., based on gender, race, and ethnicity)— that might be harmed. We should also consider groups that are relevant to a particular scenario or deployment setting. For example, in a depression screening application, gender could be relevant as a sensitive attribute due to reported gender differences in the disorder’s signals [107]. Relevant attributes can be identified either from the theoretical literature or through fairness literature related to the target application domain.

When defining the problem statement, researchers should also have in mind the generalizability of the prediction task (e.g., applicable across demographic groups). To achieve that, it is of prime importance to consider an adequate enough sample size that would enable fairness to be studied (e.g., through sub-group analyses). For example, an AFib detection system showed poor performance in people with abnormal heart rhythms other



DATA RECOMMENDATIONS	MODEL RECOMMENDATIONS
Identify potential biases and demographic groups that might be harmed	Mitigate training model bias using ML methods that adjust for group fairness by using task-specific objectives and constraints
Enrich UbiComp (sensor) datasets with diverse sensitive attributes	Consider indirect notions of fairness, such as unfair resource allocation (energy, connectivity) in federated learning, when evaluating emerging ML paradigms
Ensure annotators' diversity during ground truth data labeling	Evaluate model performance across different groups of sensitive attributes using multiple fairness metrics
Investigate potential sources of measurement error from devices used as objective inputs	Generate synthetic data covering every sensitive attribute and potential intersections
Use data validation and visualization methods across sensitive attributes to surface potential data anomalies	Monitor the performance of deployed UbiComp models and devices and adjust for data and fairness drift

Fig. 8. **Recommendations for “Fairness by design” in UbiComp.** Actions to be taken by researchers for performing fairness assessments in both data and models in UbiComp works. Fairness needs to be considered from the very start of a project.

than AFib, most likely because its data annotation scheme assigned Normal sinus rhythm (NSR) and other types of heart rhythms to the same Non-AFib category (i.e., binary classification), due to the limited number of subjects with different types of heart rhythms [153]. Additionally, datasets in UbiComp are either self-collected or well-established benchmarks (e.g., those found in the UCI repository¹¹) used to evaluate new models. For self-collected data, researchers should strive for a diverse representation of human participants in both the recruitment and the data annotation phase. Considering that the models encode the biases of the labels, they should not only be assessed by multiple people to ensure agreement but also strive for demographic diversity amongst them. For benchmark data, researchers should think carefully about the pre-processing stage. Unlike other fields where the datasets are provided out-of-the-box, in UbiComp, it is not uncommon to require further slicing or windowing in order to be used for predictive modeling. For example, applying a sliding window method can generate thousands of samples from a sensor signal that belongs to a single user. This carries the risk of providing virtually “enough” samples for training, which, however, come from a handful of users. As a result, the model does not learn generalizable patterns.

As with any data science project, data validation methods play an important role in ensuring the results’ robustness. The same holds true when it comes to fairness. Typical data validation methods, therefore, should also be applied across sensitive attributes. For example, inspecting outliers that can consistently fall into particular demographic groups, data that are not missing at random and affect certain groups, or other kinds of data anomalies (e.g., measurement error due to a device malfunction or device differences). Regarding the latter, devices such as smartwatches offer model-based estimates for many well-being features. For instance, the measurement error of a heart-rate prediction model can propagate to every downstream application. If the original device has not been validated across different groups, this can affect every possible application that

¹¹UCI Machine Learning Repository: <https://archive.ics.uci.edu/ml/index.php>

is using such data. More broadly, visualization tools (e.g., What-If Tool¹², FairLens¹³, Tensorflow's Fairness Indicators¹⁴) may help surface any potential data anomalies and help correct them before they creep into models.

At the same time, the community itself could implement a mandatory data statement policy, requiring authors to report sensitive attributes concerning their participant samples. This builds on recent quests that advocate for data excellence [112], for example, by making data statements and datasheets for datasets mandatory for authors submitting their work.

Model Training and Evaluation. Until recently, most ML-based UbiComp applications employed some sort of feature engineering in order to extract statistical summaries from sensor data. However, during the past couple of years, this step has been automated since we have witnessed a remarkable consolidation of deep-learning models and architectures such as Convolutional Neural Networks [65] and Transformers [131]. As a result, such models are being used as generic feature extractors for different data types – be it images, text, time-series, or video. A side effect of this consolidation is that recently proposed mitigation methods can be applied across a wide range of models, regardless of input data types. For example, one approach modifies the weights of training samples or changes features and labels based on these attributes [14]. Another approach learns fair representations that remove correlations between sensitive and non-sensitive attributes [149], while a third approach involves dividing the training data into subgroups and modifying them to have similar feature distributions across subgroups [32]. Additionally, some methods operate on the latent space of the models by obfuscating information about protected attributes [149]. Overall, these techniques aim to promote fairness and reduce bias by focusing on various aspects of data and model architectures.

Yet enhancing fairness in machine learning requires a means to quantify it. As UbiComp systems blend into the real world, we realize that single evaluation metrics struggle to reflect the success criteria of ML models. As such, monitoring a multitude of metrics becomes the norm, and this is where we believe that monitoring and reporting fairness metrics across different groups should become standard practice. However, we acknowledge that sometimes it might not be feasible to collect data from representative demographics, especially for smaller pilot studies. In these situations, researchers should aim for a diverse user sample based on assumptions about relevant sensitive attributes. This approach can help uncover potential biases in the data and models, which can then be addressed in later stages of development. Alternatively, researchers can leverage advances in generative models to synthesize data covering multiple sensitive attributes and potential intersections [19].

This is where the concept of intersectional fairness comes in. Intersectional fairness means designing and training algorithms to account for the complex ways that different social identities can intersect and impact a person's experiences and outcomes. UbiComp technologies for diagnosing heart disease and monitoring vital signs provide an exemplary case. As reports suggest, differences in coronary heart disease are based on gender [83], socioeconomic status [114], and race [33]. In such cases, it is important to ensure that the models do not perpetuate existing biases and inequalities by failing to account for intersectional differences in health outcomes and access to healthcare—the biases encountered by a Black woman from a low socioeconomic background may not be the same as those experienced by a White woman from a high socioeconomic background.

Beyond traditional notions of fairness, such as directly discriminating based on sensitive attributes, we should also consider indirect notions of fairness. For example, within the paradigm of distributed/federated learning, the resource allocation of participating devices may also reflect the demographic and socio-economic information of owners, which makes the exclusion of such clients unfair in terms of participation. Cheaper devices cannot support the execution of large models and are either excluded or dropped together with their unique data [22, 51].

¹²<https://pair-code.github.io/what-if-tool/>

¹³<https://www.synthesized.io/fairlens>

¹⁴<https://github.com/tensorflow/fairness-indicators>

Last, as models are being deployed in real applications, we should monitor their performance in real time and adjust for data and fairness drift [42] by ensuring that models produce fair predictions independent of changes in input data and demographics.

7 CONCLUSION

The field of mobile, wearable and ubiquitous computing (UbiComp) faces significant challenges in ensuring fairness in the development of ML-based UbiComp technologies. Although efforts have been made to address biases, only a small percentage of publications in the Proceedings of the ACM IMWUT journal focus on fairness reporting and enhancement mechanisms. Sensitive attributes such as race, nationality, and language are often overlooked, while it is evident that there is a need for more diverse sample recruitment to ensure that the benefits of these technologies are shared equally across all members of society. The lack of a universal fairness definition, metric, or “fair” threshold that applies to different applications poses a sociotechnical challenge. UbiComp researchers must be explicit and transparent about their fairness priorities, definitions, and assumptions, making trade-offs between competing priorities, ethical risks, and opportunities. Despite these challenges, the UbiComp community strives for “fairer” models by conducting and reporting ablation studies, in-the-wild vs. in-the-lab experiments, and personalized model development. Ultimately, the UbiComp community must continue to prioritize fairness to ensure that the development of these technologies leads to just and equitable outcomes.

ACKNOWLEDGMENTS

This project has received funding from the European Union’s Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No 813162. The content of this paper reflects only the authors’ view and the Agency and the Commission are not responsible for any use that may be made of the information it contains.

REFERENCES

- [1] Ashraf Abdul, Jo Vermeulen, Danding Wang, Brian Y. Lim, and Mohan Kankanhalli. 2018. Trends and Trajectories for Explainable, Accountable and Intelligible Systems: An HCI Research Agenda. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems* (Montreal QC, Canada) (CHI '18). Association for Computing Machinery, New York, NY, USA, 1–18. <https://doi.org/10.1145/3173574.3174156>
- [2] Alekh Agarwal, Miroslav Dudik, and Zhiwei Steven Wu. 2019. Fair regression: Quantitative definitions and reduction-based algorithms. In *International Conference on Machine Learning*. PMLR, 120–129.
- [3] Tousif Ahmed, Apu Kapadia, Venkatesh Potluri, and Manohar Swaminathan. 2018. Up to a Limit? Privacy Concerns of Bystanders and Their Willingness to Share Additional Information with Visually Impaired Users of Assistive Technologies. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 2, 3, Article 89 (sep 2018), 27 pages. <https://doi.org/10.1145/3264899>
- [4] Julia Angwin, Jeff Larson, Surya Mattu, and Lauren Kirchner. 2020. There’s software used across the country to predict future criminals and it’s biased against blacks. 2016.
- [5] Olivier Augereau, Charles Lima Sanches, Koichi Kise, and Kai Kunze. 2018. Wordometer Systems for Everyday Life. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 1, 4, Article 123 (jan 2018), 21 pages. <https://doi.org/10.1145/3161601>
- [6] Yazeed Awwad, Richard Fletcher, Daniel Frey, Amit Gandhi, Maryam Najafian, and Mike Teodorescu. 2020. *Exploring fairness in machine learning for international development*. Technical Report. CITE MIT D-Lab.
- [7] Liliana Barrios, Pietro Oldrati, Marc Hilty, David Lindlbauer, Christian Holz, and Andreas Lutterotti. 2021. Smartphone-Based Tapping Frequency as a Surrogate for Perceived Fatigue: An in-the-Wild Feasibility Study in Multiple Sclerosis Patients. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 5, 3, Article 89 (sep 2021), 30 pages. <https://doi.org/10.1145/3478098>
- [8] Cody Blakeney, Gentry Atkinson, Nathaniel Huish, Yan Yan, Vangelis Metris, and Ziliang Zong. 2021. Measure twice, cut once: Quantifying bias and fairness in deep neural networks. *arXiv preprint arXiv:2110.04397* (2021).
- [9] Cody Blakeney, Nathaniel Huish, Yan Yan, and Ziliang Zong. 2021. Simon says: Evaluating and mitigating bias in pruned neural networks with knowledge distillation. *arXiv preprint arXiv:2106.07849* (2021).
- [10] Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. 2016. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. *Advances in neural information processing systems* 29 (2016).
- [11] Chloë Brown, Jagmohan Chauhan, Andreas Grammenos, Jing Han, Apinan Hasthanasombat, Dimitris Spathis, Tong Xia, Pietro Cicuta, and Cecilia Mascolo. 2020. Exploring Automatic Diagnosis of COVID-19 from Crowdsourced Respiratory Sound Data. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* (Virtual Event, CA, USA) (KDD '20). Association for Computing Machinery, New York, NY, USA, 3474–3484. <https://doi.org/10.1145/3394486.3412865>
- [12] Joseph M Bumgarner, Cameron T Lambert, Ayman A Hussein, Daniel J Cantillon, Bryan Baranowski, Kathy Wolski, Bruce D Lindsay, Oussama M Wazni, and Khalidoun G Tarakji. 2018. Smartwatch algorithm for automated detection of atrial fibrillation. *Journal of the American College of Cardiology* 71, 21 (2018), 2381–2388.
- [13] Joy Buolamwini and Timnit Gebru. 2018. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on fairness, accountability and transparency*. PMLR, 77–91.
- [14] 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).
- [15] Alessandro Castelnovo, Riccardo Crupi, Greta Greco, Daniele Regoli, Ilaria Giuseppina Penco, and Andrea Claudio Cosentini. 2022. A clarification of the nuances in the fairness metrics landscape. *Scientific Reports* 12, 1 (2022), 4209.
- [16] Simon Caton and Christian Haas. 2020. Fairness in machine learning: A survey. *arXiv preprint arXiv:2010.04053* (2020).
- [17] Ann Cavoukian et al. 2009. Privacy by design: The 7 foundational principles. *Information and privacy commissioner of Ontario, Canada* 5 (2009), 12.
- [18] Hongyan Chang and Reza Shokri. 2021. On the privacy risks of algorithmic fairness. In *2021 IEEE European Symposium on Security and Privacy (EuroS&P)*. IEEE, 292–303.
- [19] Bhushan Chaudhari, Himanshu Choudhary, Aakash Agarwal, Kamna Meena, and Tanmoy Bhowmik. 2022. FairGen: Fair Synthetic Data Generation. *arXiv preprint arXiv:2210.13023* (2022).
- [20] Yang Chen, Katherine Fennedy, Anna Fogel, Shengdong Zhao, Chao Zhang, Lijuan Liu, and Chingchuan Yen. 2022. SSpoon: A Shape-Changing Spoon That Optimizes Bite Size for Eating Rate Regulation. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 6, 3,

- Article 105 (sep 2022), 32 pages. <https://doi.org/10.1145/3550312>
- [21] Anshuman Chhabra, Karina Masalkovaitė, and Prasant Mohapatra. 2021. An overview of fairness in clustering. *IEEE Access* 9 (2021), 130698–130720.
- [22] Hyunsung Cho, Akhil Mathur, and Fahim Kawsar. 2022. Flame: Federated learning across multi-device environments. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 6, 3 (2022), 1–29.
- [23] Manvi Choudhary, Charlotte Laclau, and Christine Largeron. 2022. A survey on fairness for machine learning on graphs. *arXiv preprint arXiv:2205.05396* (2022).
- [24] Marios Constantinides and Daniele Quercia. 2022. Good Intentions, Bad Inventions: How Employees Judge Pervasive Technologies in the Workplace. *IEEE Pervasive Computing* (2022).
- [25] Jean Costa, François Guimbretière, Malte F Jung, and Tanzeem Choudhury. 2019. Boostmeup: Improving cognitive performance in the moment by unobtrusively regulating emotions with a smartwatch. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 3, 2 (2019), 1–23.
- [26] Kate Crawford. 2017. The trouble with bias. keynote at neurips. (2017).
- [27] Ananyananda Dasari, Sakthi Kumar Arul Prakash, László A Jeni, and Conrad S Tucker. 2021. Evaluation of biases in remote photoplethysmography methods. *NPJ digital medicine* 4, 1 (2021), 91.
- [28] Thomas G Dietterich. 1998. Approximate statistical tests for comparing supervised classification learning algorithms. *Neural computation* 10, 7 (1998), 1895–1923.
- [29] Dheeru Dua and Casey Graff. 2017. UCI Machine Learning Repository–German Credit.
- [30] 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.
- [31] EU EU et al. 2012. Charter of fundamental rights of the European Union. *The Review of International Affairs* 63, 1147 (2012), 109–123.
- [32] 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.
- [33] Contessa Fincher, Joyce E Williams, Vicky MacLean, Jeroan J Allison, Catarina I Kiefe, and John Canto. 2004. Racial disparities in coronary heart disease. *Ethnicity & disease* 14, 3 (2004), 360–371.
- [34] Thomas B Fitzpatrick. 1988. The validity and practicality of sun-reactive skin types I through VI. *Archives of dermatology* 124, 6 (1988), 869–871.
- [35] Jessica Fjeld, Nele Achten, Hannah Hilligoss, Adam Nagy, and Madhulika Srikumar. 2020. Principled artificial intelligence: Mapping consensus in ethical and rights-based approaches to principles for AI. *Berkman Klein Center Research Publication 2020-1* (2020).
- [36] Will Fleisher. 2021. What’s fair about individual fairness?. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*. 480–490.
- [37] Luciano Floridi, Josh Cowsls, Monica Beltrametti, Raja Chatila, Patrice Chazerand, Virginia Dignum, Christoph Luetge, Robert Madelin, Ugo Pagallo, Francesca Rossi, et al. 2018. AI4People—an ethical framework for a good AI society: opportunities, risks, principles, and recommendations. *Minds and machines* 28 (2018), 689–707.
- [38] Sorelle A Friedler, Carlos Scheidegger, and Suresh Venkatasubramanian. 2021. The (im) possibility of fairness: Different value systems require different mechanisms for fair decision making. *Commun. ACM* 64, 4 (2021), 136–143.
- [39] Hana Frluckaj, Laura Dabbish, David Gray Widder, Huilian Sophie Qiu, and James D. Herbsleb. 2022. Gender and Participation in Open Source Software Development. *Proceedings of the ACM on Human-Computer Interaction* 6, CSCW2 (Nov 2022), 31 pages. <https://doi.org/10.1145/3555190>
- [40] Yang Gao, Yincheng Jin, Jagmohan Chauhan, Seokmin Choi, Jiyang Li, and Zhanpeng Jin. 2021. Voice In Ear: Spoofing-Resistant and Passphrase-Independent Body Sound Authentication. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 5, 1, Article 12 (mar 2021), 25 pages. <https://doi.org/10.1145/3448113>
- [41] Pratyush Garg, John Villasenor, and Virginia Foggo. 2020. Fairness metrics: A comparative analysis. In *2020 IEEE International Conference on Big Data (Big Data)*. IEEE, 3662–3666.
- [42] Avijit Ghosh, Aalok Shanbhag, and Christo Wilson. 2022. Faircanary: Rapid continuous explainable fairness. In *Proceedings of the 2022 AAAI/ACM Conference on AI, Ethics, and Society*. 307–316.
- [43] Yanbin Gong, Qian Zhang, Bobby H.P. NG, and Wei Li. 2022. BreathMentor: Acoustic-Based Diaphragmatic Breathing Monitor System. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 6, 2, Article 53 (jul 2022), 28 pages. <https://doi.org/10.1145/3534595>
- [44] Erin Griffiths, Salah Assana, and Kamin Whitehouse. 2018. Privacy-Preserving Image Processing with Binocular Thermal Cameras. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 1, 4, Article 133 (jan 2018), 25 pages. <https://doi.org/10.1145/3161198>
- [45] Fuqiang Gu, Mu-Huan Chung, Mark Chignell, Shahrokh Valaei, Baoding Zhou, and Xue Liu. 2021. A survey on deep learning for human activity recognition. *ACM Computing Surveys (CSUR)* 54, 8 (2021), 1–34.
- [46] Bhanu Teja Gullapalli, Stephanie Carreiro, Brittany P. Chapman, Deepak Ganesan, Jan Sjoquist, and Tauhidur Rahman. 2021. OpiTrack: A Wearable-Based Clinical Opioid Use Tracker with Temporal Convolutional Attention Networks. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 5, 3, Article 102 (sep 2021), 29 pages. <https://doi.org/10.1145/3478107>

- [47] Bhanu Teja Gullapalli, Annamalai Natarajan, Gustavo A. Angarita, Robert T. Malison, Deepak Ganesan, and Tauhidur Rahman. 2019. On-Body Sensing of Cocaine Craving, Euphoria and Drug-Seeking Behavior Using Cardiac and Respiratory Signals. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 3, 2, Article 46 (jun 2019), 31 pages. <https://doi.org/10.1145/3328917>
- [48] Adnan Hadzi and Denis Roio. 2019. Restorative Justice in Artificial Intelligence Crimes. *spheres: Journal for Digital Cultures* 5 (2019), 1–18.
- [49] Lucy Havens, Melissa Terras, Benjamin Bach, and Beatrice Alex. 2020. Situated Data, Situated Systems: A Methodology to Engage with Power Relations in Natural Language Processing Research. In *Proceedings of the Second Workshop on Gender Bias in Natural Language Processing*. Association for Computational Linguistics, 107–124. <https://aclanthology.org/2020.gebnlp-1.10>
- [50] Joseph Henrich, Steven J Heine, and Ara Norenzayan. 2010. The Weirdest People in the World? *Behavioral and Brain Sciences* 33, 2-3 (2010), 61–83.
- [51] Samuel Horvath, Stefanos Laskaridis, Mario Almeida, Ilias Leontiadis, Stylianos Venieris, and Nicholas Lane. 2021. Fjord: Fair and accurate federated learning under heterogeneous targets with ordered dropout. *Advances in Neural Information Processing Systems* 34 (2021), 12876–12889.
- [52] Wiebke Toussaint Hutiri and Aaron Yi Ding. 2022. Bias in automated speaker recognition. In *2022 ACM Conference on Fairness, Accountability, and Transparency*. 230–247.
- [53] Sinh Huynh, Rajesh Krishna Balan, and JeongGil Ko. 2022. IMon: Appearance-Based Gaze Tracking System on Mobile Devices. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 5, 4, Article 161 (dec 2022), 26 pages. <https://doi.org/10.1145/3494999>
- [54] James M Joyce. 2011. Kullback-leibler divergence. In *International encyclopedia of statistical science*. Springer, 720–722.
- [55] Manasa Kalanadhabhatta, Adrelys Mateo Santana, Zhongyang Zhang, Deepak Ganesan, Adam S. Grabell, and Tauhidur Rahman. 2022. EarlyScreen: Multi-Scale Instance Fusion for Predicting Neural Activation and Psychopathology in Preschool Children. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 6, 2, Article 60 (jul 2022), 39 pages. <https://doi.org/10.1145/3534583>
- [56] Louis Henry Kamulegeya, Mark Okello, John Mark Bwanika, Davis Musinguzi, William Lubega, Davis Rusoke, Faith Nassiwa, and Alexander Börve. 2019. Using artificial intelligence on dermatology conditions in Uganda: A case for diversity in training data sets for machine learning. *BioRxiv* (2019), 826057.
- [57] 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.
- [58] 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.
- [59] Mohammed Khwaja, Sumer S. Vaid, Sara Zannone, Gabriella M. Harari, A. Aldo Faisal, and Aleksandar Matic. 2019. Modeling Personality vs. Modeling Personalidad: In-the-Wild Mobile Data Analysis in Five Countries Suggests Cultural Impact on Personality Models. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 3, 3, Article 88 (sep 2019), 24 pages. <https://doi.org/10.1145/3351246>
- [60] Minhyung Kim, Inyeop Kim, and Uichin Lee. 2021. Beneficial Neglect: Instant Message Notification Handling Behaviors and Academic Performance. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 5, 1, Article 18 (mar 2021), 26 pages. <https://doi.org/10.1145/3448089>
- [61] Barbara Kitchenham and Stuart Charters. 2007. *Guidelines for performing systematic literature reviews in software engineering*. Technical Report. Keele University and Durham University.
- [62] Ronny Kohavi and Barry Becker. 1996. Uci machine learning repository: adult data set. *Availible*: <https://archive.ics.uci.edu/ml/machine-learning-databases/adult> (1996).
- [63] Heli Koskimäki, Hannu Kinnunen, Teemu Kurppa, and Juha Röning. 2018. How do we sleep: a case study of sleep duration and quality using data from oura ring. In *Proceedings of the 2018 ACM international joint conference and 2018 international symposium on pervasive and ubiquitous computing and wearable computers*. 714–717.
- [64] Kundan Krishna, Deepali Jain, Sanket V. Mehta, and Sunav Choudhary. 2018. An LSTM Based System for Prediction of Human Activities with Durations. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 1, 4, Article 147 (jan 2018), 31 pages. <https://doi.org/10.1145/3161201>
- [65] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. 2017. Imagenet classification with deep convolutional neural networks. *Commun. ACM* 60, 6 (2017), 84–90.
- [66] Matt J Kusner, Joshua Loftus, Chris Russell, and Ricardo Silva. 2017. Counterfactual fairness. *Advances in neural information processing systems* 30 (2017).
- [67] Julia Larson, Surya Mattu, Lauren Kirchner, and Julia Angwin. 2018. How We Analyzed the COMPAS Recidivism Algorithm. <https://www.propublica.org/article/how-we-analyzed-the-compas-recidivism-algorithm>
- [68] Benjamin Laufer, Sameer Jain, A. Feder Cooper, Jon Kleinberg, and Hoda Heidari. 2022. Four Years of FAccT: A Reflexive, Mixed-Methods Analysis of Research Contributions, Shortcomings, and Future Prospects. In *2022 ACM Conference on Fairness, Accountability, and Transparency (FAccT '22)*. ACM, 401–426. <https://doi.org/10.1145/3531146.3533107>
- [69] Tai Le Quy, Arjun Roy, Vasileios Iosifidis, Wenbin Zhang, and Eirini Ntoutsis. 2022. A survey on datasets for fairness-aware machine learning. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* (2022), e1452.

- [70] Nicol Turner Lee. 2018. Detecting racial bias in algorithms and machine learning. *Journal of Information, Communication and Ethics in Society* 16, 3 (2018), 252–260.
- [71] Dong Li, Jialin Liu, Sunghoon Ivan Lee, and Jie Xiong. 2022. LASense: Pushing the Limits of Fine-Grained Activity Sensing Using Acoustic Signals. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 6, 1, Article 21 (mar 2022), 27 pages. <https://doi.org/10.1145/3517253>
- [72] Yunqi Li, Hanxiong Chen, Shuyuan Xu, Yingqiang Ge, Juntao Tan, Shuchang Liu, and Yongfeng Zhang. 2022. Fairness in recommendation: A survey. *arXiv preprint arXiv:2205.13619* (2022).
- [73] Yunji Liang, Yuchen Qin, Qi Li, Xiaokai Yan, Zhiwen Yu, Bin Guo, Sagar Samtani, and Yanyong Zhang. 2022. AccMyrinx: Speech Synthesis with Non-Acoustic Sensor. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 6, 3, Article 127 (sep 2022), 24 pages. <https://doi.org/10.1145/3550338>
- [74] Daniyal Liaqat, Mohamed Abdalla, Pegah Abed-Esfahani, Moshe Gabel, Tatiana Son, Robert Wu, Andrea Gershon, Frank Rudzicz, and Eyal De Lara. 2019. WearBreathing: Real World Respiratory Rate Monitoring Using Smartwatches. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 3, 2, Article 56 (jun 2019), 22 pages. <https://doi.org/10.1145/3328927>
- [75] Sebastian Linxen, Christian Sturm, Florian Brühlmann, Vincent Cassau, Klaus Opwis, and Katharina Reinecke. 2021. How weird is CHI?. In *Proceedings of the 2021 chi conference on human factors in computing systems*. 1–14.
- [76] Haipeng Liu, Yuheng Wang, Anfu Zhou, Hanyue He, Wei Wang, Kunpeng Wang, Peilin Pan, Yixuan Lu, Liang Liu, and Huadong Ma. 2020. Real-Time Arm Gesture Recognition in Smart Home Scenarios via Millimeter Wave Sensing. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 4, 4, Article 140 (dec 2020), 28 pages. <https://doi.org/10.1145/3432235>
- [77] Tony Liu, Jennifer Nicholas, Max M. Theilig, Sharath C. Guntuku, Konrad Kording, David C. Mohr, and Lyle Ungar. 2020. Machine Learning for Phone-Based Relationship Estimation: The Need to Consider Population Heterogeneity. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 3, 4, Article 145 (sep 2020), 23 pages. <https://doi.org/10.1145/3369820>
- [78] Xin Liu, Yuntao Wang, Sinan Xie, Xiaoyu Zhang, Zixian Ma, Daniel McDuff, and Shwetak Patel. 2022. MobilePhys: Personalized Mobile Camera-Based Contactless Physiological Sensing. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 6, 1, Article 24 (mar 2022), 23 pages. <https://doi.org/10.1145/3517225>
- [79] Francesco Locatello, Gabriele Abbati, Thomas Rainforth, Stefan Bauer, Bernhard Schölkopf, and Olivier Bachem. 2019. On the fairness of disentangled representations. *Advances in neural information processing systems* 32 (2019).
- [80] Jie Lu, Anjin Liu, Fan Dong, Feng Gu, Joao Gama, and Guangquan Zhang. 2018. Learning under concept drift: A review. *IEEE transactions on knowledge and data engineering* 31, 12 (2018), 2346–2363.
- [81] Li Lu, Jiadi Yu, Yingying Chen, Yanmin Zhu, Minglu Li, and Xiangyu Xu. 2019. I3: Sensing Scrolling Human-Computer Interactions for Intelligent Interest Inference on Smartphones. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 3, 3, Article 97 (sep 2019), 22 pages. <https://doi.org/10.1145/3351255>
- [82] Steven A Lubitz, Anthony Z Faranesh, Caitlin Selvaggi, Steven J Atlas, David D McManus, Daniel E Singer, Sherry Pagoto, Michael V McConnell, Alexandros Pantelopoulos, and Andrea S Foulkes. 2022. Detection of atrial fibrillation in a large population using wearable devices: the Fitbit heart study. *Circulation* 146, 19 (2022), 1415–1424.
- [83] Angela HEM Maas and Yolande EA Appelman. 2010. Gender differences in coronary heart disease. *Netherlands Heart Journal* 18 (2010), 598–603.
- [84] Anna Maijala, Hannu Kinnunen, Heli Koskimäki, Timo Jämsä, and Maarit Kangas. 2019. Nocturnal finger skin temperature in menstrual cycle tracking: ambulatory pilot study using a wearable Oura ring. *BMC Women’s Health* 19, 1 (2019), 1–10.
- [85] Thomas Manzini, Yao Chong Lim, Yulia Tsvetkov, and Alan W Black. 2019. Black is to criminal as caucasian is to police: Detecting and removing multiclass bias in word embeddings. *arXiv preprint arXiv:1904.04047* (2019).
- [86] Alex Mariakakis, Edward Wang, Shwetak Patel, and Mayank Goel. 2019. Challenges in realizing smartphone-based health sensing. *IEEE Pervasive Computing* 18, 2 (2019), 76–84.
- [87] Nina Markl. 2022. Language variation and algorithmic bias: understanding algorithmic bias in British English automatic speech recognition. In *2022 ACM Conference on Fairness, Accountability, and Transparency*. 521–534.
- [88] Frank J Massey Jr. 1951. The Kolmogorov-Smirnov test for goodness of fit. *Journal of the American statistical Association* 46, 253 (1951), 68–78.
- [89] Lakmal Meegahapola, Florian Labhart, Thanh-Trung Phan, and Daniel Gatica-Perez. 2021. Examining the Social Context of Alcohol Drinking in Young Adults with Smartphone Sensing. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 5, 3, Article 121 (sep 2021), 26 pages. <https://doi.org/10.1145/3478126>
- [90] 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.
- [91] Tamir Mendel, Roei Schuster, Eran Tromer, and Eran Toch. 2022. Toward Proactive Support for Older Adults: Predicting the Right Moment for Providing Mobile Safety Help. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 6, 1, Article 25 (mar 2022), 25 pages. <https://doi.org/10.1145/3517249>
- [92] Michele Merler, Nalini Ratha, Rogerio S Feris, and John R Smith. 2019. Diversity in faces. *arXiv preprint arXiv:1901.10436* (2019).

- [93] Christian Meurisch, Cristina A. Mihale-Wilson, Adrian Hawlitschek, Florian Giger, Florian Müller, Oliver Hinz, and Max Mühlhäuser. 2020. Exploring User Expectations of Proactive AI Systems. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 4, 4, Article 146 (dec 2020), 22 pages. <https://doi.org/10.1145/3432193>
- [94] Vishwali Mhasawade, Yuan Zhao, and Rumi Chunara. 2021. Machine learning and algorithmic fairness in public and population health. *Nature Machine Intelligence* 3, 8 (2021), 659–666.
- [95] Mostafa M Mohamed and Björn W Schuller. 2022. Normalise for fairness: A simple normalisation technique for fairness in regression machine learning problems. *arXiv preprint arXiv:2202.00993* (2022).
- [96] David Moher, Alessandro Liberati, Jennifer Tetzlaff, Douglas G Altman, and the PRISMA Group*. 2009. Preferred reporting items for systematic reviews and meta-analyses: the PRISMA statement. *Annals of internal medicine* 151, 4 (2009), 264–269.
- [97] Debarghya Mukherjee, Mikhail Yurochkin, Moulinath Banerjee, and Yuekai Sun. 2020. Two Simple Ways to Learn Individual Fairness Metrics from Data. In *Proceedings of the 37th International Conference on Machine Learning (ICML'20)*. JMLR.org, Article 658, 11 pages.
- [98] Deirdre K Mulligan, Joshua A Kroll, Nitin Kohli, and Richmond Y Wong. 2019. This thing called fairness: Disciplinary confusion realizing a value in technology. *Proceedings of the ACM on Human-Computer Interaction* 3, CSCW (2019), 1–36.
- [99] Claude Nadeau and Yoshua Bengio. 1999. Inference for the generalization error. *Advances in neural information processing systems* 12 (1999).
- [100] Suraj Nair, Kiran Javkar, Jiahui Wu, and Vanessa Frias-Martinez. 2019. Understanding Cycling Trip Purpose and Route Choice Using GPS Traces and Open Data. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 3, 1, Article 20 (mar 2019), 26 pages. <https://doi.org/10.1145/3314407>
- [101] Arvind Narayanan. 21. Fairness definitions and their politics. In *Tutorial presented at the Conf. on Fairness, Accountability, and Transparency*.
- [102] United Nations. 2021. With almost half of world’s population still offline, digital divide risks becoming ‘new face of inequality’, Deputy Secretary-General warns general assembly.
- [103] Ewa M Nowara, Daniel McDuff, and Ashok Veeraraghavan. 2020. A meta-analysis of the impact of skin tone and gender on non-contact photoplethysmography measurements. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*. 284–285.
- [104] Alexandra Olteanu, Carlos Castillo, Fernando Diaz, and Emre Kıcıman. 2019. Social data: Biases, methodological pitfalls, and ethical boundaries. *Frontiers in big data* 2 (2019), 13.
- [105] Cathy O’neil. 2017. *Weapons of math destruction: How big data increases inequality and threatens democracy*. Crown.
- [106] World Health Organization. 2015. *World report on ageing and health*. World Health Organization.
- [107] Gordon Parker and Heather Brotchie. 2010. Gender differences in depression. *International review of psychiatry* 22, 5 (2010), 429–436.
- [108] Ignacio Perez-Pozuelo, Dimitris Spathis, Emma AD Clifton, and Cecilia Mascolo. 2021. Wearables, smartphones, and artificial intelligence for digital phenotyping and health. In *Digital Health*. Elsevier, 33–54.
- [109] Dana Pessach and Erez Shmueli. 2022. A Review on Fairness in Machine Learning. *ACM Computing Surveys (CSUR)* 55, 3 (2022), 1–44.
- [110] Yvonne Rogers, Kay Connolly, Lenore Tedesco, William Hazlewood, Andrew Kurtz, Robert E Hall, Josh Hursey, and Tammy Toscos Phd. 2007. Why it’s worth the hassle: The value of in-situ studies when designing ubicomp. *Springer* (2007), 336.
- [111] Sherry Ruan, Jacob O. Wobbrock, Kenny Liou, Andrew Ng, and James A. Landay. 2018. Comparing Speech and Keyboard Text Entry for Short Messages in Two Languages on Touchscreen Phones. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 1, 4, Article 159 (jan 2018), 23 pages. <https://doi.org/10.1145/3161187>
- [112] Nithya Sambasivan, Shivani Kapania, Hannah Highfill, Diana Akrong, Praveen Paritosh, and Lora M Aroyo. 2021. “Everyone wants to do the model work, not the data work”: Data Cascades in High-Stakes AI. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–15.
- [113] Ari Schlesinger, W Keith Edwards, and Rebecca E Grinter. 2017. Intersectional HCI: Engaging identity through gender, race, and class. In *Proceedings of the 2017 CHI conference on human factors in computing systems*. 5412–5427.
- [114] William M Schultz, Heval M Kelli, John C Lisko, Tina Varghese, Jia Shen, Pratik Sandesara, Arshed A Quyyumi, Herman A Taylor, Martha Gulati, John G Harold, et al. 2018. Socioeconomic status and cardiovascular outcomes: challenges and interventions. *Circulation* 137, 20 (2018), 2166–2178.
- [115] Andrew D Selbst, Danah Boyd, Sorelle A Friedler, Suresh Venkatasubramanian, and Janet Vertesi. 2019. Fairness and abstraction in sociotechnical systems. In *Proceedings of the conference on fairness, accountability, and transparency*. 59–68.
- [116] Farhana Shahid, Wasifur Rahman, M. Saifur Rahman, Sharmin Akther Purabi, Ayesha Seddiqa, Moin Mostakim, Farhan Feroz, Tanjir Rashid Soron, Fahmida Hossain, Nabila Khan, Anika Binte Islam, Nipi Paul, Ehsan Hoque, and A. B. M. Alim Al Islam. 2020. Leveraging Free-Hand Sketches for Potential Screening of PTSD. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 4, 3, Article 94 (sep 2020), 22 pages. <https://doi.org/10.1145/3411835>
- [117] Muhammad Shahzad and Shaohu Zhang. 2018. Augmenting User Identification with WiFi Based Gesture Recognition. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 2, 3, Article 134 (sep 2018), 27 pages. <https://doi.org/10.1145/3264944>

- [118] Kshitij Sharma, Evangelos Niforatos, Michail Giannakos, and Vassilis Kostakos. 2020. Assessing Cognitive Performance Using Physiological and Facial Features: Generalizing across Contexts. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 4, 3, Article 95 (sep 2020), 41 pages. <https://doi.org/10.1145/3411811>
- [119] Taoran Sheng and Manfred Huber. 2020. Weakly Supervised Multi-Task Representation Learning for Human Activity Analysis Using Wearables. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 4, 2, Article 57 (jun 2020), 18 pages. <https://doi.org/10.1145/3397330>
- [120] Chen-Hsuan (Iris) Shih, Naofumi Tomita, Yanick X. Lukic, Álvaro Hernández Reguera, Elgar Fleisch, and Tobias Kowatsch. 2020. Breeze: Smartphone-Based Acoustic Real-Time Detection of Breathing Phases for a Gamified Biofeedback Breathing Training. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 3, 4, Article 152 (sep 2020), 30 pages. <https://doi.org/10.1145/3369835>
- [121] Michael W Sjoding, Robert P Dickson, Theodore J Iwashyna, Steven E Gay, and Thomas S Valley. 2020. Racial bias in pulse oximetry measurement. *New England Journal of Medicine* 383, 25 (2020), 2477–2478.
- [122] Genevieve Smith. 2020. What does “fairness” mean for machine learning systems? https://haas.berkeley.edu/wp-content/uploads/What-is-fairness_-EGAL2.pdf
- [123] Dimitris Spathis, Ignacio Perez-Pozuelo, Soren Brage, Nicholas J Wareham, and Cecilia Mascolo. 2021. Self-supervised transfer learning of physiological representations from free-living wearable data. In *Proceedings of the Conference on Health, Inference, and Learning*. 69–78.
- [124] Jie Su, Zhenyu Wen, Tao Lin, and Yu Guan. 2022. Learning Disentangled Behaviour Patterns for Wearable-Based Human Activity Recognition. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 6, 1, Article 28 (mar 2022), 19 pages. <https://doi.org/10.1145/3517252>
- [125] Weigao Su, Daibo Liu, Taiyuan Zhang, and Hongbo Jiang. 2022. Towards Device Independent Eavesdropping on Telephone Conversations with Built-in Accelerometer. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 5, 4, Article 177 (dec 2022), 29 pages. <https://doi.org/10.1145/3494969>
- [126] Angela Sucerquia, José David López, and Jesús Francisco Vargas-Bonilla. 2017. SisFall: A fall and movement dataset. *Sensors* 17, 1 (2017), 198.
- [127] Tony Sun, Andrew Gaut, Shirlyn Tang, Yuxin Huang, Mai ElSherief, Jieyu Zhao, Diba Mirza, Elizabeth Belding, Kai-Wei Chang, and William Yang Wang. 2019. Mitigating gender bias in natural language processing: Literature review. *arXiv preprint arXiv:1906.08976* (2019).
- [128] The, U. S. Equal Employment Opportunity Commission (EEOC). 1979. Uniform Guidelines on employee selection procedures.
- [129] Wiebke Toussaint, Akhil Mathur, Aaron Yi Ding, and Fahim Kawsar. 2021. Characterising the Role of Pre-Processing Parameters in Audio-Based Embedded Machine Learning. In *Proceedings of the 19th ACM Conference on Embedded Networked Sensor Systems* (Coimbra, Portugal) (*SenSys '21*). Association for Computing Machinery, New York, NY, USA, 439–445. <https://doi.org/10.1145/3485730.3493448>
- [130] Daniel Van Ness and Karen Heetderks Strong. 2014. *Restoring justice: An introduction to restorative justice*. Routledge.
- [131] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).
- [132] Lev Velykoivanenko, Kavous Salehzadeh Niksirat, Noé Zufferey, Mathias Humbert, Kévin Huguenin, and Mauro Cherubini. 2022. Are Those Steps Worth Your Privacy? Fitness-Tracker Users’ Perceptions of Privacy and Utility. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 5, 4, Article 181 (dec 2022), 41 pages. <https://doi.org/10.1145/3494960>
- [133] Raghav H. Venkatnarayan and Muhammad Shahzad. 2018. Gesture Recognition Using Ambient Light. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 2, 1, Article 40 (mar 2018), 28 pages. <https://doi.org/10.1145/3191772>
- [134] Mingyang Wan, Daochen Zha, Ninghao Liu, and Na Zou. 2022. In-Processing Modeling Techniques for Machine Learning Fairness: A Survey. *ACM Transactions on Knowledge Discovery from Data (TKDD)* (2022).
- [135] Liang Wang, Wen Cheng, Lijia Pan, Tao Gu, Tianheng Wu, Xianping Tao, and Jian Lu. 2018. SpiderWalk: Circumstance-Aware Transportation Activity Detection Using a Novel Contact Vibration Sensor. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 2, 1, Article 42 (mar 2018), 30 pages. <https://doi.org/10.1145/3191774>
- [136] Lei Wang, Wei Li, Ke Sun, Fusang Zhang, Tao Gu, Chenren Xu, and Daqing Zhang. 2022. LoEar: Push the Range Limit of Acoustic Sensing for Vital Sign Monitoring. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 6, 3, Article 145 (sep 2022), 24 pages. <https://doi.org/10.1145/3550293>
- [137] Rui Wang, Weichen Wang, Alex daSilva, Jeremy F. Huckins, William M. Kelley, Todd F. Heatherston, and Andrew T. Campbell. 2018. Tracking Depression Dynamics in College Students Using Mobile Phone and Wearable Sensing. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 2, 1, Article 43 (mar 2018), 26 pages. <https://doi.org/10.1145/3191775>
- [138] Xiyue Wang, Kazuki Takashima, Tomoaki Adachi, Patrick Finn, Ehud Sharlin, and Yoshifumi Kitamura. 2020. AssessBlocks: Exploring Toy Block Play Features for Assessing Stress in Young Children after Natural Disasters. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 4, 1, Article 30 (mar 2020), 29 pages. <https://doi.org/10.1145/3381016>
- [139] Chatchai Wangwiwattana, Xinyi Ding, and Eric C. Larson. 2018. PupilNet, Measuring Task Evoked Pupillary Response Using Commodity RGB Tablet Cameras: Comparison to Mobile, Infrared Gaze Trackers for Inferring Cognitive Load. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 1, 4, Article 171 (jan 2018), 26 pages. <https://doi.org/10.1145/3161164>

- [140] Justin C. Wilson, Suku Nair, Sandro Scielzo, and Eric C. Larson. 2021. Objective Measures of Cognitive Load Using Deep Multi-Modal Learning: A Use-Case in Aviation. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 5, 1, Article 40 (mar 2021), 35 pages. <https://doi.org/10.1145/3448111>
- [141] Blake Woodworth, Suriya Gunasekar, Mesrob I Ohannessian, and Nathan Srebro. 2017. Learning non-discriminatory predictors. In *Conference on Learning Theory*. PMLR, 1920–1953.
- [142] Chenhan Xu, Huining Li, Zhengxiong Li, Hanbin Zhang, Aditya Singh Rathore, Xingyu Chen, Kun Wang, Ming-chun Huang, and Wenyao Xu. 2021. CardiacWave: A MmWave-Based Scheme of Non-Contact and High-Definition Heart Activity Computing. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 5, 3, Article 135 (sep 2021), 26 pages. <https://doi.org/10.1145/3478127>
- [143] Fengli Xu, Zongyu Lin, Tong Xia, Diansheng Guo, and Yong Li. 2020. SUME: Semantic-Enhanced Urban Mobility Network Embedding for User Demographic Inference. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 4, 3, Article 98 (sep 2020), 25 pages. <https://doi.org/10.1145/3411807>
- [144] Jie Xu, Yunyu Xiao, Wendy Hui Wang, Yue Ning, Elizabeth A Shenkman, Jiang Bian, and Fei Wang. 2022. Algorithmic fairness in computational medicine. *EBioMedicine* 84 (2022), 104250.
- [145] Wei Xu, ZhiWen Yu, Zhu Wang, Bin Guo, and Qi Han. 2019. AcousticID: Gait-Based Human Identification Using Acoustic Signal. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 3, 3, Article 115 (sep 2019), 25 pages. <https://doi.org/10.1145/3351273>
- [146] Joel Xue and Robert M Farrell. 2014. How can computerized interpretation algorithms adapt to gender/age differences in ECG measurements? *Journal of Electrocardiology* 47, 6 (2014), 849–855.
- [147] Samuel Yeom and Michael Carl Tschantz. 2021. Avoiding Disparity Amplification under Different Worldviews. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (Virtual Event, Canada) (FAccT '21)*. Association for Computing Machinery, New York, NY, USA, 273–283. <https://doi.org/10.1145/3442188.3445892>
- [148] Shichao Yue, Yuzhe Yang, Hao Wang, Hariharan Rahul, and Dina Katabi. 2020. BodyCompass: Monitoring Sleep Posture with Wireless Signals. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 4, 2, Article 66 (jun 2020), 25 pages. <https://doi.org/10.1145/3397311>
- [149] Rich Zemel, Yu Wu, Kevin Swersky, Toni Pitassi, and Cynthia Dwork. 2013. Learning fair representations. In *International conference on machine learning*. PMLR, 325–333.
- [150] Hanbin Zhang, Chenhan Xu, Huining Li, Aditya Singh Rathore, Chen Song, Zhisheng Yan, Dongmei Li, Feng Lin, Kun Wang, and Wenyao Xu. 2019. PDMove: Towards Passive Medication Adherence Monitoring of Parkinson’s Disease Using Smartphone-Based Gait Assessment. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 3, 3, Article 123 (sep 2019), 23 pages. <https://doi.org/10.1145/3351281>
- [151] Hanbin Zhang, Li Zhu, Viswam Nathan, Jilong Kuang, Jacob Kim, Jun Alex Gao, and Jeffrey Olgin. 2021. Towards Early Detection and Burden Estimation of Atrial Fibrillation in an Ambulatory Free-Living Environment. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 5, 2, Article 86 (jun 2021), 19 pages. <https://doi.org/10.1145/3463503>
- [152] Jingwen Zhang, Dingwen Li, Ruixuan Dai, Heidy Cos, Gregory A. Williams, Lacey Raper, Chet W. Hammill, and Chenyang Lu. 2022. Predicting Post-Operative Complications with Wearables: A Case Study with Patients Undergoing Pancreatic Surgery. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 6, 2, Article 87 (jul 2022), 27 pages. <https://doi.org/10.1145/3534578>
- [153] Shibo Zhang, Yuqi Zhao, Dzung Tri Nguyen, Runsheng Xu, Sougata Sen, Josiah Hester, and Nabil Alshurafa. 2020. NeckSense: A Multi-Sensor Necklace for Detecting Eating Activities in Free-Living Conditions. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 4, 2, Article 72 (jun 2020), 26 pages. <https://doi.org/10.1145/3397313>
- [154] Yunke Zhang, Fengli Xu, Tong Xia, and Yong Li. 2022. Quantifying the Causal Effect of Individual Mobility on Health Status in Urban Space. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 5, 4, Article 193 (dec 2022), 30 pages. <https://doi.org/10.1145/3494990>
- [155] Pengyuan Zhou, Hengwei Xu, Lik Hang Lee, Pei Fang, and Pan Hui. 2022. Are You Left Out? An Efficient and Fair Federated Learning for Personalized Profiles on Wearable Devices of Inferior Networking Conditions. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 6, 2, Article 91 (jul 2022), 25 pages. <https://doi.org/10.1145/3534585>

Received 15 February 2023