Whose Ground Truth? Accounting for Individual and Collective Identities Underlying Dataset Annotation

Emily Denton* Google Research dentone@google.com Mark Díaz*
Google Research
markdiaz@google.com

Ian Kivlichan*
 Jigsaw
kivlichan@google.com

Vinodkumar Prabhakaran* Google Research vinodkpg@google.com Rachel Rosen*
Jigsaw
rachelrosen@google.com

Abstract

Human annotations play a crucial role in machine learning (ML) research and development. However, the ethical considerations around the processes and decisions that go into building ML datasets has 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 put forth a concrete set of recommendations and considerations for dataset developers at various stages of the ML data pipeline: task formulation, selection of annotators, platform and infrastructure choices, dataset analysis and evaluation, and dataset documentation and release.

1 Introduction

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 short survey paper we explore the inherent challenges and decision points that stem from crowdsourced dataset annotation. In particular, we 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. 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 after worker engagement ends. We argue that this impacts future models built from this data, and that understanding the perspectives captured within datasets is crucial to understanding resulting models and their potential social impact.

Our work complements and extends prior scholarship examining ethical considerations relating to crowdsourcing (Vakharia and Lease, 2015; Schlagwein et al., 2019; Kocsis and de Vreede, 2016; Shmueli et al., 2021). Our work is distinct from previous scholarship in that we focus our attention on unresolved ethical problems in crowdsourcing that relate specifically to individual worker subjectivity and individual worker experiences. We start by outlining a comprehensive set of concerns regarding how annotators' individual and collective social experiences, as well as their working

^{*}Equal contribution; authors listed alphabetically.

conditions, may impact the nature of the data they provide for machine learning development, in particular the biases that may be captured in and propagated through those datasets. Based on this analysis, we offer a set of ethical considerations and recommendations for dataset developers that apply to different steps of a typical data annotation pipeline, from task formulation to dataset release.

2 Who is annotating ML datasets and why does it matter?

Recent empirical work has revealed that relatively little attention is given or documented about annotator positionality—how annotator social identity shapes their understanding of the world. (Geiger et al., 2020; Scheuerman et al., 2021). 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, since crowd-sourced annotations are often used to build datasets capturing subjective phenomena such as sentiment and hate-speech, and hence crowd workers' values and subjective judgments shape the perspectives that machine learning models learn from in a manner that is wholly unaccounted for. Indeed, crowdsourcing platforms are often explicitly designed in a manner that positions crowdworkers as *interchangeable* (Irani and Silberman, 2013).

Accounting for the socio-cultural backgrounds of dataset annotators is important for at least two reasons. First, subjective interpretations of a task can produce divergent annotations across different communities (Sen et al., 2015). As Aroyo and Welty (2015) argue, the notion of "one truth" in crowdsourcing responses is a myth; disagreement between annotators, which is often viewed as negative, can actually provide a valuable signal. Secondly, since many crowdsourced annotator pools are socio-demographically skewed, there are implications for which populations are represented in datasets as well as which populations face the challenges of crowdwork (Irani and Silberman, 2013; Gray and Suri, 2019). 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.

Accounting for lived experiences of annotators as expertise may be of great utility in some cases. 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 (Vogels, 2021). However, such lived experiences do not always fall along demographic lines. Waseem (2016) demonstrated that incorporating feminist and antiracist activists' perspectives into hate speech annotations yielded better aligned models. Similarly, Patton et al. (2019) 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.

In summary, a core question to answer in data collection is how much annotator subjectivity matters for the task at hand, and how it impacts what the resulting dataset is meant to capture. While we used relatively subjective tasks as examples above, even seemingly objective tasks such as annotating medical texts vary surprisingly with annotator backgrounds and experience (Aroyo and Welty, 2015).

3 Worker Experiences of Dataset Annotation

Another layer of considerations relate to annotators' experiences with annotation work itself and how it can impact how they do their work. These include issues related to worker compensation, imbalances in the relationship 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.

Compensation policies of the 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 (Berg, 2015), and the Fair Labor Standards Act that established the

minimum wage,² is not applicable for crowdworkers as they are independent contractors (Semuels, 2018). Moreover, for every hour of paid work, workers spend another 18 minutes on unpaid work, including searching for tasks (Berg, 2015). Time spent working is compounded by competition from other crowdworkers (Semuels, 2018), which can pressure workers to be constantly available to look for work as well as work longer hours (Berg, 2015). In addition, a large majority of crowdworkers (94% as per (Berg, 2015)) have had work that was rejected or for which they were not paid. Yet, requesters retain full rights over the data they receive regardless of whether they accept or reject it; Roberts (2016) 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 (Berg, 2015).

Power dynamics between the requesters and annotators is another major challenge. 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 (Miceli et al., 2020). 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* (Selbst et al., 2019): a "failure to understand how repurposing algorithmic solutions designed for one social context may be misleading, inaccurate, or otherwise do harm when applied to a different context." 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 (Sambasivan et al., 2021).

In summary, a core consideration for responsible data collection is whether there exist mechanisms for the workers to address these power asymmetries. The anonymous and geographically distributed nature of crowdsourced annotation work imposes significant barriers to collective action. In response, several community forums have been developed independently from crowdwork platforms to support crowd workers, e.g. TurkerNation, Turk Alert, MTurkGrind, and Reddit's /r/HITsWorthTurkingFor. Turkopticon (Irani and Silberman, 2013; tur, [n.d.]) and Dynamo (Salehi et al., 2015) have also emerged as activist tools that support and enable collective action for crowdworkers.

4 Implications for Dataset Developers

We now outline a comprehensive set of considerations for the collection, use, and dissemination of crowd-sourced ML datasets. We discuss them as they apply to different parts of a typical dataset construction pipeline, from the formulation of tasks to dissemination of datasets.

Task formulation: A core objective of constructing ML datasets is to capture the aspects of human intelligence that are of importance to a given task. 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?). It is important to be intentional about whether to lean on annotators' subjective judgements. This determination 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. Not accounting for task subjectivity may lead to inadvertent biases and misses critical insights about tasks that could benefit from the annotators' lived experiences. However, even for relatively subjective tasks such as labeling offensiveness, a dataset developer may want to restrict the annotators from relying on their lived experiences, say, if the dataset is meant to capture a set of policies defined by a platform. Clarifying such aspects of the tasks, has ramifications for how successfully the datasets capture the aspects of human intelligence they are meant to capture.

Recommendations and considerations

- Consider the subjective nature of your annotation task. Is it possible 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. What are the risks of this expertise not being reflected in the annotator pool?

²https://www.dol.gov/agencies/whd/flsa

- 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 how the final dataset annotations will relate to individual annotator responses. For instance, will you release only the aggregated labels, e.g. through a majority vote? Consider what valuable information might be lost through such aggregation.

Selecting annotators: As outlined in Section 2, the selection of an annotator pool is a highly consequential decision, especially given the subjective nature of many annotation tasks. It is important to choose annotation platforms that allow flexibility in designing custom annotator pools along various socio-demographic axes. These decisions should ideally be guided by considering which communities will be most impacted by models built from the data, and which communities could be harmed the most if they are not represented in the annotator pool.

Recommendations and considerations

- While there is no single "correct" way to assemble an annotator pool, the decisions in this stage could impact the biases captured in the resulting dataset. For instance, annotator demographics may serve as a form of expertise that is important for the task (cf. Section 2).
- 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. (Posch et al., 2018), consider how fair payment, or the lack thereof, could impact this group.

Platform and infrastructure choices: As described in Section 2, the platform policies around compensation and power asymmetries play a huge role in the quality of work the annotators produce. Some platforms offer platform-mediated channels of communication that allow task requesters to 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. Separately from the platform, task creators should be aware of worker pay per hour, since this number is not often given explicitly. 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.

Recommendations and 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.

Dataset analysis and evaluation: A common practice in building crowdsourced annotations 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 (Sabou et al., 2014). However, the disagreements between annotators may embed valuable nuances about the task (Ovesdotter Alm, 2011; Aroyo and Welty, 2013). Aggregation, in such cases may obscure such nuances, and in that process potentially exclude perspectives from minority annotators (Prabhakaran et al., 2021).

Recommendations and considerations

- Consider including the uncertainty or disagreement between individual 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.

Dataset documentation and release: Rigorous documentation of design decisions and outcomes relating to dataset annotations is an important aspect of responsible dataset development. Several dataset documentation frameworks have been proposed to contextualize a dataset and offer guidance regarding intended or unintended use, as well as facilitate accountability for development decisions.

Recommendations and considerations

- Consider adopting or adapting an existing dataset documentation framework (e.g. Gebru et al. (2020); Holland et al. (2018); Bender and Friedman (2018); Kazimzade and Miceli (2020); Hutchinson et al. (2021)) to guide your dataset documentation. Consider publishing aggregate statistics on the sociodemographic make-up of your annotator pool.
- Consider including individual annotator responses for each data points in the dataset in addition to a final aggregated ground truth label, where applicable.

References

- [n.d.]. Turkopticon. https://turkopticon.net/. Accessed: 2021-07-21.
- Lora Aroyo and Chris Welty. 2013. Crowd truth: Harnessing disagreement in crowdsourcing a relation extraction gold standard. *WebSci2013. ACM* (2013).
- Lora Aroyo and Chris Welty. 2015. Truth is a lie: Crowd truth and the seven myths of human annotation. *AI Magazine* 36, 1 (2015), 15–24.
- 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 (Dec. 2018), 587–604. https://doi.org/10.1162/tacl_a_00041
- 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.
- Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé III, and Kate Crawford. 2020. Datasheets for Datasets. *arXiv:1803.09010[cs]* (March 2020). arXiv:cs/1803.09010
- 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
- Mary L Gray and Siddharth Suri. 2019. *Ghost work: How to stop Silicon Valley from building a new global underclass*. Eamon Dolan Books.
- 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:1805.03677 [cs] (May 2018). arXiv:cs/1805.03677
- 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.
- 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.
- Gunay Kazimzade and Milagros Miceli. 2020. Biased Priorities, Biased Outcomes: Three Recommendations for Ethics-Oriented Data Annotation Practices. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society (AIES '20)*. Association for Computing Machinery, New York, NY, USA, 71. https://doi.org/10.1145/3375627.3375809
- David Kocsis and Gert-Jan de Vreede. 2016. Towards a taxonomy of ethical considerations in crowdsourcing. (2016).

- 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
- 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
- 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.
- Lisa Posch, Arnim 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).
- Vinodkumar Prabhakaran, Aida Mostafazadeh Davani, and Mark Diaz. 2021. On Releasing Annotator-Level Labels and Information in Datasets. In *Proceedings of the 15th Linguistic An*notation Workshop. Association for Computational Linguistics, Virtual.
- 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.
- 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.
- 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
- 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
- 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).
- Daniel Schlagwein, Dubravka Cecez-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 arXiv:https://onlinelibrary.wiley.com/doi/pdf/10.1111/isj.12227
- 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
- Alana Semuels. 2018. The internet is enabling a new kind of poorly paid hell. *The Atlantic* 23 (2018).
- 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

- 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).
- Donna Vakharia and Matthew Lease. 2015. Beyond Mechanical Turk: An analysis of paid crowd work platforms. *Proceedings of the iConference* (2015), 1–17.
- Emily Vogels. 2021. The state of online harassment. Pew Research Center (2021).
- 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