



CrowdWorkSheets: Accounting for Individual and Collective Identities Underlying Crowdsourced Dataset Annotation

Mark Díaz
markdiaz@google.com
Google Research
New York, New York, USA

Ian D. Kivlichan
kivlichan@google.com
Jigsaw, Google
New York, New York, USA

Rachel Rosen
rachelrosen@google.com
Jigsaw, Google
New York, New York, USA

Dylan K. Baker*
dylan@dair-institute.org
Distributed AI Research Institute
Seattle, Washington, USA

Razvan Amironesei
amironesei@gmail.com
Google Research
Mountain View, California, USA

Vinodkumar Prabhakaran
vinodkpg@google.com
Google Research
San Francisco, California, USA

Emily Denton
dentone@google.com
Google Research
New York, New York, USA

ABSTRACT

Human annotated data plays a crucial role in machine learning (ML) research and development. However, the ethical considerations around the processes and decisions that go into dataset annotation have not received nearly enough attention. In this paper, we survey an array of literature that provides insights into ethical considerations around crowdsourced dataset annotation. We synthesize these insights, and lay out the challenges in this space along two layers: (1) who the annotator is, and how the annotators' lived experiences can impact their annotations, and (2) the relationship between the annotators and the crowdsourcing platforms, and what that relationship affords them. Finally, we introduce a novel framework, CrowdWorkSheets, for dataset developers to facilitate transparent documentation of key decisions points at various stages of the data annotation pipeline: task formulation, selection of annotators, platform and infrastructure choices, dataset analysis and evaluation, and dataset release and maintenance.

CCS CONCEPTS

• **Social and professional topics** → **User characteristics**; • **Computing methodologies** → **Machine learning**.

ACM Reference Format:

Mark Díaz, Ian D. Kivlichan, Rachel Rosen, Dylan K. Baker, Razvan Amironesei, Vinodkumar Prabhakaran, and Emily Denton. 2022. CrowdWorkSheets: Accounting for Individual and Collective Identities Underlying Crowdsourced Dataset Annotation. In *2022 ACM Conference on Fairness, Accountability, and Transparency (FAccT '22)*, June 21–24, 2022, Seoul, Republic of Korea.

*Work completed while at Google.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

FAccT '22, June 21–24, 2022, Seoul, Republic of Korea

© 2022 Copyright held by the owner/author(s).

ACM ISBN 978-1-4503-9352-2/22/06.

<https://doi.org/10.1145/3531146.3534647>

Korea. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3531146.3534647>

1 INTRODUCTION

Human computation refers to the practice of tapping into human intelligence and cognition as computational elements within an information processing system design, often done on a large global scale [41]. The sheer scale of human computation that the Web enables has made possible things that were previously unimaginable, e.g., *Captchas* digitizing the entire NYTimes historical publications, global participatory platforms for human rights and crises response,¹ and large-scale data and distributed analyses enabled by citizen science projects.² In particular, human computation has played a critical role in the research, development, and deployment of modern-day artificial intelligence systems, through the creation of training datasets [10], and human-in-the-loop systems [14, 35]. By enabling efficient and scalable distribution of data labelling microtasks, crowdsourcing platforms are a natural choice for dataset developers aiming to cheaply and efficiently generate dataset annotations.

In this paper we explore the challenges and decision points inherent to crowdsourced annotation of machine learning datasets and propose a framework, CrowdWorkSheets, for reflecting on dataset annotation decisions, and documenting them in a standardized manner. At a high level, CrowdWorkSheets prompts dataset developers to ask: who is annotating the data, and why is that important? We consider how the ethical concerns of data annotation intersect with the identities of the annotators, the social structures surrounding their work, and how their individual perspectives may become encoded within the dataset labels. In doing so, we push back against the prevalent notion that crowdworkers are interchangeable and instead seek to illuminate why they are not. Data generated in crowdwork tasks is shaped by a range of social factors and the datasets that workers help to build continue to shape systems long

¹<https://www.usahidi.com/>

²<https://www.citizenscience.gov/>

after worker engagement ends. Processes of annotation thus impact future models built from this data; therefore, understanding the perspectives captured through data labeling is crucial to fully understanding these models and the potential social impact they can have.

Our work is motivated by, and extends, prior scholarship examining ethical considerations relating to crowdsourcing. For instance, Vakharia and Lease [57] outline various kinds of challenges encountered in this space by analyzing and comparing seven different crowdsourcing platforms. In addition, Schlagwein et al. [49] conducted extensive fieldwork, engaging crowdworkers, platform organizers, and requesters over the course of three years to uncover a range of ethical dilemmas relating to gig economy crowdsourcing. Shmueli et al. [54] identified risks of harm to crowdworkers engaged in NLP tasks. Attending to ethical issues more broadly, Kocsis and De Vreede [32] used value-sensitive design and transparency literature to develop a taxonomic framework of ethical considerations in crowdsourcing.

Our primary contribution is the introduction of CrowdWorkSheets, a novel framework designed to facilitate critical reflection and transparent documentation of dataset annotation decisions, processes, and outcomes. CrowdWorkSheets complements and extends dataset development and documentation frameworks that have previously been developed in service of transparency, accountability, and reproducibility [6, 9, 17, 26, 27, 40, 42], but focuses specifically on unique considerations relating to crowdsourced dataset annotation. Similar to recent dataset documentation frameworks that have tailored to specific domains (e.g. [44, 55]), our work starts from a recognition of the limitations of “one-size-fits-all” solutions to ethical issues in dataset development. More specifically, we offer CrowdWorkSheets as a targeted intervention to address unresolved ethical problems in crowdsourcing that relate specifically to worker subjectivity and worker experiences.

The remainder of this paper is structured as follows. First, we review literature relating to (1) how annotators’ individual and collective social experiences can impact their annotations, and (2) the relationship between the annotators and the crowdsourcing platforms, and what that relationship means for their ability to engage in fair work. Next, we introduce the CrowdWorkSheets considerations and documentation questions. Finally, we step through a hypothetical case study to illustrate how a dataset developer might use CrowdWorkSheets to document their decisions.

2 WHO IS ANNOTATING ML DATASETS AND WHY DOES IT MATTER?

The historical lineage of crowdsourced labor can be traced back to manufacturing innovation of piecework [2]—a form of labor that produced the “unskilled worker” as the paradigmatic *interchangeable* component and which has been credited with giving rise to the productivity and ingenuity of American manufacturing [16]. In an analogous manner, crowdwork platforms are often designed to position crowdworkers as *interchangeable* [30]. While some forms of digital work can be decomposed and distributed, the presumption that crowdsourced dataset annotators exercise near-identical capacities of perception and judgement ignores the fact that social position, identity, and experience shape how annotators apply

knowledge. Yet, recent empirical work has revealed that dataset annotators are often treated as interchangeable in practice. For example, relatively little attention is given or documented about annotator positionality—how annotator social identity shapes their understanding of the world [18, 48]. Crowd workers are often selected by task requesters based on quality metrics, rather than on any socially defining features of their knowledge or experience. This is concerning; when crowd-sourced annotations are used to build datasets capturing subjective phenomena, such as sentiment or hate speech, annotators’ values and subjective judgments shape the perspectives that machine learning models learn from in a manner that is wholly unaccounted for.

2.1 Accounting for the socio-cultural backgrounds of annotators

Understanding socio-cultural factors of an annotator pool—or even selecting annotators based on these factors—is important because annotator’s identity and lived experience can impact how annotation questions are interpreted and responded to. More generally, subjective interpretations of a task can produce divergent annotations across different communities [53]. As Aroyo and Welty [5] argue, the notion of “one truth” in crowdsourcing responses is a myth; disagreement between annotators, which is often viewed as problematic noise, can actually provide a valuable signal.

A variety of social, cultural, economic, and infrastructural factors contribute to the sociodemographic distribution of workers on any given platform. For example, as [22] points out, the remote nature of crowdwork differentially attracts workers along gender lines, such as mothers who do crowdwork because it allows for an easier balance of childcare in comparison to other work. Other work similarly notes significant gender differences among workers who report engaging in crowdwork because they are only able to conduct work from their homes [7]. This leads to a different gender balance among crowdworkers in the United States than in many other parts of the world; crowdworkers in most of the world are disproportionately male, while in contrast over 60% of U.S. annotators are female [38]. Ipeirotis [29] hypothesizes this to be due to the remote nature of the work, which attracts stay-at-home parents and unemployed or underemployed adults, who are more likely to be women. Additionally, health problems and disability are also a factor that cause many workers to only be able to work from home and motivates them to pursue crowd work [7].

Since many crowdsourced annotator pools are sociodemographically skewed, there are implications for which populations and cultural values are represented in datasets and models [20] as well as which populations face the challenges of crowdwork [22, 30]. Accounting for skews in annotator demographics is critical for contextualizing datasets and ensuring responsible downstream use. In short, there is value in acknowledging, and accounting for, worker’s socio-cultural background—both from the perspective of data quality and societal impact.

2.2 Lived experiences of annotators as expertise

Just as substantive work experience lends valuable domain expertise for a given problem (e.g., annotation of medical imagery by

a medical professional), lived experience with, and proximity to, a problem domain can provide a valuable source of expertise for dataset annotation. For example, women experience higher rates of sexual harassment online compared to men, and among those who have experienced online abuse, women are more likely to identify it as such [58]. This underscores the importance of considering raters' experience with gender-based harassment, when using crowdwork to annotate/moderate online harassment. Recent work has highlighted how the "average" rater, in terms of gender and other social characteristics, varies dramatically depending on which geographies raters are selected from [38]. Additionally, as a result of the previously mentioned sociodemographic differences among who is likely to conduct crowdwork, ratings on sexual harassment data, for example, may differ according to the geographic distribution of raters.

At the same time, relevant lived experience among annotators does not always fall along demographic lines. Waseem [59] demonstrated that incorporating feminist and antiracist activists' perspectives into hate speech annotations yielded better aligned models. Similarly, Patton et al. [37] demonstrated the importance of situated domain expertise—including contextualized knowledge of local language, concepts, and gang activity—when annotating Twitter images to detect pathways to violence among gang-involved youth in Chicago. They found that expert annotations (i.e., those from individuals situated in Chicago and with community ties) significantly diverged from those of graduate students who were scholars of social work and who were trained to perform the annotation task but who lacked this lived experience.

In summary, a core question to answer in data collection is how much annotator's identity, lived experience, and prior knowledge of a problem space matters for the task at hand, and how it impacts what the resulting dataset is intended to capture. While the aforementioned examples constitute relatively subjective tasks, even seemingly objective tasks such as annotating medical text vary surprisingly with annotator backgrounds and experience. Aroyo and Welty [5] show that medical experts are more likely to erroneously identify medical relations as being expressed in text compared with non experts because the experts already know the relation is true based on knowledge external to the task. Their work underscores a need to examine annotator experience even in tasks that appear to be unambiguous or objective.

3 WORKER EXPERIENCES OF DATASET ANNOTATION

Another series of considerations are rooted in annotators' experiences with annotation work itself and how those experiences impact how they do their work. These include issues related to worker compensation, power imbalances in between worker and requester, and the structure of annotation work itself—all of which can pose barriers to crowdworker well-being and their ability to produce quality work.

3.1 Compensation and working conditions

Compensation policies of crowdwork platforms should be a core aspect to consider when thinking about responsible data collection. For instance, in the U.S., there are currently no regulations around

worker pay for crowdwork [7], and the Fair Labor Standards Act that established the minimum wage,³ is not applicable for crowdworkers as they are independent contractors [52]. Reports on how much crowdworkers actually earn vary, but generally show an average lower than minimum wage [30]; surveys of workers from Amazon Mechanical Turk and Crowdfunder place it on average between \$1 and \$5.5 per hour [7] with a median wage of roughly \$2 / hour [24, 52]; only a small fraction of workers (4%) earn more than \$7.25 / hour [24].

Recent research has also identified how crowdworking platforms often necessitate various kinds of unpaid labor from crowdworkers, which reduces overall wages. For example, one report found that for every hour of paid work, workers spend another 18 minutes on unpaid work, including searching for tasks [7]. Another recent study found that once daily invisible labor was accounted for, the median hourly wage for crowdworkers on Amazon Mechanical Turk dropped from \$3.76 to \$2.83 [56]. Workers often invest significant labor outside the platform itself to find tasks, relying on web browsers extensions and participating in crowd work forums [23, 31]. Time spent working is compounded by competition from other crowdworkers [52], which can pressure workers to be constantly available to look for work [7]. The working conditions of crowdworkers are characterized by long working hours, partially as a result of this competition. As Berg [7] notes, this conflicts with the work flexibility motivates many workers to choose crowd work.

Worker psychological safety is a particular area of concern. Crowdworkers who work on content moderation of user generated content often need to look at content that includes violent imagery or sexual and pornographic content [43], or to transcribe conversations about trafficking children into sexual slavery [13]. In many cases, it is impossible to ascertain that a job may contain such content [13]. If crowdworkers find themselves upset or disturbed by this content, they have little recourse; often, workers need to sign non-disclosure agreements preventing them from talking to anyone about the awful things they must look at, even for support [43]. Additionally, raising concerns to their employers is quite difficult; both bureaucracy and physical distance (many of these workers are in the Global South) prohibit any direct lines of feedback or complaints. There is research available on the long-term impacts of viewing harmful user-generated content, but it is difficult to assess the full harm this causes to workers' well-being [43].

3.2 Power dynamics

Power dynamics between the requesters and annotators is another major challenge. Annotators are often heavily distanced from those leading the development of datasets are requesting tasks, which can obfuscate working conditions. Top-down organizational structures often results in the workers viewing requesters as more informed as they are the ones who provided the data and the label schema [34]. Hence, instead of resolving ambiguities, workers are more likely to try to judge from the standpoint of the requester, often with limited exposure to the goals of the annotation. This contributes to the *portability trap* [51]: a "failure to understand how repurposing algorithmic solutions designed for one social context

³<https://www.dol.gov/agencies/whd/flsa>

may be misleading, inaccurate, or otherwise do harm when applied to a different context.”

Power dynamics are also at play in the rejection of work: a large majority of crowdworkers (94% as per [7]) have had work that was rejected or for which they were not paid. Yet, some platforms give requesters full rights over the data they receive, regardless of whether they accept or reject it, and workers have no way of taking legal action if requesters use rejected work anyway [30]; Roberts [43] describes this system as one that “enables wage theft”. Moreover, rejecting work and withholding pay is painful because rejections are often caused by unclear instructions and the lack of meaningful feedback channels. Many crowdworkers report that poor communication negatively affects their work [7]. Moreover, requesters get to choose whether the work is up to their standards before choosing whether to pay for it, even though rejections are often caused by unclear instructions and the very lack of feedback channels they refuse to provide [7]. Workers also feel powerless to speak up about perceived injustices from requesters or the platform; Amazon Mechanical Turk (AMT) users have reportedly had their accounts suspended for speaking negatively about Amazon [52]. Additionally, requesters can block users who offer them feedback without consequence [7].

Power asymmetries also reflect global power dynamics. For instance, since technology development happens primarily in the West, human computation from the Global South is often relegated to the margins [47]. In particular, [47] points out that the technical, social, ethical, and physical distance between the builders of a technology and the communities it is meant to serve is large, in such settings. [8] has pointed out the potential of crowdsourcing to revolutionize civic participation in many developing countries to address complex challenges in governance around global issues such as climate change, poverty, armed conflict, and other crises. They also point out the challenges when it comes to employing crowdsourced interventions on the ground in the Global South. They note that systemic disparities endemic to local contexts are often reflected in who is represented in *crowd*; for instance, the digital crowd in the Global South tends to over-represent the elite, educated, young males who belong to the upper tiers of local social hierarchies.

Roberts [43] compares the commercial content moderation work to the practice of developed nations offloading their hazardous e-waste refuse on countries in the Global South. Her interviewees characterized this work as “akin to being immersed in ‘a cesspool’ – feeling that they are within a pit of toxic matter and waste day in and day out”. The metaphor goes further in highlighting the fact that digital content moderation when outsourced to countries in the Global South, serves to keep the digital refuse away from the field of vision of those in the Global North who are responsible for its existence, and for whom it was intended, in much the same way the rotting garbage and e-waste produced in the Global North is kept away.

On the other hand, some platforms have geographical blocking, which many non-Americans find problematic [7] since it can be used to exclude them. This reinforces the dynamic where requesters in the United States get to decide which global perspectives they want to consider for their task, and which they want to disregard.

The anonymous and geographically distributed nature of crowd-sourced annotation work imposes significant barriers to collective action on the part of dataset annotators. While some platforms offer communication spaces, such as discussion forums, for workers to communicate with one another, these platform-moderated spaces have been shown to be ineffective at supporting labor organizing or worker power [19]. In response, several tools have been developed independently from crowdwork platforms to support crowd workers. For example, *TurkerNation*, *Turk Alert*, *MTurkGrind*, and *Reddit’s /r/HITsWorthTurkingFor* offer online forums for AMT workers to share information about well-paying work and share experiences with different requesters and *Turkopticon* [1, 30] is a browser add-on that enables AMT workers to review and report requesters and view reviews from other workers. These tools can help workers overcome the information asymmetries built into the AMT platform [33]. *Dynamo* is another community platform designed specifically to support and enable collective action for AMT workers, creating “unities without unions” [46]. A 2015 study of the platform found that twenty-two ideas for action had been generated and two active campaigns had been initiated.

In summary, responsible data annotation requires careful consideration of the power dynamics that structure the working relationship between requesters, annotators, and the platforms.

4 CROWDSHEETS: A DOCUMENTATION FRAMEWORK FOR CROWDSOURCED DATASET ANNOTATION

We now introduce our framework, *CrowdWorkSheets*, which outlines a series of *considerations* designed to guide the collection, use, and dissemination of crowd-sourced annotations and *questions* designed to elicit information about various decisions and outcomes. We have decomposed the framework into sections based on different parts of a typical dataset construction pipeline, from the formulation of tasks to dissemination of datasets.

4.1 Task formulation

First, we must ask: *what are we asking annotators to do?* Our considerations and documentation questions focus on many aspects of task formulation including which assumptions we make about annotators, how we handle ambiguity and subjectivity within our task, and how our task is ultimately framed and communicated.

While some tasks tend to pose objective questions with a correct answer (*is there a human face in an image?*), oftentimes datasets aim to capture judgement on relatively subjective tasks with no universally correct answer (*is this piece of text offensive?*). Moreover, even seemingly objective tasks can still be rife with ambiguity or corner-cases and ultimately require subjective judgements to be made on the part of annotators. As such, it is important to consider how questions afford varied interpretations or may require subjective judgements on the part of the annotators. Clarifying such aspects of an annotation task as critical to ensuring a resulting dataset captures the aspects of human intelligence they are meant to capture. Moreover, as discussed in Section 3, a survey of crowdworkers on AMT found that many instances of work rejection were due to unclear instructions [7].

While we discuss the nuances of annotator selection in greater depth in Section 4.2, tasks should be formulated based on considerations regarding *who* will be annotating data and what perspectives should (or should not) be included. Determinations should be tied to the purpose of dataset creation and the downstream use cases it is meant to serve, rather than what is convenient, efficient, or scalable. Some tasks may benefit from being informed by the annotators' lived experiences and thus may be designed to explicitly seek out such expertise. On the other hand, a dataset developer may want to frame task instructions so as to restrict the annotators from relying on their lived experiences, e.g. for a dataset meant to capture a set of policies defined by a platform.

Finally, when formulating a task, it is important to consider how much information to disclose to annotators about the task in advance. Some information may be essential to disclose in order to enable annotators to make informed decision regarding whether or not to accept the task. For example, disclosure of how data will be stored, packaged, and potentially published may be particularly important when sociodemographic, or other sensitive information, about annotators is being requested. Similarly, disclosure of risks relating to psychological harm should be included where appropriate.

Considerations

- Consider the role subjectivity plays in your annotation task. Remember that individuals with different social and cultural backgrounds might differ in their judgements.
- Consider the forms of expertise that should be incorporated through data annotation, including both formal disciplinary training and lived experience with the problem domain. Remember that insufficiently capturing this expertise in the annotator pool may carry risks for downstream model usage.
- Make sure task instructions are clear and unambiguous in order to prevent annotators from wasting time on a task where their work will be rejected due to misunderstandings. Consider assessing the task instructions in a small-scale setting prior to launching your full annotation task.
- Consider the personal information you are collecting from annotators and the potential ethical or privacy risks that may accompany such collection.
- Consider the amount of information you disclose to annotators prior to engagement with the task and ensure annotators have an opportunity to make informed decisions based on any potential risks the task carries.

Documentation questions

- (1) At a high level, what are the subjective aspects of your task?
- (2) What assumptions do you make about annotators?
- (3) How did you choose the specific wording of your task instructions? What steps, if any, were taken to verify the clarity of task instructions and wording for annotators?
- (4) What, if any, risks did your task pose for annotators and were they informed of the risks prior to engagement with the task?
- (5) What are the precise instructions that were provided to annotators?

4.2 Selecting annotators

Next, we ask: *who is annotating the data?* While there is no single “correct” way to assemble an annotator pool, the selection of an annotator pool is a highly consequential decision. Since annotators from different communities can produce significantly different annotations given the same task [53], it is important to recognize that annotator selection may have a significant impact on the labels of your dataset. With this in mind, it is important to consider the intended use of the datasets—which communities will be most impacted by models built from the data, and which communities could be harmed the most by resulting biases present if they are not represented in the annotator pool?

In some cases, social identities of annotators indicate a form of expertise relevant to our task so it may be prudent to select annotators based on self-identified sociodemographic factors. In other cases, it may be important to select annotators based on other forms of expertise or experience with a problem domain. Understanding one's desired annotator pool may subsequently impact decisions regarding platform selection, as different platforms offer differing degrees of flexibility to assemble custom annotator pools.

While selecting annotators based on sociodemographic factors may help ensure a dataset reflects perspectives of certain groups, targeted data collection efforts—particularly those oriented towards the inclusion of marginalized groups—are not without risk. For example, [11] discuss how the mere inclusion of marginalized groups within a dataset, without sufficient attention to broader considerations of data capture and use, can operate as a form of “predatory inclusion”⁴. Discourses of inclusion can serve to “further rather than subvert vulnerability to what might more broadly be called ‘data violence’” [25]. From a privacy perspective, if sociodemographic information is collected and published with a dataset, developers should take extra care to mitigate risks of unintentionally making annotators identifiable.

Considerations

- While there are multiple valid ways to assemble an annotator pool, remember that annotators are not interchangeable, and that the decisions in this stage can heavily impact the final dataset.
- Consider the ways in which social identities of annotators may relate to the forms of expertise important for the task.
- Consider the intended usage contexts of the dataset, and the marginalized communities therein, when choosing which annotators to be prioritized to be included.
- Consider how labor practices intersect with the choice of who the annotators are. For example: if female annotators make up the majority, as they do in the U.S. [38], consider how fair payment, or a lack thereof, could impact this group.

Documentation questions

- (1) Are there certain perspectives that should be privileged? If so, how did you seek these perspectives out?
- (2) Are there certain perspectives that would be harmful to include? If so, how did you screen these perspectives out?
- (3) Were sociodemographic characteristics used to select annotators for your task? If so, please detail the process.

⁴The term “predatory inclusion” has been used to describes modes of inclusion that are extractive and predatory in nature in other domains (e.g. [50])

- (4) If you have any aggregated sociodemographic statistics about your annotator pool, please describe.
- (5) Do you have reason to believe that sociodemographic characteristics of annotators may have impacted how they annotated the data? Why or why not?
- (6) Consider the intended context of use of the dataset and the individuals and communities that may be impacted by a model trained on this dataset. Are these communities represented in your annotator pool?

4.3 Platform and infrastructure choices

Next, we ask, *under what conditions are data annotated?* As described in Section 3, platform policies around compensation and power asymmetries play a huge role in shaping worker experiences and the quality of work that annotators produce. Different platforms offer different affordances for communication between task requesters and annotators, which might impact the extent to which task requesters can incorporate annotator feedback into the task framing or annotator guidelines. Different platforms also impose different minimum-pay constraints; requesters may want to support platforms that uphold fair pay standards. Additionally, requesters should be mindful of potential differences between legal minimum wages and a living wage [21]. Separately from the platform, task creators should be aware of worker pay per hour; some platforms may only offer requesters the option to select pay per item for an annotation task, and the defaults may be set low. Task creators should take care when estimating work time per item to ensure they are paying workers fairly. Another thing to consider when choosing a platform for data annotation is how well that platform supports rater psychological safety. Some platforms provide more affordances than others for crowdworkers to seek out support if they are experiencing distress, or if they otherwise have questions or feedback for requesters.

Considerations

- Consider platform's underlying annotator pool and the options they provide to source specialized rater pools, and whether they enable you to curate an appropriate pool of annotators (e.g. considering sociodemographic factors or domain expertise).
- Consider comparing and contrasting the minimum pay requirements established across different platforms. You may choose to support a platform that upholds fair pay standards.
- Consider the extent to which you would like to establish a channel of communication and feedback between your team and the annotators. Platform mediated channels of communication can give annotators an opportunity to provide feedback on confusing instructions, or otherwise seek out support.

Documentation questions

- (1) What annotation platform did you utilize?
- (2) At a high level, what considerations informed your decision to choose this platform?
- (3) Did the chosen platform sufficiently meet the requirements you outlined for annotator pools? Are any aspects not covered?

- (4) What, if any, communication channels did your chosen platform offer to facilitate communication with annotators? How did this channel of communication influence the annotation process and/or resulting annotations?
- (5) How much were annotators compensated? Did you consider any particular pay standards, when determining their compensation? If so, please describe.

4.4 Dataset analysis and evaluation

Once data instances are annotated, *what do we do with the results?* This section focuses on considerations related to the process of converting the “raw” annotations into the labels that are ultimately packaged in a dataset. A common practice in building crowdsourced annotations for discrete labeling tasks is to obtain multiple annotator judgements that are then aggregated (e.g., through majority voting) to obtain a single “ground truth” that is released in the dataset [45]. However, the disagreements between annotators may embed valuable nuances about the task [4, 36]. Aggregation, in such cases may obscure such nuances, and potentially exclude perspectives from minority annotators [39]. It is thus critical to consider uncertainty and disagreement between annotators, and potentially leverage this as a signal, to avoid losing nuanced and diverse opinions in the aggregation process. It might be important to analyze how annotators disagree along sociodemographic lines in order to be able to share this information with potential users of the dataset, so they can best understand how to represent these diverse perspectives in their use of the data.

Considerations

- Consider including uncertainty or disagreement between annotations on each instance as a signal in the dataset.
- Consider analyzing systematic disagreements between annotators of different sociodemographic groups in order to better understand how diverse perspectives are represented.
- Consider how the final dataset annotations will relate to individual annotator responses. For instance, one option is to release only the aggregated labels, e.g. through a majority vote. Consider what valuable information might be lost through such aggregation.

Documentation questions

- (1) How do you define the annotation quality in your context, and how did you assess quality in your dataset?
- (2) Have you conducted any analysis on disagreement patterns? If so, what analyses did you use and what were the major findings?
- (3) Did you analyze potential sources of disagreement?
- (4) How do the individual annotator responses relate to the final labels released in the dataset?

4.5 Dataset release and maintenance

Finally, it is critical to consider *what is the future of the dataset?* Data exists within an ever-changing world, and should be viewed and used in that context. Users of the dataset now and in the future should understand the limitations of the data based on when and how it was collected. For example, a dataset may require periodic updates to remain robust to new slang or changes in language use over time. In addition, annotation tasks may be predicated upon

legal definitions or medical standards that may change according to decisions by institutions or governing bodies.

Considerations

- Consider designing and sharing a dataset maintenance plan [28].
- Consider potential conditions under which annotations may become outdated or less useful.

Documentation Questions:

- (1) Do you have reason to believe the annotations in this dataset may change over time? Do you plan to update your dataset?
- (2) Are there any conditions or definitions that, if changed, could impact the utility of your dataset?
- (3) Will you attempt to track, impose limitations on, or otherwise influence how your dataset is used? If so, how?
- (4) Were annotators informed about how the data is externalized? If changes to the dataset are made, will they be informed?
- (5) Is there a process by which annotators can later choose to withdraw their data from the dataset? Please detail.

5 CASE STUDY

We now present a *hypothetical* case study to demonstrate how our considerations outlined in Section 4 might be incorporated in practice and how dataset annotation decisions might be documented using CrowdWorkSheets. Responses to documentation questions are not intended to be prescriptive, nor are they completely comprehensive. Instead, they should be considered as one of many valid responses to this line of inquiry, and as a way to provoke further thought and discussion.

In this hypothetical case study, we take our goal to be the development of a benchmark dataset for public release to support academic research in social media content moderation. A Twitter corpus of 20,000 English-language tweets has been collected and we seek to label each tweet independently on a four-point “toxicity” scale defined in [12].

Task Formulation

At a high level, what are the subjective aspects of your task?

Judgements of toxicity of online comments is highly subjective. What makes a tweet harmful or hurtful varies greatly not only by the literal content of the tweet, but by the context surrounding it. In our task setup, tweets are presented to annotators in isolation, so they do not have access to the overall context of the online conversation. As such, we anticipate that annotators may infer surrounding context and make subjective judgements based on this inference.

What assumptions do you make about annotators?

Some of the key assumptions we make of our annotators:

- Annotators that claim proficiency in English and familiarity with social media have enough context to reasonably interpret the task.
- By giving a clear understanding of the goals of this work and explicitly indicating that this is a subjective task where disagreement is expected, we will increase the likelihood that annotators will allow their lived experiences to inform how they label toxicity.
- By paying well, we will increase the likelihood that annotators will take time to think through particularly challenging examples.

How did you choose the specific wording of your task instructions? What steps, if any, were taken to verify the clarity of task instructions and wording for annotators?

To align with existing research in the area, we’ve chosen to give annotators an existing definition of toxicity, as “rude, disrespectful or otherwise likely to make someone leave a discussion” [3]. To settle on a final task wording, our research team first completed 50 annotation tasks each to identify any obvious challenges applying this definition. We then ran several small pilot studies with slightly varying task instructions, and allowed annotators the option to give feedback on aspects that were unclear. Looking over these results, we settled on the question phrasing that yielded the least reported confusion. We intentionally chose to leave our definition of toxicity somewhat open to interpretation, operating under the understanding that being overly specific in task instructions for subjective work does not improve response quality [5]. We also explicitly informed annotators that we expect a variety of interpretations of each comment, and that we were looking for their personal best judgements in the given situation. To motivate thoughtful responses, we chose to pay well above minimum wage and gave annotators a clear idea of the ultimate purpose of their work. However, we know that it is inevitable that some annotators will simply give answers they think we want as quickly as possible. While we screen out responses below a minimum duration, it’s impossible to ensure every answer is honest and thoughtful. We assume that these responses are randomly distributed; we leave it up to dataset users to do further analysis.

What, if any, risks did your task pose for annotators and were they informed of the risks prior to engagement with the task?

Our task required annotators to read text that potentially contained hate speech, slurs, and other harmful content. As such, the task posed a risk of psychological harm to annotators. Moreover, given that we selected annotators who had previously experienced online harassment, there is a potential for the task to trigger an emotional response related to past trauma. We informed annotators about this risk prior to the start of the task. We also informed annotators that we would be requesting sociodemographic information in order to assess disagreement across different groups. We outlined our data storage policy and steps we took to prevent responses from being linked to sociodemographic information.

What are the precise instructions that were provided to annotators?

The final task instructions used for data collection reflected in the released data is available at HypotheticalTaskInstructions.com.

Selecting Annotations

Are there certain perspectives that should be privileged? If so, how did you seek these perspectives out?

We want to privilege the perspectives of annotators who have personally experienced online harassment or hold marginalized identities that are often targeted online. To this end, we included screening questions such that our annotator pool consisted of raters who have direct experience with online harassment. We intentionally defined “direct experience” very broadly to capture a wide range of experiences, intending to include annotators who’ve been personally harassed by others via online channels, who’ve encountered online content that threatened or disparaged identities they share, who have experience moderating online forums, or who have felt otherwise personally affected by harmful online content.

Are there certain perspectives that would be harmful to include? If so, how did you screen these perspectives out?

We believe that there are many harmful worldviews annotators might hold that we do not want captured by our annotations; we do not want to employ

annotators who participate in hateful online communities, for example. To attempt to account for this, we identified several tweets that we agreed were unambiguously toxic, and screened out any annotators that did not label these as toxic.

Were sociodemographic characteristics used to select annotators for your task? If so, please detail the process.

In addition to screening for annotators who have previously experienced online harassment, we selected annotators based on self-identified gender and age. We aimed for an approximately gender balanced pool and we selected for at least 10% of the annotators to be older than 65 years old. Because annotators were sourced from multiple geographic regions, we could not easily specify thresholds for racial or ethnic diversity; however, because we are screening for annotators who have experienced harassment online, we achieved decent representation among marginalized groups.

If you have any aggregated sociodemographic statistics about your annotator pool, please describe.

We first selected annotators who indicated that they had previously experienced online harassment. This resulted in a pool that is disproportionately composed of women and people of color compared with platform demographics. More specific demographic breakdowns are available with the released dataset.

Do you have reason to believe that sociodemographic characteristics of annotators may have impacted how they annotated the data? Why or why not?

Yes, we believe that annotators who have themselves experienced online harassment may be more likely to identify tweets as toxic. Based on rates of reported experience with hate speech attacks, we also expect that these annotators will disproportionately be members of marginalized social groups in their respective geographic region.

Consider the intended context of use of the dataset and the individuals and communities that may be impacted by a model trained on this dataset. Are these communities represented in your annotator pool?

Our intended audience is researchers studying English-language online content moderation, although we can anticipate that our work may have impact within industry. Content moderation has far-reaching and pervasive influence on online discourse, which impacts a wide range of individuals and communities. Not everyone is equally vulnerable to the worst impacts of toxic language online, so we specifically selected for an annotator pool where this more vulnerable population is represented.

Platform and Infrastructure Choices

What annotation platform did you utilize?

We're using HypotheticalPlatform.

At a high level, what considerations informed your decision to choose this platform?

We have selected HypotheticalPlatform for several reasons: First, they are a generally reliable platform with a history of high data quality. Second, they are able to guarantee that annotators are paid at or above a living wage. Third, their platform's interface allows annotators to easily communicate feedback and concerns. And finally, their platform allows us to make ample use of screening questions to select the annotator pool for our main body of work.

Did the chosen platform sufficiently meet the requirements you outlined for annotator pools? Are any aspects not covered?

We were able to meet all of our requirements for annotator pools through the use of many screening and demographic questions. The main trade-off we made to accomplish this is in cost; to pay annotators well, including for their time answering screening questions, we set a limit on the number of tweets we could label.

What, if any, communication channels did your chosen platform offer to facilitate communication with annotators? How did this channel of communication influence the annotation process and/or resulting annotations?

We included a free response section at the end of our survey to allow feedback from annotators. In our pilot studies, we used this to clarify our task instructions. In the full study, most annotators left this blank, so we chose to leave them out of the final dataset.

How much were annotators compensated? Did you consider any particular pay standards, when determining their compensation? If so, please describe.

Informed by the 2020 results of the MIT Living Wage Calculator [21], we aimed for annotators to take home at least \$25/hr on our work, with the goal of comfortably reaching a living wage for a single adult with no dependents, and decrease the pressure to complete tasks as quickly as possible. Annotators were paid \$6.25 for labeling a batch of 40 tweets, designed to take no more than 15 minutes, and verified over the course of the annotation job.

Dataset Analysis and Evaluation

How do you define the quality of annotations in your context, and how did you assess the quality in the dataset you constructed?

We assessed quality along several dimensions, each of which had an associated question in each 40-question batch:

- Attention: We included 1 attention check question was introduced that instructs the annotator to give a particular response so ensure annotators are reading each question;
- Self-consistency: We included 2 duplicated questions within each batch, to ensure annotators were actually reading each tweet and being self-consistent in their responses.
- Alignment with pre-defined ratings: We included several 2 tweets that the research team had pre-labeled as unoffensive and highly offensive. We chose tweets for which we would expect no disagreement from annotators.

We removed from the final dataset all batches where 2 or more of these 5 data quality questions were incorrectly answered. This ultimately accounted for 12% of our data.

Have you conducted any analysis on disagreement patterns? If so, what analyses did you use and what were the major findings?

While the main purpose of this work is data collection and not analysis, we did conduct very preliminary analyses as a starting point for dataset users. We ran standard inter-annotator agreement metrics and found a relatively low interannotator agreement across all raters (Fleiss' $\kappa = 0.25$ [15]). However, we do not believe this to be an issue of data quality—when we looked at the data aggregated along different demographic axes, we found many demographic groups with high interannotator agreement whose annotations differ significantly from the majority opinion.

Did you analyze potential sources of disagreement?

In our preliminary analysis, we looked at a few annotator demographics as a source of disagreement. There are a myriad of other factors one could

analyze with respect to disagreement—tweet topic, presence or absence of particular words, or how quickly annotators responded, for example—but as this is intended to be released as a research dataset, we have not conducted all of these analyses.

How do the individual annotator responses relate to the final labels released in the dataset?

After bucketing annotator demographics such that no annotator was uniquely identifiable, we released all responses, attached to the demographics of the annotator that gave each response. We chose not to aggregate responses into final tweet toxicity labels, and instead leave this to dataset users to aggregate in a way that's appropriate for their use case.

Dataset Release and Maintenance

Do you have reason to believe the annotations in this dataset may change over time? Do you plan to update your dataset?

The relevancy of and perceptions about tweets will certainly change over time. In an effort to remind dataset users that this data should be taken in its temporal context, we include the month and year that each tweet was (a) written and (b) annotated as meta-data. However, as a longer-term strategy, we are also open-sourcing and making public all parts of our annotation pipeline, including rater instructions, data formatting schemes, and information on how to coordinate with our data labeling partners. We will publicly extend an open invitation to future collaborators who want to reuse our pipeline to annotate more data. If this pipeline is used and our guidelines followed satisfactorily, we will append future annotations to our existing dataset.

Are there any conditions or definitions that, if changed, could impact the utility of your dataset?

Over time we expect societal views to deviate somewhat from the annotations collected. For example, it will not capture any shifts in attitude regarding language targeting social groups that may be considered marginalized in the future but that are not considered marginalized today.

Will you attempt to track, impose limitations on, or otherwise influence how your dataset is used? If so, how?

To access the data, we require dataset users indicate their affiliation, contact information, and use case. The research team will be assessing uses on a case-by-case basis, with particular attention given to risks associated with use cases that explicitly include sociodemographic data in their modeling. We also ask that any publications cite our dataset release paper so we can track academic uses of the dataset. Our full data license is available at [HypotheticalDataLicense.com](https://hypotheticalDataLicense.com).

Were annotators informed about how the data is externalized? If changes to the dataset are made, will they be informed?

Annotators were informed that this data will be released as a research dataset prior to engaging in the task. We allowed raters to opt in to an email list that with share updates about data release availability. This site will contain an automatically-updated list of papers that cite our dataset release paper.

Is there a process by which annotators can later choose to withdraw their data from the dataset? If so, please detail.

By design, we have no mechanisms of linking individual annotators to specific responses, and so have no option for annotators to withdraw their annotations from our dataset. We make this explicit to the annotators, and allow them to stop answering questions at any point if they decide they no longer want to continue.

6 CONCLUSION

In this work, we challenge the common portrayal of dataset annotators as interchangeable. Rather, we argue, their individual histories and experiences bring unique perspectives to the table that can become encoded in the overall dataset in significant ways. Therefore, it becomes imperative to consider how the process of selecting annotators, and their experience working on annotation, is documented alongside other aspects of dataset development. Towards this end, we introduced CrowdWorkSheets, a framework for reflecting on and documenting key decision points of crowdsourced dataset development, and a set of recommendations for dataset developers. While this framework is oriented towards individual dataset developers, we also recognize the role large institutions can play in shifting incentives to engage with these recommendations, e.g. incentivizing transparent dataset documentation through conference submission and reviewer guidelines.

Funding: This research was supported by Google.

REFERENCES

- [1] 2008. Turkopticon. <https://turkopticon.net/>. (2008). Accessed: 2021-07-21.
- [2] Ali Alkhatib, Michael S. Bernstein, and Margaret Levi. 2017. Examining Crowd Work and Gig Work Through The Historical Lens of Piecework. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems (CHI '17)*. Association for Computing Machinery, New York, NY, USA, 4599–4616. <https://doi.org/10.1145/3025453.3025974>
- [3] Lora Aroyo, Lucas Dixon, Nithum Thain, Olivia Redfield, and Rachel Rosen. 2019. Crowdsourcing Subjective Tasks: The Case Study of Understanding Toxicity in Online Discussions. In *Companion Proceedings of The 2019 World Wide Web Conference (WWW '19)*. Association for Computing Machinery, New York, NY, USA, 1100–1105. <https://doi.org/10.1145/3308560.3317083>
- [4] Lora Aroyo and Chris Welty. 2013. Crowd truth: Harnessing disagreement in crowdsourcing a relation extraction gold standard. *WebSci2013. ACM* (2013).
- [5] Lora Aroyo and Chris Welty. 2015. Truth Is a Lie: Crowd Truth and the Seven Myths of Human Annotation. *AI Magazine* 36, 1 (Mar. 2015), 15–24. <https://doi.org/10.1609/aimag.v36i1.2564>
- [6] Emily M. Bender and Batya Friedman. 2018. Data Statements for Natural Language Processing: Toward Mitigating System Bias and Enabling Better Science. *Transactions of the Association for Computational Linguistics* 6 (2018), 587–604. https://doi.org/10.1162/tacl_a_00041
- [7] Janine Berg. 2015. Income security in the on-demand economy: Findings and policy lessons from a survey of crowdworkers. *Comp. Lab. L. & Pol'y J.* 37 (2015), 543.
- [8] Maja Bott and Gregor Young. 2012. The role of crowdsourcing for better governance in international development. *Praxis: The Fletcher Journal of Human Security* 27, 1 (2012), 47–70.
- [9] Kasia S. Chmielinski, Sarah Newman, Matt Taylor, Josh Joseph, Kemi Thomas, Jessica Yurkofsky, and Yue Chelsea Qiu. 2020. The Dataset Nutrition Label (2nd Gen): Leveraging Context to Mitigate Harms in Artificial Intelligence. (2020). <http://securedatalol/NeurIPSWorkshoponDatasetCurationandSecurity>.
- [10] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. 2009. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR09*.
- [11] Emily L. Denton, A. Hanna, Razvan Amironesei, Andrew Smart, Hilary Nicole, and M. Scheuerman. 2020. Bringing the People Back In: Contesting Benchmark Machine Learning Datasets. *ICML Workshop on Participatory Approaches to Machine Learning* (2020).
- [12] Lucas Dixon. 2018. Annotation instructions for Toxicity with sub-attributes. https://github.com/conversationai/conversationai.github.io/blob/master/crowdsourcing_annotation_schemes/toxicity_with_subattributes.md. (2018). Accessed: 2021-01-19.
- [13] Sarah Emerson. 2019. 'I Can Hear the Suffering': Rev Exposes Freelance Transcribers to Violent, Disturbing Content. *Medium OneZero* (2019).
- [14] Anna Filippova, Connor Gilroy, Ridhi Kashyap, Antje Kirchner, Allison C. Morgan, Kivan Polimis, Adaner Usmani, and Tong Wang. 2019. Humans in the Loop: Incorporating Expert and Crowd-Sourced Knowledge for Predictions Using Survey Data. *Socius* 5 (2019), 2378023118820157.
- [15] Joseph L. Fleiss. 1971. Measuring nominal scale agreement among many raters. *Psychological bulletin* 76, 5 (1971), 378.
- [16] Joshua Benjamin Freeman. 2018. *Behemoth : a history of the factory and the making of the modern world*. W.W. Norton & Company, Inc.,.

- [17] Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé III, and Kate Crawford. 2018. Datasheets for datasets. *arXiv preprint arXiv:1803.09010* (2018).
- [18] R. Stuart Geiger, Kevin Yu, Yanlai Yang, Mindy Dai, Jie Qiu, Rebekah Tang, and Jenny Huang. 2020. Garbage in, Garbage out? Do Machine Learning Application Papers in Social Computing Report Where Human-Labeled Training Data Comes From?. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency (FAT* '20)*. Association for Computing Machinery, New York, NY, USA, 325–336. <https://doi.org/10.1145/3351095.3372862>
- [19] Christine Gerber. 2021. Community building on crowdwork platforms: Autonomy and control of online workers? *Competition & Change* 25, 2 (2021), 190–211.
- [20] Sayan Ghosh, Dylan Baker, David Jurgens, and Vinodkumar Prabhakaran. 2021. Detecting Cross-Geographic Biases in Toxicity Modeling on Social Media. In *Proceedings of the Seventh Workshop on Noisy User-generated Text (W-NUT 2021)*. Association for Computational Linguistics, Online, 313–328. <https://aclanthology.org/2021.wnut-1.35>
- [21] Amy K Glasmeier. 2020. Living Wage Calculator. (2020). livingwage.mit.edu
- [22] Mary L Gray and Siddharth Suri. 2019. *Ghost work: How to stop Silicon Valley from building a new global underclass*. Eamon Dolan Books.
- [23] Benjamin V Hanrahan, Anita Chen, JiaHua Ma, Ning F Ma, Anna Squicciarini, and Saiph Savage. 2021. The Expertise Involved in Deciding which HITs are Worth Doing on Amazon Mechanical Turk. *Proceedings of the ACM on Human-Computer Interaction* 5, CSCW1 (2021), 1–23.
- [24] Kotaro Hara, Abigail Adams, Kristy Milland, Saiph Savage, Chris Callison-Burch, and Jeffrey P Bigham. 2018. A data-driven analysis of workers' earnings on Amazon Mechanical Turk. In *Proceedings of the 2018 CHI conference on human factors in computing systems*. 1–14.
- [25] Anna Lauren Hoffmann. 2021. Terms of inclusion: Data, discourse, violence. *New Media & Society* 23, 12 (2021), 3539–3556. <https://doi.org/10.1177/1461444820958725>
- [26] Sarah Holland, Ahmed Hosny, Sarah Newman, Joshua Joseph, and Kasia Chmielinski. 2018. The dataset nutrition label: A framework to drive higher data quality standards. *arXiv preprint arXiv:1805.03677* (2018).
- [27] Ben Hutchinson, Andrew Smart, Alex Hanna, Emily Denton, Christina Greer, Oddur Kjartansson, Parker Barnes, and Margaret Mitchell. 2021. Towards Accountability for Machine Learning Datasets: Practices from Software Engineering and Infrastructure. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*.
- [28] Ben Hutchinson, Andrew Smart, Alex Hanna, Emily Denton, Christina Greer, Oddur Kjartansson, Parker Barnes, and Margaret Mitchell. 2021. Towards Accountability for Machine Learning Datasets: Practices from Software Engineering and Infrastructure. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (FAccT '21)*. Association for Computing Machinery, 560–575.
- [29] Panagiotis G Ipeirotis. 2010. Demographics of mechanical turk. (2010).
- [30] Lilly C Irani and M Six Silberman. 2013. Turkopticon: Interrupting worker invisibility in amazon mechanical turk. In *Proceedings of the SIGCHI conference on human factors in computing systems*. 611–620.
- [31] Toni Kaplan, Susumu Saito, Kotaro Hara, and Jeffrey P Bigham. 2018. Striving to earn more: a survey of work strategies and tool use among crowd workers. In *Sixth AAAI Conference on Human Computation and Crowdsourcing*.
- [32] David Kocsis and Gert Jan De Vreede. 2016. Towards a taxonomy of ethical considerations in crowdsourcing (22nd Americas Conference on Information Systems: Surfing the IT Innovation Wave, AMCIS 2016).
- [33] David Martin, Benjamin V. Hanrahan, Jacki O'Neill, and Neha Gupta. 2014. Being a Turker. In *Proceedings of the 17th ACM Conference on Computer Supported Cooperative Work & Social Computing (CSCW '14)*. Association for Computing Machinery, New York, NY, USA, 224–235. <https://doi.org/10.1145/2531602.2531663>
- [34] Milagros Miceli, Martin Schuessler, and Tianling Yang. 2020. Between Subjectivity and Imposition: Power Dynamics in Data Annotation for Computer Vision. *Proc. ACM Hum.-Comput. Interact.* 4, CSCW2, Article 115 (Oct. 2020), 25 pages. <https://doi.org/10.1145/3415186>
- [35] Robert Munro Monarch. 2021. *Human-in-the-Loop Machine Learning: Active learning and annotation for human-centered AI*. Simon and Schuster.
- [36] Cecilia Ovesdotter Alm. 2011. Subjective Natural Language Problems: Motivations, Applications, Characterizations, and Implications. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, Portland, Oregon, USA, 107–112. <https://aclanthology.org/P11-2019>
- [37] Desmond Upton Patton, Philipp Blandfort, William R Frey, Michael B Gaskell, and Svebor Karaman. 2019. Annotating twitter data from vulnerable populations: Evaluating disagreement between domain experts and graduate student annotators. (2019).
- [38] Lisa Posch, Armin Bleier, Fabian Flöck, and Markus Strohmaier. 2018. Characterizing the global crowd workforce: A cross-country comparison of crowdworker demographics. *arXiv preprint arXiv:1812.05948* (2018).
- [39] Vinodkumar Prabhakaran, Aida Mostafazadeh Davani, and Mark Diaz. 2021. On Releasing Annotator-Level Labels and Information in Datasets. In *Proceedings of the 15th Linguistic Annotation Workshop*. Association for Computational Linguistics, Virtual.
- [40] Mahima Pushkarna, Andrew Zaldivar, and Oddur Kjartansson. 2022. Data Cards: Purposeful and Transparent Dataset Documentation for Responsible AI. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*.
- [41] Alexander J. Quinn and Benjamin B. Bederson. 2011. *Human Computation: A Survey and Taxonomy of a Growing Field*. Association for Computing Machinery, New York, NY, USA, 1403–1412. <https://doi.org/10.1145/1978942.1979148>
- [42] Jorge Ramirez, Burcu Sayin, Marcos Baez, Fabio Casati, Luca Cernuzzi, Boualem Benatallah, and Gianluca Demartini. 2021. On the State of Reporting in Crowdsourcing Experiments and a Checklist to Aid Current Practices. *Proc. ACM Hum.-Comput. Interact.* 5, CSCW2, Article 387 (oct 2021), 34 pages. <https://doi.org/10.1145/3479531>
- [43] Sarah T Roberts. 2016. Digital refuse: Canadian garbage, commercial content moderation and the global circulation of social media's waste. *Wi: Journal of Mobile Media* 10, 1 (2016), 1–18.
- [44] Negar Rostamzadeh, Subhrajit Roy, Diana Mincu, Andrew Smart, Lauren Wilcox, Mahima Pushkarna, Razvan Amironesei, Jessica Schrouff, Madeleine Elish, Nyaleng Moorosi, Berk Ustun, Noah Broesti, and Katherine Heller. 2021. Specialized Healthsheet for Healthcare Datasets. In *Machine Learning for Health (ML4H)*.
- [45] Reka Marta Sabou, Kalina Bontcheva, Leon Derczynski, and A. Scharl. 2014. Corpus Annotation through Crowdsourcing: Towards Best Practice Guidelines. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*. European Language Resources Association (ELRA), Reykjavik, Iceland.
- [46] Niloufar Salehi, Lilly C. Irani, Michael S. Bernstein, Ali Alkhatib, Eva Ogbe, Kristy Milland, and Clickhappier. 2015. *We Are Dynamo: Overcoming Stalling and Friction in Collective Action for Crowd Workers*. Association for Computing Machinery, New York, NY, USA, 1621–1630. <https://doi.org/10.1145/2702123.2702508>
- [47] Nithya Sambasivan, Erin Arnesen, Ben Hutchinson, Tulsee Doshi, and Vinodkumar Prabhakaran. 2021. Re-Imagining Algorithmic Fairness in India and Beyond. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (FAccT '21)*. Association for Computing Machinery, New York, NY, USA, 315–328. <https://doi.org/10.1145/3442188.3445896>
- [48] Morgan Klaus Scheuerman, Emily Denton, and Alex Hanna. 2021. Do Datasets Have Politics? Disciplinary Values in Computer Vision Dataset Development. *Computer Supported Cooperative Work (CSCW)* (2021).
- [49] Daniel Schlagwein, Dubravka Cercec-Kecmanovic, and Benjamin Hanckel. 2019. Ethical norms and issues in crowdsourcing practices: A Habermasian analysis. *Information Systems Journal* 29, 4 (2019), 811–837. <https://doi.org/10.1111/isj.12227>
- [50] Louise Seamster and Raphaël Charron-Chénier. 2017. Predatory Inclusion and Education Debt: Rethinking the Racial Wealth Gap. *Social Currents* 4, 3 (2017), 199–207. <https://doi.org/10.1177/2329496516686620>
- [51] 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 (FAT* '19)*. Association for Computing Machinery, New York, NY, USA, 59–68. <https://doi.org/10.1145/3287560.3287598>
- [52] Alana Semuels. 2018. The internet is enabling a new kind of poorly paid hell. *The Atlantic* 23 (2018).
- [53] Shilad Sen, Margaret E. Giesel, Rebecca Gold, Benjamin Hillmann, Matt Lesicko, Samuel Naden, Jesse Russell, Zixiao (Ken) Wang, and Brent Hecht. 2015. Turkers, Scholars, "Arafat" and "Peace": Cultural Communities and Algorithmic Gold Standards. In *Proceedings of the 18th ACM Conference on Computer Supported Cooperative Work & Social Computing (CSCW '15)*. Association for Computing Machinery, New York, NY, USA, 826–838. <https://doi.org/10.1145/2675133.2675285>
- [54] Boaz Shmueli, Jan Fell, Soumya Ray, and Lun-Wei Ku. 2021. Beyond Fair Pay: Ethical Implications of NLP Crowdsourcing. *arXiv preprint arXiv:2104.10097* (2021).
- [55] Ramya Malur Srinivasan, Emily Denton, Jordan Jennifer Famularo, Negar Rostamzadeh, Fernando Diaz, and Beth Coleman. 2021. Artsheets for Art Datasets. In *Proceedings of Neural Information Processing Systems (NeurIPS), Datasets & Benchmarks Track*. https://openreview.net/pdf?id=K7ke_GZ_6N
- [56] Carlos Toxtli, Siddharth Suri, and Saiph Savage. 2021. Quantifying the Invisible Labor in Crowd Work. *Proc. ACM Hum.-Comput. Interact.* 5, CSCW2 (2021).
- [57] Donna Vakharina and Matthew Lease. 2015. Beyond Mechanical Turk: An analysis of paid crowd work platforms. *Proceedings of the iConference* (2015), 1–17.
- [58] Emily Vogels. 2021. The state of online harassment. *Pew Research Center* (2021).
- [59] Zeerak Waseem. 2016. Are You a Racist or Am I Seeing Things? Annotator Influence on Hate Speech Detection on Twitter. In *Proceedings of the First Workshop on NLP and Computational Social Science*. Association for Computational Linguistics, Austin, Texas, 138–142. <https://doi.org/10.18653/v1/W16-5618>