

Linguistically Differentiating Acts and Recalls of Racial Microaggressions on Social Media

UMA SUSHMITHA GUNTURI, Department of Computer Science, Virginia Tech, USA

ANISHA KUMAR, Department of Computer Science, Virginia Tech, USA

XIAOHAN DING, Department of Computer Science, Virginia Tech, USA

EUGENIA H. RHO, Department of Computer Science, Virginia Tech, USA

In this work, we examine the linguistic signature of online racial microaggressions (acts) and how it differs from that of personal narratives recalling experiences of such aggressions (recalls) by Black social media users. We manually curate and annotate a corpus of acts and recalls from *in-the-wild* social media discussions, and verify labels with Black workshop participants. We leverage Natural Language Processing (NLP) and qualitative analysis on this data to classify (RQ1), interpret (RQ2), and characterize (RQ3) the language underlying acts and recalls of racial microaggressions in the context of racism in the U.S. Our findings show that neural language models (LMs) can classify acts and recalls with high accuracy (RQ1) with contextual words revealing themes that associate Blacks with objects that reify negative stereotypes (RQ2). Furthermore, overlapping linguistic signatures between acts and recalls serve functionally different purposes (RQ3), providing broader implications to the current challenges in content moderation systems on social media.

CCS Concepts: • **Human Centered Computing** → **Empirical studies in collaborative and social computing**; *computer supported cooperative work*; • **Social and professional topics** → Race and Ethnicity.

Additional Key Words and Phrases: natural language processing, NLP, race, racism, microaggressions, discourse, social media

ACM Reference Format:

Uma Sushmitha Gunturi, Anisha Kumar, Xiaohan Ding, and Eugenia H. Rho. 2024. Linguistically Differentiating Acts and Recalls of Racial Microaggressions on Social Media. *Proc. ACM Hum.-Comput. Interact.* 8, CSCW1, Article 89 (April 2024), 36 pages. <https://doi.org/10.1145/3637366>

Offensive Content Warning: This paper contains offensive language and content that may cause distress for readers.

1 INTRODUCTION

Experiences of interpersonal racism - whether implicit or explicit - are still a regular part of life for most Black individuals living in the United States. While overt forms of racism against Blacks may have subsided compared to the decades past [121], race scholars argue that racism in modern society has not gone away; rather, it has morphed into more implicit and covert forms of expressions and subconscious acts that manifest in everyday life [138]. This *modern racism*, which is also referred to as symbolic racism or racial microaggressions, is often “highly disguised, invisible, and takes on subtle forms that lie outside the level of conscious awareness” [138, 156]. In this work, we examine

Authors’ addresses: [Uma Sushmitha Gunturi](mailto:umasushmitha@vt.edu), Department of Computer Science, Virginia Tech, USA, umasushmitha@vt.edu; [Anisha Kumar](mailto:anishak@vt.edu), Department of Computer Science, Virginia Tech, USA, anishak@vt.edu; [Xiaohan Ding](mailto:xiaohan@vt.edu), Department of Computer Science, Virginia Tech, USA, xiaohan@vt.edu; [Eugenia H. Rho](mailto:eugenia@vt.edu), Department of Computer Science, Virginia Tech, USA, eugenia@vt.edu.



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

© 2024 Copyright held by the owner/author(s).

ACM 2573-0142/2024/4-ART89

<https://doi.org/10.1145/3637366>

the linguistic signature of online racial microaggressions and how it differs from that of personal narratives recalling experiences of such aggressions shared by Black social media users.

Racial microaggressions are subtle acts of racism that often leave the victims questioning the intent of the aggressor, as the line of offense is often blurred, not immediately recognizable, and masked through humor or seemingly harmless intentions [156]. While microaggressions can occur in both on- and offline settings, our work focuses on racial microaggressions in online discussion communities. In recent years, it has been reported that an increasing number of Black individuals feel neither safe nor comfortable discussing race-related issues or experiences of racism and other potentially sensitive topics on social media [43, 66]. Despite the vocal prominence of public figures and influencers who are increasingly becoming more visible at the forefront of race-related conversations on the internet [107], about 43% of regular Black users report feeling anxious when it comes to discussing race-related matters publicly online [43]. Fearing risks of harassment, hate speech, and invalidation of lived experiences [67], some even choose to self-censor their views on racism or keep personal experiences of racism strictly private, especially on social media [10, 11]. Further, it is not just backlash from others, but also content moderation and hate speech policies by social media companies that seem to stifle conversations about race and racism for Black users [66, 67, 84]. For instance, Facebook and Nextdoor have been reported to remove posts shared by Black users regarding their personal experiences with racism either through automated filtering or through moderation [4, 46, 152]. Some Black users have reported being banned from the platform altogether [56] or locked out of their accounts for several hours or even days, a punishment known as being sent to “Facebook Jail” for sharing their views or experiences on race-related topics [66].

Human-Computer Interaction (HCI) scholarship in content moderation has shown that both human moderators and moderation systems disproportionately and often erroneously chastise users who often post about issues related to their marginalized identities, leading to false positives in content moderation decisions [41, 67]. At the same time, moderation systems often fail to detect race- or gender-based microaggressions that target marginalized users [64], allowing harmful content to remain online as false negatives [26, 113]. This is understandable given the significant topical and linguistic overlap between false positives and negatives, which can make it challenging for content moderation systems and human moderators to distinguish between the two. Such challenge underlies the premise of our work. In this work, we motivate the need to examine both false positives and false negatives in tandem, specifically in the context of acts and recalls of racial microaggressions. Through this work, we examine the linguistic signature of online racial microaggressions (acts) and how it differs from that of personal narratives recalling experiences of racial aggressions (recalls) shared by Black social media users, by asking the following questions:

- **RQ1:** *How can we leverage state-of-the-art language models to differentiate acts and recalls?*
- **RQ2:** *What are the similarities and differences between acts and recalls in terms of:*
 - a) **Themes:** *what themes best characterize acts vs. recalls?*
 - b) **Contexts:** *what contextual words are most predictive of acts vs. recalls?*
- **RQ3:** *What are the similarities and differences in the linguistic signature of acts and recalls?*

We answer these questions in the context of racism in the United States from the perspective of Black individuals. Through this work, we present a manually curated corpus of acts (2000 posts) and recalls (1264 posts), which were hand-annotated and iteratively verified through a workshop with Black participants. We use techniques in natural language processing (NLP) along with in-depth qualitative analyses to examine the language underlying acts and recalls of racial microaggressions with an aim to comparatively understand the lexical patterns that differentiate the two. Our findings show that state-of-the-art neural classifiers are able to distinguish acts and recalls with relatively high accuracy (95.4%) while more traditional language models based on n-grams features do so less

efficiently (RQ1). While acts and recalls are thematically, contextually (RQ2), and linguistically (RQ3) distinct, they also share certain themes (appearance, criminality, ability, personality, and sexual exoticism) and key linguistic signatures (use of first person pronouns and out-group language) that can make it challenging for platforms and human moderators to distinguish between acts and recalls.

Our findings represent an initial step towards better understanding semantic differences between acts and recalls of racial microaggressions on social media platforms and a re-evaluation of whether and how current socio-technical systems are able to differentiate between false positives and false negatives. Together, we aim to understand how Computer-Supported Cooperative Work (CSCW) research can best support members of marginalized groups. We argue that it is crucial to differentiate between acts (false negatives) and recalls (false positives) and utilize this knowledge to build and enhance online systems. By doing so, we can create constructive and safe online spaces where users can discuss, share, and learn from dialogues on race and racism with others. This understanding is vital in enacting inclusive and supportive online environments. Our contributions are as follows:

- We provide an in-depth characterization of the themes and linguistic signatures that underlie acts and recalls of racial microaggressions on social media communities, which has not been empirically examined at-scale by prior research to the best of our knowledge. Our insights show both distinct and shared themes and language patterns across acts and recalls, highlighting key challenges faced by content moderation systems and human moderators in their effort to distinguish between false positives and false negatives.
- Unlike the publicly available, generic off-the-shelf toxicity detectors that merely provide a score without an underlying explanation, in our comparative analysis of acts and recalls, we go beyond just classifying between the two, by leveraging interpretation techniques in deep learning (DL) to explain our classification results. Specifically, we address the lack of interpretive insight, typically associated with large pretrained language models by using Integrated Gradients [141] to identify contextual words that are most predictive of acts and recalls of racial microaggressions associated with Black users. By doing so, we overcome the black-box nature of DL language models by making our classification results explainable – a practice we believe is contextually crucial when working with textual corpora such as ours that contain semantically nuanced and subjective content.
- We complement recent efforts in the NLP community to capture and surface implicit hate speech and microaggressions online [21, 49], by providing a dataset that is more specific to the context of race and racism in the U.S. Our data is hand-annotated and validated by Black participants whose labels we used as gold-standard truths to resolve any discrepancies between non-Black annotators. We make our dataset publicly available for the wider research community in hopes that it would serve as a benchmark in examining language associated with racial biases and microaggressions.
- We provide insights from our workshop discussions with Black participants to further inform and validate our findings. Given their racial identity as Black individuals and the contextual familiarity with the content of our data, we corroborate our findings with the rich insights from our participants. By doing so, we aim to strive towards the goal of directly incorporating the experiences and views of the marginalized in CSCW and HCI research.

2 RELATED WORK

2.1 Racial Microaggressions and Implicit Biases

Microaggressions are often subconscious [53, 138] or even unintentional [105, 139, 156], meaning that they are driven by an individual's implicit biases toward people who are not members of one's own in-group [92, 134]. In the context of racial microaggressions, social psychologists argue that most people generally do not deliberately exhibit or act on racial biases all the time [73, 93, 115]. Instead, most racial biases today often take implicit forms and manifest through social conditions to which people are exposed to, or through which they interact with others [115]. Nonetheless, implicit racial biases can have consequential damages. For example, researchers found that teachers were more likely to discipline Black students on their first rather than second offense, implying that instructors were quicker to see so-called "patterns" of bad behavior in Black children compared to those of other races [111, 126]. In another study, people who saw images of Black families tended to associate them with poorer and less safe neighborhoods, despite how middle-class those families appeared in the pictures [18]. Even when people generally do not consciously harbor racist views, research has shown that they tend to subconsciously link criminality and primitiveness with Blackness [22, 62, 118, 131]. In essence, people's cognition can subconsciously interact with the conditions they are exposed to when determining responses to other people, especially in the context of race [126]. Such a premise suggests that social conditions shape the nature of interactions, and this is not exclusive to offline realms.

2.2 Racial Microaggressions in Offline vs. Online Settings

While racial microaggressions that occur face-to-face may be difficult to respond to, people can still immediately call out the transgression as it presently occurs (or has just happened), especially when the offender is visibly identified and present. In fact, there are a countless number of guidelines providing recommendations on how to respond to racial microaggressions in various contexts [50, 99, 148]. However, such guidelines are specifically tailored to in-person, offline interactions, which occur in settings much different from those of online environments. According to [138], an example of an offline racial microaggression would be a White person checking their wallet or clutching their bag as a Black man approaches or passes them, insinuating a sense of fear that Black people are most likely to be criminals. Online, racial microaggressions often surface on social media as posts and comments reacting to content posted by Black users (e.g., "You're too pretty for a Black girl").

The different affordances of on- vs. offline settings in which racial microaggression occur can potentially shape how victims experience or react to such transgressions. For example, online, offenders can hide behind the anonymity afforded by throwaway accounts, which can sometimes disinhibit bad behavior [5]. While de-identified settings are necessary and crucial in circumstances where users disclose sensitive personal experiences [5] or exchange support in stigmatized contexts [121], the affordance of online anonymity can make hate speech and harassment effortless [9, 109]. Further, throughout networked publics on social media, single users can easily connect to a congregation of thousands of others, meaning the scale of interaction and exposure can be one-to-many [1, 20]. While there are advantages to such scalability [46, 85, 150], this also means that users sharing or recalling personal experiences of racism online, can potentially face an army of aggressors who can instantaneously pile on their post by flagging or downvoting content or trolling in the comment threads. This can lead to an uptick in engagement metrics that may trigger content moderation systems to automatically flag the post for further review or removal for potentially violating platform policies [41]. Hence, scaled interactions in online settings can make it difficult for users to share recalls of racial microaggressions, or even personally respond to acts

without the burden of risking oneself against mob harassment [108, 129]. Furthermore, unlike racial microaggressions that occur in person, online acts of racial microaggressions are rarely directed at individuals. Instead, the stereotyping language of acts often targets racial groups as a whole in generalized expressions [40, 41, 160], making it harder for users to directly call out the aggressor's offense based on personal grounds beyond the context of one's race.

2.3 Language as a Condition of Online Discourse

Scholars in CSCW and HCI studying online discourse have shown early on that the conditions through which social interactions take place inevitably shape the nature of such interactions [122–124]. It is therefore not just the design of platform features or algorithmic content-ranking, but also the language users regularly encounter through others that characterize the conditions of how people come across and talk about certain topics [124]. For example, when people process contentious issues online, negatively charged affective words induce more negative conclusions in ensuing discussions [122]. As such, linguistic patterns can effectuate negativity biases toward the subject of discourse among users [13, 24, 102, 136, 149]; language that characterizes online discussions on race or racism is not an exception [123, 124]. Further, while it is important to recognize that racism expressed through language can and does take extreme and overt forms, our present study focuses primarily on the more subtle and implicit expressions of racial biases in the context of online racial microaggressions targeting Black users. We do so for several reasons. First, research has shown that racial microaggressions can profoundly impact people's physical and mental health, self-esteem, and academic performance [16, 68, 82, 105]. However, the implicit nature of racial microaggressions can be camouflaged across everyday life [139, 156], such that the impact of offense and harm is often unrecognizable or perceived as insignificant [147], especially by the offenders (and bystanders), while receivers are relegated to self-doubt and distress [2, 137]. Second, while automated detection of explicit racial profanity (e.g., via customized lexicons or regular expressions) is currently possible and widely implemented across online platforms, detecting the more nuanced and inconspicuous language around racial microaggressions masked through everyday language is not yet systematically feasible at-scale [92, 136]. Finally, human content-moderation too can be predisposed to the moderator's own unconscious biases and subjective understandings of racism and race-related matters, which makes drawing the line between acts and recalls of racial microaggressions difficult. As a result, racially marginalized users repeatedly face the burden of navigating and resolving situations where they are censored, locked-out, or banned from their accounts for sharing personal views or experiences of racism or race-related issues [66]. Such experiences aggregated over time can invalidate or elicit self-doubt towards the user's lived experience as a racial minority [7, 69]. Such issues are precisely the challenges we bring attention to and aim to address through this work.

2.4 Understanding False Positives and False Negatives in Tandem in Content Moderation Decision-Making

Content moderation research in CSCW and HCI has shown that people from marginalized communities are disproportionately affected compared to other users [41, 55, 67, 88]. For example, trans and Black users are more likely to become victims of content moderation false positives, meaning that their comments are mis-classified and censored as harmful even when they do not violate platform policies [67]. Likewise, false negatives (toxic content that violate platform policies, but are left undetected) and their impact on users has been well documented in prior HCI literature [26, 112, 119]. For example, [113] shows that content moderators removed only one in 20 comments violating macro Reddit moderation norms in 2016, and one in 10 violating comments in 2020, highlighting that some categories of violation were more likely to slip through the cracks,

leaving most anti-social behaviors unmoderated. This is because false negative content, which often includes microaggressions or implicit hate speech, are rarely explicit [50, 51, 71] and tend to be subtly disguised through humor [52], insider expressions, and neologisms [19]. As a result, this makes it harder for language models, let alone even human moderators, to identify false negatives [23].

Although content moderation researchers in CSCW and HCI have examined both false positives [41, 67, 75] and negatives [113, 119], as well as the impact they have on users [55, 60], such studies have generally examined the two discretely rather than in tandem. Our research motivates the need to examine false positives and negatives in conjunction. The significant topical and linguistic overlap between recalls and acts of racial microaggressions [67, 88] makes it difficult for content moderation systems and human moderators to detect the two apart [41, 159], often leading the former to be censored as a false positive and the latter to remain on the platform as a false negative [67, 89]. Incorrect content flagging [28, 143], or the inability of moderation systems and human moderators to differentiate when someone is critiquing racism (false positive) versus being racist (false negative) [27, 77], can be a gate-keeping practice that can evolve into a form of digital gentrification [55, 94], further exacerbating disparities between platform members and content moderators [127]. Our work aims to address such challenges highlighted by prior CSCW work in content moderation by examining the intertwined relationship between false positives and false negatives, specifically in the context of race-related social media posts. In so doing, this work takes a step towards building content moderation practices that aims to distinguish between false positives (recalls of racial microaggressions) and false negatives (acts of racial microaggressions) on social media.

Furthermore, most commercial toxicity models (e.g., Perspective API¹) used in social media content moderation systems, are not designed to distinguish false positives and false negatives in the first place, but to merely assign toxicity scores. As a result, content moderation algorithms that rely on these toxicity scores have difficulty differentiating online acts of racial microaggressions (false negatives) and discussions recalling experiences of racial microaggressions (false positives), often penalizing users they are supposed to protect [61, 79]. We explicitly address this gap by training a language model to comparatively learn the linguistic and topical nuances that are intertwined between false positives and the false negatives.

2.5 Lack of Explanations in Content Removal Decisions on Social Media

One major concern highlighted by prior research examining marginalized user's experiences with online content moderation is that most moderation systems fail to explain content removal decisions [77, 159]. Users who frequently experience content removals are often given vague (e.g., violated terms of service) or no explanations at all [67, 92]. As a result, they are left with very little knowledge as to what part of their language may have caused their post to be removed in the first place [155, 158]. Such lack of transparency and context around content removals [80, 117] make users feel frustrated and helpless as they are left unsure of how to subsequently engage on the platform [155]. Prior research has argued that enhancing the explainability of moderation decisions can not only empower users [48, 140], but also moderators as well [135, 142]. Moderators who rely on semi-automated crowdsourcing [58, 70, 95] or AI-led moderation [87, 89] experience less cognitive burden, stress, and reduced symptoms of Post Traumatic Stress Disorder (PTSD) [31, 34, 132] compared to those who do not [135]. However, the semantic and topical overlap between recalls and acts of racial microaggressions can make it difficult for content moderation systems

¹<https://perspectiveapi.com/>

and human moderators to detect the two apart, leading to false positives and negatives. Decision-making around grey content areas [67, 151], such as microaggressions, are also heavily subject to the moderator's personal and often limited understanding of what is and is not microaggressions [67], which can often lead to arbitrary and inconsistent removal justifications [145].

We address this challenge by providing a computational approach that makes classification decisions around false positives (recalls) and false negatives (acts) interpretable. In our work, we not only build a model that can distinguish between false negatives (acts) and false positives (recalls), but also provide a computational approach that helps contextualize as to what contributed most to the model's classifications decisions. Such insight can not only help moderators learn and understand the nuanced differences between acts and recalls, but also inform their decision-making as well. In our work, we use a deep learning interpretation technique to demonstrate which input feature (word) of a given post contributes most to the output decision as to why that post is most likely to be a false positive or a false negative. In so doing, we provide a scaled approach in help contextualizing and explaining model decisions in classifying between acts and recalls.

2.6 Challenges in Detecting Toxicity Through Language Models

CSCW and HCI scholars have demonstrated early on, both the effectiveness and limitations of prior approaches in identifying and curbing online toxicity, such as crowdsourcing [76, 88, 91], nudging user behavior (e.g turning off comments) [3, 130], and human moderation [25, 89, 157]. However, the challenge of countering the growing volume of toxic language has yet to be resolved [100, 161]. This is perhaps due to the diversification of toxic language [15] and online hate speech that increasingly contain neologisms [103], coded expressions, or subtle and indirect phrases that mask harmful language [65].

Both microaggressions and hate speech are similar in that they can both stem from underlying biases and perceptions. However, they differ significantly in their manifestations. Compared to hatespeech, microaggressions, are often unintentional and usually emerge as stereotypes, micro-invalidations, and micro-insults [138, 139]. By contrast, hatespeech is characterized by deliberate use of explicit language, often expressed as profanity, name-calling, bullying, or accusations [6, 8]. Moreover, while microaggressions tend to be subtle and context sensitive, hatespeech is generally less nuanced, context-independent, and less subject to interpretation. As a result, online hatespeech is often regulated by clearer definitions and guidelines [114, 153], while there is a noticeable lack of comparable guidelines and definitions for managing online microaggressions. Consequently, detecting online microaggressions is difficult for both humans and machines, making content moderation decisions more challenging [146, 162]. Recent scholarship in NLP has aimed to uncover linguistic patterns in implicit hate speech by utilizing neural models [21, 49, 71, 163]. However, such studies use data that focus on a broad array of topics. Our work adds to this existing body of research by providing and analyzing a dataset that is more contextually targeted to the topic of race and racism in the U.S.

Furthermore, a growing body of NLP research in language and fairness has highlighted how language models unintentionally capture, reflect, or even amplify various social biases that manifest in the data they are trained on [17, 128]. For example, linguistic models that power YouTube's automatically generated captions tend to identify the language spoken by male and white users more efficiently than they do of female and minority users [39]. Scholars have also demonstrated racial disparities in NLP systems by showing how widely used commercial, off-the-shelf language models fail to recognize African-American English compared to other dialects [14]. Furthermore, state-of-the-art language models pretrained on large amounts of data from the internet collected at specific points in time are susceptible to learning unintended biases towards real-world entities [116]. For example, even though the phrases, "*I hate Justin Timberlake*" and "*I hate Rihanna*" both

express the same semantics based on identical constructions, language models tend to classify the former as significantly more toxic than the latter [116], exhibiting gender disparity in toxicity scoring. Such issues may arise from disagreements in annotation labels, especially when it comes to annotating texts associated with gender and race-related content – a task that is immensely difficult for human annotators to reach consensus around what is considered ground truth [45, 104]. Throughout our analyses, we strive to be aware of such aforementioned biases. Hence, we host a workshop with Black participants through which we iteratively discuss, learn and validate our labels of acts, and use the annotations from the Black participants as the gold standard to resolve any discrepancies between non-Black annotators.

Finally, another problem with neural classifiers is that the results from such models are difficult to understand due to their lack of interpretability [54, 81]. Therefore, users across online platforms that rely on toxicity classifiers powered by neural networks might question how their posts were evaluated [77, 158]. Therefore, as a step towards encouraging linguistic tools that can provide end-users an explanation as to why their post was (mis)classified as toxic, we interpret the classification result from our best performing language model that predicts acts apart from recalls of online microaggressions.

3 METHODS

3.1 Data Collection

Type	List of subreddits
Acts	askblackpeople, askScience, BlackPeopleTwitter, casualconversation, circlejerk, confessions, darkjokes, darkjokeunlocked, explainlikeimfive, Forwardsfromgrandma, gatekeeping, hiphopcirclejerk, insanepeoplefacebook, Jokes, offensivejokes, outoftheloop, pewdiepiesubmissions, relationship_advice, shitliberalsays, shitredditsays, shitaskscience, shittylifeprotips, showerthoughts, subredditdrama, Tankiejerk, TrueOffMychest, TrueUnpopularopinion, unpopularOpinion, WhitePeopleTwitter
Recalls	Blackladies, Cptsd_bipoc, datingadvice, interracialdating, Mixedrace, TwoXChromosomes
Both	askReddit, nostupidquestions, Teenagers, TooAfraidToAsk

Table 1. List of Subreddits Used to Retrieve the Posts and Comments for Acts and Recalls.

Type	List of search keywords
Acts	'African American people', 'African American men', 'African American women', 'African American individual', 'African American individuals', 'African American person', 'African American man', 'African American woman', 'African American girl', 'African American boy', 'African American ladies', 'African American guy', 'African American gal', 'African American lady', 'African American dude', 'African American kid', 'African American chick', 'African American parent', 'African American student', 'black people', 'black boy', 'black girl', 'black men', 'black women', 'black man', 'black woman', 'black person', 'black individual', 'black individuals', 'black dude', 'black guy', 'black gal', 'black lady', 'black ladies', 'black chick', 'black kid', 'black parent', 'black student'
Recalls	"Microaggressions"; "Microaggressions, I face as a Black" + man, woman, kid, girl, gal, individual, person, guy, lady, dude, parent, student; "As a black" + man, woman, kid, girl, gal, individual, person, guy, lady, dude, parent, student; "I'm a black" + man, woman, kid, girl, gal, individual, person, guy, lady, dude, parent, student; "I'm an African American" + man, woman, kid, girl, gal, individual, person, guy, lady, dude, parent, student

Table 2. List of Search Terms Used to Collect Acts and Recalls From Reddit.

We introduce a new corpus called Recalls and Acts of Racial Microaggressions (RAMA) containing 2,000 instances of acts and 1,264 recalls of racial microaggressions from posts and comments from Reddit and Tumblr. Our data consists of acts and recalls of racial microaggressions, specifically

against Black people. For Reddit, two authors first manually examined acts and recalls across posts and comments on subreddits known to contain racial microaggressions (e.g., r/showerthoughts, r/TooAfriadToAsk, r/unpopularopinion)². As subreddit profile pages display a list of other similar subreddits, we used this information and adopted a snowball approach to expand our list of subreddits to manually look for posts and comments containing acts and recalls of racial microaggressions targeting Black people. Upon examining approximately 150 subreddit pages, we finalized our list of subreddits to those shown in Table 1. Using these subreddits, three researchers then manually verified 300 acts and 200 recalls of racial microaggression posts and comments across a diverse array of topics. We then used these 500 posts and comments to identify an initial set of common keywords to be used in the API search as shown in Table 2. We iteratively expanded on the search keywords through multiple discussions to ensure their relevance and search strength. These keywords are similar to those used in prior work examining posts containing or recalling experiences of microaggressions [21]. We then used these keywords to collect posts, comments using Reddit's official API PRAW [96] and pushshift.io³. A randomized selection of posts and comments collected from the API search were then inspected by the authors and annotated and verified by workshop participants (Refer to Section 3.2).

Furthermore, we also added posts and comments from a Tumblr website⁴ that contains a collection of self-reported accounts of acts of microaggressions across various topics [21]. We scraped all posts and comments topically pertaining to racial microaggressions and manually verified for acts and recalls related to Black people. Data from Tumblr were also annotated and verified by workshop participants.

3.2 Participant Workshop: Verifying Labels

In order to verify and annotate our RAMA corpus, we hosted workshop sessions with a total of 15 participants (6 Black and 9 Non-Black, see Table 3 for demographic information). The purpose of the workshop was to obtain high-quality human labels as to whether or not a given social media post or comment contained an act of racial microaggression against a Black person/people.

We sent out workshop flyers through campus mailing lists and contacted respondents through email. Participants were then invited to a 90 minute workshop session, which was held over lunch that was provided and paid for. We conducted a total of two separate, but procedurally identical sessions - one with Black and another with non-Black participants. We used the labels from the Black participants as the gold standard to resolve any discrepancies between non-Black annotators. We began each workshop by introducing ourselves, the project background, key research motivations behind identifying acts vs. recalls of racial microaggressions, and the take-home annotation task. We aimed to minimize participant fatigue and encourage a safe discussion environment given the difficult nature of the topic and sensitive content associated with racial microaggressions. Hence, participants were informed that they could take a break or leave at any point during the workshop and that there was no time constraints for the take-home task. In addition, participants had the opportunity to receive research credits by obtaining approval from their respective course instructors. This meant that their active participation in the workshop could be recognized and counted towards fulfilling their research credit requirements.

Once participants introduced themselves to one another, we used a guided PowerPoint presentation to explain and walk through multiple examples of comments and posts that contained acts

²We identified that a large share of the acts in our data came from comments reacting to recalls on r/BlackPeopleTwitter.

³<https://pushshift.io/>

⁴<https://www.microaggressions.com/>

and recalls of racial microaggressions across various themes (Table 9), and invited participants to reflect on these examples and to share their perspectives.

Subsequently, we collectively annotated 20 carefully chosen examples, representing distinct types of microaggressions as defined by Sue et al. (2009), through in-depth group discussions. During the discussions, participants shared why they did or did not choose to classify a post/comment as an act of racial microaggression. Participants explained what part of the comment or the post/comment specifically contributed to their decision. At the end of the workshop, we invited participants to annotate approximately 300 random samples as a take-home task. There was no requirement to further participate nor a deadline imposed on the take-home tasks. All workshop material (presentation guidelines and annotation examples) were shared with the participants at the conclusion of the workshop.

All workshop sessions were audio-recorded for transcription with participants' consent. Once we received the completed take-home tasks from participants, we checked for inter-annotator agreement. Annotation agreement across all acts was substantial considering the difficulty of understanding the subtle nature of acts of microaggressions ($k=0.77$) [90]. Each of our workshop participants annotated approximately 300 posts, such that every post in our dataset received at least 3 annotations to ensure reliability and robustness in annotation verification. Labels from Black participants were used as ground truth values to resolve discrepancies among annotations from non-Black participants. [Refer to Section 3.4]

ID	Racial-Ethnic group	Age	Gender	Ongoing/ Highest Degree
P1	Black/ African American	19	M	BA
P2	Black/ African American	20	M	BA
P3	Black/ African American	20	M	BA
P4	Black/ Caribbean American	20	F	BA
P5	Black/ Caribbean American	21	F	MS
P6	Black/ Caribbean American	25	M	PhD
P7	South East Asian	24	M	PhD
P8	South Asian	24	F	MS
P9	Middle Eastern	26	F	PhD
P10	Middle Eastern	25	F	PhD
P11	South Asian	25	F	MS
P12	South Asian	25	M	MS
P13	Asian American	26	F	MS
P14	South Asian	21	F	MS
P15	South Asian	25	F	BA

Table 3. Participant's Demographic Data.

3.3 Analysis

RQ1: *How can we leverage state-of-the-art language models to differentiate acts and recalls?*

To answer RQ1, we initially tested traditional ML-based classifiers known to work well with small amounts of data. To further understand the challenges of implicit hate detection and achieve compositional understanding beyond simple keyword-matching, we fine-tuned large language models (LLMs) such as BERT, RoBERTa and XLNet to distinguish acts and recalls of racial microaggressions. First, using a 80-20 split on the RAMA corpus, we balanced the distribution of the target class in our data in both train (80%, $N=1680$) and test (20%, $N=420$) sets. For the traditional ML models, we used the Naive Bayes (NB), Support Vector Machine (SVM) and Logistic Regression (LR) with standard unigrams, Term Frequency-Inverse Document Frequency (TF-IDF), and GloVe embedding (Pennington et al., 2014) features. We used scikit-learn's *feature_extraction* attribute to extract

features and *CountVectorizer* module with a (default) minimum word frequency = 2, (2, 2) n-gram range. We used k-fold cross validation with $k = 5$ to avoid overfitting of the models. For neural models, we fine-tuned BERT, RoBERTa and XLNet and set the batch size to 16 with 8 training epochs and used AdamW for optimization with $2e-05$ learning rate. All baseline ML models were implemented using *sklearn*⁵ and LLMs using *PyTorch*⁶.

RQ2: *What are the similarities and differences between acts and recalls in terms of:*

- a) **Contexts:** *what contextual words are most predictive of acts vs. recalls?*
- b) **Themes:** *what themes best characterize acts vs. recalls?*

Identification of Influential Tokens using Integrated Gradients (IG). While state-of-the-art neural models are effective at high-level hate speech classification, they are not effective at spelling out more fine-grained categories with detailed explanations of the implied message [159]. First, to address this, we employ Integrated Gradients (IG) [141], a model interpretability technique in deep learning (DL) that helps identify key input features that contribute most to the model's predictive decision [98] and are calculated by computing the gradient of the model's prediction output to its input features. IG (1) provides intuitive explanation for output decisions from transformer based models like BERT that often lack interpretive insight [29, 38]. We used the best performing classification model in RQ1 along with IG to derive words predictive of the model's decision (RQ2a). In IG, a feature's contribution to the output of a neural classifier is calculated by considering the gradient of the model prediction with respect to that of the input feature. Integrated Gradients calculate the average of gradients at all points along a straight line from the baseline x' , which is set to zero vector for text-based models to input x [29]. Formally, if: $R_n \rightarrow [0, 1]$ represents BERT, then the integrated gradient of the i -th dimension is:

$$\text{Integrated Gradients } (x; F) = (x - x'_i) x \int_{\Theta=0}^1 \partial F(x' + \Theta(x - x')) \div \partial x' d\Theta \quad (1)$$

The final attribution score of a particular word is the sum of the integrated gradients for each dimension of that word's embedded vector [42].

Uncovering Themes and Contexts using Qualitative Analysis. Once we identified all the tokens that were most predictive of acts vs. recalls, we manually examined every occurrence of each IG token and its contextual use across the entire dataset of acts (2,000) and recalls (1,264) – this resulted in a total of 490 acts and 626 recalls that contained at least one of the top 40 tokens that were most predictive of their respective class. We then manually examined every single instance of these posts (490 acts, 626 recalls) by examining the themes in which these tokens most frequently occurred in context. This process enabled us to identify key themes and contexts within the posts, offering thematic insights into the nature of acts and recalls of racial microaggressions.

RQ3: *What are the similarities and differences between the linguistic signature of acts and recalls?* The methodological strength and nature of Integrated Gradients (IG) is such that it identifies words that contribute most to why a post is classified as one class versus another, by focusing on factors that differentiate the two classes. Hence, as with most deep learning (DL) interpretation techniques used in explaining binary classification decisions, IG may overlook overlapping characteristics shared between the two classes. Furthermore, IG as with most NLP methods, involves preprocessing of textual data (e.g., lemmatization) which may bypass detecting verb tenses or other linguistic aspects that might be characteristic of acts and recalls. Hence, to ensure a more robust and comprehensive

⁵<https://scikit-learn.org/>

⁶<https://pytorch.org/>

examination of acts and recalls in RQ3, we overcome this limitation by analyzing a randomly selected set of posts from our entire data— not just those that included IG tokens, but those that did not as well.

To analyze the rhetorical styles of acts and recalls, we followed an iterative open coding procedure. We first approached the transcribed text through axial [32], thematic coding [59] and then performed discourse analysis [57]. Two of the authors independently coded a test sample of 100 randomly selected acts and recalls drawn from the larger racial microaggressions dataset, then discussed each post together with assigned codes to establish a shared set of categories for acts and recalls respectively. The authors then coded another 100 randomly selected recalls and similarly discussed them one by one in detail. Applying the resulting codebook, the authors coded a final set of randomly selected 100 acts and recalls.

To analyze linguistic patterns across acts and recalls, we then proceeded to perform discourse analysis. We used discourse analysis given its methodological focus on understanding how people contextually use language to convey underlying intent [57]. Discourse analysis has also been widely used across CSCW and HCI scholarship examining social media comments [47, 122, 124]. Using this method, we identified differences and similarities in the linguistic patterns between acts and recalls with regards to how posters conveyed underlying intent. Further, we validated our qualitative findings with the help of workshop discussions with the Black participants.

3.4 Understanding Participant Rationale for Labeling Posts as Acts of Microaggressions

In this section, we highlight key instances where Black and non-Black participants differed in labeling posts as racial microaggressions, aiming to elucidate their respective reasoning behind such annotation decisions. Black participants, informed by their personal encounters with racial bias, exhibited heightened sensitivity to language nuances in the posts, often missed by non-Black participants. Black participants adeptly identified posts perpetuating harmful stereotypes, discerned underlying intent, and identified subtle racial microaggressions, elements often missed by non-Black participants. Non-Black participants' labeling errors predominantly originated from a lack of awareness regarding these stereotypes and a mistaken sense of color neutrality, leading them to overlook microaggressions subtly embedded in seemingly innocuous questions or statistical assertions.

Implicit Bias. For example, all Black participants labeled the post, "*Why are all Black people so good at sports? They must be naturally athletic*" as an act. They observed that the language used in such posts subtly reinforced implicit biases about Black people. By contrast, two non-Black participants did not label this post as an act due to the seemingly positive sentiment attached to phrases like "*naturally athletic*". According to Black participants, while such expressions might initially present positive stereotypes, they carry harmful implications. They argued that such statements can diminish the individuality and diversity of Black people, attributing their achievements in sports solely to inherent physical abilities rather than hard work, dedication, and skill. Black participants labeled such posts under the racial microaggression themes, such as "*Assumptions about intelligence, competence, or status*" and "*Pathologizing minority culture or appearance*" as shown in A.1.

Lack of Familiarity with Stereotypes. Black participants displayed more profound knowledge of stereotypes than their non-Black counterparts, as exemplified in several annotation discrepancies. Notably, non-Black participants overlooked the microaggressions embedded in acts like "*Do black people find it hard to swim?*" and "*Why do Black people like watermelon?*" which subtly perpetuate stereotypes about Black individuals that were immediately recognizable to our Black participants. This disparity in awareness underlines the essential role of cultural insight and firsthand lived experiences in identifying and understanding microaggressions associated with one's social group.

Failure to recognize perceived intent or broader implications. On certain occasions, non-Black participants fell short in understanding the underlying intent of posts or the broader implications they carried, aspects that Black participants consistently discerned. For example, with the post, *"Do Black people even tan?"*, non-Black participants saw it as an innocent question, likely stemming from curiosity or a lack of knowledge. Conversely, Black participants questioned the poster's underlying intent, recognizing that the question seemed to subtly propagate belittling stereotypes through curiosity and humor.

This pattern was also evident when participants encountered the post - *"90% of violent outbreaks were caused by Black people"*. Three non-Black participants interpreted this as a mere statistical observation, failing to consider the broader implications. In contrast, all Black participants pointed out the potential harm of such assertions, underscoring that statistics can be misleading and that the intent behind their use is critically important. As one Black participant noted, *"Definitely heard that one before, it is trying to use statistics to conclude that black people are dangerous; however, there is a lot more to that statistic because the original intent of the police was to monitor minority groups. If the intent is to monitor minority groups, it is almost obvious that this statistic is going to be there"* -P2. Another, Black participant stated, *"The context matters - who is saying it..."* -P4, arguing that the identity of the speaker—their race, relationship to the listener, etc.—and the context in which the statement is made—whether in jest or during a serious conversation—can significantly sway its interpretation as a microaggression.

False Colorblindness. In some instances, non-Black participants failed to identify acts due to the absence of explicit racial references. In contrast, Black participants, acknowledging the role of race as an inseparable component of their experiences, leveraged this understanding to detect even subtly embedded microaggressions in posts. This difference was particularly evident in labeling posts like *"I don't care if they are black, gay, green, alien, or inanimate objects."* While several non-Black participants perceived this as an attempt to equate all individuals irrespective of their race, Black participants highlighted the problem of such seemingly neutral statements. For example, one Black participant stated, *"They are trying to see everyone the same, but race plays a factor in our lives - cut and dry. Just saying you don't see color doesn't change the fact that there is systemic racism. It is more important to acknowledge the oppression that people face."* -P5.

Recognizing Hedging Language as a Red Flag. The use of hedging expressions like *"Not to sound racist"* or *"I don't mean to be racist"* was adeptly identified by Black participants as a potential marker for microaggressions. Such phrases, rather than mitigating the impact of a poster's statements, served as warnings to Black participants, hinting at microaggressions. One illustrative example is the post: *"Why do some Black people have large yellow eyes? Ok, I know this probably sounds really racist, but I don't mean to be one."* Black participants classified this post as an act, recognizing that the hedging phrases did not negate or neutralize the racially charged comments that followed. This understanding allowed Black participants to navigate microaggressions more adeptly than non-Black participants. As one participant stated, *"Just because you say 'I don't mean to be such and such' does not automatically absolve you of what you said before or after"* -P6. This perception, however, was often missed by non-Black participants, emphasizing the importance of sensitivity and awareness in detecting and understanding microaggressions.

4 FINDINGS

In this section, we first provide our findings related to classifier performance for differentiating acts and recalls of racial microaggressions (RQ1). Further, our findings revealed higher level differences between acts and recalls across various themes identified via Integrated Gradients (RQ2), as well

Binary Classification Result				
Model	Precision	Recall	F-1	Accuracy
SVM (n-grams)	0.752	0.747	0.749	0.716
SVM (TF-IDF)	0.639	0.653	0.645	0.618
SVM (GloVe)	0.709	0.735	0.721	0.708
Naive Bayes (n-grams)	0.717	0.721	0.718	0.722
Naive Bayes (TF-IDF)	0.622	0.608	0.614	0.621
Naive Bayes (GloVe)	0.667	0.525	0.587	0.545
Log. Regression (n-grams)	0.633	0.692	0.661	0.693
Log. Regression (TF-IDF)	0.769	0.724	0.746	0.719
Log. Regression (GloVe)	0.719	0.72	0.719	0.733
BERT	0.864	0.925	0.893	0.906
XLNET	0.891	0.968	0.927	0.917
RoBERTa	0.921	0.950	0.934	0.954

Table 4. Classification Metrics for Acts vs. Recalls (Best Performance Is Bolded)

as more granular differences in their linguistic signatures via discourse analysis and focus group interactions (RQ3).

4.1 RQ1: Classifying Acts vs. Recalls

To answer RQ1, we first built natural language based classification models to classify acts vs. recalls of racial microaggressions. We experimented with several traditional models such as Support Vector Machine (SVM), Naive Bayes, and Logistic Regression as baseline classifiers along with three different feature extraction methods (GloVe, TF-IDF, n-grams), as well as state-of-the-art neural models such as XLNet, BERT, and RoBERTa, to classify a user post/comment into an act or recall. We observed that RoBERTa and XLNet achieved a high level of performance, with accuracies at 0.954 and 0.917, respectively; in contrast, SVM and BERT had far lower accuracy scores at 0.760 and 0.906, respectively. The modeling results are presented in Table 4.

4.2 RQ2: Interpreting Acts vs. Recalls

To explain the predictive decisions of our best-performing neural language classifier, we leverage Integrated Gradients (IG), a model interpretability technique, to extract key words that the model considered most predictive of acts vs. recalls. We use IG to identify specific tokens that contribute most to the model's predictive decision by computing attribution scores for each of the tokens. Attribution scores indicate how much a specific token correlates to the model's prediction judgment. We categorized the top 40 tokens with the highest and lowest attribution scores ranging from positive (predictive of recalls) to negative (predictive of acts) values by each class (act/recall), with the magnitude indicating the predictive strength for each class. We then manually examined every occurrence of each IG token and its contextual use across the entire dataset of acts and recalls and identified prevalent themes in which the top 40 tokens most frequently occurred in context. Table 5 details the themes that emerged in the acts and recalls respectively. In total, we identified three themes unique to acts ("*Questions*", "*Ethnicity*", "*Evolution*" and "*Human Race*"), four themes associated with recalls ("*Relationships*", "*Workplace*", "*Everyday Life*", "*Geographical Location*"), and five overlapping themes ("*Appearance*", "*Criminality*", "*Ability*", "*Personality*", "*Sexual Exoticism*"). Tokens are listed in the descending order of their predictive strength and grouped by salient themes based on the most common contexts in which these tokens appeared across sentences.

4.2.1 Themes of Acts.

Questions. Among posters of acts of racial microaggression, the tokens, *why* and *question*, are often used to pose a question along the lines of why a certain stereotype exists. Consider this example of an act: “*Why do black men have such dry hands?...I have just realized that while making this question that I’ve never [shook] the hand of a black woman.*” The poster asks a question (“*Why do black men have such dry hands?*”) and then proceeds to acknowledge that they posed a question. By questioning the validity of their own question while writing the post, the poster appears to be thinking out loud, thereby not filtering their thoughts.

Ethnicity. The tokens *ethnic*, *indian*, *hispanic*, and *latino* are often used by posters of acts to compare Black people to other races, especially racial and ethnic minorities, such as Indian and Hispanic people. For example, the poster of the following act creates a hierarchy of intelligence based on race in order to highlight their belief that Black people are less intelligent than other people, using different racial majorities (‘*caucasians*’) as well as racial and ethnic minorities (‘*asians,*’ ‘*indians,*’ ‘*latinos/hispanics*’) to elucidate their point:

“*Some say that you can order races on their intelligence with asians on top, indians and caucasians a little below, latinos/hispanics about 1/4 std. deviation below whites, and blacks about 2/3-1 standard deviation below whites.*”

Evolution and Human Race. Tokens such as *earth* and *evolution* are frequently used by posters of acts to discuss where Black people came from as well as what led them to develop their distinct appearance and abilities. In addition to using the token *earth* to imply a point of origin, the poster also uses the word to casually place Black people outside of the human race: “*If Adam and Eve are the first people in the Earth and they are white, why are there Black people?*” Moreover, the following act uses the word *evolution* as a way to justify why Black people are faster than other races by connecting a common stereotype to evolution: “*Black people are faster because of evolution.*”

4.2.2 Themes of Recalls.

Relationships. Posters of recalls commonly use the word *friend* to describe the perpetrator of an act of racial microaggression or the person that experiences a racial microaggression alongside them. The poster of this recall uses the word *friend* to describe the person that they were with when they both experienced a racial microaggression: “*I’m Dominican-American and one day me and my friend who’s Bengali went to the mall. We walked into a MAC store and a White lady approached us and asked us: Where are you guys from? You guys look exotic.*” Our findings suggest that this *friend* is often times Black or another racial minority, such as Bengali. On the other hand, the word *friend* is also often used to describe the perpetrator of an act of racial microaggression: “*I had a white male friend of mine in high school tell me that I’m the darkest he would go.*” Our findings also suggest that the token, *husband*, is often used to describe the person that the victim of the racial microaggression is compared to by the perpetrator of the act: “*After expressing his shock upon seeing my husband’s last name (whose family is German) he said, “Oh, so you’ve got a good Jewish boy, huh?” You must feel lucky!*”

Workplace. Recalls often contain tokens such as *job* and *interviews* to describe acts of racial microaggression that occur in the workplace: “*At job interviews, I tell them where I’m from, born and raised in the Dominican Republic, and they say “Oooh” with a tone of disappointment.*” Here, the words *job* and *interviews* serve to specify where in particular the poster has experienced a racial microaggression.

Attributed to Acts				Attributed to Recalls			
Themes	Token	Attribution Score	Freq.	Themes	Token	Attribution Score	Freq.
Questions (247)	question	-0.077	29	Relationships (176)	friend	0.1820	76
	why	-0.169	218		husband	0.1597	12
Ethnicity (23)	ethnic	-0.0147	9		mother	0.1419	13
	indian	-0.020	11		partner	0.1004	7
	hispanic	-0.029	6		father	0.0987	14
	latino	-0.044	3		family	0.0848	37
Evolution and Human Race (22)	earth	-0.009	2		boyfriend	0.0697	17
	existence	-0.017	5		job	0.1751	5
	evolution	-0.017	2		manager	0.1665	31
	population	-0.017	6		company	0.1596	6
	civilization	-0.038	3		office	0.1508	75
	humans	-0.066	2		teacher	0.1221	13
	primates	-0.015	2		interview	0.1122	14
Appearance (42)	fat	-0.004	2	employment	0.104	54	
	monkeys	-0.006	2	internship	0.1025	19	
	palms	-0.007	6	cafe	0.1577	1	
	gorilla	-0.009	4	salon	0.1464	1	
Criminality (52)	skinned	-0.024	28	shopping	0.1340	3	
	police	-0.003	6	church	0.1195	4	
	robbed	-0.014	2	grocery	0.0941	4	
	dangerous	-0.014	11	gym	0.0776	2	
	commit	-0.016	14	california	0.2480	1	
	violent	-0.023	11	london	0.153584	1	
Ability (17)	streets	-0.056	5	midwest	0.2480	4	
	attacking	-0.058	3	europa	0.1228	3	
	sports	-0.003	7	chicago	0.1185	2	
	intelligent	-0.020	3	hair	0.0843	143	
	IQ	-0.048	7	straightened	0.0113	6	
Personality (54)	ignorance	-1.300	5	ape	0.0003	2	
	scary	-0.001	3	criminals	0.0078	3	
	names	-0.004	15	neighborhood	0.0017	9	
	dumb	-0.004	5	basketball	0.0040	4	
	disgusting	-0.004	2	athletic	0.0283	4	
	funny	-0.007	6	smelled	0.0115	2	
	ghetto	-0.009	7	lazy	0.0110	3	
	loud	-0.018	11	stronger	0.0034	4	
Sexual Exoticism (33)	sex	-0.017	10	ignorant	0.0029	6	
	attracted	-0.037	12	stupid	0.0004	3	
	hotter	-0.005	7	exotic	0.0466	6	
	sexual	-0.065	4	Sexual Exoticism (18)	attractive	0.0322	12

Table 5. Themes Emerged From Attributive Words of Acts and Recalls With Scores in Descending Order. In addition to the attribution scores, we provide the number of times each token occurred in the corpus (Freq.)

Everyday Life. Recalls often contain tokens relating to everyday activities or places that people frequent on a day to day basis such as a *café* in order to convey what they were doing when experiencing an act of racial microaggression: “*On Friday morning, as I walked to the café between classes at my predominantly white university, the school appointed photographer offered me a free coffee if I agreed to play the role of the cheerful token black woman in a group of strangers.*”

Geographical Location. Posters of recalls often use tokens relating to geographical location (countries, cities, regions, etc.), such as *California*, to communicate where they experienced a racial microaggression: “*I walk into a gas station market in California with about ten of my Latina/o and black high school students to buy snacks for our college road trip, and within five minutes, we hear, “SECURITY CHECK ON ALL AISLES.”*”

4.2.3 Overlapping Themes Across Acts and Recalls.

Appearance. Both acts and recalls contain words such as *gorilla* and *ape*, respectively. While posters of acts commonly use the word *gorilla* in order to imply that Black people do not deserve to be treated as human beings: “*Black people appear more closely related to gorillas than human,*” posters of recalls commonly use the word *ape* in order to describe instances in which they have been compared to such an animal: “*He was also the worst of the people making these jokes in high-school*

and shortly after, making hundreds of “all black people are criminals” jokes and “comparing me to an ape” one time.” As it relates to the topic of appearance, given that ape is an umbrella term that includes several species, one of which is a gorilla, it follows that acts within this theme are often more specific than recalls. Similarly, while acts utilize a variety of words such as *fat*, *palms*, and *skinned* with the purpose of negatively stereotyping Black peoples’ appearances: “Black people are fat because McDonalds is all they can afford”, recalls frequently use words relating to a single feature, hair (*hair*, *straightened*) to describe their experience of being the victim of microaggressions about the texture of their hair: “This is why I don’t ever straighten my hair anymore, even though it’s something I used to like to do for fun on occasion, because the compliments always seem to insinuate that my normal hair is unprofessional, unruly, or otherwise socially unacceptable.”

Criminality. Posters of acts often use anecdotal evidence and a variety of words related to criminality, such as *police*, *robbed*, and *dangerous* in order to make stereotypical statements about Black people: “Today I almost I got my car robbed from me. I’ve gotten robbed twice by a black person and this is the 3rd time but this time I was able to get away.” In contrast, posters of recalls frequently use a single word, *criminals*, in order to describe being associated with criminals: “He was also the worst of the people making these jokes in high-school and shortly after, making hundreds of “all black people are criminals” jokes and “comparing me to an ape” one time.”

Ability. While acts commonly contain tokens such as *sports* and *IQ*, recalls frequently contain tokens related to sports, such as *basketball* and *athletic*. For example, in this act the poster uses the word *sports* to highlight that Black people are only good at rap and sports: “Black people of Reddit, How does it feel to be inferior and only good at rap and sports?”. In addition, the acronym *IQ* is used in order to demonstrate that Black people are naturally less intelligent than other races, as *IQ* is seen as more of an inherent intelligence as opposed to intelligence derived from hard work: “Black people have a way lower IQ across the board compared to their white counterparts.” In contrast, recalls focus solely on sports: “The normal stereotypical things like us liking fried chicken and watermelon, every black person knows how to dance (my brothers are proof that’s a lie) we’re all good at basketball (I’m proof that’s a lie).” Finally, similar to criminality and appearance, recalls within this theme contain a fewer variety of words, indicating that the content of acts is more dispersed as compared to recalls.

Personality. Posters of both acts and recalls commonly make use of tokens related to negative personality characteristics such as *dumb* and *stupid*. While acts frequently use the word *dumb* to criticize Black people: “They never know what I’m talking about because they are dumb,” recalls frequently make use of the word *stupid* to describe how they are perceived by others: “So that’s why I’m not only likely to be a thief, I’m likely to be a stupid thief!”

Sexual Exoticism. Posters of acts frequently use the tokens, *attracted* and *sex*, to highlight their perception of Black people as highly sexually desirable just because of their race: “I’m extremely sexually attracted to Black men and women.” In addition, this act uses the word *sex* to underscore that being Black is a condition that must be satisfied for them to agree to have sex with someone: “I will only have sex with Black men.” In contrast, posters of recalls frequently use the word *attractive* to describe that other races typically find their race unattractive: “Like when people tell you aren’t black because you’re attractive.”

While Integrated Gradient (IG) offers valuable insights into the contexts and themes of racial microaggressions, as with most deep learning interpretation techniques used to explain binary classification decisions [97, 125], IG too has limitations in capturing complex linguistic patterns. For instance, it may struggle to capture the *overlapping* characteristics between acts and recalls, which are essential in understanding the nuances in posts that may result in FPs and FNs. Additionally,

IG relies on preprocessing steps like lemmatization and stop word removal, which may not fully capture linguistic cues such as verb tense or the use of absolute terminology in acts and recalls. To overcome these limitations, we adopted a more comprehensive approach in RQ3 by analyzing a randomly selected set of posts from our entire data, including those that did not contain the top 40 IG tokens. This allowed us to be more exhaustive in our examination of linguistic patterns that characterize acts and recalls of microaggressions.

4.3 RQ3: Characterizing Acts vs. Recalls

To further understand the nature of acts and recalls of racial microaggressions in Reddit posts and comments, we examined the linguistic attributes manifest in their content using discourse analysis to better understand the social purpose underlying the linguistic patterns observed in acts and recalls of racial microaggressions. Tables 6, 7, 8 represent the linguistic pattern that emerged from the analysis and examples for each linguistic pattern unique to acts and recalls respectively. Our findings suggest three linguistic patterns in both acts and recalls of racial microaggression from our data. We utilize linguistic analysis as well as data gathered from our workshop participants to better understand the functional purposes of the similarities and differences we observed in the linguistic patterns of acts and recalls.

4.3.1 Linguistic Signature of Acts.

Questions. Our findings revealed questions to be a key linguistic pattern in acts of racial microaggression. Consider this example of an act: *“Why are black people so athletic? Not to sound racist, but I have recently noticed that black students excel at all the sports in my school. I am White and most people in my school are white, However the black players are the best in every single sport for our school. Why is that? I don’t mean to be racist or offend anyone. Just wondering.”* This statement serves to create a broad generalization of all Black people as being athletic, falling under several themes of racial microaggressions such as racial categorization and sameness, assumptions about intelligence, competence, or status, and connecting via stereotypes [156]. Consider the statement below from one of our workshop participants, P4:

“It is still a microaggression in the form of a question. You are still generalizing black people and you kind of believe that statement which is why you are curious about it” —P4

This statement is a racial microaggression masked in the form of a question, which can give the statement more of a tone of curiosity than aggressiveness [Refer to Table 8]. Nevertheless, just like P4 highlighted, this curiosity is still a form of generalization. The poster of this statement seems to believe that all Black people are athletic, which is why they are curious about why that is the case. Consider the statement below from another one of our workshop participants, P5:

“There are only a few black players that are the best but when they are asking this question, they clearly generalize. If the star athletes weren’t black, you wouldn’t be wondering, “why are these people so athletic?” It would just be obvious that they probably practice a lot” —P5

Despite the subtlety of this statement, P4 further highlights how this statement is a generalization, emphasizing that if the star players were white, the posters would likely attribute their talent to practice and hard work instead of race. On the other hand, our findings also suggest that posters commonly disguise microaggressions using a combination of curiosity and humor. Consider this example of an act: *“Black people of Reddit, which one of you stole my bike?”* This statement uses a common stereotype of Black people, criminality, to make a joke in the form of a question. By directly addressing members of the Black community on Reddit (*“Black people of Reddit”*) and

Linguistic Pattern	Examples of acts of Racial Microaggression
Questions	<ul style="list-style-type: none"> • <i>Why are black people so bad at swimming?</i> • <i>Why do so many black people litter? I don't think a day goes by that I don't see someone throw trash out of their car, 9 times out of 10 that person is black. Why?</i> • <i>Black people of Reddit, which one of you stole my bike?</i> • <i>Black people of Reddit, can I touch your hair?</i>
Use of absolute terminology ('all', 'never', 'ever', 'should', 'absolutely', 'only')	<ul style="list-style-type: none"> • <i>[Black people] are all mentally handicapped and physically incapable of supporting themselves.</i> • <i>I would never, ever hire someone with a 'black' name on their resume. I wouldn't even interview them.</i>
Use of statistics	<ul style="list-style-type: none"> • <i>Black people are not oppressed and if they want to be in prison less, they should not be committing 53% of all homicides while only being 12% of the population.</i> • <i>12% of the population is black people, yet they commit so much more crimes.</i>
Use of modifying adverbs or adjectives ('most', 'usually', 'consistently')	<ul style="list-style-type: none"> • <i>Black people usually name their kids after stuff they can't afford. Like Mercedes, Diamond, Hope, and Insurance.</i> • <i>Black people are consistently the most rude, demanding, ignorant, of what want and shady.</i> • <i>I fail black students way more often because, objectively, they make the most mistakes on driving test.</i>

Table 6. Linguistic Patterns Observed Using Discourse Analysis for Acts With Examples.

assuming one of these members stole their bike (“*which one of you*”), the poster of this act utilizes humor to soften a harsh stereotype.

Use of Absolute terminology and statistics. In addition to questions, our findings show that acts utilize absolute terminology and statistics in order to justify making a racial microaggression. Consider this example of an act: “*Black people are not oppressed and if they want to be in prison less they should not be committing 53% of all homicides while only being 12% of the population.*” This statement uses the problematic “13/50” argument [83], which is commonly used to stereotype Black crime in order to make statements appear factual as opposed to stereotypical [72]. The 13/50 argument is an overused and often misleading talking point that poses that black people make up only 13% of the population but commit 50% of all known crimes [154]. Consider the statement below from one of our workshop participants, P3:

“There is a lot more to that statistic because the original intent of the police was to monitor Black people”-P3

According to P4, such statistical reference is biased because the United States history of systemic racism has resulted in Black people often being the target of police [12, 101]. Moreover, another

Linguistic Pattern	Examples of Recalls of Racial Microaggression
'White' and 'White people'	<ul style="list-style-type: none"> • <i>I wish white people in general would stop commenting on my appearance unless it's to compliment me or to tell me that I have something stuck in my teeth.</i> • <i>White people singling me out at social events to make small talk with me about race/politics.</i>
Paste Tense Verbs	<ul style="list-style-type: none"> • <i>I was at work and the topics of racism came up with my boss who is Italian...</i> • <i>I felt irritated at having to explain that yes, I am a REAL programmer.</i>
'Only Black'	<ul style="list-style-type: none"> • <i>I had an English teacher who loved to talk, and whenever she'd say anything about race or Black culture, she's turn to me(the only black kid in the room) as if to validate/confirm the statement.</i> • <i>I was the only black girl in the room with him and ten other coworkers</i>

Table 7. Linguistic Patterns Observed Using Discourse Analysis for Recalls With Examples.

participant highlights how the use of statistics matters in determining whether a statement should be considered a racial microaggression:

“There is a difference between using a statistic to prove that there is a problem with the prison system versus using it to say something about Black people.”-P4

In this statement, P4 highlights that she believes that using a statistic to make a statement about an institution, such as the prison system, is different from using a statistic to make a statement about a certain race, such as Black people. Given that institutions such as the prison system are not human, she implies that using a statistic to make a seemingly “*factual*” statement that is negative about a particular group of humans is bound to be hurtful to people. The poster of this act uses the statistic in order to justify making a racial microaggression on the grounds that they are unbiased and are just disseminating a fact.

Use of modifying adverbs or adjectives. The last key linguistic pattern distinct to acts is the use of modifying adverbs/adjectives. Consider this example of an act: “*Black people usually name their kids after things they can’t afford. Like Mercedes, Diamond, Hope, or Insurance.*” In response to this statement, one of our workshop participants highlights why this statement is a stereotypical/generalized statement, stating that she doesn’t know anyone with those names:

“This statement is too generalizing. I don’t know a single person with those names. It is not as common as you think.”-P5

The poster of this act uses the modifying adverb, ‘*usually*’, in order to characterize this behavior of Black people as frequent. Unlike the posters of acts that use statistics to justify making a racist remark, the poster of this act uses a colloquial term, ‘*usually*’. By using the modifying adverb, ‘*usually*’, the poster of this act seeks to give the impression that their personal knowledge is sufficient to justify such a claim. Thus, similar to posters of acts that utilize statistics, this poster tries to justify making a stereotypical racist remark.

4.3.2 *Linguistic signature of Recalls.*

‘White’ and ‘White people’. One notable linguistic feature of recalls is the use of the phrases ‘White’ and ‘White people.’ Consider this example of a recall: *“White people singling me out at social events to make small talk with me about race/politics. I think they want to see me get impassioned or educate them. No I’m tired, I came out to have fun.”* In this recall, the poster is expressing his/her discomfort of being ‘singled out’ at social events by ‘white people’. Clearly stating the subject (‘white people’) that is causing his/her discomfort, the poster uses this phrase to highlight the source of their tiredness. By making the source of their tiredness very clear, the poster seeks to have his/her experiences validated.

Past Tense Verbs. Another common linguistic feature we noticed in recalls was the use of past tense verbs. Since recalls are recounts of acts of racial microaggressions, it follows that most recalls describe events that took place in the past. Consider this example of a recall: *“I was at work when the topic of racism came up with my boss who is Italian...”*. The poster of this recall is describing the setting of when he/she experienced an act of racial microaggression. The use of past tense verbs, such as ‘was,’ and ‘came’ aids in disclosing the setting in which the act took place. This is a key element in the poster’s recount of their experience. Prior work in human communications research suggests that lying individuals use fewer words and fewer past tense verb forms [44]. These verbs help provide readers with a confidence that the poster must be telling the truth because they seem to be recounting their past experience very clearly. This type of disclosure helps create an overall tone of honesty and authenticity in the statement, as the poster hopes to have their experience validated by others.

‘Only Black’. The use of the phrase ‘only black’ was the last prominent linguistic pattern distinct to recalls that our findings revealed. Consider this example of a recall: *“I was the only black girl in the room with him and ten other coworkers.”* The poster’s use of the phrase ‘only Black’ in contrast to ‘him and ten other coworkers’ underscores her discomfort at being the only Black girl in the room.

4.3.3 Linguistic Similarities between Acts and Recalls.

Use of First-Person. One notable difference between acts and recalls of racial microaggression is the role that first person voice serves in context. Consider the example: *“I am not a racist, but it seems whenever I sit near a group of black people I can’t hear the movie over all the noise they make.”* Here, the poster uses first person to preemptively defend himself/herself from being called a racist. One of our workshop participants, P3, states that this defensiveness does nothing to absolve the poster of what he/she said:

“Just because you say I don’t be mean to be such and such does not automatically resolve you of what you said before or after. Maybe they don’t mean to be racist but words are words and the implication is still going to be there regardless of intent.”-P3

While this type of linguistic pattern appears to mask a racial microaggression, P1 points out that while the poster may not have bad intentions, the statement is still an act of racial microaggression, and therefore, the hedging does nothing to reduce the severity of the statement. Unlike posters of acts that seek anonymity, posters of recalls typically use first person voice to thoroughly describe themselves- *“I am a 22 year old, brown skinned African American girl. In school in Maryland. I felt out of place and isolated.”* Here, the poster utilizes the first person voice twice to describe her age and feelings. Based on LIWC analyses [144] and manual inspection, our results indicate that roughly 40% of recalls that utilize first person do so to describe themselves. This self-disclosure serves to create an overall tone of authenticity and honesty, which the poster hopes will allow his/her experiences to be validated.

Moreover, our findings suggest that posters of acts utilize first person voice in even subtler ways to hedge the aggressiveness of their comments. Consider this example of an act from our dataset:

“Speaking as a capital “C” Conservative. I totally agree. I enjoy good, well developed characters, I don’t care if they are black, gay, green, alien, or inanimate objects. Scandal is a great show, orange is the new black was great for the first couple seasons, and Jesus is a great character in the walking dead. (That’s coming from an orthodox catholic conservative). Now to be honest, I hated black panther, just didn’t think it was a good movie.”

Workshop participants agreed that the last sentence of the post contains the microaggression; nevertheless, the poster makes several statements before in order to take attention away from it. The repetitive use of first person voice (“I”) prior to the last statement serves to highlight the poster’s desire to portray themselves as an objective critic of entertainment. Similar to our previous findings, one of our workshop participants highlights that this superfluous build up is just another method of hedging a racist remark:

“They try to use color blindness to remove themselves from what they are saying”-P6

The poster seeks to show that they see everyone as equal by saying, *“I don’t care if they are black, gay, green, alien, or inanimate objects.”* By making a somewhat extreme statement, that they don’t see color, the poster seeks to curb the aggressiveness of his/her statement. Moreover, one of our workshop participants, P3, highlights this type of hedging has consequences beyond just being a way to remove oneself from the consequences of making a racist remark:

“I definitely can take offense to this since you are not validating who people are-you are not acknowledging where they came from or what they experienced”-P2

This type of false color blindness fails to acknowledge the systemic racism that has existed in American culture for many decades [156]. Nevertheless, like P3 points out, these types of statements are often invalidating for Black people if they are proud of their identity or have suffered because of it [for more details on the definition and examples of color blindness, refer to Table 9, Appendix]. While acts utilize superfluous build up and false color blindness to hedge the severity of microaggressions, our findings suggest that posters of recalls also incorporate build up prior to the point they are trying to make.

“I will ALWAYS be one to want to expand my viewpoints, appreciate history, and under the social climates, but I’m fucking tired... I 10000% care about race and gender issues within our communities, but I’m tired of including “others” in the conversation. How can we navigate things like institutionalized racism without begging “others” to see us as human? Unfortunately, it’s never going to change. Maybe I’m just being a pessimist, but...”

By emphasizing the word *always* and using the first two sentences to prove that he/she is an open minded individual, the poster of this recall is trying to establish themselves as someone that is legitimate and trustworthy in the Black community. Therefore, the poster uses superfluous build up to gain peoples’ trust and have his/her thoughts and feelings heard.

Linguistic Pattern	Examples of acts of Racial Microaggression	Examples of recalls of Racial Microaggression
Use of First Person	<ul style="list-style-type: none"> • <i>I'm not a racist, but it seems whenever I see sit near a group of Black people, I can't hear the movie over all the noise they make.</i> • <i>I don't hate them, I don't bully them, but I'm careful with them, as if they were criminals.</i> • <i>I know I'll be downvoted, but anyone in the U.S. who works for tips knows that Black people are far less likely to tip...</i> 	<ul style="list-style-type: none"> • <i>I'm an African American graduate student, and I teach at a large university.</i> • <i>It angers me when people measure my race by the way I talk, dress and carry myself.</i>
"Us" vs. "Them" Language	<ul style="list-style-type: none"> • <i>Black redditors, what is your take on having a white friend? What do you see us as? Or any other race?</i> • <i>Maybe we shouldn't let them (Black people) vote.</i> 	<ul style="list-style-type: none"> • <i>How can we navigate things like institutionalized racism without begging "others" to see us as human?</i> • <i>They always said it was a joke but they kept doing it over and over.</i>

Table 8. Linguistic Patterns Observed for Acts and Recalls Using Discourse Analysis With Examples.

Use of "Us" vs. "Them" Language. Another notable characteristic of acts and recalls of racial microaggression is the use of Us vs. Them language. Consider the example from our discourse analysis: *"Maybe we shouldn't let them (Black people) vote"*. By utilizing both *we* and *them*, this statement is characteristic of Us vs. Them language. The juxtaposition of these words serves to highlight the presence of two different groups, *we* or *us* and *them* or *"Black people."* This creation of an in-group and out-group serves to portray Black people as second class citizens. On the other hand, recalls utilize *Us* vs. *Them* language in order or to emphasize their feelings of being discriminated against. Consider this recall from our dataset: *"They always said it was a joke but they kept doing it over and over."* The use of the word *"they"* in this statement is characteristic of *"them"* language in *Us* vs. *Them* language. The sentence highlights that the out-group is hurting the in-group's feelings by making what they considered to be a joke about the victim's race. In context, the use of *Us* vs. *Them* language helps underscore the victim's discomfort due to the "joke" by highlighting the actions of the out-group.

As we conclude our findings for RQ3, it's important to reflect on how they differ from our findings from RQ2. While both research questions aimed to strengthen our understanding of racial microaggressions in online content, they each focused on different aspects of these interactions. RQ2 centered on **'what'** was being discussed in the posts, identifying themes and contexts of racial microaggressions through the analysis of influential tokens and their contextual use. This analysis revealed themes such as *'Questions'*, *'Evolution'*, and *'Criminality'* for acts (among others), and *'Relationships'*, *'Ability'*, and *'Workplace'* for recalls (among others), offering thematic insights into the content of racial microaggressions. On the other hand, RQ3 delved into **'how'** these discussions were being framed and expressed. We manually examined a randomized selection of posts from the entire dataset, identifying linguistic patterns and discourse structures (such as the "Use of modifying adverbs or adjectives", "Past tense verbs", and "Use of first-person" etc) that characterize acts versus recalls.

To this end, our findings from RQ2 and RQ3 complement each other, with RQ2 providing thematic insights into the content of racial microaggressions, and RQ3 offering linguistic insights into their expression. Together, they contribute to a more holistic understanding of racial microaggressions in online content, each illuminating different aspects of these complex interactions.

5 DISCUSSION

5.1 Enhancing Content Moderation Classifiers with Context Awareness of FPs and FNs

Current hate speech classifiers often rely on predetermined thresholds, such as classifying a post as "offensive" if its toxicity score exceeds 0.8 (a benchmark recommended by Perspective API for determining harmful posts). As a result, these classifiers may struggle to adequately capture implicit microaggressions, leading to false negatives (FNs) [49]. This raises the question around whether toxicity thresholds alone are reliable measures for evaluating implicit forms of harmful content.

To improve the precision of these classification systems, our research findings highlight the necessity for a deeper examination into the specific linguistic patterns that trigger FPs and FNs in the first place. For example, our study reveals that FPs and FNs often exhibit overlapping, yet distinct linguistic patterns and themes, which can be overlooked by threshold-based classifiers. For instance, posters of acts often use a variety of words (e.g., 'police', 'robbed', and 'dangerous') thematically associated with *criminality* to stereotype Black people as "*criminals*" in their posts, while posters of recalls tend to use the word '*criminal*' as a singular word to recount their experience of being wrongfully perceived as a potential lawbreaker. This strong thematic overlap in words that characterize acts and recalls can make it extremely difficult for classifiers, from simple keyword-matching algorithms (which often overlook the context of the flagged words) to even the more context-aware ones (e.g., BERT), to distinguish the subtle nuances shared between FPs and FNs. Similarly, the presence of common semantic characteristics between recalls and acts can also potentially mislead classifiers. For instance, both posters of acts and recalls use the first person voice for different purposes. While the former tend to use the first person alongside hedging language to defend themselves (e.g., "*I am not a racist*"), the latter use it to depict or narrate past experiences (e.g., "*As a Black woman, I've always been ignored by my coworkers.*"). This highlights that classifiers, despite the different contexts in which the language is used, can fail to discern FPs apart from FNs.

Hence, we foresee the need for future content moderation systems to incorporate the nuanced contexts around acts and recalls into toxicity classifiers, as decision thresholds alone cannot fully capture the complex linguistic subtleties shared by recalls and acts, leading them to be misclassified as FPs and FNs, respectively. To address this issue, we suggest that threshold-based classifiers could be improved by incorporating *contextual features* that incorporate the thematic and linguistic patterns of acts and recalls observed in our study. These features can capture the semantic meaning of each linguistic pattern within its specific context. For instance, the word "*I*" would generate different features when used in varying contexts, thereby capturing the nuances of its use in acts versus recalls. By extracting such contextual features, we can improve the accuracy of toxicity classifiers and overcome the limitations of simple decision thresholds, thereby reducing the likelihood of misclassifying posts as FPs or FNs.

Furthermore, insights from our workshop demonstrate that capturing the true essence of context when distinguishing acts from recalls depends on one's familiarity with stereotypes or being able to discern the role of hedging language in an argument. This involves understanding underlying implications that extend beyond the mere text—a nuance most classifiers overlook. Encapsulating such context more holistically, as shown in this work, is crucial for fostering more accurate and equitable content moderation systems.

5.2 Increasing Awareness of Acts and Recalls Across Moderators and Users

Improving Explanations Around Moderation Decisions. Prior work on transparency in content moderation calls attention to key challenges in the moderation process [78]: moderators often fail to articulate *what* aspect of the content prompted moderation or *why* such moderation was

necessary. This begs the question: how can we inform users about why their posts were flagged or removed, and which aspects of their posts led to such moderation outcomes? What kind of strategies can moderators undertake to effectively communicate such moderation decisions to users? Our study provides insights to help address these questions in the context of improving moderation explanations. Our analysis reveals key themes and linguistic patterns that could potentially equip moderators with a deeper understanding around acts versus recalls, thereby enabling more informed decisions. For instance, posters of acts often pose questions that are seemingly driven by curiosity, or statements incorporating statistics about Black people. It is possible that individuals committing acts of microaggressions are not fully cognizant of their behavior in their discussions of certain racial groups [139]. For example, posters frequently pose questions, such as "*Why are black people so bad at swimming?*" or "*Why did evolution turn us white people white when dark skinned Africans have no problem surviving in places like Northern Europe?*". These questions, while appearing innocent and curiosity-driven at first glance, contain offensive undertones and implications that the person posing the question might overlook. As our findings show, acts such as these examples, frequently feature themes around Black people's *ability* or their status in the *evolution* or *human race*. Moderators could employ the themes and linguistic patterns (e.g., 'Questions', etc.) identified in our research that commonly characterize acts to articulate their content moderation decisions to users. They can explain that posts like the examples above, containing seemingly harmless questions, are in fact associated with common themes of racial microaggressions. By leveraging these insights, moderators could bring a new level of transparency and understanding to their moderation process, potentially minimizing confusion and dispute within the community. Similarly, moderators can leverage the detailed characterizations of recalls from our research to identify potential misclassifications of recalls in their content moderation decisions. This would enable moderators to effectively communicate to posters of recalls who are mistakenly caught up in false positives, elaborating on reasons why their post was erroneously flagged as harmful.

Improving Moderation and Community Guidelines. Prior work in content moderation has highlighted the issue of vague and unclear guidelines on social media platforms [75, 80, 133]. These rules often lack explicit and specific wording, making their operationalization and enforcement processes opaque and non-transparent [80]. As a result, it becomes challenging for both users and moderators to understand how these rules are applied and why certain content is flagged or removed [75]. To solve this, we can incorporate clear illustrative examples of posts analyzed from our study into the moderation and community guidelines. Specifically, we can update moderation guidelines to include illustrative posts (e.g., "*Why are all black people so loud?*") highlighted with keywords derived from IG (e.g., 'Why', 'loud') as well as clear descriptions of the themes (e.g., 'Questions', 'Personality') and linguistic patterns (e.g., 'Questions', 'Use of absolute terminology') associated with acts and recalls in our study. Such comprehensive guidelines can empower moderators to make more informed and consistent decisions when evaluating user-generated content. Moreover, prior research suggests that "explicit rules and guidelines increase the ability for community members to know the norms" [75, 86]. Our findings can be used to establish clear and explicit community guidelines that informs users what type of content is acceptable to post on the community while minimizing the risk of unintentionally perpetuating harmful narratives (e.g., *acts*) or disregarding the experiences of marginalized groups (e.g., *recalls*).

5.3 Language Mimicry and Relational Dynamics in Online Discussion Communities

Online discussion communities, such as Reddit, naturally embody relational dynamics across users based on community roles (e.g., moderators/ admins vs. regular users) and membership statuses (e.g., old vs. new members). Research shows that such social structures and relational hierarchies

within online discussion groups can potentially play a role in language mimicry and adoption. Language coordination is a phenomenon in which people tend to unconsciously mimic the language of others by responding with similar words or phrases [106]. Research in computational linguistics has demonstrated how language coordination persists across conversations in ways that reflect power differentials between people. For example, in Supreme Court case settings, lawyers tend to linguistically mimic the language of the Supreme Court justices rather than vice versa [33]. Such language coordination also occurs online: Wikipedians tend to echo the linguistic style of admins significantly more than that of non-admins who are perceived to have a lower status within the community [33]. Our findings show that acts of racial microaggressions on social media embody persistent linguistic patterns, such as absolutist expressions (e.g., *never, ever hire someone with a Black name*) and modifying adverbs (e.g., *Black people are consistently the rudest*) that generalize or racially discriminate against Black people. Acts are also frequently masked in the form of questions disguised as genuine curiosity, or conveyed with statistics that tend to factualize selective information as broader truths. Given the presence of relational dynamics in online communities on top of platform affordances (e.g., up/down votes, likes, volume of comments) that interplay with such dynamics, linguistic patterns of acts can be mimicked and adopted across users, potentially amplifying racial biases, and endorsing harmful assumptions that underly acts of racial microaggressions. For future work, we intend to empirically capture how membership statuses and power dynamics within online discussion groups are associated with the adoption and spread of linguistic patterns of acts and recalls.

5.4 Critical Race Theory in Language and the Importance of Counter-Storytelling

According to Critical Race Theory, social conceptions of race and racism shape, and are shaped by laws, social movements, politics, and the media [110]. Such an argument is well-reflected across the theoretical premise of several anthropological research studies on race and language. Anthropological linguists have long recognized the importance of treating racial categories and concepts "not as objective facts about the world, but as the outcome of discursive processes that operate across intersecting scales of space and time" [30]. That is because, the process through which language itself is racialized, or the way language racializes certain groups of people over time, inevitably involves linking certain objects, ideas, and themes to a racial group [63, 74, 120], thereby concretizing stereotypes about a particular race, as shown in our findings. For example, the dominant themes that emerge across acts tend to link Blackness with **crime, sexual exoticism, and questionable belonging to the human race**, which falsely perpetuate racial tropes about Black people's **personality** (e.g., *funny, loud, dumb, creepy, ghetto, etc.*), **ability** (*sports, intelligent, IQ, etc.*), and **appearance** (*fat, hair, etc.*). Interestingly, many of these identical themes appear in recalls, wherein Black users engage in what critical race scholars describe as counter-storytelling. Counter-storytelling is the act of recounting an individual's experience with racism, typically through language that operates as a discursive tool for challenging majoritarian perspectives in culturally dominant discourses on race [37]. Stereotypes perpetuated through racial attitudes, and conceptions of race and racism that have persisted across centuries, tend to become normalized into culturally dominant narratives [36]. As a result, implicit racism as observed through online acts of racial microaggressions in our data can falsely appear as race neutral. Black users, as shown in our findings, in essence, call-out such biases through counter-storytelling, through which they directly challenge racial stereotypes and attitudes by conversing on the same topics and themes present in the acts through autobiographical language. Sociotechnical systems that fail to distinguish acts and recalls risk suppressing these counter-stories shared by Black users. Both critical race scholars and historians argue that sharing personal stories has always been essential to the survival and liberation of racially oppressed groups [35]. Ensuring sociotechnical systems that safeguard rather

than impede important conversations and experiential knowledge shared through counter-stories, such as the ones shown in this work, are critical to establishing more inclusive and enriching environments for online discourse.

6 CONCLUSION

Through this work, we call for a deeper understanding of the semantic differences between acts and recalls of racial microaggressions on social media, and a re-evaluation of how users and current socio-technical systems differentiate the two. As an initial step towards this effort, we manually curated a corpus of acts and recalls, which were discussed, hand-annotated, and verified by Black participants through a workshop session. We then used this data to classify, interpret and characterize the language underlying acts vs. recalls of racial microaggressions associated with Black people. By doing so, we provide an empirical characterization of the underlying themes, contexts, and the linguistic signature between acts and recalls.

7 LIMITATIONS

While our research is the first to systematically investigate acts and recalls of racial microaggressions comprehensively, our work is not without limitations. First, our analysis is limited to the context of racism in the U.S. Hence, implications around our findings may not be generalized to foreign contexts of racism against Blacks in other countries. Second, since Reddit is a global site, we have limited understanding of whether all the posters are from the U.S or not, which skews the earlier assumption of our analysis being limited to the context of racism in the U.S. Third, given the limited Black population in our area, we were only able to recruit Black college students for our workshop discussions. We plan to extend this study to a broader group in the future. Further, obtaining “ground-truth” labels for discursive data such as ours, are often fraught with subjective interpretations of race and social values linked with race-related matters, which are subject to the annotator’s own perspectives on racism, personal experiences, identity, and social background. Hence, while we endeavored towards obtaining “ground-truth” labels by discerning insights from discussions with Black participants through our workshop, we acknowledge that this process too, can be subject to biases. Finally, in our data collection, we excluded posts with neutral references to race to maintain focus on acts and recalls of microaggressions. This decision, while necessary for our study, may limit the generalizability of our machine learning models to broader or different types of contexts.

8 ACKNOWLEDGMENTS

We thank our workshop participants for providing helpful discussions and sharing their personal experiences with us.

REFERENCES

- [1] 2010-09-10. Social Network Sites as Networked Publics: Affordances, Dynamics, and Implications. In *A Networked Self* (0 ed.), Zizi Papacharissi (Ed.). Routledge, 47–66. <https://doi.org/10.4324/9780203876527-8>
- [2] Kupiri Ackerman-Barger, Dowin Boatright, Rosana Gonzalez-Colaso, Regina Orozco, and Darin Latimore. 2020. Seeking Inclusion Excellence. *Academic Medicine* 95, 5 (May 2020), 758–763. <https://doi.org/10.1097/acm.0000000000003077>
- [3] Shiza Ali, Mohammad Hammas Saeed, Esraa Aldreabi, Jeremy Blackburn, Emiliano De Cristofaro, Savvas Zannettou, and Gianluca Stringhini. 2021-06-21. Understanding the Effect of Deplatforming on Social Networks. In *13th ACM Web Science Conference 2021* (New York, NY, USA) (*WebSci '21*). Association for Computing Machinery, 187–195. <https://doi.org/10.1145/3447535.3462637>
- [4] Bobby Allyn. 2020. It’s ‘our fault’: Nextdoor CEO takes blame for deleting of black lives matter posts. <https://tinyurl.com/3353yeb4>
- [5] Nazanin Andalibi, Oliver L. Haimson, Munmun De Choudhury, and Andrea Forte. 2016-05-07. Understanding Social Media Disclosures of Sexual Abuse Through the Lenses of Support Seeking and Anonymity. In *Proceedings*

of the 2016 CHI Conference on Human Factors in Computing Systems (San Jose California USA). ACM, 3906–3918. <https://doi.org/10.1145/2858036.2858096>

- [6] Zahra Ashktorab and Jessica Vitak. 2016. Designing Cyberbullying Mitigation and Prevention Solutions through Participatory Design With Teenagers. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems* (San Jose, California, USA) (CHI '16). Association for Computing Machinery, New York, NY, USA, 3895–3905. <https://doi.org/10.1145/2858036.2858548>
- [7] Shervin Assari and Maryam Moghani Lankarani. 2018. Depressive Symptoms and Self-Esteem in White and Black Older Adults in the United States. *Brain sciences* 8, 6 (2018), 105.
- [8] Agathe Balayn, Jie Yang, Zoltan Szlavik, and Alessandro Bozzon. 2021. Automatic Identification of Harmful, Aggressive, Abusive, and Offensive Language on the Web: A Survey of Technical Biases Informed by Psychology Literature. *Trans. Soc. Comput.* 4, 3, Article 11 (oct 2021), 56 pages. <https://doi.org/10.1145/3479158>
- [9] James Banks. 2010. Regulating hate speech online. *International Review of Law, Computers & Technology* 24, 3 (2010), 233–239.
- [10] Daniel Bar-Tal. 2017. Self-censorship as a socio-political-psychological phenomenon: Conception and research. *Political Psychology* 38 (2017), 37–65.
- [11] Daniel Bar-Tal. 2017. Self-censorship: The conceptual framework. *Self-censorship in contexts of conflict: Theory and research* (2017), 1–18.
- [12] Mario L. Barnes. 2015. *Law & Society Review* 49, 1 (2015), 279–282. <http://www.jstor.org/stable/43670469>
- [13] Roy F Baumeister, Ellen Bratslavsky, Catrin Finkenauer, and Kathleen D Vohs. 2001. Bad is stronger than good. *Review of general psychology* 5, 4 (2001), 323–370.
- [14] Su Lin Blodgett, Lisa Green, and Brendan O'Connor. 2016. Demographic dialectal variation in social media: A case study of African-American English. *arXiv preprint arXiv:1608.08868* (2016).
- [15] Su Lin Blodgett, Q. Vera Liao, Alexandra Olteanu, Rada Mihalcea, Michael Muller, Morgan Klaus Scheuerman, Chenhao Tan, and Qian Yang. 2022. Responsible Language Technologies: Foreseeing and Mitigating Harms. In *Extended Abstracts of the 2022 CHI Conference on Human Factors in Computing Systems* (New Orleans, LA, USA) (CHI EA '22). Association for Computing Machinery, New York, NY, USA, Article 152, 3 pages. <https://doi.org/10.1145/3491101.3516502>
- [16] Arthur W Blume, Laura V Lovato, Bryan N Thyken, and Natasha Denny. 2012. The relationship of microaggressions with alcohol use and anxiety among ethnic minority college students in a historically White institution. *Cultural diversity and ethnic minority psychology* 18, 1 (2012), 45.
- [17] Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. 2016. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. *Advances in neural information processing systems* 29 (2016).
- [18] Courtney M. Bonam, Hilary B. Bergsieker, and Jennifer L. Eberhardt. 2016. Polluting Black space. 145 (2016), 1561–1582. <https://doi.org/10.1037/xge0000226> Place: US Publisher: American Psychological Association.
- [19] Valeria Borsotti and Pernille Bjorn. 2022. Humor and Stereotypes in Computing: An Equity-focused Approach to Institutional Accountability. *Computer Supported Cooperative Work (CSCW)* 31 (07 2022). <https://doi.org/10.1007/s10606-022-09440-9>
- [20] Danah Boyd. 2010. Social Network Sites as Networked Publics: Affordances, Dynamics, and Implications. <https://api.semanticscholar.org/CorpusID:158379198>
- [21] Luke Breitfeller, Emily Ahn, David Jurgens, and Yulia Tsvetkov. 2019-11. Finding Microaggressions in the Wild: A Case for Locating Elusive Phenomena in Social Media Posts. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)* (Hong Kong, China). Association for Computational Linguistics, 1664–1674. <https://doi.org/10.18653/v1/D19-1176>
- [22] Jerome S Bruner. 1957. On perceptual readiness. *Psychological review* 64, 2 (1957), 123.
- [23] Tommaso Caselli, Valerio Basile, Jelena Mitrović, Inga Kartoziya, and Michael Granitzer. 2020. I Feel Offended, Don't Be Abusive! Implicit/Explicit Messages in Offensive and Abusive Language. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*. European Language Resources Association, Marseille, France, 6193–6202. <https://aclanthology.org/2020.lrec-1.760>
- [24] Nicoletta Cavazza and Margherita Guidetti. 2014. Swearing in political discourse: Why vulgarity works. *Journal of Language and Social Psychology* 33, 5 (2014), 537–547.
- [25] Stevie Chancellor, Yannis Kalantidis, Jessica A. Pater, Munmun De Choudhury, and David A. Shamma. 2017. Multimodal Classification of Moderated Online Pro-Eating Disorder Content. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems* (Denver, Colorado, USA) (CHI '17). Association for Computing Machinery, New York, NY, USA, 3213–3226. <https://doi.org/10.1145/3025453.3025985>
- [26] Stevie Chancellor, Zhiyuan (Jerry) Lin, and Munmun De Choudhury. 2016. “This Post Will Just Get Taken Down”: Characterizing Removed Pro-Eating Disorder Social Media Content. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems* (San Jose, California, USA) (CHI '16). Association for Computing Machinery,

- New York, NY, USA, 1157–1162. <https://doi.org/10.1145/2858036.2858248>
- [27] Eshwar Chandrasekharan, Chaitrali Gandhi, Matthew Wortley Mustelier, and Eric Gilbert. 2019. Crossmod: A Cross-Community Learning-based System to Assist Reddit Moderators. *Proceedings of the ACM on human-computer interaction* 3, CSCW (2019), 1–30.
- [28] Eshwar Chandrasekharan, Mattia Samory, Shagun Jhaver, Hunter Charvat, Amy Bruckman, Cliff Lampe, Jacob Eisenstein, and Eric Gilbert. 2018. The Internet’s hidden rules: An empirical study of Reddit norm violations at micro, meso, and macro scales. *Proceedings of the ACM on Human-Computer Interaction* 2, CSCW (2018), 1–25.
- [29] Lijuan Chen, Yanyan Lan, Liang Pang, Jiafeng Guo, and Xueqi Cheng. 2021. Toward the Understanding of Deep Text Matching Models for Information Retrieval. *arXiv preprint arXiv:2108.07081* (2021).
- [30] Elaine W Chun and Adrienne Lo. 2015. Language and racialization. In *The Routledge handbook of linguistic anthropology*. Routledge, 220–233.
- [31] Chrissy Cook, Jie Cai, and Donghee Yvette Wohn. 2022. Awe Versus Aww: The Effectiveness of Two Kinds of Positive Emotional Stimulation on Stress Reduction for Online Content Moderators. *Proceedings of the ACM on Human-Computer Interaction* 6 (2022), 1 – 19.
- [32] Juliet M. Corbin and Anselm Strauss. 2008. Basics of Qualitative Research (3rd ed.): Techniques and Procedures for Developing Grounded Theory.
- [33] Cristian Danescu-Niculescu-Mizil, Lillian Lee, Bo Pang, and Jon Kleinberg. 2012. Echoes of Power: Language Effects and Power Differences in Social Interaction. In *Proceedings of the 21st International Conference on World Wide Web* (Lyon, France) (WWW ’12). Association for Computing Machinery, New York, NY, USA, 699–708. <https://doi.org/10.1145/2187836.2187931>
- [34] Anubrata Das, Brandon Dang, and Matthew Lease. 2020. Fast, Accurate, and Healthier: Interactive Blurring Helps Moderators Reduce Exposure to Harmful Content. In *AAAI Conference on Human Computation & Crowdsourcing*.
- [35] Richard Delgado and Jean Stefancic. 1989. Why do we tell the same stories?: Law reform, critical librarianship, and the triple helix dilemma. *Stanford Law Review* (1989), 207–225.
- [36] Richard Delgado and Jean Stefancic. 1991. Images of the outsider in American law and culture: Can free expression remedy systemic social ills. *Cornell L. Rev.* 77 (1991), 1258.
- [37] Richard Delgado and Jean Stefancic. 1993. Critical Race Theory: An Annotated Bibliography. *Virginia Law Review* 79, 2 (1993), 461–516. <http://www.jstor.org/stable/1073418>
- [38] Xiaohan Ding, Michael Horning, and Eugenia H. Rho. 2023. Same Words, Different Meanings: Semantic Polarization in Broadcast Media Language Forecasts Polarity in Online Public Discourse. *Proceedings of the International AAAI Conference on Web and Social Media* 17, 1 (Jun. 2023), 161–172. <https://doi.org/10.1609/icwsm.v17i1.22135>
- [39] Nicola Döring and M Rohangis Mohseni. 2019. Male dominance and sexism on YouTube: results of three content analyses. *Feminist Media Studies* 19, 4 (2019), 512–524.
- [40] Bryan Dosono. 2018. AAPI Identity Work on Reddit: Toward Social Support and Collective Action. In *Proceedings of the 2018 ACM Conference on Supporting Groupwork* (Sanibel Island, Florida, USA) (GROUP ’18). Association for Computing Machinery, New York, NY, USA, 373–378. <https://doi.org/10.1145/3148330.3152697>
- [41] Bryan Dosono and Bryan Semaan. 2019. Moderation Practices as Emotional Labor in Sustaining Online Communities: The Case of AAPI Identity Work on Reddit. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems* (Glasgow, Scotland Uk) (CHI ’19). Association for Computing Machinery, New York, NY, USA, 1–13. <https://doi.org/10.1145/3290605.3300372>
- [42] Mengnan Du, Varun Manjunatha, Rajiv Jain, Ruchi Deshpande, Franck Dernoncourt, Jiuxiang Gu, Tong Sun, and Xia Hu. 2021. Towards interpreting and mitigating shortcut learning behavior of NLU models. *arXiv preprint arXiv:2103.06922* (2021).
- [43] Maeve Duggan. 2017. 1 in 4 Black Americans have faced online harassment because of their race or ethnicity. <https://tinyurl.com/5n8n845e>
- [44] Jr. Dulaney, Earl F. 2006. Changes in Language Behavior as a Function of Veracity. *Human Communication Research* 9, 1 (03 2006), 75–82. <https://doi.org/10.1111/j.1468-2958.1982.tb00684.x> arXiv:<https://academic.oup.com/hcr/article-PDF/9/1/75/22343175/jhumcom0075.PDF>
- [45] Anca Dumitrache, Oana Inel, Benjamin Timmermans, Carlos Ortiz, Robert-Jan Sips, Lora Aroyo, and Chris Welty. 2021. Empirical methodology for crowdsourcing ground truth. *Semantic Web* 12, 3 (2021), 403–421.
- [46] Elizabeth Dvoskin, Nitasha Tiku, and Craig Timberg. 2021. Facebook’s race-blind practices around hate speech came at the expense of black users, new documents show. <https://www.washingtonpost.com/technology/2021/11/21/facebook-algorithm-biased-race/>
- [47] Sara E Shaw and Julia Bailey. 2009. Discourse analysis: what is it and why is it relevant to family practice? *Family Practice* 26, 5 (06 2009), 413–419. <https://doi.org/10.1093/fampra/cmp038> arXiv:<https://academic.oup.com/fampra/article-PDF/26/5/413/1779242/cmp038.PDF>

- [48] Upol Ehsan, Q. Vera Liao, Michael Muller, Mark O. Riedl, and Justin D. Weisz. 2021. Expanding Explainability: Towards Social Transparency in AI Systems. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems* (Yokohama, Japan) (CHI '21). Association for Computing Machinery, New York, NY, USA, Article 82, 19 pages. <https://doi.org/10.1145/3411764.3445188>
- [49] Mai ElSherief, Caleb Ziems, David Muchlinski, Vaishnavi Anupindi, Jordyn Seybolt, Munmun De Choudhury, and Diyi Yang. 2021-11. Latent Hatred: A Benchmark for Understanding Implicit Hate Speech. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing* (Online and Punta Cana, Dominican Republic). Association for Computational Linguistics, 345–363. <https://doi.org/10.18653/v1/2021.emnlp-main.29>
- [50] Rob Eschmann. 2021. Digital Resistance: How Online Communication Facilitates Responses to Racial Microaggressions. *Sociology of Race and Ethnicity* 7, 2 (2021), 264–277. <https://doi.org/10.1177/2332649220933307> arXiv:<https://doi.org/10.1177/2332649220933307>
- [51] Rob Eschmann, Jacob Groshek, Rachel Chanderdatt, Khea Chang, and Maysa Whyte. 2020. Making a Microaggression: Using Big Data and Qualitative Analysis to Map the Reproduction and Disruption of Microaggressions through Social Media. *Social Media + Society* 6, 4 (2020), 2056305120975716. <https://doi.org/10.1177/2056305120975716> arXiv:<https://doi.org/10.1177/2056305120975716>
- [52] Andre Espaillet, Danielle Panna, Dianne L. Goede, Matthew J. Gurka, Maureen A Novak, and Zareen Zaidi. 2019. An exploratory study on microaggressions in medical school: What are they and why should we care? *Perspectives on Medical Education* 8 (2019), 143 – 151.
- [53] Philomena Essed. 1991. *Understanding everyday racism: An interdisciplinary theory*. Sage Publications, Inc. Pages: x, 322.
- [54] Allyson Ettinger. 2020. What BERT is not: Lessons from a new suite of psycholinguistic diagnostics for language models. *Transactions of the Association for Computational Linguistics* 8 (2020), 34–48.
- [55] Jessica L. Feuston, Alex S. Taylor, and Anne Marie Piper. 2020. Conformity of Eating Disorders through Content Moderation. *Proc. ACM Hum.-Comput. Interact.* 4, CSCW1, Article 40 (may 2020), 28 pages. <https://doi.org/10.1145/3392845>
- [56] Meira Gebel. 2020. Black creators say Tiktok still secretly hides their content. <https://www.digitaltrends.com/social-media/black-creators-claim-tiktok-still-secretly-blocking-content/>
- [57] James Paul Gee. 2014. *An Introduction to Discourse Analysis: Theory and Method*. routledge.
- [58] Arpita Ghosh, Satyen Kale, and Preston McAfee. 2011. Who Moderates the Moderators? Crowdsourcing Abuse Detection in User-Generated Content. In *Proceedings of the 12th ACM Conference on Electronic Commerce* (San Jose, California, USA) (EC '11). Association for Computing Machinery, New York, NY, USA, 167–176. <https://doi.org/10.1145/1993574.1993599>
- [59] Graham R Gibbs. 2018. *Analyzing qualitative data*. Vol. 6. Sage.
- [60] Tarleton Gillespie. 2018. *Custodians of the Internet: Platforms, content moderation, and the hidden decisions that shape social media*. Yale University Press.
- [61] Tarleton Gillespie. 2022. Do Not Recommend? Reduction as a Form of Content Moderation. *Social Media + Society* 8, 3 (2022), 20563051221117552.
- [62] Phillip Atiba Goff, Jennifer L Eberhardt, Melissa J Williams, and Matthew Christian Jackson. 2008. Not yet human: implicit knowledge, historical dehumanization, and contemporary consequences. *Journal of personality and social psychology* 94, 2 (2008), 292.
- [63] Joseph P Goodwin. 2002. *Stories in the time of cholera: Racial profiling during a medical nightmare*. JSTOR.
- [64] Robert Gorwa, Reuben Binns, and Christian Katzenbach. 2020. Algorithmic content moderation: Technical and political challenges in the automation of platform governance. *Big Data & Society* 7, 1 (2020), 2053951719897945. <https://doi.org/10.1177/2053951719897945> arXiv:<https://doi.org/10.1177/2053951719897945>
- [65] Jonathan Grudin. 1988. Why CSCW applications fail: problems in the design and evaluation of organizational interfaces. In *Proceedings of the 1988 ACM conference on Computer-supported cooperative work*. 85–93.
- [66] Jessica Guynn. July 2019. Facebook while black: Users call it getting 'Zucked' say talking about racism is censored as hate speech. *Usa today* 24 (July 2019).
- [67] Oliver L. Haimson, Daniel Delmonaco, Peipei Nie, and Andrea Wegner. 2021. Disproportionate Removals and Differing Content Moderation Experiences for Conservative, Transgender, and Black Social Media Users: Marginalization and Moderation Gray Areas. *Proc. ACM Hum.-Comput. Interact.* 5, CSCW2, Article 466 (oct 2021), 35 pages. <https://doi.org/10.1145/3479610>
- [68] Joanne M Hall and Becky Fields. 2015. "It's killing us!" Narratives of Black adults about microaggression experiences and related health stress. *Global qualitative nursing research* 2 (2015), 2333393615591569.
- [69] Chayla Haynes, Saran Stewart, and Evette Allen. 2016. Three paths, one struggle: Black women and girls battling invisibility in U.S classrooms. *Journal of Negro Education* 85, 3 (2016), 380–391.

- [70] Danula Hettiachchi and Jorge Goncalves. 2020. Towards Effective Crowd-Powered Online Content Moderation. In *Proceedings of the 31st Australian Conference on Human-Computer-Interaction* (Fremantle, WA, Australia) (OZCHI'19). Association for Computing Machinery, New York, NY, USA, 342–346. <https://doi.org/10.1145/3369457.3369491>
- [71] Sharon Heung, Mahika Phutane, Shiri Azenkot, Megh Marathe, and Aditya Vashistha. 2022. Nothing Micro About It: Examining Ableist Microaggressions on Social Media. In *Proceedings of the 24th International ACM SIGACCESS Conference on Computers and Accessibility* (Athens, Greece) (ASSETS '22). Association for Computing Machinery, New York, NY, USA, Article 27, 14 pages. <https://doi.org/10.1145/3517428.3544801>
- [72] Jon Hurwitz and Mark Peffley. 1997. Public Perceptions of Race and Crime: The Role of Racial Stereotypes. *American Journal of Political Science* 41 (04 1997), 375. <https://doi.org/10.2307/2111769>
- [73] Jevan A. Hutson, Jessie G. Taft, Solon Barocas, and Karen Levy. 2018. Debiasing Desire: Addressing Bias & Discrimination on Intimate Platforms. *Proc. ACM Hum.-Comput. Interact.* 2, CSCW, Article 73 (Nov 2018), 18 pages. <https://doi.org/10.1145/3274342>
- [74] Miyako Inoue. 2018. Word for word: Verbatim as political technologies. *Annual Review of Anthropology* 47 (2018), 217–232.
- [75] Shagun Jhaver, Darren Scott Appling, Eric Gilbert, and Amy Bruckman. 2019. "Did You Suspect the Post Would Be Removed?": Understanding User Reactions to Content Removals on Reddit. *Proc. ACM Hum.-Comput. Interact.* 3, CSCW, Article 192 (nov 2019), 33 pages. <https://doi.org/10.1145/3359294>
- [76] Shagun Jhaver, Christian Boylston, Diyi Yang, and Amy Bruckman. 2021. Evaluating the effectiveness of deplatforming as a moderation strategy on Twitter. *Proceedings of the ACM on Human-Computer Interaction* 5, CSCW2 (2021), 1–30.
- [77] Shagun Jhaver, Amy Bruckman, and Eric Gilbert. 2019. Does transparency in moderation really matter? User behavior after content removal explanations on reddit. *Proceedings of the ACM on Human-Computer Interaction* 3, CSCW (2019), 1–27.
- [78] Shagun Jhaver, Amy Bruckman, and Eric Gilbert. 2019. Does Transparency in Moderation Really Matter? User Behavior After Content Removal Explanations on Reddit. *Proc. ACM Hum.-Comput. Interact.* 3, CSCW, Article 150 (nov 2019), 27 pages. <https://doi.org/10.1145/3359252>
- [79] Shagun Jhaver, Sucheta Ghoshal, Amy Bruckman, and Eric Gilbert. 2018. Online Harassment and Content Moderation: The Case of Blocklists. *ACM Trans. Comput.-Hum. Interact.* 25, 2, Article 12 (mar 2018), 33 pages. <https://doi.org/10.1145/3185593>
- [80] Perna Juneja, Deepika Rama Subramanian, and Tanushree Mitra. 2020. Through the Looking Glass: Study of Transparency in Reddit's Moderation Practices. *Proc. ACM Hum.-Comput. Interact.* 4, GROUP, Article 17 (jan 2020), 35 pages. <https://doi.org/10.1145/3375197>
- [81] Harmanpreet Kaur, Harsha Nori, Samuel Jenkins, Rich Caruana, Hanna Wallach, and Jennifer Wortman Vaughan. 2020. Interpreting Interpretability: Understanding Data Scientists' Use of Interpretability Tools for Machine Learning. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) (CHI '20). Association for Computing Machinery, New York, NY, USA, 1–14. <https://doi.org/10.1145/3313831.3376219>
- [82] Micere Keels, Myles Durkee, and Elan Hope. 2017. The psychological and academic costs of school-based racial and ethnic microaggressions. *American Educational Research Journal* 54, 6 (2017), 1316–1344.
- [83] Kennedy Kelis. 2021. Deconstructing The "13/50" Argument. <https://detester.org/publications/162kennedy>
- [84] Ben Kew. 2018. Poll: Two-thirds of Conservatives don't trust Facebook, believe social media censors conservatives. *Breitbart* (29 Aug 2018). <https://tinyurl.com/pzavznpb>
- [85] Gohar Feroz Khan, Bobby Swar, and Sang Kon Lee. 2014-10-01. Social Media Risks and Benefits: A Public Sector Perspective. 32, 5 (2014-10-01), 606–627. <https://doi.org/10.1177/0894439314524701> Publisher: SAGE Publications Inc.
- [86] Sara B. Kiesler, Robert E. Kraut, Paul Resnick, and Aniket Kittur. 2010. Regulating Behavior in Online Communities.
- [87] Yubo Kou and Xinning Gui. 2020. Mediating Community-AI Interaction through Situated Explanation: The Case of AI-Led Moderation. *Proc. ACM Hum.-Comput. Interact.* 4, CSCW2, Article 102 (oct 2020), 27 pages. <https://doi.org/10.1145/3415173>
- [88] Yubo Kou and Xinning Gui. 2021. Flag and Flagability in Automated Moderation: The Case of Reporting Toxic Behavior in an Online Game Community. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems* (Yokohama, Japan) (CHI '21). Association for Computing Machinery, New York, NY, USA, Article 437, 12 pages. <https://doi.org/10.1145/3411764.3445279>
- [89] Vivian Lai, Samuel Carton, Rajat Bhatnagar, Q. Vera Liao, Yunfeng Zhang, and Chenhao Tan. 2022. Human-AI Collaboration via Conditional Delegation: A Case Study of Content Moderation. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems* (New Orleans, LA, USA) (CHI '22). Association for Computing Machinery, New York, NY, USA, Article 54, 18 pages. <https://doi.org/10.1145/3491102.3501999>
- [90] J Richard Landis and Gary G Koch. 1977. An Application of Hierarchical Kappa-type Statistics in the Assessment of Majority Agreement among Multiple Observers. *Biometrics* (1977), 363–374.

- [91] Alyssa Lees, Daniel Borkan, Ian Kivlichan, Jorge Nario, and Tesh Goyal. 2021-04. Capturing Covertly Toxic Speech via Crowdsourcing. In *Proceedings of the First Workshop on Bridging Human–Computer Interaction and Natural Language Processing* (Online). Association for Computational Linguistics, 14–20.
- [92] Bruno Lepri, Nuria Oliver, Emmanuel Letouzé, Alex Pentland, and Patrick Vinck. 2018. Fair, Transparent, and Accountable Algorithmic Decision-making Processes: The Premise, the Proposed Solutions, and the Open Challenges. *Philosophy & Technology* 31 (12 2018). <https://doi.org/10.1007/s13347-017-0279-x>
- [93] Weiwen Leung, Zheng Zhang, Daviti Jibuti, Jinhao Zhao, Maximilian Klein, Casey Pierce, Lionel Robert, and Haiyi Zhu. 2020. Race, Gender and Beauty: The Effect of Information Provision on Online Hiring Biases. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) (*CHI '20*). Association for Computing Machinery, New York, NY, USA, 1–11. <https://doi.org/10.1145/3313831.3376874>
- [94] Jessa Lingel. 2019. The gentrification of the internet. <https://culturedigitally.org/2019/03/the-gentrification-of-the-internet/>
- [95] Daniel Link, Bernd Hellgrath, and Jie Ling. 2016. A Human-in-the-Loop Approach for Semi-Automated Content Moderation. In *International Conference on Information Systems for Crisis Response and Management*.
- [96] Yang Liu, Christopher Whitfield, Tianyang Zhang, Amanda Hauser, Taeyonn Reynolds, and Mohd Anwar. 2021. Monitoring COVID-19 pandemic through the lens of social media using natural language processing and machine learning. *Health Information Science and Systems* 9, 1 (June 2021). <https://doi.org/10.1007/s13755-021-00158-4>
- [97] Scott M Lundberg and Su-In Lee. 2017. A Unified Approach to Interpreting Model Predictions. In *Advances in Neural Information Processing Systems*, I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.), Vol. 30. Curran Associates, Inc. https://proceedings.neurips.cc/paper_files/paper/2017/file/8a20a8621978632d76c43dfd28b67767-Paper.pdf
- [98] Andreas Madsen, Siva Reddy, and Sarath Chandar. 2022-04-29. Post-hoc Interpretability for Neural NLP: A Survey. arXiv:2108.04840 [cs]
- [99] Andrea Marshall, Angela D. Pack, Sarah A. Owusu, Rainbo Hultman, David Drake, Florentine U. N. Rutaganira, Maria Namwanje, Chantell S. Evans, Edgar Garza-Lopez, Samantha C. Lewis, Cristina Termini, Salma AshShareef, Innes Hicsasmaz, Brittany L. Taylor, Melanie R. McReynolds, Haysetta D. Shuler, and Antentor O. Hinton. 2021. Responding and navigating racialized microaggressions in STEM. *Pathogens and Disease* 79 (2021).
- [100] Binny Mathew, Anurag Illendula, Punyajoy Saha, Soumya Sarkar, Pawan Goyal, and Animesh Mukherjee. 2020. Hate Begets Hate: A Temporal Study of Hate Speech. *Proc. ACM Hum.-Comput. Interact.* 4, CSCW2, Article 92 (oct 2020), 24 pages. <https://doi.org/10.1145/3415163>
- [101] Ryan Moore. 2017. *The New Jim Crow*. Macat Library. <https://doi.org/10.4324/9781912282586>
- [102] Samantha L Moore-Berg, Boaz Hameiri, and Emile G Bruneau. 2022. Empathy, dehumanization, and misperceptions: A media intervention humanizes migrants and increases empathy for their plight but only if misinformation about migrants is also corrected. *Social Psychological and Personality Science* 13, 2 (2022), 645–655.
- [103] Michael Muller and Thomas Erickson. 2018. In the Data Kitchen: A Review (a Design Fiction on Data Science). In *Extended Abstracts of the 2018 CHI Conference on Human Factors in Computing Systems* (Montreal QC, Canada) (*CHI EA '18*). Association for Computing Machinery, New York, NY, USA, 1–10. <https://doi.org/10.1145/3170427.3188407>
- [104] Michael Muller, Christine T. Wolf, Josh Andres, Michael Desmond, Narendra Nath Joshi, Zahra Ashktorab, Aabhas Sharma, Kristina Brimijoin, Qian Pan, Evelyn Duesterwald, and Casey Dugan. 2021. Designing Ground Truth and the Social Life of Labels. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems* (Yokohama, Japan) (*CHI '21*). Association for Computing Machinery, New York, NY, USA, Article 94, 16 pages. <https://doi.org/10.1145/3411764.3445402>
- [105] Kevin L Nadal, Katie E Griffin, Yinglee Wong, Sahran Hamit, and Morgan Rasmus. 2014. The impact of racial microaggressions on mental health: Counseling implications for clients of color. *Journal of Counseling & Development* 92, 1 (2014), 57–66.
- [106] Kate G Niederhoffer and James W Pennebaker. 2002. Linguistic style matching in social interaction. *Journal of Language and Social Psychology* 21, 4 (2002), 337–360.
- [107] Poppy Noor. 2020. The celebrities who are doing anti-racism right. <https://www.theguardian.com/world/2020/jun/15/selena-gomez-anti-racism-celebrities-getting-it-right>
- [108] Fayika Farhat Nova, Michael Ann DeVito, Pratyasha Saha, Kazi Shohanur Rashid, Shashwata Roy Turzo, Sadia Afrin, and Shion Guha. 2021. “Facebook Promotes More Harassment” Social Media Ecosystem, Skill and Marginalized Hijra Identity in Bangladesh. *Proceedings of the ACM on Human-Computer Interaction* 5, CSCW1 (2021), 1–35.
- [109] Fayika Farhat Nova, MD. Rashidujjaman Rifat, Pratyasha Saha, Syed Ishtiaque Ahmed, and Shion Guha. 2019. Online Sexual Harassment over Anonymous Social Media in Bangladesh. In *Proceedings of the Tenth International Conference on Information and Communication Technologies and Development* (Ahmedabad, India) (*ICTD '19*). Association for Computing Machinery, New York, NY, USA, Article 1, 12 pages. <https://doi.org/10.1145/3287098.3287107>

- [110] Ihudiya Finda Ogbonnaya-Ogburu, Angela D.R. Smith, Alexandra To, and Kentaro Toyama. 2020-04-21. Critical Race Theory for HCI. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems* (Honolulu HI USA). ACM, 1–16. <https://doi.org/10.1145/3313831.3376392>
- [111] Jason A Okonofua, Gregory M Walton, and Jennifer L Eberhardt. 2016. A vicious cycle: A social–psychological account of extreme racial disparities in school discipline. *Perspectives on Psychological Science* 11, 3 (2016), 381–398.
- [112] Joon Sung Park, Lindsay Popowski, Carrie Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2022. Social Simulacra: Creating Populated Prototypes for Social Computing Systems. In *Proceedings of the 35th Annual ACM Symposium on User Interface Software and Technology* (Bend, OR, USA) (UIST '22). Association for Computing Machinery, New York, NY, USA, Article 74, 18 pages. <https://doi.org/10.1145/3526113.3545616>
- [113] Joon Sung Park, Joseph Seering, and Michael S. Bernstein. 2022. Measuring the Prevalence of Anti-Social Behavior in Online Communities. *Proceedings of the ACM on Human-Computer Interaction* 6 (2022), 1 – 29.
- [114] John Pavlopoulos, Jeffrey Sorensen, Lucas Dixon, Nithum Thain, and Ion Androutsopoulos. 2020. Toxicity Detection: Does Context Really Matter?. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Online, 4296–4305. <https://doi.org/10.18653/v1/2020.acl-main.396>
- [115] Jennifer L. Eberhardt PhD. 2020-03-03. *Biased: Uncovering the Hidden Prejudice That Shapes What We See, Think, and Do*. Penguin. Google-Books-ID: hpDODwAAQBAJ.
- [116] Vinodkumar Prabhakaran, Ben Hutchinson, and Margaret Mitchell. 2019-11. Perturbation Sensitivity Analysis to Detect Unintended Model Biases. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)* (Hong Kong, China). Association for Computational Linguistics, 5740–5745. <https://doi.org/10.18653/v1/D19-1578>
- [117] Emilee Rader, Kelley Cotter, and Janghee Cho. 2018. Explanations as Mechanisms for Supporting Algorithmic Transparency. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems* (Montreal QC, Canada) (CHI '18). Association for Computing Machinery, New York, NY, USA, 1–13. <https://doi.org/10.1145/3173574.3173677>
- [118] Aneeta Rattan and Jennifer L Eberhardt. 2010. The role of social meaning in inattentive blindness: When the gorillas in our midst do not go unseen. *Journal of Experimental Social Psychology* 46, 6 (2010), 1085–1088.
- [119] Elizabeth Reid, Regan L. Mandryk, Nicole A. Beres, Madison Klarkowski, and Julian Frommel. 2022. Feeling Good and In Control: In-Game Tools to Support Targets of Toxicity. *Proc. ACM Hum.-Comput. Interact.* 6, CHI PLAY, Article 235 (oct 2022), 27 pages. <https://doi.org/10.1145/3549498>
- [120] Angela Reyes. 2017. *Language, identity, and stereotype among Southeast Asian American youth: The other Asian*. Routledge.
- [121] Eugenia Ha Rim Rho, Oliver L. Haimson, Nazanin Andalibi, Melissa Mazmanian, and Gillian R. Hayes. 2017-05-02. Class Confessions: Restorative Properties in Online Experiences of Socioeconomic Stigma. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems* (New York, NY, USA) (CHI '17). Association for Computing Machinery, 3377–3389. <https://doi.org/10.1145/3025453.3025921>
- [122] Eugenia Ha Rim Rho, Gloria Mark, and Melissa Mazmanian. 2018. Fostering civil discourse online: Linguistic behavior in comments of# metoo articles across political perspectives. *Proceedings of the ACM on human-computer interaction* 2, CSCW (2018), 1–28.
- [123] Eugenia Ha Rim Rho and Melissa Mazmanian. 2019. Hashtag burnout? a control experiment investigating how political hashtags shape reactions to news content. *Proceedings of the ACM on human-computer interaction* 3, CSCW (2019), 1–25.
- [124] Eugenia Ha Rim Rho and Melissa Mazmanian. 2020-04-21. Political Hashtags & the Lost Art of Democratic Discourse. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems* (Honolulu HI USA). ACM, 1–13. <https://doi.org/10.1145/3313831.3376542>
- [125] Marco Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. “Why Should I Trust You?”: Explaining the Predictions of Any Classifier. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Demonstrations*. Association for Computational Linguistics, San Diego, California, 97–101. <https://doi.org/10.18653/v1/N16-3020>
- [126] Travis Riddle and Stacey Sinclair. 2019. Racial disparities in school-based disciplinary actions are associated with county-level rates of racial bias. *Proceedings of the National Academy of Sciences* 116, 17 (2019), 8255–8260.
- [127] Sarah Roberts. 2018. Digital detritus: ‘Error’ and the logic of opacity in social media content moderation. *First Monday* 23 (03 2018). <https://doi.org/10.5210/firstmonday.v23i3.8283>
- [128] Maarten Sap, Dallas Card, Saadia Gabriel, Yejin Choi, and Noah A. Smith. 2019-07. The Risk of Racial Bias in Hate Speech Detection. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics* (Florence, Italy). Association for Computational Linguistics, 1668–1678. <https://doi.org/10.18653/v1/P19-1163>
- [129] Morgan Klaus Scheuerman, Stacy M. Branham, and Foad Hamidi. 2018. Safe Spaces and Safe Places: Unpacking Technology-Mediated Experiences of Safety and Harm with Transgender People. *Proc. ACM Hum.-Comput. Interact.* 2, CSCW, Article 155 (nov 2018), 27 pages. <https://doi.org/10.1145/3274424>

- [130] Christoph Schneider, Markus Weinmann, and Jan Vom Brocke. 2018. Digital nudging: guiding online user choices through interface design. *Commun. ACM* 61, 7 (2018), 67–73.
- [131] Daniel J Simons and Christopher F Chabris. 1999. Gorillas in our midst: Sustained inattention blindness for dynamic events. *perception* 28, 9 (1999), 1059–1074.
- [132] C. Estelle Smith, William Lane, Hannah Miller Hillberg, Daniel Kluver, Loren Terveen, and Svetlana Yarosh. 2021. Effective Strategies for Crowd-Powered Cognitive Reappraisal Systems: A Field Deployment of the Flip*²Doubt Web Application for Mental Health. *Proc. ACM Hum.-Comput. Interact.* 5, CSCW2, Article 417 (oct 2021), 37 pages. <https://doi.org/10.1145/3479561>
- [133] Jean Y. Song, Sangwook Lee, Jisoo Lee, Mina Kim, and Juho Kim. 2023. ModSandbox: Facilitating Online Community Moderation Through Error Prediction and Improvement of Automated Rules. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems* (Hamburg, Germany) (CHI '23). Association for Computing Machinery, New York, NY, USA, Article 107, 20 pages. <https://doi.org/10.1145/3544548.3581057>
- [134] Michael Spencer. 2017. Microaggressions and Social Work Practice, Education, and Research. *Journal of Ethnic & Cultural Diversity in Social Work* 26 (02 2017), 1–5. <https://doi.org/10.1080/15313204.2016.1268989>
- [135] Miriah Steiger, Timir J Bharucha, Sukrit Venkatagiri, Martin J. Riedl, and Matthew Lease. 2021. The Psychological Well-Being of Content Moderators: The Emotional Labor of Commercial Moderation and Avenues for Improving Support. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems* (Yokohama, Japan) (CHI '21). Association for Computing Machinery, New York, NY, USA, Article 341, 14 pages. <https://doi.org/10.1145/3411764.3445092>
- [136] Stefan Stieglitz and Linh Dang-Xuan. 2013. Emotions and information diffusion in social media—sentiment of microblogs and sharing behavior. *Journal of management information systems* 29, 4 (2013), 217–248.
- [137] Derald Wing Sue, Sarah Alsaidi, Michael N Awad, Elizabeth Glaeser, Cassandra Z Calle, and Narolyn Mendez. 2019. Disarming racial microaggressions: Microintervention strategies for targets, White allies, and bystanders. *American Psychologist* 74, 1 (2019), 128.
- [138] Derald Wing Sue, Christina M Capodilupo, Gina C Torino, Jennifer M Bucceri, Aisha Holder, Kevin L Nadal, and Marta Esquilin. 2007. Racial microaggressions in everyday life: implications for clinical practice. *American psychologist* 62, 4 (2007), 271.
- [139] Derald Wing Sue and Lisa Beth Spanierman. 2020. *Microaggressions in everyday life, 2nd ed.* John Wiley & Sons, Inc. Pages: xx, 349.
- [140] S Shyam Sundar, Haiyan Jia, T Franklin Waddell, and Yan Huang. 2015. Toward a theory of interactive media effects (TIME) four models for explaining how interface features affect user psychology. *The handbook of the psychology of communication technology* (2015), 47–86.
- [141] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017. Axiomatic Attribution for Deep Networks. In *International Conference on Machine Learning*. <https://api.semanticscholar.org/CorpusID:16747630>
- [142] Nicolas Suzor, Sarah Myers West, Andrew Quodling, and Jillian York. 2019. What do we mean when we talk about transparency? Towards meaningful transparency in commercial content moderation. *International Journal of Communication* 13 (2019), 1526–1543. <https://eprints.qut.edu.au/126386/>
- [143] Michael Swart, Ylana Lopez, Arunesh Mathur, and Marshini Chetty. 2020. Is This An Ad?: Automatically Disclosing Online Endorsements On YouTube With AdIntuition. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) (CHI '20). Association for Computing Machinery, New York, NY, USA, 1–12. <https://doi.org/10.1145/3313831.3376178>
- [144] Yla R. Tausczik and James W. Pennebaker. 2010. The Psychological Meaning of Words: LIWC and Computerized Text Analysis Methods. *Journal of Language and Social Psychology* 29, 1 (2010), 24–54. <https://doi.org/10.1177/0261927X09351676> arXiv:<https://doi.org/10.1177/0261927X09351676>
- [145] Hibby Thach, Samuel Mayworm, Daniel Delmonaco, and Oliver Haimson. 0. (In)visible moderation: A digital ethnography of marginalized users and content moderation on Twitch and Reddit. *New Media & Society* 0, 0 (0), 14614448221109804. <https://doi.org/10.1177/14614448221109804> arXiv:<https://doi.org/10.1177/14614448221109804>
- [146] Roelien C. Timmer, David Liebowitz, Surya Nepal, and Salil S. Kanhere. 2021-12. Can pre-trained Transformers be used in detecting complex sensitive sentences? - A Monsanto case study. In *2021 Third IEEE International Conference on Trust, Privacy and Security in Intelligent Systems and Applications (TPS-ISA)*. 90–97. <https://doi.org/10.1109/TPSISA52974.2021.00010>
- [147] Alexandra To, Hillary Carey, Geoff Kaufman, and Jessica Hammer. 2021. Reducing Uncertainty and Offering Comfort: Designing Technology for Coping with Interpersonal Racism. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems* (Yokohama, Japan) (CHI '21). Association for Computing Machinery, New York, NY, USA, Article 398, 17 pages. <https://doi.org/10.1145/3411764.3445590>
- [148] Alexandra To, Wenxia Sweeney, Jessica Hammer, and Geoff Kaufman. 2020. “They Just Don’t Get It”: Towards Social Technologies for Coping with Interpersonal Racism. *Proceedings of the ACM on Human-Computer Interaction* 4,

CSCW1 (2020), 1–29.

- [149] Stephen M Utych. 2018. Negative affective language in politics. *American Politics Research* 46, 1 (2018), 77–102.
- [150] Sonja Utz and Johannes Breuer. 2016. Informational benefits from social media use for professional purposes: Results from a longitudinal study. *Cyberpsychology: Journal of Psychosocial Research on Cyberspace* 10, 4 (2016).
- [151] Michael Veale, Max Van Kleek, and Reuben Binns. 2018. Fairness and Accountability Design Needs for Algorithmic Support in High-Stakes Public Sector Decision-Making. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems* (Montreal QC, Canada) (CHI '18). Association for Computing Machinery, New York, NY, USA, 1–14. <https://doi.org/10.1145/3173574.3174014>
- [152] James Vincent. 2020. Nextdoor CEO says it's 'our fault' moderators deleted black lives matter posts. <https://www.theverge.com/2020/7/2/21311046/nextdoor-ceo-admits-fault-moderators-racial-bias-black-lives-matter>
- [153] Zeerak Waseem, Thomas Davidson, Dana Warmusley, and Ingmar Weber. 2017. Understanding Abuse: A Typology of Abusive Language Detection Subtasks. In *Proceedings of the First Workshop on Abusive Language Online*. Association for Computational Linguistics, Vancouver, BC, Canada, 78–84. <https://doi.org/10.18653/v1/W17-3012>
- [154] Kelly Welch. 2007. Black Criminal Stereotypes and Racial Profiling. *Journal of Contemporary Criminal Justice* 23, 3 (2007), 276–288. <https://doi.org/10.1177/1043986207306870> arXiv:<https://doi.org/10.1177/1043986207306870>
- [155] Sarah Myers West. 2018. Censored, suspended, shadowbanned: User interpretations of content moderation on social media platforms. *New Media & Society* 20, 11 (2018), 4366–4383. <https://doi.org/10.1177/1461444818773059> arXiv:<https://doi.org/10.1177/1461444818773059>
- [156] Monnica T Williams, Matthew D Skinta, and Renée Martin-Willett. 2021. After Pierce and Sue: A revised racial microaggressions taxonomy. *Perspectives on Psychological Science* 16, 5 (2021), 991–1007.
- [157] Donghee Yvette Wohn. 2019. Volunteer Moderators in Twitch Micro Communities: How They Get Involved, the Roles They Play, and the Emotional Labor They Experience. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems* (Glasgow, Scotland Uk) (CHI '19). Association for Computing Machinery, New York, NY, USA, 1–13. <https://doi.org/10.1145/3290605.3300390>
- [158] Austin Wright, Omar Shaikh, Haekyu Park, Will Epperson, Muhammed Ahmed, Stephane Pinel, Duen Chau, and Diyi Yang. 2021. RECAST: Enabling User Recourse and Interpretability of Toxicity Detection Models with Interactive Visualization. *Proceedings of the ACM on Human-Computer Interaction* 5 (04 2021), 1–26. <https://doi.org/10.1145/3449280>
- [159] Austin P Wright, Omar Shaikh, Haekyu Park, Will Epperson, Muhammed Ahmed, Stephane Pinel, Diyi Yang, and Duen Hornng Chau. 2020. RECAST: Interactive Auditing of Automatic Toxicity Detection Models. In *The Eighth International Workshop of Chinese CHI* (Honolulu, HI, USA) (*Chinese CHI 2020*). Association for Computing Machinery, New York, NY, USA, 80–82. <https://doi.org/10.1145/3403676.3403691>
- [160] Qunfang Wu, Louisa Kayah Williams, Ellen Simpson, and Bryan Semaan. 2022. Conversations About Crime: Re-Enforcing and Fighting Against Platformed Racism on Reddit. *Proc. ACM Hum.-Comput. Interact.* 6, CSCW1, Article 54 (apr 2022), 38 pages. <https://doi.org/10.1145/3512901>
- [161] Yan Xia, Haiyi Zhu, Tun Lu, Peng Zhang, and Ning Gu. 2020. Exploring Antecedents and Consequences of Toxicity in Online Discussions: A Case Study on Reddit. *Proc. ACM Hum.-Comput. Interact.* 4, CSCW2, Article 108 (oct 2020), 23 pages. <https://doi.org/10.1145/3415179>
- [162] Qian Yang, Aaron Steinfeld, Carolyn Rosé, and John Zimmerman. 2020. Re-Examining Whether, Why, and How Human-AI Interaction Is Uniquely Difficult to Design. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) (CHI '20). Association for Computing Machinery, New York, NY, USA, 1–13. <https://doi.org/10.1145/3313831.3376301>
- [163] Savvas Zannettou, Mai ElSherief, Elizabeth Belding, Shirin Nilizadeh, and Gianluca Stringhini. 2020. Measuring and characterizing hate speech on news websites. In *12th ACM Conference on Web Science*. 125–134.

A APPENDIX

A.1 Annotation Guidelines

We analyze 13 of 16 themes of racial microaggressions from the revised Sue et al. (2009)'s taxonomy [156]. Three of the 16 themes (tokenism, environmental exclusion, and environmental attacks) were excluded as they were not applicable to online contexts: the three themes pertain to situations in which people are present in a physical environment. For each of the 13 themes, we provide the definition from [156] and examples of racial microaggressions from our RAMA corpus in Table 9. We discussed these examples and definitions with workshop participants as an annotation guideline for labeling instances of acts of online racial microaggressions.

Theme	Definition	Examples
THEME-1: Alien in own land / Not a True Citizen	When a question, statement, or behavior indicates that a person of color is not a real citizen or a meaningful part of society because they are not White; Questioning the legitimacy of their identity.	<ul style="list-style-type: none"> •“If Adam and Eve are the first people in the Earth and they are white, why are there Black people?” •“God gave black people rights in all corners of the globe and then he made the Earth round.”
THEME-2: Racial Categorization and Sameness	When a person is compelled to disclose their racial group to enable others to attach pathological racial stereotypes to the person; includes the assumption that all people from a particular group are alike	<ul style="list-style-type: none"> •“All Black people look alike.”
THEME-3: Assumptions about intelligence, competence, or status	When behavior or statements are based on an assumption about a person’s intelligence, competence, education, income, or social status derived from racial stereotypes.	<ul style="list-style-type: none"> •“How did Black people get so good at science? and why are they so athletic?”
THEME-4: Connecting via stereotypes	When a person tries to communicate or connect with a person through the use of stereotyped speech or behavior to be accepted or understood; can include racist jokes and epithets as terms of endearment.	<ul style="list-style-type: none"> •“Why do Black people love fried chicken and watermelon?” •“Why don’t Black people tip?”
THEME-5: False color blindness/ invalidating racial or ethnic identity	Expressing that an individual’s racial or ethnic identity should not be acknowledged, which can be invalidating for people who are proud of their identity or who have suffered because of it.	<ul style="list-style-type: none"> •“I don’t care if they are black, gay, green objects.” •“I’m white and I don’t care, we’re all the same human race.”
THEME-6: Myth of meritocracy/ race is irrelevant for success	When someone makes statements about success being rooted in personal efforts and denial of the existence of racism/White privilege; Statements which assert that race does not play a role in succeeding in career advancement or education.	<ul style="list-style-type: none"> •“Role should go to the best performer regardless of race.” •“Rich black people don’t face any form of systematic racism.”
THEME-7: Reverse-racism hostility	Expressions of jealousy or hostility surrounding the notion that POC get unfair advantages and benefits because of their race.	<ul style="list-style-type: none"> •“I was fully qualified for the job, but they gave it to a Black girl.” •“Oh wait you’re Black; they practically guarantee you’d get into that college.”
THEME-8: Criminality or dangerousness	Demonstrating belief in stereotypes that POC are dangerous, untrustworthy, and likely to commit crimes or cause bodily harm; A person of color is presumed to be dangerous, criminal, or deviant on the basis of their race.	<ul style="list-style-type: none"> •“I held back because he was Black” [user is speaking in the context of avoiding conflict with a Black person out of fear of physical retaliation]. •“Black men are dangerous.”
THEME-9: Avoidance and distancing	When POC are avoided or measures are taken to prevent physical contact or close proximity.	<ul style="list-style-type: none"> •“When I see a Black person approaching me, I cross the street.”
THEME-10: Denial of individual racism	When a person tries to make a case that they are not biased, often by talking about antiracist things they have done to deflect perceived scrutiny of their own biased behaviors; A statement made when Whites renounce their racial biases.	<ul style="list-style-type: none"> •“I’m not a racist. I have several Black friends.” •“I’m not racist but Black people make me uncomfortable.”
THEME-11: Pathologizing minority culture or appearance	When people criticize others on the basis of perceived or real cultural differences in appearance, traditions, behaviors, or preferences; The notion that the values and communication styles of the dominant culture are ideal.	<ul style="list-style-type: none"> •“Black kids shouldn’t dress that way.” •You’re pretty for a Black girl.
THEME-12: Exoticization and eroticization	When a person of color is treated according to sexualized stereotypes or attention to differences that are characterized as exotic in some way.	<ul style="list-style-type: none"> •“Black women are exotic. I have a fetish for Black women, am I racist?”
THEME-13: Second-class citizen/ ignored and invisible	When POC are treated with less respect, consideration, or care than is normally expected or customary; may include being ignored or being unseen/ invisible; Occurs when a White person is given preferential treatment as a consumer over a person of color.	<ul style="list-style-type: none"> •“Oh, sorry we kept you waiting so long. From your surname on the form, we thought you were Black!” •“Black people and LGBT are untouchable.”

Table 9. Themes and examples of acts of racial microaggressions from our RAMA corpus.

Received January 2023; revised July 2023; accepted November 2023