

An Investigation into Black and Brown Communities' Engagement with Data & Technology

Ebtesam Al-Haque^{1*}, Gabriella Thompson^{2*}, Angela D. R. Smith², Brittany Johnson¹

¹George Mason University

²University of Texas at Austin

ehaque4@gmu.edu, gabriella.thompson@utexas.edu, adrsmith@utexas.edu, johnsonb@gmu.edu

Abstract

Over the years, we have witnessed significant biases in datasets and AI-driven systems. While these biases can impact anyone, there is a heightened risk for disproportionate harm to Black and Brown communities. Despite efforts to address these inequities, a critical gap remains in understanding how Black and Brown communities perceive and interact with data-centric innovations. In this paper, we present findings from a survey of 60 technology users with diverse racial, ethnic, and gender identities. Our findings reveal that while Black and Brown users may frequently contribute data through social media and research participation, discomfort arises when data is used without explicit consent, particularly by for-profit organizations. Transparency, trust, and familiarity with data collection entities and outcomes emerged as key factors influencing engagement. These insights inform our ongoing efforts, as well as the development of ethical, inclusive approaches to data-driven innovation that center marginalized voices and foster equitable outcomes.

Introduction & Background

Data plays a pivotal role in contemporary research and innovation, offering valuable insights into societal trends and individual behaviors (Wu et al. 2013). However, the utility of data to drive equitable and inclusive technological outcomes is limited by representation gaps that disproportionately affect marginalized communities. These gaps often exacerbate existing inequalities, as evidenced by significant disparities in data-driven technologies, such as facial recognition and voice-based systems, which frequently underperform for Black and Brown users (Buolamwini and Gebru 2018; Koenecke et al. 2020).

The persistence of these inequities highlights how biased datasets and non-inclusive design processes contribute to high error rates for historically marginalized populations. Researchers have documented the effects of algorithmic bias, with studies showing that systems trained on datasets lacking racial diversity reproduce and amplify structural inequities (Sweeney 2013; Hern 2016; Mengesha et al. 2021). Efforts to address these issues have emphasized the importance of diversifying datasets and employing inclusive design practices, such as enabling broader self-identification

on demographic forms and adopting methods that prioritize diversity in data collection (Slade et al. 2021; Stasaski, Yang, and Hearst 2020).

Despite these initiatives, scholarship on the lived experiences and perceptions of Black, Indigenous, and People of Color (BIPOC) regarding data-driven technologies remains limited. Existing studies have largely focused on technical solutions without adequately addressing the social and cultural dimensions that shape how these communities engage with data and technology (Harrington and Egede 2023; Kim et al. 2022). We understand BIPOC communities are not monolithic and can vary greatly in their experiences, however our research focuses on Black and Brown communities as populations that are consistently overlooked and underserved in technology. This gap in perspective underscores the need for research that centers the voices and experiences of Black and Brown end users to understand their concerns, priorities, and expectations regarding data-centric innovations.

In this paper, we aim to address this gap by presenting preliminary findings from a survey that explores the factors influencing Black and Brown communities' engagement with data and technology. Our study investigates not only the barriers to equitable representation in datasets but also the potential for fostering trust and meaningful engagement in these communities.

Negative Impacts in Black and Brown Communities

Data-driven technologies have the capacity to facilitate targeted violence against marginalized communities, especially as identity becomes a form of quantifiable data that can be collected and analyzed. Centralized databases exemplify this harm when operating as a surveillance technology to monitor and control specific groups of people. Risk-identification systems demonstrate such usage by explicitly targeting and criminalizing unconsenting individuals, such as the Automated Targeting System from the U.S. Homeland Security that gathers biometric and behavioral data from travel records to profile an individual based on their "risk" to national security (Browne 2010). The California Gang Database, which was similarly built to target and harm marginalized communities, collects data based on racist contextual markers that disproportionately target Black and Lat-

*Both authors contributed equally to this research.

inx people, who make up around 87% of the total names (Benjamin 2019). Data surveillance technology can inflict equal harm when operating outside of centralized databases, as evidenced by the use of EBT data to track and stigmatize recipients of a benefits program in a political scheme against public assistance initiatives (O'neil 2017). In these technologies, the datafication of personal information enabled the explicit oppression of marginalized groups.

While overt efforts to inflict harm through data do exist, many forms of technology reproduce and exacerbate structures of oppression without direct intent. Algorithms may perpetuate harm against non-white communities when trained on data imbued with racism, as shown by Facebook's algorithm that discriminated against non-white users when targeting housing ads (Hern 2016) and Google's display of discriminatory ads when users searched for racially-coded names (Sweeney 2013). Though these algorithms are often confronted and subsequently removed or reformed, the rapid growth of Artificial Intelligence technologies has amplified the call to address bias in data technologies (Schlesinger, O'Hara, and Taylor 2018).

An absence of data from BIPOC (Black, Indigenous, and people of color) and marginalized communities may also contribute to the manifestation of real-world inequity and discrimination. A study on racial bias in machine learning from Kostick-Quenet et al. revealed that a lack of data on Black patients in need of cardiology care obscured racial disparities in the health outcomes of a specific heart assistance treatment (Kostick-Quenet et al. 2022). The authors attribute the gap in data to Black patients having difficulty accessing trustworthy and local cardiologists, a lack of assistance in being directed to specialized clinics, and a selection bias against Black patients in healthcare that results in delayed or inadequate treatment. The implications of this study suggest that algorithms trained on these datasets may misinterpret the role of race in the device's outcomes and can further perpetuate health inequity in Black communities. Our work aims to explore these gaps in data by engaging Black and Brown communities' to understand their perspectives and experiences when contributing data.

Black and Brown Perspectives on Engagement with Data & Technology

To further contextualize the inadequate representation of BIPOC communities in datasets, we must understand their lived realities as they engage with technology. Racism can permeate even minor interactions in digital spaces, such as experiencing cultural devaluation from anglo-biased auto-correction algorithms (Dyal-Chand 2021). Many researchers have sought to understand the effects that interacting with existing AI technology, such as Automated-Speech Recognition (ASR) software, has on BIPOC users' perceptions of usability. These systems have demonstrated notable performance gaps for Black users, likely due to a lack of training data on African-American Vernacular English (AAVE) (Koenecke et al. 2020). Mengesha et al. found that when Black participants used ASR, they felt othered by the tech and attributed the errors to the system, perceiving it as not made for them (Mengesha et al. 2021). Conversely, a com-

parative study on white and Black participants interacting with a voice assistant found that Black participants were more likely to internalize the technology by blaming the system's errors on themselves (Wenzel et al. 2023). Both studies found that interacting with the AI elicited negative emotional responses, such as heightened self-consciousness, frustration, and disappointment. Wenzel et al.'s study on voice assistant interactions with BIPOC users echoed these findings, pointing to the identity, cultural, and emotional harms that arise when this technology fails to recognize their speech (Wenzel and Kaufman 2024).

Technological systems that incorporate culturally-informed data by working with and for marginalized communities offer alternative experiences for BIPOC users. Several HCI scholars have recognized the potential for AI technology to support access to health information for Black and Brown adults and uncovered challenges that may inhibit their use (Mendu et al. 2018; Harrington and Egede 2023; Kim et al. 2022). In a pilot study with 66 Hispanic women, Mendu et al. deployed a culturally-informed virtual agent to educate on cervical cancer and found that the majority of participants reported high levels of satisfaction, usability, and trust after interacting with the technology (Mendu et al. 2018). Harrington et al. found that older Black adults initially distrusted health chatbots and questioned the source and credibility of their information, a pattern that the authors attributed to the history of distrust Black patients have towards the healthcare system (Harrington and Egede 2023). These findings were echoed by a design study from Kim et al., in which Black participants were wary of a health chatbot and sought transparency from the information it provided (Kim et al. 2022). Researchers from this study found that participants were more receptive to the chatbot when other Black community members acknowledged it, indicating the need for widespread community trust of a technology before it is adopted. The importance of community trust in relation to Black users' acceptance of new technology is supported by a study on the effects of a racial mirroring chatbot agent, in which Black participants demonstrated a greater preference for the same-race chatbot agent than any other racial group in the study (Liao and He 2020). We hope to provide novel insight to this area of research by identifying the conditions needed for Black and Brown communities to feel comfortable engaging various technologies.

User Attitude Towards Data Sharing

Understanding user attitudes toward data sharing requires recognizing that privacy expectations are contextual rather than universal. Nissenbaum's theory of privacy as contextual integrity argues that privacy violations occur when information flows violate established contextual norms, which fall into two categories: norms of appropriateness (determining what information is suitable to share) and norms of distribution (governing how information should flow between parties) (Nissenbaum 2004). This framework emphasizes that privacy expectations are shaped by the specific roles, relationships, purposes, and cultural norms operating within particular social settings. Rather than applying

universal privacy principles, contextual integrity recognizes that every social setting operates under its own informational guidelines. Whether in healthcare settings, personal relationships, or research contexts, each environment has developed specific expectations about what information sharing is acceptable and how that sharing should occur. These expectations emerge from the historical, cultural, and social experiences of the communities operating within those contexts. This theoretical perspective is crucial for understanding research participation attitudes because it indicates that different communities may have developed distinct expectations about appropriate information sharing based on their unique social contexts and cultural histories. Communities cannot be assumed to share identical privacy preferences, and research findings from one demographic group may not apply to others without careful consideration of the contextual factors that shape each community's information-sharing norms.

Building on this idea, recent scholarship has examined user attitudes toward uses of their social media data, revealing relationships between contextual factors and user comfort with data collection. Gilbert et al. conducted factorial vignette surveys to understand how users perceive researchers' use of their social media data, first focusing on Facebook users (Gilbert, Vitak, and Shilton 2021) and later expanding to examine cross-platform differences across dating apps, Instagram, and Reddit (Gilbert, Shilton, and Vitak 2023). Their research found that factors such as researcher domain, content type, purpose of research, and participant awareness significantly impact user comfort, with consent being the most influential factor across all contexts. Users consistently rated research scenarios as more appropriate and less concerning when informed consent was obtained prior to data collection.

Research on trust in research participation has found that communities with histories of exploitation develop different engagement patterns. Guillemin et al. (Guillemin et al. 2018) found that Indigenous participants required researchers to personally earn trust due to historical research exploitation, rather than trusting institutional reputation like non-Indigenous participants. This suggests that marginalized communities may have fundamentally different expectations for research relationships, highlighting the importance of understanding community-specific perspectives on research participation.

Bessenyei et al. conducted a large-scale survey investigating user comfort with smartphone data collection for mental health monitoring (Bessenyei et al. 2021). Their findings revealed that participants with mental health treatment history showed higher comfort levels than those without treatment, challenging assumptions about vulnerable populations being more wary of data collection. The authors speculated that this increased comfort could stem from stronger relationships that foster trust in healthcare providers, though they also raised concerns that cognitive and emotional aspects of mental health conditions might affect risk perception and decision-making capacity.

However, research suggests that Black and Brown populations may be less willing to share data for research and

monitoring purposes. Corbie-Smith et al. documented significant racial differences in research distrust, with African Americans demonstrating nearly five times the odds of high distrust scores compared to whites, even after controlling for socioeconomic factors (Corbie-Smith, Thomas, and George 2002). They speculate that this distrust stems from documented histories of medical and research exploitation, including unethical studies like Tuskegee (Brandt 1978). A study on Hispanic and Latinx community members' perceptions of wearable health technologies revealed a similar systemic distrust of healthcare and technology companies as participants were reluctant to share data with healthcare providers, despite being interested in the health benefits of wearable tech (Cruz et al. 2024). Some exploratory evidence suggests these patterns may extend to digital data sharing, with psychiatric outpatients who were non-white showing lower willingness to share passive smartphone data compared to white patients (Rieger et al. 2019).

While research identifies trust as a key factor influencing Black and Brown communities' research participation, it provides limited insight into the conditions and factors that might drive engagement or willingness to share data within these communities. Simultaneously, studies that do examine the factors driving engagement and comfort with data sharing either do not examine race as a factor or lack sufficient diverse participation to understand how these factors might operate differently across racial groups. Kyi and et al. focused their efforts on perceptions of Europeans regarding data collection and use, citing the importance of detailed and intentional descriptions of the research purposed (Kyi et al. 2024). Gilbert et al.'s studies predominantly sampled white participants, with their cross-platform study including only 8.5%-11.7% Black participants across different platforms (Gilbert, Shilton, and Vitak 2023) and no detailed reporting of race/ethnicity in their Facebook study (Gilbert, Vitak, and Shilton 2021). Similarly, while Bessenyei et al. acknowledged the role race plays in willingness to share data, they did not report detailed racial/ethnic demographics for their own sample (Bessenyei et al. 2021). Even research that has documented racial differences has been limited by small sample sizes and exploratory analyses that cannot adequately examine the specific factors that influence different racial communities' engagement with research (Rieger et al. 2019).

This results in a lack understanding both of what drives engagement within Black and Brown communities specifically and of whether the factors identified as important in predominantly white samples operate similarly across different racial groups. Understanding how Black and Brown populations specifically perceive and engage with research data collection, and what drives their comfort with participation, becomes crucial for developing more inclusive and equitable research practices. Our work addresses this gap by centering the voices and experiences of Black and Brown communities in understanding factors that influence their comfort with and participation in data-driven innovation.

Research Agenda

The goal of our ongoing research efforts is to investigate and improve the lack of diverse, representative datasets and insights in the development and use of technology. This includes investigations into experiences, community-engaged innovation, and the development of novel methods for engaging historically marginalized groups in data-centric computing research and innovation. The following overarching question guides our research agenda:

What does it mean for researchers and technologists to meaningfully and sustainably engage in data-centric innovation with historically marginalized communities?

Our efforts began with an investigation into the experiences of technologists seeking and using representative datasets for research and development. In parallel, we began collecting data from end users. In this paper, we report on the beginning of this effort and the plans for building on these efforts and insights.

Survey Methods

We designed our survey to support our ability to answer the following research questions:

RQ1: *In what ways and to what extent are Black and Brown tech users contributing data?*

RQ2: *To what extent are Black and Brown tech users aware of how their data is being used?*

RQ3: *What makes Black and Brown tech users feel comfortable when contributing data?*

Survey Design

As a part of our overall research agenda (page 4), we designed a survey to collect data from technologists, researchers, and end users on their experiences *contributing data*, *seeking data*, and *collecting data*. The **Full Version** of our survey included logic that would help respondents navigate the questions such that they only saw questions relevant to their role (e.g., a technologist seeking data vs. an end-user contributing data). However, during our efforts we discovered that it may be easier to collect end user data if we had a separate version that only included questions relevant to end users. Therefore, we created the **Tech User Version** of our survey with only *contributing data* questions. The data and insights we report on in this paper are from this version of the survey.

The **Tech User Version** survey included questions such as “*What types of practices do you have as it relates to contributing or providing your data?*” along with questions regarding factors that impact their comfort with their data being accessed and use with and without consent. We also asked demographic questions to increase our ability to characterize our sample. The survey was comprised mostly of multiple choice questions regarding respondents background and experiences with data and technology. We also included likert scaled questions that gauged respondents comfort with (e.g., very comfortable to very uncomfortable)

and awareness of (e.g., aware or not aware) regarding data collection and use practices by researchers and for-profit organizations. We were also intentional about survey question ordering, where we put questions regarding experiences and background before questions regarding comfort or awareness. This helps situate respondents in their experiences before providing experience-based opinions about the more specific considerations we were interested in (Hjortskov 2017; Schiff, Montagnes, and Peskowitz 2022). In addition, some questions were open-ended free-response items, allowing respondents to elaborate on their experiences, reasoning, or concerns in their own words.

Both the **Full Version** and **Tech User Version** of our survey are available online for reuse.¹

Data Collection

Studies have shown that the typical approach of advertising and recruiting via social media for research may not be effective when attempting to engage niche groups (Zindel 2023). Therefore, to engage our audience of interest we leveraged more intimate methods that allowed for direct, sometimes face-to-face, opportunities to find respondents. More specifically, for our survey, we leveraged our internal networks in collaboration with our community partner, who actively works with local Black and Brown populations to empower these communities through data equity. This community-based organization distributed our recruitment materials to their local networks and collaborated with our research team to engage Black and Brown communities at a community event in our area. This also helped decrease the potential for survey timing to impact responses or outcomes. At the event, we provided information on our research efforts that attendees could take with them and share with others. We also had iPad devices and seating at the event so attendees could take the survey on the spot. This effort yielded of our total responses (31). We acquired the remaining 29 responses through our community partner’s efforts advertising among their contacts. We created a unique distribution link for each of our recruitment efforts and did not include any of these links in any mass distribution or advertising.

Respondent Demographics

Our efforts thus far have yielded a total of 60 responses. Within that sample, respondents’ ages range from 18–44 years old, with most respondents being 25 or older. Our respondents cover a range of non-White racial identities 1, including Black (38), Latin American (9), and South Asian (3). Some respondents identified themselves with mixed racial identities, most of whom identified at least partially as Black (4) or Latin American (3). Two respondents preferred not to specify. Though our current sample consists of mostly cisgender female respondents (38), we also engaged cisgender males (10), a transgender male, and a non-binary respondent. Nine respondents opted to not specify their gender.

¹https://anonymous.4open.science/r/ETM_EndUser-8BAF/

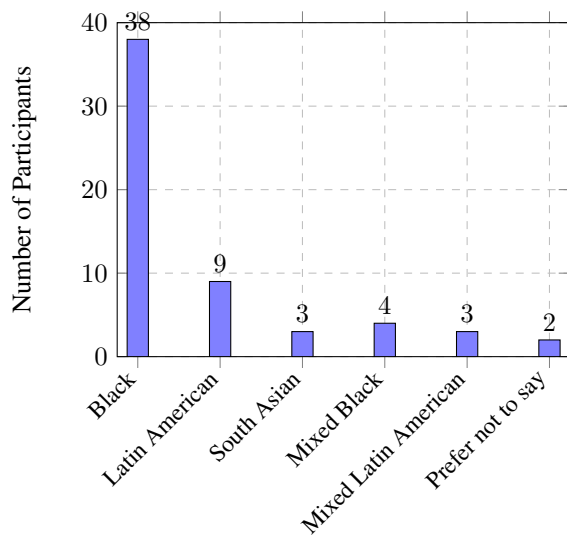


Figure 1: Self-reported ethnicity distribution of participants

Data Analysis

We conducted preliminary analyses on our current set of responses to glean initial insights into our research questions. Our analyses at this phase focused on descriptive statistics and statistical correlations using the Chi-Square independence test. To answer **RQ1**, we analyzed responses to questions regarding respondent data contribution practices, the types of data they contribute, the ways in which they share data and engage online, forms of engagement where they most often contribute data, and the devices they most often use. To answer **RQ2**, we analyzed responses to regarding rationale for current data sharing practices and awareness regarding the use of their data. To answer **RQ3**, we analyzed respondents' comfort with the ways in which their data is used (*with* and *without* consent), factors that contribute to their comfort with their data being used (with and without consent), factors that influence their decision to contribute data, and things it would be helpful to know about their data and its use.

Preliminary Findings

Data Contribution Patterns (RQ1)

The most common data contribution practice reported by our respondents is reading consent forms (51), followed by reading disclosure forms (42) and researching the platform they may be contributing data to (39). For 32 respondents, asking others was another way to make decisions when considering contributing data.

The kind of data most commonly (and recently) contributed among our respondents are *social media posts* (47), *social media engagement* (44), and *demographic data* (43). Our respondents also contributed data via public reviews, research participation, and short internet polls, though less frequently.

The least frequent kinds of data contributions included engaging with online advertisements (28), online knowledge

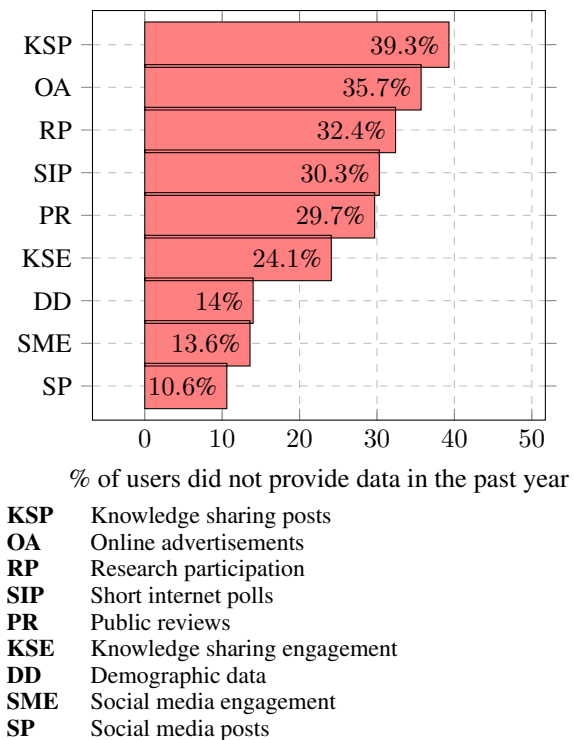


Figure 2: Decline in forms of online engagement

sharing/community posts (28) and engagement (29). According to our survey data, data contributions in most forms had decreased for our respondents. In the past year, while for most participants engagement on social media continued, many decreased other forms of engagement (Figure 2).

When analyzing the forms of online engagement where our respondents contribute data, we found that social media was the most common form of engagement across respondents. Most contributed data on Facebook and Instagram; many also contributed data on LinkedIn, TikTok, and X (formally known as Twitter). Additionally, we found comparable engagement in survey or interview research. Respondents cited contributing data on Stack Overflow, Reddit, and HackerNews much less. When analyzing the frequency of data contributions, most respondents report Facebook and Instagram were where they contributed data the most frequently.

Data Use Awareness (RQ2)

We asked respondents what experiences informed their current data sharing practices. For those who provided a response, we found two dominant themes: *professional or educational background* (9) and *media or informational sources* (9). For most responses that related to professional or educational background, data practices stemmed from having a technical or research background. However, for one respondent, the practices in a previous employment setting influenced his practice of reading terms and conditions, stating “, I used to work for a company that had some insidious clauses in their terms and conditions. Since then I do the best I can

to read terms and conditions as it pertains to EULA (End User License Agreement).”

For many respondents, social media and the news informed their data sharing practices. More specifically, respondents reported seeing information on news and media platforms regarding the things to be aware of regarding their data. One respondent noted that watching “*The Social Dilemma*”, a documentary that sheds light on the unintended consequences of social media and modern technology.²

Another theme that emerged was the impact of *negative experiences* on data contribution practices (6). One of the most common negative experiences was with online scams. One respondent elaborated on how experience with scams impact her practices:

“I have had about 5 instances of my bank account or social media being hacked. I’ve received notifications from a software I pay for that my personal data is on the dark web. Knowing where sensitive pieces of information not only live but have the potential of showing up is extremely important to me. I’ve also seen what happens when someone’s identity is stolen; it’s terrifying.”

We found that most of our respondents were aware that their data may be used without their permission or consent. In fact, we observed a high correlation between awareness of data use without explicit consent by both researchers and for-profits. As put by one respondent “, I have been scammed a lot online, so I generally assume that use of the internet comes with a risk to my data and privacy.”

Comfort with Contributing Data (RQ3)

Most respondents reported being comfortable with or indifferent towards their data being accessed or used for research, as long as consent or permission was obtained (36). Even with consent, many of our respondents reported being uncomfortable with for-profit companies using their data (31). This finding directly contradicts conclusion from prior work (Gilbert, Vitak, and Shilton 2021) that users are more wary of data collected for research purposes. Our study focused specifically on Black and Brown populations, while Gilbert et al. did not report race or ethnicity in their demographics, making it impossible to determine the racial composition of their sample. This limits our ability to fully understand whether the contrasting findings reflect methodological differences, population-specific attitudes toward research participation, or other factors. Our findings suggest that Black and Brown communities may demonstrate comfort with research participation when proper consent processes are followed.

When asked about data access and use without explicit consent or permission, an overwhelming majority of our respondents reported being uncomfortable (46). Even more reported being uncomfortable with for-profit companies accessing and using their data without consent (51).

Along with consent, respondents indicated other factors that contributed to their comfort with use of their data.

²<https://thesocialdilemma.com/the-dilemma/>

Most common among respondents was the desire for **transparency regarding data collection, analysis, and outcomes** (51), followed by **trust in the source or organization** (45), **familiarity with the source or organization** (40), and **efforts to explicitly acquire consent, including how specifically the data will be used, at some point prior to using your data** (39). Despite concerns regarding data collection and use, one respondent noted its necessity stating, “*I dislike this entire system of gathering information but when done well it can help unlock funding for vulnerable groups so it is a necessary evil even if imperfect. I support improving the designs of such projects for my community, many of us do.*” As implied by our findings regarding consent, most respondents are not comfortable with their data being collected and used without consent. However, for some respondents *familiarity with the source or organization* (20), **trust in the source or organization** (19), **transparency regarding data collection, analysis, and outcomes** (19), and **efforts to explicitly acquire consent, including how specifically the data will be used, at some point prior to using your data** (17) would increase their comfort without consent being obtained.

While some respondents would just “prefer it not being used at all,” many shared what it would be helpful to know when it comes to their data being accessed and used. The most common concerns among respondents were **how data is protected** and **what data is being used and why**. While it is not uncommon for data collection efforts to provide information on **how data is protected**, this was still a common consideration among respondents. This may be due to the fact that many users “*often don’t review that closely because the information is dense and hard to access.*” One respondent noted that they would like a “warning label” that summarized any concerns they should have regarding their data being access and used. Another respondent expanded on this notion, stating, “Having a quick easy way to review my data safety like a check up and a recommendation for who to quickly ‘unsubscribe’ from my data...if that exists.” Another respondent noted that it would be helpful to have “*publicized reports that are centralized on what information you have given permission to be collected.*”

Most respondents provided blanket responses regarding the desire to know **what data is being used and why**, making statements such as “*knowing the scope of where my data will be used*” and “*why they need it*”. However, some respondents provided more specific considerations for building their trust in data collection and use. One respondent preferred specific insights into “*how [the data] is being used to effect my demographic.*” Another stated that they “*would like to know where to find the results*” of when and how their data has been used.

Discussion & Future Work

Our efforts thus far have provided insights into the experiences and perceptions of Black and Brown tech users when it comes to contributing their data for use in research and technology. While our findings overlap with insights from prior works, to our knowledge, ours is the first effort to engage

and acquire insights from Black and Brown tech users regarding their data. Below we discuss some of these insights and our plans for building on this foundation.

Appropriateness & Distribution Norms

One of the most dominant themes in our efforts thus far is that Black and Brown communities may be more willing to contribute data for research than data for-profit. While prior work has examined and provided insights into end user perceptions of research data collection and use more broadly (Tenopir et al. 2011; Kyi et al. 2024), our findings provide empirical evidence for specific contextual privacy norms that shape research participation among Black and Brown communities, offering actionable insights for researchers seeking to conduct more inclusive studies. Building on Nissenbaum’s theory of privacy as contextual integrity (Nissenbaum 2004), we identified distinct **appropriateness norms** and **distribution norms** that can guide future research practices.

Our participants articulated specific expectations regarding what information is suitable to share for research purposes, revealing clear appropriateness norms for research contexts. Most prominently, they require clear articulation of research purposes, with particular emphasis on understanding what data is being used, why and how that data would be used to impact their own communities. This suggests that research appropriateness norms in these communities prioritize **community relevance and benefit**, and that it is extremely important for researchers to explicitly and intentionally connect their work to community outcomes and demonstrate tangible value for participants’ communities.

Participants also expressed expectations about how information should flow and be managed once shared, revealing important distribution norms for ongoing engagement. These norms emphasize the significance of **transparency** regarding how data is protected, **accessible communication** through simplified summaries or “warning labels,” **ongoing visibility** into data use through centralized reporting mechanisms, and **access to research outcomes**. These expectations reflect a preference for **sustained relationships** rather than one-time data transactions, suggesting that effective research partnerships require continuous engagement and accountability.

These findings suggest concrete strategies for researchers seeking to respect these contextual norms. Appropriateness norms can be honored by explicitly articulating community benefits, connecting research questions to community-identified priorities, and demonstrating how findings will be used to address relevant issues. Distribution norms can be respected through iterative engagement practices, accessible reporting mechanisms, and transparency tools that allow ongoing monitoring of data use. Our work demonstrates how contextual integrity theory can be operationalized in research practice by identifying the specific contextual norms operating in Black and Brown communities around research participation, providing a framework that other researchers can adapt and apply.

Trust Building & Iterative Engagement

Trust emerges as a critical factor shaping data contribution behaviors among Black and Brown communities, as evidenced by our findings. Respondents emphasized that transparency, familiarity with data-collecting organizations, and explicit consent mechanisms are essential to their willingness to share data. These observations align with existing research in human-computer interaction (HCI) that highlights trust-building as central to fostering engagement, particularly through culturally-informed design and clear communication about data use (Harrington and Egede 2023; Kim et al. 2022). Trust is especially vital given the historical exploitation of historically marginalized communities through biased surveillance systems and discriminatory technologies, such as facial recognition and algorithmic profiling, which perpetuate inequities (O’neil 2017; Buolamwini and Gebru 2018; Markl 2022; Baack 2024; Chen, Johansson, and Sontag 2018; Lin et al. 2020; Asudeh, Jin, and Jagadish 2019). Without mechanisms to demonstrate accountability and mitigate concerns about misuse, these communities remain wary of data contribution efforts.

Our findings further underscore the importance of **iterative engagement** in building and maintaining trust while encouraging sustained data contributions. Many respondents expressed dissatisfaction with the “one-time consent” model, which provides little visibility into how their data is used over time. Instead, they advocated for ongoing updates about the impact and outcomes of their contributions, reinforcing the value of continuous feedback as a trust-building mechanism in HCI and information sciences (Schlesinger, O’Hara, and Taylor 2018; Wenzel and Kaufman 2024). Iterative engagement ensures transparency and empowers participants to make informed decisions. This aligns with user-centered design principles, which prioritize accessible communication methods, such as visual dashboards, summarized reports, or automated alerts that are culturally relevant and easy to understand (Liao and He 2020; Stasaski, Yang, and Hearst 2020). These approaches address the opacity of algorithmic systems, which disproportionately marginalize communities by obscuring data practices and decision-making processes (Mengesha et al. 2021; Kim et al. 2022).

Beyond transparency, iterative engagement fosters prolonged participation by allowing community members to provide input and observe how their data influences research and technological development. This feedback loop adheres to participatory design principles, which have been shown to improve ethical alignment and promote trust in data-centric innovation (Stasaski, Yang, and Hearst 2020; Harrington and Egede 2023). Respondents noted that receiving regular updates on how their data is utilized and seeing visible accountability for ethical data use increased their confidence in contributing data. Such practices are particularly critical for improving the diversity of datasets, as representative data is essential to reducing algorithmic bias and advancing equitable technological outcomes (Kostick-Quenet et al. 2022; Sweeney 2013).

These insights have significant implications for researchers, technologists, and policymakers. Researchers must incorporate iterative updates into their study proto-

cols to maintain transparency and trust with participants. To implement these practices effectively with Black and Brown communities, adopting transparency-enhancing design patterns (Rossi and Lenzini 2020) that emphasize user-centric communication, appropriate timing of information provision, and iterative engagement throughout the research process can provide concrete tools for implementing the transparency and trust-building principles our participants identified as essential. Technologists should design systems that facilitate two-way communication and allow users to monitor data usage in real-time, while also incorporating culturally-sensitive design features to enhance engagement. Policymakers, in turn, can establish accountability frameworks that mandate iterative engagement as a standard for ethical data practices. Such measures not only address historical inequities and rebuild trust but also ensure that Black and Brown communities are active stakeholders in shaping inclusive and equitable technological advancements (Koencke et al. 2020; Wenzel et al. 2023). By implementing these approaches, organizations can build trust, address historical inequities, and ensure that Black and Brown communities are active and valued contributors in shaping equitable, inclusive technological advancements.

Next Steps

The overarching aim of our research agenda is to investigate the causes of underrepresentation and misrepresentation of Black and Brown identities and experiences in data-centric innovation and to develop strategies to address these inequities, fostering more equitable outcomes. Thus far, our efforts have yielded valuable insights into the perspectives of technologists and historically marginalized end users regarding data and technology.

Building on these foundational findings, we have refined our research agenda to better understand and operationalize the principles required for responsible innovation that ensures equitable benefits for Black and Brown users. Our next steps involve collecting additional survey data to develop intersectional personas that capture the nuanced privacy norms and data contribution preferences we have identified. While our current findings reveal general patterns among Black and Brown communities, we recognize that contextual privacy norms likely vary across intersections of race, gender, age, education, socioeconomic status, and other identity factors. To conduct meaningful intersectional analysis, we need to expand our survey reach to gather sufficient data across these various intersectional categories. This larger-scale data collection will help move beyond broad demographic categories toward understanding how multiple identity factors shape data contribution attitudes and trust-building requirements, providing researchers and technologists with more guidance for improving diverse representation in datasets. As part of this agenda, we are also conducting semi-structured interviews with technologists to examine their experiences in seeking and using representative BIPOC data. In parallel, we plan to conduct follow-up interviews with survey participants who opted in to gather deeper insights into the perceptions and experiences of Black and Brown communities regarding data contribution.

The insights from these interviews will inform the design of a series of workshops that will bring together technologists and end users in both homogeneous and heterogeneous group discussions. Findings from our current work have already revealed potential interventions to enhance trust and engagement within Black and Brown communities. These workshops, supported by insights from this study, will serve as a platform to further refine these interventions and explore co-design activities. Through this iterative process, we will synthesize potential solutions from the end-user perspective with feasibility considerations from the technologist perspective, laying a foundation for actionable, inclusive practices in data-centric innovation.

Acknowledgments

This material is based upon work supported by the National Science Foundation under Award Nos. 2224674 and 2224675. The authors also extend their sincere gratitude to the interviewees for this paper, without whom these insights would not be possible.

References

- Asudeh, A.; Jin, Z.; and Jagadish, H. 2019. Assessing and remedying coverage for a given dataset. In *2019 IEEE 35th International Conference on Data Engineering (ICDE)*, 554–565. IEEE.
- Baach, S. 2024. A Critical Analysis of the Largest Source for Generative AI Training Data: Common Crawl. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '24, 2199–2208. New York, NY, USA: Association for Computing Machinery. ISBN 9798400704505.
- Benjamin, R. 2019. *Race after Technology: Abolitionist Tools for the New Jim Code*. Cambridge and Medford: Polity Press.
- Bessenyei, K.; Suruliraj, B.; Bagnell, A.; McGrath, P.; Wozney, L.; Huguet, A.; Elger, B. S.; Meier, S.; and Orji, R. 2021. Comfortability with the passive collection of smartphone data for monitoring of mental health: an online survey. *Computers in Human Behavior Reports*, 4: 100134.
- Brandt, A. M. 1978. Racism and research: the case of the Tuskegee Syphilis Study. *Hastings center report*, 21–29.
- Browne, S. 2010. Digital epidermalization: Race, identity and biometrics. *Critical Sociology*, 36(1): 131–150.
- Buolamwini, J.; and Gebru, T. 2018. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on fairness, accountability and transparency*, 77–91. PMLR.
- Chen, I.; Johansson, F. D.; and Sontag, D. 2018. Why is my classifier discriminatory? *Advances in neural information processing systems*, 31.
- Corbie-Smith, G.; Thomas, S. B.; and George, D. M. M. S. 2002. Distrust, race, and research. *Archives of internal medicine*, 162(21): 2458–2463.
- Cruz, S.; Lu, C.; Ulloa, M.; Redding, A.; Hester, J.; and Jacobs, M. 2024. Perceptions of wearable health tools post the

- COVID-19 emergency in low-income Latin communities: qualitative study. *JMIR mHealth and uHealth*, 12: e50826.
- Dyal-Chand, R. 2021. Autocorrecting for Whiteness. *BUL Rev.*, 101: 191.
- Gilbert, S.; Shilton, K.; and Vitak, J. 2023. When research is the context: Cross-platform user expectations for social media data reuse. *Big Data & Society*, 10(1): 20539517231164108.
- Gilbert, S.; Vitak, J.; and Shilton, K. 2021. Measuring Americans' comfort with research uses of their social media data. *Social Media+ Society*, 7(3): 20563051211033824.
- Guillemin, M.; Barnard, E.; Allen, A.; Stewart, P.; Walker, H.; Rosenthal, D.; and Gillam, L. 2018. Do research participants trust researchers or their institution? *Journal of Empirical Research on Human Research Ethics*, 13(3): 285–294.
- Harrington, C. N.; and Egede, L. 2023. Trust, Comfort and Relatability: Understanding Black Older Adults' Perceptions of Chatbot Design for Health Information Seeking. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, CHI '23. New York, NY, USA: Association for Computing Machinery. ISBN 9781450394215.
- Hern, A. 2016. Facebook's "ethnic affinity" advertising sparks concerns of racial profiling.
- Hjortskov, M. 2017. Priming and context effects in citizen satisfaction surveys. *Public Administration*, 95(4): 912–926.
- Kim, J.; Muhic, J.; Robert, L. P.; and Park, S. Y. 2022. Designing Chatbots with Black Americans with Chronic Conditions: Overcoming Challenges against COVID-19. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, CHI '22. New York, NY, USA: Association for Computing Machinery. ISBN 9781450391573.
- Koenecke, A.; Nam, A.; Lake, E.; Nudell, J.; Quartey, M.; Mengesha, Z.; Toups, C.; Rickford, J. R.; Jurafsky, D.; and Goel, S. 2020. Racial disparities in automated speech recognition. *Proceedings of the national academy of sciences*, 117(14): 7684–7689.
- Kostick-Quenet, K. M.; Cohen, I. G.; Gerke, S.; Lo, B.; Antaki, J.; Movahedi, F.; Njah, H.; Schoen, L.; Estep, J. E.; and Blumenthal-Barby, J. 2022. Mitigating racial bias in machine learning. *Journal of Law, Medicine & Ethics*, 50(1): 92–100.
- Kyi, L.; Mhaidli, A.; Santos, C. T.; Roesner, F.; and Biega, A. J. 2024. "It doesn't tell me anything about how my data is used": User Perceptions of Data Collection Purposes. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, 1–12.
- Liao, Y.; and He, J. 2020. Racial mirroring effects on human-agent interaction in psychotherapeutic conversations. In *Proceedings of the 25th International Conference on Intelligent User Interfaces*, IUI '20, 430–442. New York, NY, USA: Association for Computing Machinery. ISBN 9781450371186.
- Lin, Y.; Guan, Y.; Asudeh, A.; and Jagadish, H. 2020. Identifying insufficient data coverage in databases with multiple relations. *Proceedings of the VLDB Endowment*, 13(11).
- Markl, N. 2022. Language variation and algorithmic bias: understanding algorithmic bias in British English automatic speech recognition. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '22, 521–534. New York, NY, USA: Association for Computing Machinery. ISBN 9781450393522.
- Mendu, S.; Boukhechba, M.; Gordon, J. R.; Datta, D.; Molina, E.; Arroyo, G.; Proctor, S. K.; Wells, K. J.; and Barnes, L. E. 2018. Design of a culturally-informed virtual human for educating hispanic women about cervical cancer. In *Proceedings of the 12th EAI International Conference on Pervasive Computing Technologies for Healthcare*, 360–366.
- Mengesha, Z.; Heldreth, C.; Lahav, M.; Sublewski, J.; and Tuennerman, E. 2021. "I don't think these devices are very culturally sensitive."—Impact of automated speech recognition errors on African Americans. *Frontiers in Artificial Intelligence*, 4: 725911.
- Nissenbaum, H. 2004. Privacy as contextual integrity. *Wash. L. Rev.*, 79: 119.
- O'neil, C. 2017. *Weapons of math destruction: How big data increases inequality and threatens democracy*. Crown.
- Rieger, A.; Gaines, A.; Barnett, I.; Baldassano, C. F.; Gibbons, M. B. C.; Crits-Christoph, P.; et al. 2019. Psychiatry outpatients' willingness to share social media posts and smartphone data for research and clinical purposes: Survey study. *JMIR formative research*, 3(3): e14329.
- Rossi, A.; and Lenzini, G. 2020. Transparency by design in data-informed research: A collection of information design patterns. *Computer Law & Security Review*, 37: 105402.
- Schiff, K. J.; Montagnes, B. P.; and Peskowitz, Z. 2022. Priming self-reported partisanship: implications for survey design and analysis. *Public Opinion Quarterly*, 86(3): 643–667.
- Schlesinger, A.; O'Hara, K. P.; and Taylor, A. S. 2018. Let's Talk About Race: Identity, Chatbots, and AI. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*, CHI '18, 1–14. New York, NY, USA: Association for Computing Machinery. ISBN 9781450356206.
- Slade, T.; Gross, D. P.; Niwa, L.; McKillop, A. B.; and Guptill, C. 2021. Sex and gender demographic questions: improving methodological quality, inclusivity, and ethical administration. *International Journal of Social Research Methodology*, 24(6): 727–738.
- Stasaski, K.; Yang, G. H.; and Hearst, M. A. 2020. More Diverse Dialogue Datasets via Diversity-Informed Data Collection. In Jurafsky, D.; Chai, J.; Schluter, N.; and Tetreault, J., eds., *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 4958–4968. Online: Association for Computational Linguistics.
- Sweeney, L. 2013. Discrimination in Online Ad Delivery: Google ads, black names and white names, racial discrimination, and click advertising. *Queue*, 11(3): 10–29.
- Tenopir, C.; Allard, S.; Douglass, K.; Aydinoglu, A. U.; Wu, L.; Read, E.; Manoff, M.; and Frame, M. 2011. Data sharing by scientists: practices and perceptions. *PLoS one*, 6(6): e21101.

Wenzel, K.; Devireddy, N.; Davison, C.; and Kaufman, G. 2023. Can voice assistants be microaggressors? Cross-race psychological responses to failures of automatic speech recognition. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, 1–14.

Wenzel, K.; and Kaufman, G. 2024. Designing for Harm Reduction: Communication Repair for Multicultural Users' Voice Interactions. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems, CHI '24*. New York, NY, USA: Association for Computing Machinery. ISBN 9798400703300.

Wu, X.; Zhu, X.; Wu, G.-Q.; and Ding, W. 2013. Data mining with big data. *IEEE transactions on knowledge and data engineering*, 26(1): 97–107.

Zindel, Z. 2023. Social media recruitment in online survey research: a systematic literature review. *Methods, data, analyses: a journal for quantitative methods and survey methodology (mda)*, 17(2): 207–248.