



# Is Your Toxicity My Toxicity? Exploring the Impact of Rater Identity on Toxicity Annotation

NITESH GOYAL, Google Research, Google, USA

IAN D. KIVLICHAN, Jigsaw, Google, USA

RACHEL ROSEN, Jigsaw, Google, USA

LUCY VASSERMAN, Jigsaw, Google, USA

Machine learning models are commonly used to detect toxicity in online conversations. These models are trained on datasets annotated by human raters. We explore how raters' self-described identities impact how they annotate toxicity in online comments. We first define the concept of specialized rater pools: rater pools formed based on raters' self-described identities, rather than at random. We formed three such rater pools for this study—specialized rater pools of raters from the U.S. who identify as African American, LGBTQ, and those who identify as neither. Each of these rater pools annotated the same set of comments, which contains many references to these identity groups. We found that rater identity is a statistically significant factor in how raters will annotate toxicity for identity-related annotations. Using preliminary content analysis, we examined the comments with the most disagreement between rater pools and found nuanced differences in the toxicity annotations. Next, we trained models on the annotations from each of the different rater pools, and compared the scores of these models on comments from several test sets. Finally, we discuss how using raters that self-identify with the subjects of comments can create more inclusive machine learning models, and provide more nuanced ratings than those by random raters.

*Please be advised that this work contains examples of toxic and offensive content.*

CCS Concepts: • **Human-centered computing** → **Empirical studies in HCI**; • **Computing methodologies** → **Machine learning**; **Natural language processing**; • **Social and professional topics** → **User characteristics**; **Hate speech**.

Additional Key Words and Phrases: Human Annotations, Identity, Toxicity, Harassment, Machine Learning, Data Annotation, Moderation, Subjectivity, Raters, LGBTQ, African-American

## ACM Reference Format:

Nitesh Goyal, Ian D. Kivlichan, Rachel Rosen, and Lucy Vasserman. 2022. Is Your Toxicity My Toxicity? Exploring the Impact of Rater Identity on Toxicity Annotation. *Proc. ACM Hum.-Comput. Interact.* 6, CSCW2, Article 363 (November 2022), 28 pages. <https://doi.org/10.1145/3555088>

## 1 INTRODUCTION

Toxic language, defined as rude, disrespectful, or unreasonable language that is likely to make someone leave a discussion [14], is a pervasive problem online. Detecting toxic language in online conversations using machine learning (ML) models is an active area of interest. However, models built to detect such online toxicity in conversations can be biased. Recent work has shown that some of the classifiers based on these models are more likely to label non-toxic language from minority communities as toxic compared to equivalent language from non-minority communities [18, 33].

Authors' addresses: Nitesh Goyal, [teshg@google.com](mailto:teshg@google.com), Google Research, Google, 111 8th Ave, New York, NY, USA, 11201; Ian D. Kivlichan, [kivlichan@google.com](mailto:kivlichan@google.com), Jigsaw, Google, 111 8th Ave, New York, NY, USA, 11201; Rachel Rosen, [rachelrosen@google.com](mailto:rachelrosen@google.com), Jigsaw, Google, 111 8th Ave, New York, NY, USA, 11201; Lucy Vasserman, [lucyvasserman@google.com](mailto:lucyvasserman@google.com), Jigsaw, Google, 111 8th Ave, New York, NY, USA, 11201.



This work is licensed under a Creative Commons Attribution International 4.0 License.

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

2573-0142/2022/11-ART363

<https://doi.org/10.1145/3555088>

For example, it has been shown that the Perspective API<sup>1</sup>, a publicly available API to detect toxicity in text, is more likely to predict high toxicity scores for comments written in African American English than for other comments [33]. Likewise, Gomes et al. [18] has shown that non-toxic tweets from drag queens are more likely to be scored as toxic by the Perspective API than tweets by known white supremacists.

As of today, ML models that define online content as toxic or not fall short of reaching the goal of being free from bias. As is evident from the examples above, when models fail, they do not fail equally for all identities. Some identity groups, already disenfranchised, are hurt even more than others because the prevalence of toxic comments directed towards all identity groups is not equal. Similarly, the presence of harder to detect microaggressions is likely unequal across groups. While this is not intentional, it does create inequity. ML models are based on huge datasets that have been trained on data labeled by human raters or annotators, also known as crowd workers. However, every human annotator cannot be fully aware of the intricacies of how different communities are being hurt using toxic language. Though certain phrases and words might seem innocuous to those outside an identity group, they may be known to be toxic to people who self-identify as that identity through their shared and lived experiences—or vice versa: some comments are perceived within identity groups to be non-toxic, but may be perceived by outsiders as toxic. What if we could ask those who self-identify with an identity group to rate content that was related to their community? This way, those who are likely to be targeted, and who would be best equipped to label the data, would be the ones to determine the ground truth for models that classify toxicity online. This paper continues to build upon research in this space of creating groups of annotators based on some differentiating factor(s) [3, 16, 22, 41]. More specifically, we explore how raters from two relevant identity groups, African American and LGBTQ, label data that represents those identities, and whether their ratings vary from those provided by a randomly selected pool of raters who do not self-identify with these identity groups.

## 2 RELATED WORK

Our work fits within a broader picture of trying to understand annotator identity or lived experience as a source of expertise, particularly within the context of abusive language and toxicity. We use the term “annotator” and “rater” interchangeably throughout this work to acknowledge that past literature has used both terms.

We are not the first researchers who have been interested in understanding demographics of crowd workers and how that may impact crowdwork. For example, Berg [8] and Posch et al. [29] have asked questions about gender, age, and education level, among others. Further, Posch et al. [29] found that in the U.S., female crowdworkers are the majority, but in most other countries they are the minority. With respect to age, the authors found that most crowdworkers are between 18 and 34 years old; that remains consistent with previous findings [8].

We are also not the first to explore the impact of identity in building models for online conversation. Some papers have explored the identity of comment authors [20, 34]. In particular, Halevy et al. [19] explored author dialect as a proxy for author identity, fine-tuning a toxic comment classifier on African American English to reduce racial bias in the model. And by looking at identity on the community level, Wich et al. [40] explored how conversations in different political communities have different norms. Our work expands upon this previous work by focusing on the identity of the annotators, instead of the authors.

Others have explored identity from the perspective of who the comment is targeting. Kurrek et al. [23] looked at the usage of online slurs, in particular pejorative terms used against three distinct

<sup>1</sup>[www.perspectiveapi.com](http://www.perspectiveapi.com)

groups of people—gay men, Black people, and transgender people. The authors have explicitly mentioned that due to the limitations of their research environment, understanding relationships between annotator identity and annotations is a limitation of their work, something we explore in this paper. Others have also identified annotator bias as an area for future work [7, 20]. Our work expands Kurrek et al. [23] in three ways: first, by going beyond slurs and considering impact on multiple facets: toxicity, identity attacks, profanity, insults, and threats. Second, by annotating an entirely different dataset (Civil Comments instead of Reddit) by hundreds of carefully selected raters that belong to a specific identity group instead of 20 raters of mixed identity groups. Third, Kurrek et al. [23] point out in Section 3.3 of their work that they could only perform limited analysis on the relationship between annotator demographics and annotations. They left this as an area for future work, which our work expands upon.

Other papers have considered the identity of annotators in groups, but along different lines than this work. For example, Sap et al. [33] grouped annotators based on whether they have been primed to think about dialect and race. Authors found that annotators' different understandings of the same word in the same language could lead to racial biases in machine learning models when detecting hate speech. Our work expands on these previous works by focusing on lived experiences of people instead of asking annotators to imagine those experiences. Furthermore, Waseem [37] grouped annotators based on domain expertise and found that amateur annotators (recruited without selection criteria) were more likely than experts (feminist and antiracist activists) to mislabel content as racist or sexist. This points to the importance of expertise, as expressed through the lived experiences and identity of the feminist and antiracist activists. Our work builds upon and differentiates in a few ways: we use and contribute an annotated version of the Civil Comments dataset instead of annotating Twitter, and we focus on lived identities of crowd workers and not the expertise of activists, who can be challenging to recruit, making the work difficult to scale. Annotations by specialized raters, however, can be easily scaled up using various platforms that provide the ability to hire panels of raters based on several identity metrics. Perhaps most importantly, Waseem [37] grouped feminist and antiracist activists together into one group, while we separate the identities into separate groups explicitly, showing nuanced differences between identities and annotations, i.e. that different identities are not the same and can not generically be grouped together.

Further, even when using experts who have deeper knowledge, annotating data appropriately is a challenge. In particular, Davidson et al. [11] found biases even in the expert-annotated data. Authors found that that a classifier trained on expert-annotated data flagged Black-aligned tweets as sexist at nearly twice the rate of White-aligned tweets. Based on this, they cautioned that experts, like activists, may also hold biases similar to those of academics and crowdworkers, and emphasized the need for further work on exploring expertise and its role in the annotation process. In this work, we do not focus on expertise of trained experts because such experts are not always available and can be expensive to find and recruit. Instead we focus on crowd worker annotators and use their identities and lived experiences as a proxy for expertise.

Yet other works have grouped annotators into pools by forming clusters of annotators who rate similarly [2, 3, 6] or by using community-detection algorithms on a graph of annotators [39]. Other have incorporated a diversity of annotators in their rater pools but not explicitly grouped annotators according to identity and drawn conclusions based on identity groups specifically [44].

Identity effects on annotation have also been studied from the perspective of how gender can impact inter-annotator agreement, toxicity labels, and the resultant classifier performance [9]. Wulczyn et al. [42] found that female annotators gave lower toxicity ratings on average than male annotators on the Wikipedia Detox dataset (note that there are different possible reasons for this), and analyzed which lexical features were weighted most heavily by classifiers trained on data

labeled only by male or female annotators. More recently, broad comparisons of classifiers trained on data from demographically different annotator groups based on gender, first language, age, and education performed differently as measured by the classifier F1 score (harmonic mean of precision and recall) [4]. They found that these features correlate with significant differences in classifier performance.

Identity has also been recently explored by Larimore et al. [24]. Using a relatively small dataset, the authors found that White and non-White annotators rate tweets differently, especially in the context of certain topics (e.g. police brutality, antiracist politics, or empowering history), suggesting that the standards for evaluating racist language annotations should reflect the interpretations of those who are impacted or are at the receiving end. This lays the foundation for further research with larger data sets that can focus on particular identities and unpack their impact on ratings that power machine learning models that in turn are used to moderate content on the internet.

To summarize, understanding how demographics of annotators can impact annotations has been studied in seven different ways so far in the literature:

- (1) Creating rater pools *a posteriori* by clustering raters based on their ratings and maximizing distances between these clusters instead of on the basis of their identities [2, 6, 39].
- (2) Creating identity-based pools on pre-existing datasets that looks for differences based on markers like age, gender, ESL, education e.g. on the Wikipedia Detox Dataset [22].
- (3) Creating small, expert-based pools that perform annotations based on certain markers e.g. 3 annotators annotating immigrant/native status [3, 16, 41].
- (4) As non-grouped random raters representing diversity to understand impact on annotations [26, 44]
- (5) As identity-trained classifier-evaluation where classifiers trained on different pools of raters based on markers like gender, first language, age and education are evaluated as to how they perform in terms of precision and recall [4].
- (6) As effects of comment authors' demographic differences and not the annotators' demographics as we study [20, 25, 27, 43], or creating new data sets based on political sub-communities of authors [40].
- (7) As focusing on model-robustness-evaluation to detect identity markers like race and gender correctly [36], multi-lingual hate speech classification [5], and model-fairness-evaluation [17, 19].

Our work was motivated by three pieces of recent research and these works are perhaps the closest to ours. Work by Basile et al. [7] acts as the motivation for our work since it provides theoretical grounding of why data annotations matter for machine learning. Another close work is that by Sap et al. [33], where annotators are primed to think as if they might belong to specific ethnicities. While there is indeed value in asking annotators to consider putting themselves in the shoes of others, our work builds upon Sap et al. [33]'s work by actually recruiting annotators that belong to these communities. Third, Kurrek et al. [23] presents work where the focus is on detecting slurs in Reddit data by a team of 20 annotators that are a mix of genders, ethnicities, and sexual orientations. The authors created a team of 20 annotators who represented inter-sectional identities across the three identity factors and they annotated slurs together as a team. Our work is a direct next step to this where we provide annotations on a different data set (Civil Comments dataset instead of Reddit), and we focus on toxicity, identity attacks, insults, profanity, and threats instead of slurs, while showing relationships between annotators' ethnicities and sexual orientations and their annotations. One of the goals of doing such annotations has also been to contribute this annotated dataset to the wider research community, who can use it to build models that can be used for classification and evaluation, or perform deeper qualitative analysis.

This work has the following three primary contributions:

- (1) We explore the impact on toxicity, insult, threat, identity attack, and profanity ratings of online conversations as perceived by two groups of annotators (referred to as specialized rater pools): African American and LGBTQ American raters
- (2) We are open-sourcing a large corpora of Civil Comments dataset and raters' annotations of conversations in this corpora across toxicity, insult, threat, identity attack, and profanity. We created and used this dataset to answer the above question and are sharing this dataset to encourage future research
- (3) We found that while specialized rater pools do create a statistically significant difference in annotations for online conversations as perceived by the two specialized rater pools, there are nuances to these differences, highlighting pros and cons of using specialized rater pools.

### 3 METHODS

#### 3.1 Research Questions

Based on prior research, it is clear that individual raters will rate content differently. However, it is unclear whether raters will rate toxicity differently in comments written by or about their own identities than a randomly selected pool of raters who do not self-identify with these identities. Since we are focused on two particular communities, African Americans and LGBTQ Americans, we pose the following two Research Questions:

- *RQ1: Annotations by African American-identified raters on data related to the African American community will be measurably different from data annotated by raters who do not identify as African American*
- *RQ2: Annotations by LGBTQ identified raters on data related to the LGBTQ community will be measurably different from data annotated by raters who do not identify as LGBTQ*

#### 3.2 Specialized and Control Rater Pools

We define **specialized rater pools** as groups of raters that self-identify as specific identities, instead of a randomly selected group of raters chosen irrespective of their identities. Thus, for the purposes of this paper, we have two specialized rater pools: an African American specialized rater pool and an LGBTQ American specialized rater pool. We chose to focus on these two identities because repeated bias with respect to these groups has been demonstrated in toxicity models previously [18, 33]. So, to investigate our research questions, we ran three crowdrating tasks on a sample set of data with the following U.S.-based rater pools:

- (1) A specialized rater pool with raters who identify as African American.
- (2) A specialized rater pool with raters who identify as LGBTQ.
- (3) A control rater pool with raters who identify as neither African American nor LGBTQ.

One thing to note is the language used to describe the raters in rater pool (1); we are using the community description “African American” instead of “Black” to refer to this group. This is an intentional distinction; these group names, while often used interchangeably, are not identical—African American refers to a specific ethnicity within the Black community [31]. Since previous work focused on African American English, we wanted to choose a group description that would be more likely to encompass speakers of that dialect, in order to build upon past work.

For the control rater pool (3), we alternatively considered having a rater pool consisting of raters selected at random irrespective of self-described identities. However, this would have led to interaction effects owing to identity groups present across both the control and specialized rater pools. Instead, the control rater pool consists only of raters who do not self-identify with either of

the groups. This ensures that the rater pools remain disjoint. While we recognize that in a normal crowd-rating situation, raters of all identities are included at random, the goal of this paper is to understand if rater identity impacts toxicity annotation, and so we separate raters into distinct groups to test this question.

### 3.3 Designing Rater Pools: Crowd Contributor and Identity Considerations

We worked with a third-party company to provide raters that perform annotation jobs. They recruited participants, on our behalf, to participate in this experiment. For the sake of this experiment, they managed constraints for recruiting raters, for which we specified age 18-35, English speaking, U.S. based raters.

Further, the third-party company provided a screener survey to select participants. We provided them with 50 screener questions and explanations for this purpose. These screener questions covered all aspects of toxicity and its subtypes, including identity attack, but we made sure to not include any ambiguous examples related to the experiment's identity groups in the screener questions to avoid asserting correct answers for these cases. Candidates taking the screener task saw the explanations when they got an item wrong, and were required to maintain an accuracy of 75% or higher to be considered for further participation. The third party company chose participants with the highest accuracy for the task and placed them into rater pools based on their self-described identities.

Raters with intersectional identities were chosen for a single rater pool rather than multiple rater pools so that the experimental groups remained disjoint. We did request that the third-party company keep the percentages of intersectional identities within each rater group within a tolerance of 1.5× the U.S. population averages (e.g. at the time the task was run, the LGBTQ percentage of the U.S. population was 4.5% [28], and the African American population percentage was 13.4% [35]. So we requested the percentage of the LGBTQ rater pool that also identifies as African American should be no higher than 20.1%, and the percentage of the African American rater pool that also identifies as LGBTQ should be no higher than 6.75%, 1.5× the population rates). This study has been approved by all internal review processes.

### 3.4 Ethical Considerations when Working with Rater Identity

There is a growing recognition of the importance of considering the ethical implications of how researchers perform work, especially with crowdworkers. In our context, raters perform the crowdwork of annotating potentially toxic text. We discuss the ethical considerations we applied throughout this work.

- (1) **Minimize Data Leakage:** As our third-party partner professionally manages raters, all of the three rater pools were constructed by them. They interfaced with the raters and had access to rater level information. We only received pseudonymized annotations that we analyzed. This was done to minimize leaking any identity-related information about the raters.
- (2) **Wage Fairness:** We wanted to make sure that raters were paid fairly [8]. While the raters are employed by the third-party vendor that we partnered with for the study, the vendor ensured us that raters were paid at least minimum wage for the jurisdiction.
- (3) **Limit Toxic Content Exposure:** Additionally, we restricted the toxicity of the data we asked for annotations, to limit participants' exposure to toxic content around their own identities. We did this by sampling 30% of comments above a Perspective API toxicity threshold of 0.6, and the rest of the comments below that threshold, so that according to Perspective, they would be exposed to a maximum of 30% toxic comments.



- (4) **Psychological Safety:** Additionally, we considered raters' psychological safety. We gave raters the option to use a crowdrating platform-provided chatroom to communicate about the task so that they would have support, if needed. While this work involved limited engagement with toxic content, we still wanted to provide extra support.
- (5) **Upfront Transparency:** We committed to transparency as recommended by Kazimzade and Miceli [21]. Since rater identity was being used to place raters into specialized pools, we informed raters about this. We also asked the vendor to communicate to raters that this is a study, so that we had rater consent to use the annotations for research purposes, since raters might otherwise expect the data was only for other common data annotation use cases, such as model building.

### 3.5 Dataset

The dataset used for this study is the Civil Comments dataset, which has been used previously for similar research about annotating toxic comments [10]. The full Civil Comments dataset consists of approximately 2 million news comments from a now-defunct commenting platform. It has crowdsourced labels for toxicity and toxicity-subtypes, with 22% of comments also labeled for identities. This data is public<sup>2</sup> and the dataset itself is released under a CC0 license. It contains many references to identities, which is why we chose this dataset for this study.

From this dataset of 2 million news comments, we randomly sampled identity-neutral comments, as well as comments that mention each of the two identity groups in the study (African American and LGBTQ). The identity-neutral comments were sampled by excluding comments with identity labels (as provided by the dataset itself) for these identities. Next, we controlled for the rate of toxicity using Perspective API to mitigate the negative effects of toxicity exposure on raters, as we discussed in the Ethical Considerations section above. We recognize that Perspective API will make some errors here—this will affect our efforts to mitigate toxicity exposure, but will not otherwise affect the final results of the experiment, as raters will see each comment without knowing what the Perspective API score is or how the comment was sampled.

When sampling comments that mentioned each identity, we used a combination of identity labels as provided by the dataset as well as machine learning models that detect identity mentions in the text for the same identity labels used in the Civil Comments dataset. The comments that mentioned the identities were meant to serve as a proxy for both community-specific discussions and comments written about identities. Our sampling strategy means both are included, though it is challenging to distinguish between the two.

Overall, we created a dataset that contained a total of 25,500 comments from the Civil Comments dataset, with 8500 comments sampled to be identity agnostic, and 8500 comments sampled for each identity group (8500 for LGBTQ and 8500 for AA). The complete annotated data (including all individual annotations) reflecting 382,500 annotations is available on Kaggle in CSV and TSV formats as Google Specialized Rater Pools Dataset at <https://www.kaggle.com/datasets/google/jigsaw-specialized-rater-pools-dataset><sup>3</sup>. We hope that by releasing this data, we will enable the broader research community to build models using it to better understand the differences between the way annotators in the different pools annotate. Besides building new models, the community can also use this opportunity to dig deeper and perform content analysis on this dataset, which is beyond the goals of this paper.

We next discuss the Likert scales and annotation template used for the task.

<sup>2</sup><https://www.kaggle.com/c/jigsaw-unintended-bias-in-toxicity-classification/data>

<sup>3</sup>Nitesh Goyal, Ian Kivlichan, Rachel Rosen, & Lucy Vasserman. (2022). [Data set]. Kaggle. <https://doi.org/10.34740/KAGGLE/DSV/3533200>

### 3.6 Task Design

All three rater pools were presented with the same full set of 25,500 comments, in a pre-sorted randomized order, for which we receive 5 ratings per annotator, a standard practice for similar work as shown by Larimore et al. [24]. Hence, the resulting dataset contained 15 ratings per comment: 5 ratings from annotators from each rater pool in the study. The task as seen by the raters uses the same template that has been used in previous works, for example by Dixon [14]. The task asks raters to rate comment toxicity, as well as other components of Toxicity on a Likert scale as first defined in Dixon [14] and defined inline below.

- (1) Toxicity is defined as “a rude, disrespectful, or unreasonable comment that is likely to make people leave a discussion”. This is measured on a 4-point Likert scale with values between  $-2$  and  $1$ , where  $-2$  = Very toxic,  $-1$  = Toxic,  $0$  = Unsure, and  $1$  = Not toxic.
- (2) Identity Attack is defined as “negative or hateful comments targeting someone because of their identity”. This is measured on a 3-point Likert scale from  $-1$  to  $1$ .
- (3) Insult is defined as “insulting, inflammatory, or negative comment towards a person or a group of people”. This is measured on a 3-point Likert scale from  $-1$  to  $1$ .
- (4) Profanity is defined as “swear words, curse words, or other obscene or profane language”. This is measured on a 3-point Likert scale from  $-1$  to  $1$ .
- (5) Threat describes “an intention to inflict pain, injury, or violence against an individual or group”. This is measured on a 3-point Likert scale from  $-1$  to  $1$ .

More details on the task itself are included in the Appendix: the full instructions are included in Figure 2, and sample questions in Figure 3.

### 3.7 Measures

#### 3.7.1 Descriptive Statistics.

- (1) Toxicity Mean Difference: For each annotated comment, Toxicity Mean Difference is the difference in means of scores between each of the specialized rater pools and the control group, when the annotators rated the comment on a 4-point Likert scale as “a rude, disrespectful, or unreasonable comment that is likely to make people leave a discussion”.
- (2) Identity Comments with High Agreement: This is defined as the percentage of the comments that contain identity information and raters between control and specialized rater pools agree highly such that the Toxicity Mean Difference =  $0$ .
- (3) Identity Comments with Low Agreement: This is defined as the percentage of the comments that contain identity information and raters between control and specialized rater pools disagree highly such that the Toxicity Mean Difference  $\geq 1$ .

#### 3.7.2 Regression Analysis.

- (1) Toxicity Odds Ratio: This is the proportional odds for specialized rater pools to rate annotations as more likely to be toxic
- (2) Identity Attack Odds Ratio: This is the proportional odds for specialized rater pools to rate annotations indicating higher likelihood of involving an identity attack
- (3) Insult Odds Ratio: This is the proportional odds for specialized rater pools to rate annotations as more likely to be insulting
- (4) Profanity Odds Ratio: This is the proportional odds for specialized rater pools to rate annotations as more likely to include profanity
- (5) Threat Odds Ratio: This is the proportional odds for specialized rater pools to rate annotations as more likely to be threatening



## 4 RESULTS AND ANALYSIS

### 4.1 Descriptive Statistics

**4.1.1 Rater disagreement.** For each label, we consider the histogram (counts) where raters disagree. In Figure 1 we show the probability distribution of the mean differences:  $\text{mean}(\text{specialized rater pool}) - \text{mean}(\text{control rater pool})$ , with negative differences meaning that the specialized pool rated the comment as more likely to be toxic (or another label), and positive differences meaning the specialized pool rated the comment as less likely to be toxic. Notably for all labels, the distributions of the histograms are similar on the negative and positive sides, indicating that there is not a trend towards more or less toxicity (or other labels) among the specialized pools; the disagreements go in both directions.

In Table 1 we also explore the overall amount of disagreement between the rater pools, by showing the percentage of comments where the absolute value of the mean difference is greater than or equal to 1. We find that toxicity has the largest proportion of comments with disagreement (>12% for both African American and LGBTQ rater pools), whereas the threat and profanity attributes have the least amount of disagreement, with <1%.

We also consider how toxicity itself interacts with agreement between the rater pools by looking at the percentage of comments that are toxic (mean score < 0) among high and low disagreement comments. In Table 2 we see these differences for the African American and control rater pools, and in Table 3 we see these differences for the LGBTQ and control rater pools. From this data we can see that according to all the rater pools, high disagreement comments have a higher percentage of toxicity than low disagreement comments, with the control rater pools overall finding more of these comments toxic than the specialized rater pools.

Table 1. The percentage of comments where the specialized rater pool groups disagreed with the control rater pools. We see that toxicity has the largest proportion of comments with disagreement (>12% for both African American and LGBTQ rater pools), whereas the threat and profanity attributes have the least amount of disagreement, with <1%. These are not to be read as stat. significant differences, but to highlight that some labels had more differences (Toxic Score, Identity Attack, and Insult) than others (Threat, and Profanity).

Label	% Comments with high AA-control disagree- ment	% Comments with high LGBTQ-control disagree- ment
Toxic score	12.4%	12.5%
Identity attack	5.7%	6.2%
Insult	8.5%	8.5%
Threat	0.5%	0.4%
Profanity	0.8%	0.8%

**4.1.2 Comments mentioning identity.** In Table 4 we examine the percentage of comments that mention identities for different levels of agreement between each of the African American and LGBTQ specialized rater pools and the control pool. We can see that comments with high disagreement between the specialized rater pools and the control pools tend to mention identities more often (71% of the time) compared to comments with no disagreement (50% and 49% of the time). This supports the hypothesis that specialized rater pools are critical for understanding comments that reference identity groups, since annotators in specialized rater pools are more likely to disagree with control group raters on these comments.

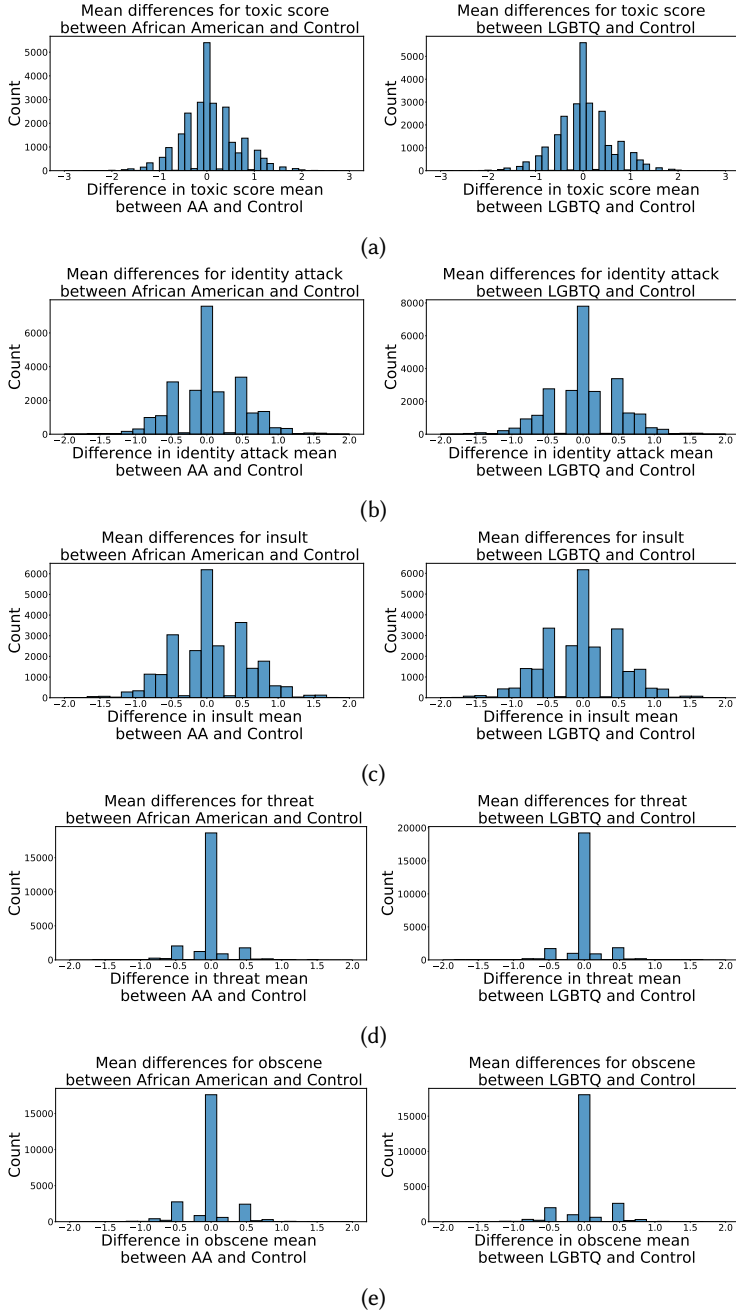


Fig. 1. Histograms of differences in scores, mean(specialized rater pool) - mean(control rater pool). The differences range from  $-3$  (for primary toxicity label;  $-2$  for other labels), meaning that the specialized group rated the comment as much more likely to be toxic (or other label) than the control, to  $3$  ( $2$  for other labels), meaning that the specialized group felt the comment was less likely to be toxic (or other label). Histograms for mean(African American) - mean(control) are on the left, and histograms for mean(LGBTQ) - mean(control) are on the right. (a) Represents the histograms for toxic score, (b) identity attack, (c) insult, (d) threat, and (e) profanity. Notably for all labels, histograms are similar on the negative and positive sides.

Table 2. The percent toxicity according to the African American and control rater pools. From the data, we can see that comments are more likely to be toxic according to both groups if there is high disagreement between the African American and control rater pools, with the control rater pool finding a higher percentage of comments to be toxic overall (7.2%) than the African American rater pool (4.5%). These are not to be read as stat. significant differences, but to highlight that AA and control pools think of Toxicity differently.

	High AA-control disagreement	Low AA-control disagreement
% of Comments that are Toxic according to African American pool	4.5%	2.3%
% of Comments that are Toxic according to control pool	7.2%	2.3%

Table 3. Percentage of comments that are toxic according to the LGBTQ and control rater pools. From the data, we can see that comments are more likely to be toxic according to both groups if there is high disagreement between the LGBTQ and control rater pools, with the control rater pool finding a higher percentage of comments to be toxic overall (6.6%) than the LGBTQ rater pool (5.3%). These are not to be read as stat. significant differences, but to highlight that LGBTQ and control pools think of Toxicity differently.

	High LGBTQ-control disagreement	Low LGBTQ-control disagreement
% of Comments that are Toxic according to LGBTQ pool	5.3%	2.4%
% of Comments that are Toxic according to control pool	6.6%	2.4%

Table 4. Percentage of comments that mention identities for different levels of agreement between African American/LGBTQ rater pools and the control rater pool. Comments with no disagreement (mean difference = 0) have lower percentages of identity mentions than comments with higher disagreement between groups (mean difference  $\geq 1$ ).

	High AA-control disagreement	No AA-control disagreement	High LGBTQ-control disagreement	No LGBTQ-control disagreement
% identity-mentioning comments	71.0	49.5	71.0	49.4

## 4.2 Regression Analysis

Since the outcome variables for Regression Analysis are not binary (4-point Likert scale or 3-point Likert scale), we ran an ordinal logistic regression analysis to evaluate RQ1 and RQ2 [1, 30]. We used the SPSS Advanced Statistics module to create dummy variables for recoding the dependent variables (Toxicity Score, Identity Attack, Insult, Profanity and Threat) and used the PLUM and GENLIN commands to perform the analysis. This test has 4 assumptions and all of them were met: (a) dependent variable (Toxicity Score, Identity Attack, Insult, Profanity and Threat) should be ordinal; (b) one or more independent variables that are continuous or categorical (Rater Pool

category: control, LGBTQ, African American); (c) There should be no multicollinearity if there are two or more continuous independent variables (we only have one non-continuous variable: Rater Pool category); (d) There should be proportional odds. Using separate Binomial Logistic Regressions on the dummy variables, we found that variables in the equation had similar  $\exp(B)$  values across the 3 dummy variable tests for all the measures, where  $\exp$  denotes the exponential and  $B$  is the coefficient estimate. Subsequently, a cumulative odds ordinal logistic regression with proportional odds was run as a final model to determine the effect of rater pool as an independent variable, with comment id and rater id as covariates, on identifying dependent variables Toxicity Score, Identity Attack, Insult, Profanity and Threat. The deviance goodness-of-fit test indicated that the model was a good fit to the observed data but most cells were sparse with zero frequencies.

**4.2.1 Toxicity Odds Ratio.** The final model statistically significantly predicted the dependent variable over and above the intercept-only model,  $\chi^2(3) = 624.02, p < .001$ . The final model didn't code control as the intercept since the control involved random participants. The odds of the control pool rating comments to be toxic was 0.957 (95% CI, 0.939 to 0.976) times that of LGBTQ rater pool, a statistically significant effect,  $\chi^2(1) = 19.88, p < .001$ . The odds of the control pool rating comments to be toxic was 0.986 (95% CI, 0.968 to 1.004) times that of the African American rater pool, not a statistically significant effect,  $\chi^2(1) = 2.27, p = .132$ . This indicates that the control pool was slightly less likely to rate comments as toxic than the LGBTQ pool, with no statistically significant difference between the control and the African American rater pool.

**4.2.2 Identity Attack Odds Ratio.** The final model statistically significantly predicted the dependent variable over and above the intercept-only model,  $\chi^2(3) = 371.51, p < .001$ . The odds of the control pool considering comments to be identity attacks was 0.905 (95% CI, 0.886 to 0.925) times that of the LGBTQ rater pool, a statistically significant effect,  $\chi^2(1) = 80.15, p < .001$ . The odds of the control pool considering comments to be identity attacks was 0.942 (95% CI, 0.923 to 0.962) times that of the African American rater pool, a statistically significant effect,  $\chi^2(1) = 31.333, p < .001$ . Similar to toxicity, this indicates that the control pool was slightly less likely to rate comments as identity attacks than the specialized pools.

**4.2.3 Insult Odds Ratio.** The final model statistically significantly predicted the dependent variable over and above the intercept-only model,  $\chi^2(3) = 715.37, p < .001$ . The odds of the control pool considering comments to have to be insults was 0.930 (95% CI, 0.912 to 0.948) times that of the LGBTQ rater pool, a statistically significant effect,  $\chi^2(1) = 53.07, p < .001$ . The odds of the control pool considering comments to be insults was 1.066 (95% CI, 1.046 to 1.086) times that of the African American rater pool, a statistically significant effect,  $\chi^2(1) = 45.113, p < .001$ . Here we see that the control pool was again less likely to rate comments as insults than the LGBTQ rater pool, but was more slightly likely to rate comments as insults than the African American pool.

**4.2.4 Profanity Odds Ratio.** The final model statistically significantly predicted the dependent variable over and above the intercept-only model,  $\chi^2(3) = 116.539, p < 0.001$ . The odds of the control pool considering comments to be profanity was 0.957 (95% CI, 0.932 to 0.982) times that of the LGBTQ rater pool, a statistically significant effect,  $\chi^2(1) = 10.74, p = .001$ . The odds of the control pool considering comments to be profaned was 0.954 (95% CI, 0.930 to 0.978) times that of the African American rater pool, a statistically significant effect,  $\chi^2(1) = 13.92, p < .001$ . Again, the control pool is less slightly likely to rate comments as profanity than both specialized pools.

**4.2.5 Threat Odds Ratio.** The final model statistically significantly predicted the dependent variable over and above the intercept-only model,  $\chi^2(3) = 229.044, p < 0.001$ . The odds of the control pool considering comments to be threats was 0.820 (95% CI, 0.784 to 0.858) times that of the LGBTQ

rater pool, a statistically significant effect,  $\chi^2(1) = 75.325, p = .0001$ . The odds of the control pool considering comments to be threats was 0.714 (95% CI, 0.684 to 0.745) times that of the African American rater pool, a statistically significant effect,  $\chi^2(1) = 13.92, p < .001$ . Again the control pool is less likely to rate comments as threats than both specialized pools.

To summarize, these results mean:

(1) RQ1:

- Despite relatively small differences, ratings by the control pool were statistically significantly less likely to be tagged as identity attacks, profanity, and threats compared to ratings by the African American rater pool
- Despite relatively small differences, ratings by the control pool were statistically significantly more likely to be tagged as insults compared to ratings by the African American rater pool
- Ratings by the control pool were not statistically significantly less likely to be tagged as toxic compared to ratings by the African American rater pool

(2) RQ2:

- Despite relatively small differences, ratings by the control pool were statistically significantly less likely to be tagged as toxic, having identity attacks, including insults, profanity, and threats compared to ratings by the LGBTQ rater pool

(3) Disagreements:

- The control pool disagrees more with the LGBTQ pool about identity attacks as compared to disagreeing about identity attacks with the African American pool
- The control pool rates more comments to be toxic than perceived by both the African American and LGBTQ pools, when there is significant disagreement

### 4.3 Preliminary Content Analysis

While there might be variability within every rater pool along such lines as political affiliation, gender, educational level, etc, the statistically significant difference shown in the previous section indicates that there are reasons other than chance for certain differences. While a full qualitative analysis of 25000+ comments is beyond the scope of this paper, we performed content analysis [38] of a sample of these comments to understand why some comments are rated differently between control pools and specialized rater pools.

We created a sample dataset by identifying the top 100 comments with high mean differences ( $\geq 1$ ) between the control pool and specialized rater pools. One of the authors next categorized these comments at very high-level domain related themes, like political, religious, gender-based etc. Given the source dataset, it is not surprising that many of the comments are political in nature. So, many of the comments that raters disagreed on were political. When multiple themes were discovered in a comment—the comment was assigned to both the domains.

By highlighting the comments that were rated significantly differently between control and specialized rater pools in the sections below, we want to share some examples of what kind of content is disagreed upon significantly. We also share how likely Perspective API based ML models consider these comments to be toxic to draw comparison against a well established ML model being used in industry and for academic research.

**4.3.1 Mean Differences between Specialized Rater Pools and Control Rater Pool.** We present in Table 5 and Table 6 a few examples of comments where the rater's African American or LGBTQ identity appears closely tied to the themes presented in the comment. We chose these comments because of their high differences in means between the control and African American/LGBTQ rater groups (they are among the top 100 with the highest differences in means between rater groups).

Table 5. Sample comments with themes closely tied to identity where the African American and control rater pools had high disagreement (high differences in mean scores). We present here the fraction of raters who rated the comment as “toxic” or “very toxic”, –1 or –2 on the Likert scale respectively. These fractions are what we use for model training, which makes them comparable to Perspective’s scores for these comments.

Comment	Perspective Toxicity score	Fraction of control Pool Raters who rated this toxic	Fraction of AA Pool Raters who rated this toxic
“In other words, he’s not brown enough for ya to be labeled a terrorist.”	0.56	1	0
“I’m sure it was merely an oversight but...not mentioned in the story is that the killer was black and the victims were white. Jus’ sayin’.”	0.38	0	1
“This guy waved a flag of treason and then shot two cops. Let me guess, he’s “mentally ill,” not a terrorist, because he’s white.”	0.78	1	0
“You apparently can say whatever you want about Mexicans, Hispanics & Black people, but the Republican Party draws the line on white women.”	0.5	1	0
“South Africa today is the sad reality of the result of failed colonialism. White South Africans would rather die than share power.”	0.53	0.8	0
“Good for the IRR. Maybe those fascist middle class UCT black students - especially that idiotic one with the ‘Decolonised Science’ rubbish - might actually learn something from this. But I won’t be holding my breath - those children wouldn’t understand democracy if it jumped up and threw a burning tyre around their necks.”	0.85	0.4	1

We also consider here the fraction of the rater scores who marked the comment as toxic. The fraction of raters who mark the comment as toxic is the score we use as input for model training, as discussed in the next section. As a point of comparison, we also include Perspective API’s scores for these comments.

This matters for two particular cases:

- (1) FALSE NEGATIVE CASE: First, where the specialized rater pool rated the comment as more toxic than the control rater pool. If we were only considering ratings from the control rater pool for model building, this could result in models having false negatives—where the comment is actually more toxic than what the model score indicates.
- (2) FALSE POSITIVE CASE: Second, where the specialized rater pool rated the comment as less toxic than the control rater pool. If we were only considering ratings from the control rater



Table 6. Sample comments with themes closely tied to identity where the LGBTQ and control rater pools had high disagreement (high differences in mean scores). We present here the fraction of raters who rated the comment as “toxic” or “very toxic”, -1 or -2 on the Likert scale respectively. This fraction of raters is the value we use for model training, which makes them comparable to Perspective’s scores for these comments.

Comment	Perspective Toxicity score	Fraction of control Pool Raters who rated this toxic	Fraction of LGBTQ Pool Raters who rated this toxic
“Here’s hoping that their cats’ <i>Toxoplasma gondii</i> brain-hijacking parasite causing women to commit suicide nearly two times more often than all other women (whether having been previously diagnosed with emotional or mental problems or not, that does not matter to all <i>T. gondii</i> caused suicides), that their very own cats will be the tipping-point to make them all finally practice what they preach. The fewer of them on the planet the sooner we can bring sustainable ecosystems back into balance for all the sane, responsible, wise, and respectable humans.”	0.43	0	1
“I Homosexuality, transgender sex-all are deviations from normal human behavior as well.”	0.74	0	1
“Because a homosexual orientation is supposed to not be acted upon whether or not one enters the priesthood? Since a priest sacrifices having a wife and family, how does one sacrifice something they don’t desire to begin with?”	0.44	0	0.8
“Yeah. Remember that Jesus told us ‘Love your neighbor as yourself. Unless, of course, if your neighbor is gay; him you shall reject as intrinsically disordered.’”	0.65	0.8	0
“I am NOT protecting or condoning the alleged behavior! I’m pointing out the ‘he did this because he is gay’ bigotry.”	0.77	0.8	0
“I feel the same fear for the gay members of my family.”	0.5	0.8	0

pool for model building, this could result in false positives—where the comment is actually less toxic than what the model score indicates.

In both cases, we see that it would be problematic to consider only ratings from the control group, where voices from the minority communities are not heard.

For most comments, the Perspective scores are closer in distance to the control rater pool scores, indicating that the current models of Perspective API could benefit from the inclusion of specialized rater pools for annotations.

Table 7. AUC scores for toxicity classifiers trained on data annotated by the three different rater groups (control, LGBTQ, and African American), compared with the Perspective API. The Perspective model is trained on much more data and outperforms our models in most cases; however, the LGBTQ data model is sometimes competitive with it even with much less data. This suggests that training models on data from specialized rater pools can yield surprising performance benefits on test sets considering identities.

Test set	Control model toxicity AUC	LGBTQ data model toxicity AUC	AA data model toxicity AUC	Perspective toxicity AUC
Synthetic test set [15]	0.973	0.987	0.972	0.994
HateCheck [32]	0.746	0.792	0.709	0.664
Civil Comments [10]	0.903	0.869	0.914	0.974

#### 4.4 Models

We also wanted to consider how models would perform when trained on the data from the study. We want to understand this because while it's interesting to see differences in the underlying data, the biases observed in prior work were in models, which get trained from the data. Therefore, it is important to test not only the data, but the models themselves.

To test this, we trained three BERT [13] models (pretrained on data from online conversations) on the data from each of the rater pools. In Table 7, we show the area under the receiver operating characteristic curve (AUC) on three different test sets, a synthetic test set [15], the HateCheck dataset from [32], and the Civil Comments test data [10]. We compare these models to Perspective API's current models [12].

Unsurprisingly, the Perspective models usually outperform any of the models trained on data from the experiment, as Perspective has far more training data. However, the LGBTQ data model from the study in this paper is sometimes competitive with this model, indicating that even with less training data, models trained on data from specialized rater pools can have performance benefits on test sets that consider identities.

We also want to consider whether the models trained on the African American and LGBTQ specialized rater pool data perform differently on identity-related comments than the model trained on the control data. To test this, for the Civil Comments test data, we analyzed those comments that are tagged with the "homosexual\_gay\_or\_lesbian" identity label for which the model trained on the LGBTQ data and the model trained on the control data most disagree. We also analyzed those comments that are tagged with the "black" identity label for which the model trained on the African American data and the model trained on the control data most disagree. These represent examples where model score differences are high, and identities are also mentioned. See some examples in Table 8 and Table 9.

We also considered the bias metrics on a synthetic based test set [15]. For more details please refer to the Appendix, where these are shown in Figure 4 for the control data model, in Figure 5 for the LGBTQ data model, and in Figure 6 for the African American data model.

While deeper qualitative analysis is beyond the goals of this current work, we noticed that when the model results deviated significantly from each other for some comments, certain patterns emerged. The BPSN metric—background-positive, subgroup-negative—was first defined in [10] and represents the ROC AUC calculated on a set of comments containing toxic comments that don't contain the identity term in question (background positive), and non-toxic comments containing the identity term (subgroup negative). As seen in Figure 4, the control model performs less well for

Table 8. Sample comments where scores from models trained on data annotated by the control and LGBTQ rater groups differ most.

Comment	control model score	LGBTQ model score
“As you know, I support the LGBT cause and object to the notion that gay marriage is a sin.”	0.520	0.115
“You and pretty much every person that is in favor of the baker discrimination do not get it. Political thought is not protected. In Colorado sexual orientation is protected. Bottom line? You do not have to make a cake for a democrat. You have to make a cake for a gay person.”	0.477	0.129
“How is heavens name does homosexuals doing what for them is perfectly natural ‘show the failure to love God?’”	0.385	0.733
“How about the LGBTQH and P’s that should cover all the bases and only the animals can complain that they are left out. I am sure there is some group that feels slighted so lets abolish sports all together.”	0.337	0.656

Table 9. Sample comments where scores from models trained on data annotated by the control and African American rater groups differ most.

Comment	control model score	African American model score
“Since all black people are the same, then Michelle agrees with all rap lyrics. Of course!”	0.561	0.247
“Democrats renounced racism. I wish I could say that of Trump’s Republicans. More white people get government assistance than black people. Some people develop dependency, but most use it as the safety net it was meant to be. It helps people, and the economy, survive hard times.”	0.504	0.202
“In 2013, whites committed 4,396 homicides, While Blacks committed 5,375 homicides, even though Blacks were only 13.2% of the population, and Whites being 64% of the population. - FBI crime statistics 2013.”	0.178	0.393
“In the 500+ days between the Tray von shooting and the Zimmerman trial, more than 11,000 blacks were killed by other blacks, an inconvenient truth for sure?”	0.193	0.377

the identity terms “gay”, “queer”, “homosexual”, and “muslim” based on BPSN scores. As seen in Figure 6, the African American annotated model improves the BPSN bias performance for the term “muslim” but performance drops for the terms “gay” and “queer”. As seen in Figure 5, the LGBTQ annotated model improves the BPSN bias performance for “queer”, “homosexual”, and “muslim”,

and slightly for the term “gay”. So, in summary, the specialized rater pool data trained models show improvements in BPSN bias scores for several terms over the control data trained models.

## 5 DISCUSSION

A large amount of machine learning research involves the use of crowdworkers, but we as researchers are only now starting to think about how the raters self-identify, and how that might affect the way they annotate toxicity. We discussed earlier the example of African American English; raters who do not identify as African American and are not trained linguists may not have a nuanced understanding of African American English, and that can affect how they rate it for toxicity [33].

### 5.1 Specialized Rater Pools for Inclusive ML Models

In this work, we utilize “specialized rater pools”: pools of raters crafted based on a particular dimension. In our case, we focus on one part of their identity, either ethnicity (African American) or sexual orientation/gender-identity (LGBTQ). With specialized rater pools, we challenge the status quo for who gets to decide what is toxic. This adds an additional question: who is it that gets to decide which models for toxicity are state of the art? As it stands now, published datasets such as Civil Comments [10] and HateCheck [32] are the gold standard, and so the annotators for the datasets are in effect the deciders. Yet specialized rater pools were not used for these annotations, so it is a majority vote from an average sampling of the population that determines which models are more effective at detecting toxicity. But what if members of the affected communities were to be able to make these decisions? After all, they are the ones who are impacted the most by toxicity directed towards their identities.

Also, we would like to imagine a world in which the gold standard for toxicity datasets was data annotated by specialized rater pools—a concept recently championed by Basile et al. [7]. Then we would not only be able to measure which models perform best, but also according to *whom*. This would force researchers and the industry to reckon with the fact that not all identity groups are served equally by existing toxicity classifiers, and would encourage development of models that work well for more communities, as determined by specialized rater pool annotations.

### 5.2 When to Use Specialized Rater Pools?

In our Results section, we found that raters from specialized rater pools rated comments statistically significantly differently than those from control pools in multiple measures. In particular, while LGBTQ specialized rater pools rated significantly differently than the control pool in their toxicity score, the African American and control pool showed no significant difference in their toxicity scores. This highlights that there is something about toxicity score and LGBTQ rater pools that is worth further investigation. Is it that the toxicity score itself is not the best metric to test for identity related differences? Perhaps we should be paying closer attention to other metrics like identity attacks where there were statistically significant differences across all the pools. For example, we found that there were higher percentages of identity mentions in comments with high disagreement between the control and specialized rater pools than in comments with low disagreement. More specifically, this means that if specialized and control rater pools disagree on a comment, the comment is more likely to contain an identity reference than if the two rater pools had agreed on the comment. Therefore, we recommend that for comments containing identity references, specialized rater pools should be used to confirm whether or not there is a difference in opinion between groups, which if not addressed could lead to model bias due to the control group’s lack of understanding of in-group language around an identity group, or lack of experience being on the receiving end of toxicity towards an identity group.

Often, datasets focused on toxicity and bias contain many identity references [10, 32, 37]. Therefore, we hypothesize that the use of specialized rater pools to annotate these datasets would produce differing annotations than what is currently published. Given that these datasets are used to determine what is the state of the art for toxicity modeling, we hypothesize that the “best” models (as determined today) may not continue to be seen as the best if data was relabeled with specialized rater pools. Further, our work provides a new recommendation for other researchers. Researchers should pay attention to determine if their datasets have identity mentions or attacks. If they do, they should consider using specialized rater pools for data annotation.

### 5.3 Specialized Rater Pools beyond Ethnicity and Sexual Orientation

In Section 3, we presented a model of how such pools can be created, and what considerations are necessary when crafting such pools and doing research related to the identity of crowdworkers. We have shown that identity can make a significant impact on the way annotators annotate toxicity, identity attacks, insults, threats, and profanity in text. This opens up the design and empirical space for ML engineers and researchers to use our work as an example and consider all the different kinds of specialized rater pools we should be imagining, considering, designing, and testing with. It is equally important to perform this work as participatory research in collaboration with all the different rater pools, which we strived for.

Our goal in this research is to ask researchers and industry professionals to think more critically about who the annotators are for their data and how their identities are impacting the annotations, since this in turn impacts the models they build and models that other researchers build that depend on their data. Using specialized rater pools may enable researchers to build models that work better for more identity groups. While this does not resolve all bias issues (because, for example, members of marginalized communities may still show bias towards their own communities), it still shows a good faith effort towards reducing model bias by directly consulting the experts—the communities themselves.

### 5.4 Specialized Rater Pools Dataset and Future Directions

One of the contributions of this work is the creation and publication of the specialized rater pool dataset, as discussed in Section 3.5. This consists of 382,500 annotations of 25,500 comments in the Civil Comments dataset by three carefully curated pools of raters that self-identify as LGBTQ, African American, and neither LGBTQ nor African American. We aim to encourage the research community to dig more deeply into Google Specialized Rater Pools Dataset at <https://www.kaggle.com/datasets/google/jigsaw-specialized-rater-pools-dataset><sup>4</sup>, and perform subsequent analysis into when and how identity-based specialized rater pools might be integrated into ML development. As a starting point, this could include replicating our results. We additionally hope that the community will use this dataset as an opportunity to expand upon this work by imagining and creating their own specialized rater pools across other identities.

Further, we performed preliminary content analysis in the paper. Our initial findings suggest that the identity attack, threat, and profanity categories are the differentiating reasons for different ratings between the control pool and specialized rater pools. Qualitative researchers can choose to perform deeper linguistic and content analysis to identify the specific markers in language that lead to these differences. Practitioners can use this dataset and approach to compare how their existing datasets will perform when rated by specialized rater pools.

<sup>4</sup>Nitesh Goyal, Ian Kivlichan, Rachel Rosen, & Lucy Vasserman. (2022). [Data set]. Kaggle. <https://doi.org/10.34740/KAGGLE/DSV/3533200>

Finally, we expect this dataset to be useful for the machine learning community. In lieu of identity group information, several recent works have constructed simulated groups using methods like graph clustering on the basis of annotators' responses [2, 6, 39]. The dataset we are releasing here makes this information accessible, enabling further evaluation and comparison of these methods with the underlying identity groups. Additionally, ML practitioners could further analyze the predictions generated by models trained on different identity pools' annotations, expanding upon our work in Section 4.4. For example, the data in Sap et al. [33] could be further evaluated using models trained on data from these three specialized rater pools so as to further understand different sources of potential bias. Overall, we hope that this specialized rater pool dataset will create new opportunities and research directions for academics and practitioners alike.

## 6 LIMITATIONS

The specialized rater pools in this study were limited to African Americans and LGBTQ Americans, but future research should be expanded to include other groups as well to explore if these findings would hold true for other identity-related rater pools. Similarly, even though the comments should be rated by the group that is targeted, we need more research to understand what happens when different pools disagree. This becomes increasingly important when working with intersectional identities. This work focused on binary identities, and future work should go beyond binary identities.

Future work could also include more community-specific data sources to look at community-specific discussion, instead of discussions about communities, to be more reflective of the relationship between conversations and the identity-related linguistic markers. Future work should also consider deeper qualitative research methods to fully understand how and why the data differs.

## 7 CONCLUSION

The intent of the research was to determine if specialized rater pools, made up of individuals who identify as African American or LGBTQ American, rate comments differently than control rater pools made up of individuals who don't identify as these identities. We show this to be true in several cases. We also look at the performance on models trained on data from the study and show that even models trained on smaller datasets labeled by specialized rater pools can perform better than models trained on larger datasets labeled by randomly assigned annotators.

## ACKNOWLEDGMENTS

We thank Olivia Redfield and Raquel Saxe for contributions to the early stages of this work. We also would like to thank Alyssa Whitlock Lees, Jeffrey Sorensen and other anonymous reviewers for suggestions on improving this work.

## REFERENCES

- [1] Alan Agresti. 2002. *Categorical Data Analysis*. Wiley Series in Probability and Statistics, Vol. 482. Wiley, Hoboken, New Jersey. <https://doi.org/10.1002/0471249688>
- [2] Sohail Akhtar, Valerio Basile, and Viviana Patti. 2020. Modeling Annotator Perspective and Polarized Opinions to Improve Hate Speech Detection. *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing* 8, 1 (Oct. 2020), 151–154. <https://ojs.aaai.org/index.php/HCOMP/article/view/7473>
- [3] Sohail Akhtar, Valerio Basile, and Viviana Patti. 2021. Whose Opinions Matter? Perspective-aware Models to Identify Opinions of Hate Speech Victims in Abusive Language Detection. *arXiv preprint arXiv:2106.15896* (2021).
- [4] Hala Al Kuwatly, Maximilian Wich, and Georg Groh. 2020. Identifying and Measuring Annotator Bias Based on Annotators' Demographic Characteristics. In *Proceedings of the Fourth Workshop on Online Abuse and Harms*. Association for Computational Linguistics, Online, 184–190. <https://doi.org/10.18653/v1/2020.alw-1.21>
- [5] Sai Saketh Aluru, Binny Mathew, Punyajoy Saha, and Animesh Mukherjee. 2020. Deep learning models for multilingual hate speech detection. *arXiv preprint arXiv:2004.06465* (2020).



- [6] Valerio Basile. 2021. It's the End of the Gold Standard as We Know It. In *AIxIA 2020 – Advances in Artificial Intelligence*, Matteo Baldoni and Stefania Bandini (Eds.). Springer International Publishing, Cham, 441–453.
- [7] Valerio Basile, Federico Cabitza, Andrea Campagner, and Michael Fell. 2021. Toward a Perspectivist Turn in Ground Truthing for Predictive Computing. *arXiv preprint arXiv:2109.04270* (2021).
- [8] 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.
- [9] Reuben Binns, Michael Veale, Max Van Kleek, and Nigel Shadbolt. 2017. Like Trainer, Like Bot? Inheritance of Bias in Algorithmic Content Moderation. In *Social Informatics*, Giovanni Luca Ciampaglia, Afra Mashhadi, and Taha Yasseri (Eds.). Springer International Publishing, Cham, 405–415.
- [10] Daniel Borkan, Lucas Dixon, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. 2019. Nuanced Metrics for Measuring Unintended Bias with Real Data for Text Classification. In *Companion Proceedings of The 2019 World Wide Web Conference* (San Francisco, USA) (WWW '19). Association for Computing Machinery, New York, NY, USA, 491–500. <https://doi.org/10.1145/3308560.3317593>
- [11] Thomas Davidson, Debasmita Bhattacharya, and Ingmar Weber. 2019. Racial Bias in Hate Speech and Abusive Language Detection Datasets. In *Proceedings of the Third Workshop on Abusive Language Online*. Association for Computational Linguistics, Florence, Italy, 25–35. <https://doi.org/10.18653/v1/W19-3504>
- [12] The Perspective Developers. 2020. Perspective: Model Cards. <http://developers.perspectiveapi.com/s/about-the-api-model-cards>. Accessed: 2021-07-13.
- [13] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Association for Computational Linguistics, Minneapolis, Minnesota, 4171–4186. <https://doi.org/10.18653/v1/N19-1423>
- [14] 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](https://github.com/conversationai/conversationai.github.io/blob/master/crowdsourcing_annotation_schemes/toxicity_with_subattributes.md). Accessed: 2020-10-12.
- [15] Lucas Dixon, John Li, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. 2018. Measuring and Mitigating Unintended Bias in Text Classification. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society* (New Orleans, LA, USA) (AI/ES '18). Association for Computing Machinery, New York, NY, USA, 67–73. <https://doi.org/10.1145/3278721.3278729>
- [16] Armend Duzha, Cristiano Casadei, Michael Tosi, and Fabio Celli. 2021. Hate versus politics: detection of hate against policy makers in Italian tweets. *SN Social Sciences* 1, 9 (2021), 1–15.
- [17] Oguzhan Gencoglu. 2021. Cyberbullying Detection With Fairness Constraints. *IEEE Internet Computing* 25, 1 (2021), 20–29. <https://doi.org/10.1109/MIC.2020.3032461>
- [18] Alessandra Gomes, Dennys Antonialli, and Thiago Dias Oliva. 2019. Drag queens and Artificial Intelligence: should computers decide what is 'toxic' on the internet? <https://www.internetlab.org.br/en/freedom-of-expression/drag-queens-and-artificial-intelligence-should-computers-decide-what-is-toxic-on-the-internet>. Accessed: 2020-09-01.
- [19] Matan Halevy, Camille Harris, Amy Bruckman, Diyi Yang, and Ayanna Howard. 2021. Mitigating Racial Biases in Toxic Language Detection with an Equity-Based Ensemble Framework. In *Equity and Access in Algorithms, Mechanisms, and Optimization* (New York, NY, USA) (EAAMO '21). Association for Computing Machinery, New York, NY, USA, Article 7, 11 pages. <https://doi.org/10.1145/3465416.3483299>
- [20] Xiaolei Huang, Linzi Xing, Franck Dernoncourt, and Michael J. Paul. 2020. Multilingual Twitter Corpus and Baselines for Evaluating Demographic Bias in Hate Speech Recognition. In *Proceedings of the 12th Language Resources and Evaluation Conference*. European Language Resources Association, Marseille, France, 1440–1448. <https://aclanthology.org/2020.lrec-1.180>
- [21] 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* (New York, NY, USA) (AI/ES '20). Association for Computing Machinery, New York, NY, USA, 71. <https://doi.org/10.1145/3375627.3375809>
- [22] Jan Kocoń, Alicja Figas, Marcin Gruza, Daria Puchalska, Tomasz Kajdanowicz, and Przemysław Kazienko. 2021. Offensive, aggressive, and hate speech analysis: From data-centric to human-centered approach. *Information Processing & Management* 58, 5 (2021), 102643.
- [23] Jana Kurrek, Haji Mohammad Saleem, and Derek Ruths. 2020. Towards a Comprehensive Taxonomy and Large-Scale Annotated Corpus for Online Slur Usage. In *Proceedings of the Fourth Workshop on Online Abuse and Harms*. Association for Computational Linguistics, Online, 138–149. <https://doi.org/10.18653/v1/2020.alw-1.17>
- [24] Savannah Larimore, Ian Kennedy, Breon Haskett, and Alina Arseniev-Koehler. 2021. Reconsidering Annotator Disagreement about Racist Language: Noise or Signal?. In *Proceedings of the Ninth International Workshop on Natural Language Processing for Social Media*. Association for Computational Linguistics, Online, 81–90. <https://doi.org/10.1145/3465416.3483299>

18653/v1/2021.socialnlp-1.7

- [25] Haochen Liu, Wei Jin, Hamid Karimi, Zitao Liu, and Jiliang Tang. 2021. The Authors Matter: Understanding and Mitigating Implicit Bias in Deep Text Classification. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*. Association for Computational Linguistics, Online, 74–85. <https://doi.org/10.18653/v1/2021.findings-acl.7>
- [26] Saif M Mohammad. 2021. Ethics Sheet for Automatic Emotion Recognition and Sentiment Analysis. *arXiv preprint arXiv:2109.08256* (2021).
- [27] Edoardo Mosca, Maximilian Wich, and Georg Groh. 2021. Understanding and Interpreting the Impact of User Context in Hate Speech Detection. In *Proceedings of the Ninth International Workshop on Natural Language Processing for Social Media*. Association for Computational Linguistics, Online, 91–102. <https://doi.org/10.18653/v1/2021.socialnlp-1.8>
- [28] Frank Newport. 2018. In US, estimate of LGBT population rises to 4.5%. *Gallup News* (2018).
- [29] 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).
- [30] Daryl Pregibon. 1981. Logistic regression diagnostics. *The annals of statistics* 9, 4 (1981), 705–724.
- [31] Michael Quander and Lauryn Froneberger. 2019. Black vs. African-American: The complex conversation Black Americans are having about identity #ForTheCulture. WUSA (2019). <http://www.wusa9.com/article/news/local/black-history/black-vs-african-american-the-complex-conversation-black-americans-are-having-about-identity-fortheculture/65-80dde243-23be-4cfb-9b0f-bf5898bcf069>
- [32] Paul Röttger, Bertie Vidgen, Dong Nguyen, Zeerak Waseem, Helen Margetts, and Janet Pierrehumbert. 2021. HateCheck: Functional Tests for Hate Speech Detection Models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. Association for Computational Linguistics, Online, 41–58. <https://doi.org/10.18653/v1/2021.acl-long.4>
- [33] Maarten Sap, Dallas Card, Saadia Gabriel, Yejin Choi, and Noah A. Smith. 2019. The Risk of Racial Bias in Hate Speech Detection. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Florence, Italy, 1668–1678. <https://doi.org/10.18653/v1/P19-1163>
- [34] Shivashankar Subramanian, Xudong Han, Timothy Baldwin, Trevor Cohn, and Lea Frermann. 2021. Evaluating Debiasing Techniques for Intersectional Biases. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Online and Punta Cana, Dominican Republic, 2492–2498. <https://doi.org/10.18653/v1/2021.emnlp-main.193>
- [35] U.S. Census Bureau. 2019. QuickFacts: United States. <https://www.census.gov/quickfacts/fact/table/US/RHI125219>. Accessed: 2021-07-13.
- [36] Jialu Wang, Yang Liu, and Xin Eric Wang. 2021. Assessing Multilingual Fairness in Pre-trained Multimodal Representations. *arXiv preprint arXiv:2106.06683* (2021).
- [37] 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>
- [38] Robert Philip Weber. 1990. *Basic content analysis* (2 ed.). Sage, Thousand Oaks, California. <https://doi.org/10.4135/9781412983488>
- [39] Maximilian Wich, Hala Al Kuwatly, and Georg Groh. 2020. Investigating Annotator Bias with a Graph-Based Approach. In *Proceedings of the Fourth Workshop on Online Abuse and Harms*. Association for Computational Linguistics, Online, 191–199. <https://doi.org/10.18653/v1/2020.alw-1.22>
- [40] Maximilian Wich, Jan Bauer, and Georg Groh. 2020. Impact of Politically Biased Data on Hate Speech Classification. In *Proceedings of the Fourth Workshop on Online Abuse and Harms*. Association for Computational Linguistics, Online, 54–64. <https://doi.org/10.18653/v1/2020.alw-1.7>
- [41] Maximilian Wich, Adrian Gorniak, Tobias Eder, Daniel Bartmann, Burak Enes Cakici, and Georg Groh. 2021. Introducing an Abusive Language Classification Framework for Telegram to Investigate the German Hater Community. *arXiv preprint arXiv:2109.07346* (2021).
- [42] Ellery Wulczyn, Nithum Thain, and Lucas Dixon. 2017. Ex Machina: Personal Attacks Seen at Scale. In *Proceedings of the 26th International Conference on World Wide Web (Perth, Australia) (WWW '17)*. International World Wide Web Conferences Steering Committee, Republic and Canton of Geneva, CHE, 1391–1399. <https://doi.org/10.1145/3038912.3052591>
- [43] Michael Miller Yoder. 2021. *Computational Models of Identity Presentation in Language*. Ph.D. Dissertation. Carnegie Mellon University.
- [44] Philine Zeinert, Nanna Inie, and Leon Derczynski. 2021. Annotating Online Misogyny. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. Association for Computational Linguistics, Online, 3181–3197. <https://doi.org/10.18653/v1/2021.acl-long.247>

A RATER TEMPLATE

In this appendix, we include the rating template (Figure 2) and examples for the rating instructions (Figure 3) used for this work.

In this task, you will be given a list of online comments. For each comment, your job is to:

- 1. Read the comment.
- 2. If the comment is in a **foreign language** or **not comprehensible** for another reason (e.g. gibberish, different dialect, etc.), please indicate that by selecting the checkbox.
- 3. Answer a set of questions about the comment.
- 4. Choose the level of toxicity in the comment.

To answer the questions about each comment, please use the following definitions as a guide:

Category	Definition
Profanity/Obscenity	Swear words, curse words, or other obscene or profane language.
Identity based negativity	A negative, discriminatory, or hateful comment against a group of people based on criteria including (but not limited to) race or ethnicity, religion, gender, nationality or citizenship, disability, age, or sexual orientation.
Insults	Inflammatory, stereotyping, insulting, or negative language towards a person or a group of people.
Threatening	Language that is threatening or encouraging violence or harm, including self-harm.

To evaluate the **toxicity of language** in each comment, please use the following definitions as a guide:

Toxicity Level	Description of level
Very Toxic	A comment that is very hateful, aggressive, disrespectful, or otherwise very likely to make a user leave a discussion or give up on sharing their perspective.
Toxic	A comment that is rude, disrespectful, unreasonable, or otherwise somewhat likely to make a user leave a discussion or give up on sharing their perspective.
I'm not sure	The comment could be interpreted as toxic depending on the context but you are not sure.
Not toxic	A neutral, civil, or even nice comment very unlikely to discourage the conversation.

☐ This comment is in a foreign language or not comprehensible for another reason (e.g. gibberish, different dialect, etc.)

Does this comment contain swear words, curse words, or other obscene or profane language? (required)

- ☐ Yes
- ☐ No
- ☐ Maybe - I'm not sure

Is the comment inflammatory, stereotyping, insulting, or negative towards a person or a group of people? If you are in doubt, err on the side of Yes or Maybe - I'm not sure. (required)

- ☐ Yes
- ☐ No
- ☐ Maybe - I'm not sure

Does the comment contain threatening language? Does the comment encourage violence or harm, including self-harm? (required)

- ☐ Yes
- ☐ No
- ☐ Maybe - I'm not sure

Is this a negative, discriminatory, or hateful comment against a group of people based on criteria including (but not limited to) race or ethnicity, religion, gender, nationality or citizenship, disability, age, or sexual orientation? (required)

- ☐ Yes
- ☐ No
- ☐ I'm not sure

Please rate the toxicity level of this comment. (required)

- ☐ Very Toxic
- ☐ Toxic
- ☐ I'm not sure
- ☐ Not Toxic

Fig. 2. Rating template used for this work.

## Examples

All women are gold-digging parasites looking to prostitute themselves to the highest bidder their beauty can attract.	Is the comment in a <b>foreign language</b> or <b>not comprehensible</b> for another reason (e.g. gibberish, different dialect, etc.)?	No	The comment contains dehumanizing language about women by calling them "parasites". It also invokes a demeaning, sexist stereotype.
	Does this comment contain <b>swear words</b> , <b>curse words</b> , or other <b>obscene</b> or <b>profane</b> language?	No	
	Is the comment <b>inflammatory</b> , <b>stereotyping</b> , <b>insulting</b> , or <b>negative</b> towards a person or a group of people? If you are in doubt, err on the side of Yes or I'm not sure.	Yes	
	Does the comment contain <b>threatening</b> language? Does the comment <b>encourage violence</b> or harm, including self-harm?	No	
	Is this a <b>negative</b> , <b>discriminatory</b> , or <b>hateful comment</b> against a group of people based on criteria including (but not limited to) <b>race or ethnicity</b> , <b>religion</b> , <b>gender</b> , <b>nationality</b> or <b>citizenship</b> , <b>disability</b> , <b>age</b> , or <b>sexual orientation</b> ?	Yes	
	Please rate the <b>toxicity level</b> of this comment.	Very Toxic	
Good fucking lord.	Is the comment in a <b>foreign language</b> or <b>not comprehensible</b> for another reason (e.g. gibberish, different dialect, etc.)?	No	This is a more subtle version of toxicity perhaps more passive aggressive but it is likely to have an impact on the conversation. In addition, many will consider "fucking" to be obscene.
	Does this comment contain <b>swear words</b> , <b>curse words</b> , or other <b>obscene</b> or <b>profane</b> language?	Yes	
	Is the comment <b>inflammatory</b> , <b>stereotyping</b> , <b>insulting</b> , or <b>negative</b> towards a person or a group of people? If you are in doubt, err on the side of Yes or I'm not sure.	I'm not sure	
	Does the comment contain <b>threatening</b> language? Does the comment <b>encourage violence</b> or harm, including self-harm?	No	
	Is this a <b>negative</b> , <b>discriminatory</b> , or <b>hateful comment</b> against a group of people based on criteria including (but not limited to) <b>race or ethnicity</b> , <b>religion</b> , <b>gender</b> , <b>nationality</b> or <b>citizenship</b> , <b>disability</b> , <b>age</b> , or <b>sexual orientation</b> ?	No	
	Please rate the <b>toxicity level</b> of this comment.	I'm not sure	
ouch bad move mother fucker. i'm coming for you now,, going to rape you in your sleep.	Is the comment in a <b>foreign language</b> or <b>not comprehensible</b> for another reason (e.g. gibberish, different dialect, etc.)?	No	Threatening comment and containing obscene content "mother fucker". Would count as Very toxic.
	Does this comment contain <b>swear words</b> , <b>curse words</b> , or other <b>obscene</b> or <b>profane</b> language?	Yes	
	Is the comment <b>inflammatory</b> , <b>stereotyping</b> , <b>insulting</b> , or <b>negative</b> towards a person or a group of people? If you are in doubt, err on the side of Yes or I'm not sure.	No	
	Does the comment contain <b>threatening</b> language? Does the comment <b>encourage violence</b> or harm, including self-harm?	Yes	
	Is this a <b>negative</b> , <b>discriminatory</b> , or <b>hateful comment</b> against a group of people based on criteria including (but not limited to) <b>race or ethnicity</b> , <b>religion</b> , <b>gender</b> , <b>nationality</b> or <b>citizenship</b> , <b>disability</b> , <b>age</b> , or <b>sexual orientation</b> ?	No	
	Please rate the <b>toxicity level</b> of this comment.	Very Toxic	
That's fucking amazing ! thanks for sharing.	Is the comment in a <b>foreign language</b> or <b>not comprehensible</b> for another reason (e.g. gibberish, different dialect, etc.)?	No	Contains "fucking" which would count as profane language however it is used in a "friendly" way so would not count as toxic.
	Does this comment contain <b>swear words</b> , <b>curse words</b> , or other <b>obscene</b> or <b>profane</b> language?	Yes	
	Is the comment <b>inflammatory</b> , <b>stereotyping</b> , <b>insulting</b> , or <b>negative</b> towards a person or a group of people? If you are in doubt, err on the side of Yes or I'm not sure.	No	
	Does the comment contain <b>threatening</b> language? Does the comment <b>encourage violence</b> or harm, including self-harm?	No	
	Is this a <b>negative</b> , <b>discriminatory</b> , or <b>hateful comment</b> against a group of people based on criteria including (but not limited to) <b>race or ethnicity</b> , <b>religion</b> , <b>gender</b> , <b>nationality</b> or <b>citizenship</b> , <b>disability</b> , <b>age</b> , or <b>sexual orientation</b> ?	No	
	Please rate the <b>toxicity level</b> of this comment.	Not Toxic	

Most importantly, don't forget that we're all humans, who make mistakes. And take it easy !	Is the comment in a <b>foreign language</b> or <b>not comprehensible</b> for another reason (e.g. gibberish, different dialect, etc.)?	No	Nothing negative about particular comment.
	Does this comment contain <b>swear words</b> , <b>curse words</b> , or other <b>obscene</b> or <b>profane</b> language?	No	
	Is the comment <b>inflammatory</b> , <b>stereotyping</b> , <b>insulting</b> , or <b>negative</b> towards a person or a group of people? If you are in doubt, err on the side of Yes or I'm not sure.	No	
	Does the comment contain <b>threatening</b> language? Does the comment <b>encourage violence</b> or harm, including self-harm?	No	
	Is this a <b>negative</b> , <b>discriminatory</b> , or <b>hateful</b> comment against a group of people based on criteria including (but not limited to) <b>race</b> or <b>ethnicity</b> , <b>religion</b> , <b>gender</b> , <b>nationality</b> or <b>citizenship</b> , <b>disability</b> , <b>age</b> , or <b>sexual orientation</b> ?	No	
You're a real idiot, you know that.	Please rate the <b>toxicity level</b> of this comment.	Not Toxic	Insulting and toxic comment towards other speaker. Likely to lead to conversation stopping or escalating.
	Is the comment in a <b>foreign language</b> or <b>not comprehensible</b> for another reason (e.g. gibberish, different dialect, etc.)?	No	
	Does this comment contain <b>swear words</b> , <b>curse words</b> , or other <b>obscene</b> or <b>profane</b> language?	No	
	Is the comment <b>inflammatory</b> , <b>stereotyping</b> , <b>insulting</b> , or <b>negative</b> towards a person or a group of people? If you are in doubt, err on the side of Yes or I'm not sure.	Yes	
	Does the comment contain <b>threatening</b> language? Does the comment <b>encourage violence</b> or harm, including self-harm?	No	
Peter is an idiot.	Is this a <b>negative</b> , <b>discriminatory</b> , or <b>hateful</b> comment against a group of people based on criteria including (but not limited to) <b>race</b> or <b>ethnicity</b> , <b>religion</b> , <b>gender</b> , <b>nationality</b> or <b>citizenship</b> , <b>disability</b> , <b>age</b> , or <b>sexual orientation</b> ?	No	Insulting and toxic comment towards other speaker. Likely to lead to conversation stopping or escalating. However, even if this comment is about a specific person it is NOT identity based negativity in the sense that we define it because it does not have to do with a criteria such as race or ethnicity, religion, gender, nationality or citizenship, disability, age, or sexual orientation.
	Please rate the <b>toxicity level</b> of this comment.	Toxic	
	Is the comment in a <b>foreign language</b> or <b>not comprehensible</b> for another reason (e.g. gibberish, different dialect, etc.)?	No	
	Does this comment contain <b>swear words</b> , <b>curse words</b> , or other <b>obscene</b> or <b>profane</b> language?	No	
	Is the comment <b>inflammatory</b> , <b>stereotyping</b> , <b>insulting</b> , or <b>negative</b> towards a person or a group of people? If you are in doubt, err on the side of Yes or I'm not sure.	Yes	
French people drink wine every evening.	Does the comment contain <b>threatening</b> language? Does the comment <b>encourage violence</b> or harm, including self-harm?	No	This is certainly stereotypical but NOT identity based negativity in the sense that we define it because it is not negative or hateful or otherwise promoting discrimination. If the comment had been "French people are raging alcoholics" then it would be considered as identity based negativity.
	Is this a <b>negative</b> , <b>discriminatory</b> , or <b>hateful</b> comment against a group of people based on criteria including (but not limited to) <b>race</b> or <b>ethnicity</b> , <b>religion</b> , <b>gender</b> , <b>nationality</b> or <b>citizenship</b> , <b>disability</b> , <b>age</b> , or <b>sexual orientation</b> ?	No	
	Please rate the <b>toxicity level</b> of this comment.	Not Toxic	
	Is the comment in a <b>foreign language</b> or <b>not comprehensible</b> for another reason (e.g. gibberish, different dialect, etc.)?	No	
	Does this comment contain <b>swear words</b> , <b>curse words</b> , or other <b>obscene</b> or <b>profane</b> language?	No	

Fig. 3. Examples given for the rating instructions for this work.

## B MODEL BIAS METRICS

We list bias metrics for the three models we trained on data from the three different specialized rater pools. [Figure 4](#) lists metrics for the control pool model, [Figure 5](#) for the LGBTQ pool model, and [Figure 6](#) for the African-American pool model.

	bpsn_score_subgroup_pos	bpsn_score_subgroup_neg	bpsn_score_subgroup_avg
lesbian	0.99	0.97	0.99
gay	0.98	0.87	0.99
bisexual	1.0	0.99	0.99
transgender	1.0	0.97	0.99
trans	1.0	0.99	0.98
queer	0.97	0.65	1.0
light	1.0	0.99	0.98
lightq	1.0	1.0	0.97
homosexual	0.96	0.55	1.0
straight	1.0	1.0	0.96
heterosexual	1.0	0.99	0.99
male	1.0	1.0	0.98
female	1.0	1.0	0.96
nonbinary	1.0	1.0	0.94
african	1.0	0.98	0.99
african american	1.0	0.98	0.99
black	1.0	0.97	0.99
white	0.99	0.95	0.99
european	1.0	1.0	0.98
hispanic	1.0	0.99	0.98
latino	1.0	1.0	0.98
latina	1.0	1.0	0.98
latina	1.0	1.0	0.98
mexican	1.0	0.98	0.99
canadian	1.0	1.0	0.96
american	1.0	0.99	0.99
asian	1.0	0.99	0.96
indian	1.0	1.0	0.98
middle eastern	1.0	0.99	0.99
chinese	1.0	0.99	0.98
japanese	1.0	1.0	0.97
christian	1.0	1.0	0.99
muslim	1.0	0.84	1.0
jewish	1.0	0.98	0.99
buddhist	1.0	1.0	0.97
catholic	1.0	1.0	0.99
protestant	1.0	1.0	0.97
sikh	1.0	1.0	0.97
taoist	1.0	1.0	0.96
old	1.0	1.0	0.95
elder	1.0	1.0	0.91
young	1.0	1.0	0.92
younger	1.0	1.0	0.9
teenage	1.0	1.0	0.94
millennial	1.0	1.0	0.98
middle aged	1.0	1.0	0.96
elderly	1.0	1.0	0.94
blind	0.99	0.99	0.97
deaf	1.0	1.0	0.96
paralyzed	1.0	1.0	0.94

Fig. 4. Bias metrics for the control data model. The BPSN metric—background-positive, subgroup-negative—represents the ROC AUC calculated on a set of comments containing toxic comments that don’t contain the identity term in question (background positive), and non-toxic comments containing the identity term (subgroup negative). As seen here, the control model has lower BPSN scores for the identity terms “gay”, “queer”, “homosexual”, and “muslim”.



	bpsn_score_subgroup_auc	bpsn_score_bgm_auc	bpsn_score_bgm_auc
lesbian	1.0	0.99	0.99
gay	0.97	0.89	1.0
bisexual	1.0	0.99	1.0
transgender	1.0	0.98	1.0
trans	1.0	0.99	1.0
queer	0.98	0.9	1.0
light	1.0	0.99	0.99
light	1.0	0.99	0.99
homosexual	0.99	0.96	1.0
straight	0.99	0.99	0.99
heterosexual	1.0	1.0	0.99
male	1.0	1.0	0.99
female	1.0	1.0	0.99
nonbinary	1.0	0.99	0.99
african	1.0	0.99	1.0
african american	1.0	0.99	0.99
black	0.99	0.98	1.0
white	0.99	0.97	1.0
european	1.0	0.99	1.0
hispanic	1.0	0.99	0.99
latino	1.0	0.99	0.99
latina	1.0	0.99	1.0
latina	0.98	0.99	0.97
mexican	1.0	0.98	1.0
canadian	1.0	1.0	0.99
american	1.0	0.99	0.99
asian	1.0	0.99	0.99
indian	1.0	1.0	1.0
middle eastern	0.99	1.0	0.96
chinese	1.0	0.99	0.99
japanese	0.99	0.99	0.98
christian	1.0	1.0	0.99
muslim	1.0	0.98	1.0
jewish	1.0	0.98	1.0
buddhist	1.0	1.0	0.98
catholic	1.0	1.0	0.99
protestant	1.0	1.0	0.99
sikh	1.0	0.99	0.99
baptist	0.98	1.0	0.92
old	0.98	0.99	0.96
elder	0.99	1.0	0.97
young	0.99	1.0	0.94
younger	0.99	1.0	0.96
teenage	1.0	1.0	0.98
millennial	1.0	0.99	0.99
middle aged	1.0	1.0	0.97
elderly	1.0	0.99	0.99
blind	0.98	0.96	0.99
deaf	0.99	0.98	0.99
paralyzed	0.99	0.99	0.98

Fig. 5. Bias metrics for the LGBTQ data model. The BPSN metric—background-positive, subgroup-negative—represents the ROC AUC calculated on a set of comments containing toxic comments that don’t contain the identity term in question (background positive), and non-toxic comments containing the identity term (subgroup negative). As seen here, the LGBTQ annotated model improves the BPSN bias performance for “queer”, “homosexual”, and “muslim”, and slightly for the term “gay” compared to the control model.

	bpsn_score_subgroup_auc	bpsn_score_bpos_auc	bpsn_score_bneg_auc
lesbian	0.99	0.94	0.99
gay	0.98	0.81	1.0
bisexual	1.0	0.99	0.99
transgender	1.0	0.96	0.99
trans	1.0	1.0	0.98
queer	0.94	0.54	1.0
lgbt	1.0	1.0	0.98
lgbtq	1.0	1.0	0.96
homosexual	0.97	0.54	1.0
straight	1.0	1.0	0.95
heterosexual	1.0	1.0	0.98
male	1.0	1.0	0.97
female	1.0	1.0	0.98
nonbinary	1.0	1.0	0.92
african	1.0	0.99	0.99
african american	1.0	0.99	0.99
black	1.0	0.98	0.99
white	1.0	0.98	0.99
european	1.0	1.0	0.98
hispanic	1.0	1.0	0.98
latino	1.0	1.0	0.98
latina	1.0	0.99	0.99
latinx	1.0	1.0	0.98
mexican	1.0	0.98	0.99
canadian	1.0	1.0	0.96
american	1.0	0.99	0.98
asian	1.0	0.99	0.98
indian	1.0	1.0	0.98
middle eastern	1.0	1.0	0.98
chinese	1.0	1.0	0.98
japanese	1.0	1.0	0.97
christian	1.0	1.0	0.98
muslim	1.0	0.94	1.0
jewish	1.0	0.99	0.98
buddhist	1.0	1.0	0.97
catholic	1.0	1.0	0.98
protestant	1.0	1.0	0.96
sikh	1.0	1.0	0.96
taoist	1.0	1.0	0.96
old	1.0	1.0	0.94
older	1.0	1.0	0.92
young	1.0	1.0	0.94
younger	1.0	1.0	0.92
teenage	1.0	1.0	0.95
millennial	1.0	1.0	0.98
middle aged	1.0	1.0	0.96
elderly	1.0	1.0	0.94
blind	1.0	0.99	0.97
deaf	1.0	1.0	0.97
paralyzed	1.0	1.0	0.92

Fig. 6. Bias metrics for the African American data model. The BPSN metric—background-positive, subgroup-negative—represents the ROC AUC calculated on a set of comments containing toxic comments that don’t contain the identity term in question (background positive), and non-toxic comments containing the identity term (subgroup negative). As seen here, the African American annotated model improves the BPSN bias performance for the term “muslim” but sees a drop in performance for the terms “gay” and “queer” compared to the control model.

Received July 2021; revised November 2021; accepted February 2022