# Trust and Friction: Negotiating How Information Flows Through Decentralized Social Media

SOHYEON HWANG, Princeton University, USA
PRIYANKA NANAYAKKARA, Harvard University, USA
YAN SHVARTZSHNAIDER, York University, Canada

Decentralized social media protocols enable users in independent, user-hosted servers (i.e., instances) to interact with each other while they self-govern. This community-based model of social media governance opens up new opportunities for tailored decision-making about information flows—i.e., what user data is shared to whom and when—and in turn, for protecting user privacy. To better understand how community governance shapes privacy expectations on decentralized social media, we conducted a semi-structured interview with 23 users of the Fediverse, a decentralized social media network. Our findings illustrate important factors that shape a community's understandings of information flows, such as rules and proactive efforts from admins who are perceived as trustworthy. We also highlight "governance frictions" between communities that raise new privacy risks due to incompatibilities in values, security practices, and software. Our findings highlight the unique challenges of decentralized social media, suggest design opportunities to address frictions, and outline the role of participatory decision-making to realize the full potential of decentralization.

CCS Concepts: • **Human-centered computing → Collaborative and social computing theory, concepts and paradigms**; • **Security and privacy → Social aspects of security and privacy**.

Additional Key Words and Phrases: privacy, contextual integrity, online communities, decentralization, social media, moderation, trust and safety, online harassment, community governance

## 1 INTRODUCTION

Amidst ongoing debate about centralized, privately-owned social media platforms [31, 80], decentralized open-source alternatives are gaining attention [23, 39, 51, 94]. Users on decentralized social media like the Fediverse[1] join or set up independently-run instances (servers) that host their content and "federate" by connecting to a shared protocol which allows users across instances to interact. Although the interfaces visually resemble those of centralized platforms like Facebook or Twitter (𝕏), this decentralized design enables a form of community governance, where users (in particular, those acting as admins) can make key social and technical decisions in accordance with their values, shaping the flow of information on social media (Figure 1).

Because people have diverse privacy concerns [86, 95, 98], decentralization can enable groups to make more nuanced decisions [94] that meet their privacy needs as a community. For example,

---

[1]The Fediverse is a decentralized social media network for microblogging that has experienced rapid growth following the exodus from Twitter (now 𝕏) in 2022 [see 23].
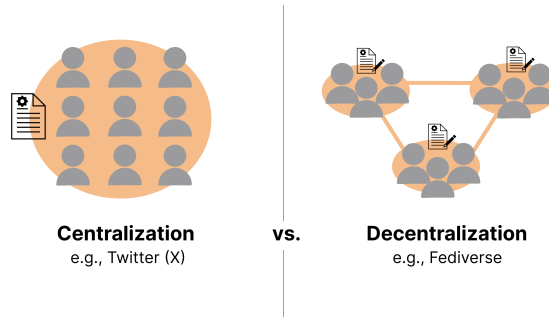
Fig. 1. Users must rely on top-down decisions in centralized social media settings (e.g., Twitter). In comparison, decentralized alternatives (e.g., the Fediverse) allow users to form self-governed communities, where members can make decisions about how information flows across the system at the community level.

for instances using Mastodon, the most widely-used software to run instances on the Fediverse, admins can toggle on or off a series of configurations. One notable option is "whitelist mode," which changes federation with other instances such that an instance will only interact with a list of approved instances — unlike the default, wherein instances interact with each other except for those that are explicitly blocked (on blocklists). Instances can decide not only whether to turn whitelist mode on, but also who is included in the list and who can make decisions about it. All of these decisions shape how easily user data can be accessed by potential bad actors on other instances who might doxx or harass users, thus impacting privacy risks users must anticipate. At the same time, privacy-preserving decisions may be difficult for instances to enact, as they require active maintenance of tools like white/blocklists or routinely ensuring back-end servers (i.e, where private messages are stored) are robustly secured. Prior work has shown that community governance requires users to have substantial time, resources, and expertise to govern effectively [48, 61, 119]; these challenges could discourage many from running an instance or rely on a default but ill-suited strategy. Meanwhile, how community members can anticipate the impact of governance practices on their privacy is unclear. In this work, we seek to better understand how the governance practices and processes of communities shape expectations of privacy on decentralized social media, toward identifying opportunities and challenges in making privacy-preserving decisions.

We use the theory of contextual integrity [77] to conceptualize privacy as *contextually-appropriate information flows*, i.e., whether the way user data (a user's profile, content, or other information about them) moves between actors, and when, is understood as appropriate in its social context. Drawing on 23 semi-structured interviews with community admins and members on the Fediverse, we investigate how communities self-govern to cultivate an understanding of what constitutes appropriate information flows for their particular community as part of a *decentralized* social media network. We contribute a rich description of how communities manage the bounds and technical understanding of information flows among users. Our findings highlight how challenges in making privacy-preserving decisions by communities are exacerbated by "governance frictions" which stem from three core incompatibilities between federated communities along: values, security practices, and software. These incompatibilities highlight key challenges that will need to be addressed to advance privacy with values such as safety in mind on decentralized social media.

## 2 BACKGROUND

### 2.1 Privacy, agency, and safety through the lens of contextual integrity

We follow growing theoretical consensus of privacy as a collective construct that is relational, contextual, and dynamic [10, 67, 77, 82, 83, 104, 107], grounding our understanding of privacy in the framework of contextual integrity (CI) [77]. CI defines privacy as the appropriate flow of information in a given context, with appropriateness determined by the established social norms of that context. The CI framework captures *information flows* using five parameters: the data type (what kind of information is being shared), subject (who or what the information is about), sender (who is sharing the data), recipient (who is receiving the data), and transmission principle (conditions imposed on the flow). Per CI, privacy is *prima facie* violated when a piece of information flows in a way that deviates from the established norms along one of these parameters, such as when information is shared with an unintended audience (recipient) or under unexpected conditions (transmission principle). CI has been a useful framework for making sense of privacy concerns around online participation [e.g., 14, 96], particularly in identifying *when* privacy is violated and *what* constitutes a violation. As such, the analytical advantages of CI have drawn calls in the CSCW community to leverage it in future empirical work [5, 58].

Although the CI framework can help identify potentially norm-violating information flows, it does not take a normative stance on the legitimacy of an information flow beyond established norms (and not all norms are desireable to all [70]). However, recognizing that contexts are social spheres "comprise[ed of] characteristic activities and practices, functions (or roles), aims, purposes, institutional structures, values, and action-governing norms" [78], CI as a framework can also bring attention to underlying structures that shape expectations about privacy in a given space. Nissenbaum [78] details heuristics to evaluate the actors, values, and potential outcomes (benefits, harms, consequences, etc.) of contexts that information flows are a part of, in order to analyze the ethical implications of how privacy norms emerge.

A key consideration is which actors have *agency* over information flows because this can change expectations about how information flows are managed. McDonald and Forte [71] note that companies running centralized platforms such as Facebook, Instagram, and Snapchat exert their power over policy and design to establish privacy norms. However, when agency is shifted towards community decision-making over information flows on social media, norms about what are acceptable information flows are likely to align more closely with specific community needs. In short, user agency can radically change what is understood as an actual privacy violation, particularly because different people have different values that shape norms around privacy. For example, in the context of social media, *safety* is a core value that may or may not be prioritized when determining what are contextually-appropriate flows of information, i.e., safety can shape perspectives on whether privacy has been violated (for example, did an information flow expose someone to harassment?). Taking the lens of CI allows us to see how key heuristics like agency and safety guide how we make normative judgments about privacy in settings like social media, where sharing information about oneself is generally "the point."

### 2.2 Decentralizing social media governance to communities

Continued ethical and regulatory challenges surrounding popular social media platforms [31, 45, 80] have rekindled interest in decentralized models of governing as promising ways forward [e.g., 23, 51, 62, 108]. Recent calls for decentralization have often been framed as efforts to recoup user agency [51] over their online privacy as centralized platforms collate and sell information about users at scale and with fine granularity [6, 13, 18, 49, 50, 69, 79, 120]. Viljoen [107] argues that such datafication (i.e., the extraction and quantification of peoples' lives into data) by centralized

platforms "materializes unjust social relations" that exacerbate inequalities. Addressing these relational harms thus requires *collective* institutional forms rather than just a reassertion of individual control over information, an argument echoed by related work calling for empowering bottom-up organizing around data by end users [26, 62, 108, 116].

As a decentralized model of decision-making, *community* governance on social media platforms offers one potential response by re-arranging the relations of data governance around groups of users. We define communities as *self-defined groups of users—having voluntarily and freely joined— interacting towards some common interests in a shared virtual space within a platform*, drawing on a more discrete conceptualization of "community" [16, 34]. Because of our interest in community governance as a form of decentralized governance, we focus on communities as units of meaningful decision-making, distinguishing the concept from the "sense of community" [11] often associated with the term. Communities are founded by users, who govern themselves by managing membership (i.e., granting different access privileges), setting rules that guide participation, and encouraging engagement among members. Community governance refers to how a community enacts practices and processes that "[create] the conditions for ordered rule" [99] in itself.

Community governance can enhance user agency in social media by placing key decisions over information flows — about moderation, data access, etc. — in the hands of decentralized groups of users rather than a centralized entity. CSCW scholars have long noted how online communities can lead to more nuanced and contextually-appropriate governance decisions in large-scale social computing systems [e.g., 56, 61, 94], by each representing a distinct subset of norms and goals [19, 101, 114]. This is particularly salient in a community's rules, i.e., the written normative prescriptions set by the community about how one should participate in it (sometimes also called codes of conduct, covenants, policies, or guidelines; for example, a community for sharing privacy tips might have a rule about citing credible sources). Differences in community governance can thus emerge along many dimensions, such as the processes they set for decision-making (making new rules, recruiting new admins, communicating issues, etc.), their organizational and participatory structure [90], or the tools they use to facilitate governance [118]. In truly decentralized systems, communities may even have say over the software and technical infrastructure they are hosted on although the need for technical expertise may pose high barriers of entry.

Given the plurality of privacy needs across different socio-demographic groups [86, 95, 98], communities are poised to enable more meaningful decision-making that advance the privacy interests of their particular members through community rules, processes, and tools. As communities can vary widely in purpose, views about what information flows are appropriate may similarly vary; communities can choose to enforce these views by allowing or restricting particular information flows. Understanding a community's acceptable costs and trade-offs can guide privacy-preserving governance processes and institutions [86, 97]. For example, many communities may want have a broad, active network, such as professional communities seeking wide engagement to advance their careers [111]. However, we can intuit that there also certain kinds of non-work content their members would not want to interact with. In general, growth is assumed to be an imperative for communities to obtain critical mass [81] and take advantage of network effects. However, a community focused on sharing intimate, highly personal, or marginalized experiences may want their content to stay within a defined safe space [44, 117]. Marginalized communities especially face heightened risks of online harassment [9, 46], which has made the ability to prioritize safety in defining and maintaining privacy crucial [32, 70]. Contemporaneous work calls for tools that strengthen "the ability of group participants to articulate and enforce privacy norms online" [21].

## 2.3 Shaping privacy expectations in social media communities

A central privacy concern in any social media context is being able to anticipate the audience one is sharing information with [7, 63, 109, 110], especially as changes in technical design can unexpectedly shift the terms of interaction among users [22, 66, 75, 105, 109]. To this end, a rich body of prior work in CSCW and related fields has examined how communities on social media can shape audience by enforcing rules around access, identity, and anonymity.

One notable perspective in prior CSCW research looks to communities as *safe spaces* where people — especially those from historically marginalized groups — can disclose information typically considered sensitive or personal without fear of anti-social responses or offline repercussions [1, 4, 8, 8, 27, 43, 44, 74, 92]. Safe spaces may focus on limiting access to the community from potential bad actors by increasing barriers to membership (e.g., requiring people apply to join the community, blocking certain users or groups [40]), and enforcing boundary-regulating norms more strictly. This work has typically focused on the experiences of individuals hoping to interact with a narrower, more sympathetic audience as they seek information, contribute their own content for others, and find new means of support online [29, 47, 52, 85, 103]. In many cases, people participate in these spaces because they can be "anonymous" (i.e., dissociate content from oneself as the source [2]), enabling them to safely share "private sentiments in the public sphere" [3].

However, enabling "anonymity" or similar strategies is not a silver bullet solution to creating privacy-respecting and safe communities. Allowing users to be anonymous can also undermine a community's ability to provide healthy and safe interactions [42, 100, 103] by creating room for harmful behaviors such as vandalism, spam, scams, toxic behavior. Cases like these warrant a certain level of information collection about individual users to enable content moderation [72]. At the same time, such information collection can also be perceived as a kind of surveillance [12, 36, 87], especially as it is not clear whether disallowing anonymity is necessarily effective. For example, Forte et al. [35] found that Wikipedians contributing via Tor — a network which enables anonymous communication online [41, 72, 103] — did so out of fear of violence, harassment, reputation loss, and fear for loved ones. But because anonymous contributions on Wikipedia are associated with spam, vandalism, trolling, and abuse, Wikipedia has attempted to block contributions coming from Tor users for nearly two decades [41, 103]. Yet, researchers have shown that Tor users that manage to bypass blocks contribute revisions of similar quality to non-Tor users [103]. McDonald et al. [72] note that community policies around contribution quality and quantity shaped whether people felt anonymous contributions should be allowed.

These tensions demonstrate that community decisions rooted in values likes safety or quality can impact how people perceive the appropriateness of information flows. However, they are not straightforward in their effects and may require substantial iteration and deliberation. Community governance can be highly labor-intensive while voluntary [30, 61]: burnout is well-documented in prior research [93]. Configuring settings that have more direct implications for privacy may require technical know-how (e.g., setting up back-end servers, selecting secure services, etc.) or legal expertise (e.g., privacy policies and terms) that not all community leaders have. Tosch et al. [102] show that, in a sample of 351 privacy policies on the Fediverse, only about 10% had been customized, echoing broader patterns of community rule isomorphism across platforms [55].

Amidst these challenges, we know precious little about which governance practices of communities shape privacy expectations, and how. In *decentralized* social media, how to guide privacy-preserving governance decisions through community governance is both pressing and relatively understudied given the recent revitalization of interest in decentralization. Compared to large, privately-owned platforms — which run on proprietary servers and code, and centrally determine and enforce a broad set of rules through design choices, algorithms, and moderation — decentralized

social media embodies critically different organizing and governing principles. Although prior research indicates that social media users rely on implicit rules (norms) that become agreed upon over time and then become taken for granted [21, 28], how those norms become communicated and established in decentralized systems is unclear. Studies on *self-disclosure* examine patterns of information sharing that provide signals of peoples' expectations of privacy [e.g., 112, 115]. However, these results are primarily on centralized social media platforms and provide limited insight into how community governance might tie to privacy expectations on decentralized social media—as well as how decentralization can mitigate or produce new challenges.

## 3 EMPIRICAL SETTING: THE FEDIVERSE

We examine the Fediverse, a decentralized social media network made up of independent *instances* (also called *servers*) that can communicate with each other, or "federate," through a shared protocol.[2] A user on the Fediverse signs up on an instance (not an overarching platform), which lets them interact with other users of their instance as well as users of other instances that their instance is federated with. Often called a decentralized alternative to Twitter, the Fediverse usually consists of micro-blogging interactions like short text status updates, photos and videos, or comments to others' posts. The exact interface of the instance depends on the software it is running to connect to the Fediverse, and the most popular software used by community-run servers is a free and open-source one called Mastodon. Following our definition of "community" in §2.2, we use "community" interchangeably with "instance" throughout this work.

On the Fediverse, any person can set up an instance, although technical barriers and hosting costs mean many join existing ones instead. Instance owners have significant governing power, shaping normative expectations for the members of their instance: what kinds of behaviors are acceptable, how content will be moderated, how data is maintained in the back-end, and which instances they will federate with. The configurations of an instance (typically, set by the admin who operates it) can shape information flows originating from and traversing its community on social media through community governance decisions. As in many online contexts, how instances make these decisions as a group varies. Some follow a "feudal" system where leaders act as "benevolent dictators" who make executive decisions [89], while others take extra effort to cultivate community engagement in decision-making [see work supporting such efforts like 118]. Depending on how the community operates, some decisions impacting privacy — such as how well the server's database is secured — may not be visible or obvious to users.

The Fediverse is also notable as a home built by and for traditionally marginalized groups, in particular LGBTQ+ communities [64]. The ActivityPub protocol the Fediverse operates on was primarily developed by queer and trans developers [38, 53]. Likewise, many conversations around safety and governance on the Fediverse have been led and advanced by these users. #Fediblock, a well-known hashtag tool for tracking bad actors to block, was started by a group of queer femmes who wanted to call out a sexual harasser [54]. Recent conversations have also highlighted concerns about racism on the Fediverse [54, 84]. This sociocultural context is crucial in understanding privacy concerns on the Fediverse, as they are thus closely tied to questions of safety.

## 4 STUDY DESIGN

Between October 2023—March 2024, the first author conducted semi-structured interviews via video calls with 23 users of the Fediverse. We recruited participants through a general call posted on the Fediverse and by directly contacting users. We took a statistically non-representative sampling strategy [106] in recruitment, intentionally aiming for diversity in the pool along dimensions of

---

[2]At the time of writing, the ActivityPub protocol is the main protocol of the Fediverse.

user type (member vs. admin), user gender, server topic, and server size. Participants were offered an Amazon e-gift card valued at 20 USD. Interviews were on average 77.5 minutes long. Participants were assigned pseudo-random numbers as participant IDs (P#s) and names were removed from transcripts. We also obfuscated details in our reporting to prevent identification. The study was approved by the [university ethics review boards].

## 4.1 Interview protocol

The interview protocol began with asking participants about their perceptions of information flows on the Fediverse, using a survey as a structured think-aloud exercise (see Appendix). Each survey item served as a vignette describing a possible information flow (drawing on CI), and participants noted the level of acceptability while explaining their choice to the interviewer. We devised these vignettes based on the language of a widely-used privacy policy on the Fediverse: in a sample of 803 servers drawn from a service listing Fediverse instances (`instances.social`), 82% of servers used this text, which was the default privacy policy text automatically generated if an instance used the Mastodon software. Thus, we expected the policy to be representative of policies on the Fediverse broadly. We followed up with open questions about any privacy concerns on the instance (and the Fediverse) and reflections on the instance's practices in protecting user privacy. We also asked about how participants joined (or started) their instances for background.

## 4.2 Interview pool

Table 1 shows the interview pool of 23 Fediverse users that included 11 men, 7 women, and 5 agender/non-binary/genderfluid people falling in the age range of 18-69. Participants were all located in the Global North (United States, Canada, Germany, and Portugal), except for one person (P77) based in India. As the majority of participants were in the United States, we aimed to gather geographic diversity within the United States; the pool includes people from 11 states.

The interview pool included 14 admins and 9 members across 19 unique instances. While most instances were running the Mastodon software, four instances used other software like Friendica. At least two more participants mentioned prior experiences with an instance that used non-Mastodon software (but not the focus of the interviews). Including these experiences expanded our analysis to examine a broader view on privacy on the Fediverse, beyond Mastodon-specific concerns.

Table 1 shows brief descriptors of the main topic/purpose of participants' instances. Instances were general (i.e., not focused on a particular topic but on socializing on the Fediverse broadly), personal (for oneself and one's friends), and/or for a specific professional, identity-based, or subculture group. The thematically-driven groups were focused on technology issues, academia, writing, queerness, and digital media and tech subcultures (open and independent web, online roleplaying). The use of one descriptor does not necessarily preclude the relevance of others. For example, the community of P90 was mostly a personal network of friends but also considers itself a space for people of color (predominantly a group of Southeast Asians) and queer people on the Fediverse.

## 4.3 Interview coding and qualitative analysis

Our analysis followed Braun and Clarke [15]'s thematic analysis approach. After the first round of interviews (11), the first author summarized the interviews to identify general themes and directions in the style of memos. Following discussion with all authors, we refined the questions in the interview protocol, to balance maintaining consistency in questions asked and addressing apparent gaps in the second round (12). We determined data saturation by a heuristic of redundancy in content, themes, and codes in our interview data that suggests replicability [25, 37].

The first author conducted line-by-line inductive coding of all transcripts after interviews were completed. This practice followed strategies recommended by Charmaz [20], in which initial codes

| P# | Gender | Age | Locale | Role | Instance software | Instance topic | Instance size |
|---|---|---|---|---|---|---|---|
| P43 | Male | 65+ | Western USA | Member | Mastodon | General | 1k-10k |
| P67 | Male | 35-44 | Eastern USA | Member | Mastodon | Writing | 1k-10k |
| P16 | Male | 25-34 | Western USA | Member | Mastodon | Tech | 10k+ |
| P22 | Male | 35-44 | Western USA | Member | Mastodon | Tech | 10k+ |
| P25 | Male | 18-24 | Midwest USA | Member | Mastodon | Academia | 501-1k |
| P77 | Male | 35-44 | India | Admin | Mastodon | Personal | 1-10 |
| P12 | Male | 35-44 | Canada | Admin | Friendica | Personal | 1-10 |
| P61 | Male | 35-44 | Eastern USA | Member | Mastodon | Academia | 501-1k |
| P08 | Male | 25-34 | Germany | Member | Mastodon | Academia | 1k-10k |
| P33 | Agender | 35-44 | Southern USA | Admin | Mastodon | Queer community | 1k-10k |
| P52 | Female | 25-34 | Midwest USA | Admin | GoToSocial | Role-play | 1-10 |
| P50 | Non-binary | 35-44 | Southern USA | Member | Mastodon | Tech | 1k-10k |
| P73 | Female | 25-34 | Western USA | Admin | Mastodon | Subculture | 51-100 |
| P11 | Non-binary | 25-34 | Western USA | Admin | Mastodon | Subculture | 11-50 |
| P45 | Female | 35-44 | Portugal | Admin | Mastodon | Tech | 1-10 |
| P17 | Female | 35-34 | Western USA | Admin | Mastodon | Personal | 51-100 |
| P51 | Genderfluid | 18-24 | Eastern USA | Admin | Sharkey | Personal | 1-10 |
| P90 | Male | 25-34 | Western USA | Admin | Sharkey | Personal | 11-50 |
| P84 | Female | 25-34 | Germany | Admin | Mastodon | Queer community | 51-100 |
| P47 | Male | 18-24 | Canada | Admin | Mastodon | General | 1k-10k |
| P70 | Female | 55-64 | Canada | Admin | Mastodon | Writing | 1k-10k |
| P31 | Female | 25-35 | Eastern USA | Admin | Mastodon | Academia | 501-1k |
| P54 | Non-binary | 55-64 | Western USA | Member | Mastodon | Tech | 1k-10k |

Table 1. Summary information about interview participants. Some details have been obfuscated to protect participant privacy. We replaced the instance name with a brief description of the topic.
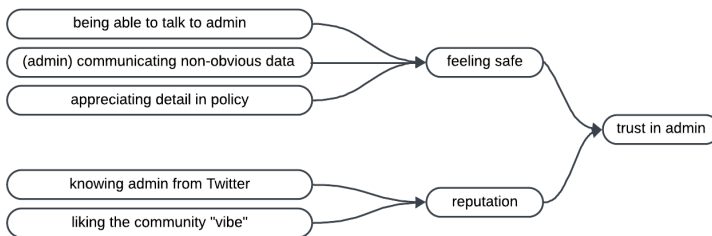


Fig. 2. Example of a sample of our inductive codes (left-most) being clustered into broader themes.

closely reflect what is being described in the data, such as actions or feelings of the interviewee, and then clustered into broader axial themes (see Figure 2). Through discussions with other authors

who had also reviewed interview transcripts, the first author generated memos — some based on themes identified in the first round of interviews — that were guided by our initial motivations to understand privacy concerns in communities, and how they are addressed or not through community governance practices. Thus, our attention leaned to the dimensions of community governance participants identified as relevant to privacy. We also focused on how people were conceptualizing and reasoning about privacy in the context of their instances and decentralized social media more broadly. We developed and refined a set of findings until reaching agreement.

## 5 ESTABLISHING INFORMATION FLOWS WITHIN A COMMUNITY

We first discuss how the governance of a community shapes the bounds and technical understanding of information flows, specifically through community rules and leadership.

### 5.1 Rules shape value-laden expectations around information flows

Like many online communities, instances on the Fediverse typically had a set of rules for participating in the community, such as those around harassment or content moderation. When asked about privacy risks or concerns, participants frequently referenced these instance rules.

Participants described reviewing rules as "*philosophies that [people] have written out or borrowed from*" (P11) to more deeply consider the shared social and political *values* they conveyed. Rules operated as proxies of what kind of interactions would be allowed. They also signified what needs and concerns were being prioritized, especially by stating what information flows would be cut off through moderation actions (such as removing harmful content, warning users about staying respectful, and blocking instances or individuals who harassed others). To this end, the quality and detail of rules were important. P54 recalled feeling so impressed by a community's rules that they decided to sign up for an account on it. The concreteness of the instance's rules (e.g., enumerating specific behaviors that would moderated as transphobic like "*no targeted dead-naming, no targeted misgendering, cis is not a slur*", P54) offered reassurance of what to expect:

> *Whether or not people have put in effort into writing server rules is an indication of whether or not this is something that the admins are seriously concerned about [...] If people are going down to that more specific level, that says to me they're probably more experienced with this kind of moderation and spend more time on it.* (P54)

Concrete, tailored rules signal care and experience in ensuring the safety of members, while vague rules could create space for bad faith interpretations. At an extreme, without rules as safeguards, bad actors could feel emboldened to "*pull up information that is not widely circulated [...] and share it freely [even if] just the amplification of that information could be dangerous [to someone]*" (P77).

Rules mattered because they offered value-laden *justifications* (or what P22 called "*legitimate use cases*") for information flows, which impacted users' perceptions of whether they found it acceptable for an instance to handle that data in particular ways. For example, participants perceived information flows such as the collection of IP addresses for moderation purposes to be relatively more acceptable than the collection of IP addresses without any stated reason, as they felt moderation was necessary if they wanted to have functional, working social media.

Rules also reflected a range of community goals. For example, participants in profession-oriented instances like those about tech-related careers (e.g., the instances of P16, P22, P50) and academia (e.g., the instances of P25, P61, P08, P31) noted that the point of joining the community was to increase the exposure of their community members to relevant people and content as widely as possible. Meanwhile, some communities aimed to heavily constrain the potential reach of their content, drawing stricter lines to proactively defend against potential harassment.

## 5.2 Informal and formal strategies for revising community members' technical understandings of information flows

Several participants felt that the Fediverse's underpinning technology — specifically, the ActivityPub protocol — did not prioritize privacy. Even if the software an instance used had different privacy options, interviews pointed to persistent misunderstandings around these options and how information flows on the Fediverse more broadly:

> *People make posts every now and again that [make me go:] Hey, listen, you know, all these different posting settings and all these different post privacy settings [that we have] are basically hacks on top of the ActivityPub protocol [...] So just be safe out there. People make posts [that reveal] quite a bit, especially when there's like big waves coming from Twitter or elsewhere.* — P73

The majority of our participants described themselves as having the knowledge and technical expertise to understand the Fediverse because of their professional and educational background, which allowed them to anticipate and identify violations of their privacy. As indicated by the quote above, they also felt that the average user would struggle to do the same, particularly newcomers — most of whom were coming from Twitter and unfamiliar with the decentralized, federated model of the Fediverse. To avoid inappropriate uses and misaligned expectations, P33 argued that the most important step in protecting privacy was therefore "*educating people, because [federation] is a different paradigm [...] from other social media we're used to.*" As the Fediverse lacked the centralized overhead that mainstream social media platforms had, the task of "educating" fell to the instance level via various informal and formal strategies.

*5.2.1 Informal strategies.* Participants noted that the Fediverse had a high volume of self-referential and meta-discussions. As, P54 jokingly remarked, the Fediverse was "*a great place to talk about decentralized social [media]!*" The popularity of talking *about* the Fediverse on the Fediverse meant that people could develop an awareness of common privacy concerns simply by being on the network, following specific hashtags and content from admins.

Admins in our interviews mentioned maintaining a separate account for posting updates and announcements to their community. For example, in the wake of corporate-owned social platforms (like Threads) attempting to add bridges to the ActivityPub protocol, P90 described posting an announcement to his community about his decision to preemptively defederate from these instances to prevent the community's content and information from being collected by corporate actors. Some admins described actively checking in with instance members about potential issues. For example, P73 reached out to users to explain some of the finer technical details about their instance (and the Fediverse in general) if users posted something that seemed rather sensitive for a general audience. Similarly, P70 described how her instance made extra efforts during newcomer onboarding:

> *We've certainly got a few users who have spoken to us before joining and said: I have concerns with my public account.* Can you support this? *And we're able to say:* Well, yes, we can — here's the settings you should use. *So we're a little more rigorous about onboarding and explaining to people how they might be vulnerable and what they can do about it.* — P70

P70 posited that such extra effort was not necessarily common, particularly for instances operating at large-scale like `mastodon.social` or with open registrations (enabling user account sign-ups that didn't require any approvals). Nevertheless, the admins in our interview pool (mostly with instances of under ten thousand users) all expressed willingness to engage with users' questions.

*5.2.2 Formal strategies.* As part of more formal ways of shaping technical understandings in the community, participants suggested that a community's privacy policy could be useful in describing

how user information got collected, shared, and stored on the instance and beyond. Functioning as a kind of information sheet, a community's privacy policy gave instances opportunities to flag potential privacy risks, such as the fact that direct messages could be viewed by admins or what country's laws the server was subject to. As a simple example, P33 recalled that they had specifically "*bold[ed] the line that if people get messages, they can make copies of them*" to warn people not to share sensitive information. P22 saw their privacy policy as making a plea to users to be careful before posting anything personal, especially the personal data of anyone else: "*This is just, like, your dad telling you to make good decisions, right?*"

A privacy policy could also be a message to admins and users of other communities that the instance was operating in good faith, demonstrating a degree of "*professionalism in running a service*" (P47) by outlining relevant processes, technical constraints, and local laws it would be beholden to. P17 described the policy as a written "*commitment in some way or another, that you aren't going to be unnecessarily endangering people who use the instance, whether it's you snooping around or doing something that allows others to snoop around.*"

However, privacy policies seemed to have limited impact because people do not read privacy policies [18, 79]. This pattern was evident in our interviews: admins in our interview pool noted that they rarely customized the privacy policy (at times, not realizing they had one), while members conceded that they did not read the policy beyond skimming it at sign-up. Overwhelmingly, participants described their privacy policies as boilerplate. The lack of engagement seemed to render privacy policies to exercises in compliance and legal "*hygiene*" (P54) — even as interviewees noted that they may help convey important knowledge.

### 5.3 Leadership builds trust that information flows will adhere to expectations

We explore how community governance created a sense of trust that information flows traversing the community will adhere to expectations. Specifically, our interviews emphasized the role of community leadership, tying trust in the community to trust in the admin.

*5.3.1 Taking risks in trusting an admin.* Interviews described the inherent risks in trusting an instance admin. Participants noted that admins had back-end access to all the information people inputted as users, from sign-up information to direct message history. This access meant admins needed "*to be very judicious with who they share information with*" (P12), including law enforcement, third-party service providers, and other instances. Participants like P17 also felt admins should be mindful of how and when admins themselves accessed back-end information.

P11 acknowledged that such access was important to what admins did to keep the server running: "*I trust them to be able to understand when it would be appropriate and wouldn't be inappropriate to review that information.*" However, some participants reported cases in which they questioned the admin's conduct. In one such case, P47 recalled that the admin of their instance had enrolled all users to a third-party marketing service without their explicit consent, to email out donation requests. P47 felt their trust had been betrayed, particularly because the marketing service's design had made the unsubscribe button difficult for them to find (using black text on a black background). Ultimately they left the instance to join a different one.

Another important consideration was whether an admin had sufficient technical expertise and knowledge to avoid security issues while maintaining the server. P16 appreciated that the admin of their instance "had been operating the server for quite a while" and was a "really technically competent person." He had seen how the admin was able to adjust the server, scale it, and keep it up to date. As admins were *de facto* fiduciaries of community members' data, good decision-making in the back-end "*that wouldn't compromise users too much*" (P50) was crucial. Several participants referenced an infamous incident in 2023 where the admin of a leftist instance had kept the instance's

data in a server in their home, unencrypted. When the United States federal government raided this admin's home, the server and all of the unencrypted messages between instance members on it were seized. Such incidents raised concerns among our participants that instances could serve as "honey pots" of users' data that could be breached or subpoenaed.

*5.3.2    Trusting the reputation and incentives of admins.* The majority of participants expressed high levels of trust in the admin of the instance they were currently on. We identified a few factors which cultivated trust in the admin.

Several participants noted that they knew the admin personally or through their network (or if they were an admin, that they knew their community members personally). P17, who was an instance owner, recalled that "*[people joined the instance] just from me telling my friends [...] and it has just sort of steadily gone from there.*" On the other side, P22 recalled seeing well-known users they liked from Twitter join the instance that he was considering:

> *Getting that kind of endorsement from [User A] and [B] and from everybody else saying that:* Yeah, we know this guy is great; he's super friendly, very level-headed, a very reasonable person. He's got a server that's been running since 2017 or something like that. *So there was also some history there to refer back to. And then taking a look [at the instance] for myself and being like, alright, this looks - this isn't a cesspit.*

As suggested by P22, participants went through the admin's public profile and scrolling through the recent posts on the instance as a way to gauge their trustworthiness further.

Another important source of trust came from having well-configured, informative documentation about the server including (1) a privacy policy (even if generic), (2) lists of who was on the admin team, and (3) community rules. The combination of these showed that one had "*covered all their bases*" (P47). For P84, having such information offered transparency: "*I trust my friends, and I give them access to the server. Because they are listed as moderators or administrators, it's transparent [to others] that those two or three people have access to the server.*

Some participants trusted their admin simply because the instance was not corporate-run and instead, voluntary and community-run. For example, P43 noted that although he had no idea who their admin was ("*I have no reason to trust that person, and I have no reason to distrust that person.*"), he felt the admin seemed reliable. When probed further, participants emphasized the lack of financial incentives behind instances as a key part of their preconceived trust, usually placed in contrast to corporate social media that people saw as "*using the data for their own good and not necessarily for the users*" (P08). P43 repeatedly compared his instance to Twitter, while other participants like P25 said that they trusted their instance precisely because it "*feels less corporate.*" As an admin of a larger, fairly active instance, P33 addressed questions of profit directly:

> *We do get donations for the instance. I started collecting donations last November because the cost [of running the server] has gotten tremendously high, [...] about $1,000 a month [...] It is fully paid for [by donations]. So I really appreciate the community for that. But there's not really — there's not a profit incentive here.*

As financial considerations were about covering costs rather than profit, P22 likened an instance to public radio: "*[...] no one is trying to wring every dollar out of it, and it's more that we need to get the lights on and pay the salaries.*" Moreover, the fairly small size of any given instance compared to centralized platforms suggested that an instance couldn't make a profit from selling data even if they wanted to. This became salient when participants described ongoing debates about whether instances should federate with emerging services developed by companies (primarily Threads) once they became compatible with the Fediverse, as these services could re-aggregate data in corporate hands. Finally, participants hypothesized that constraints on the admin bandwidth made it unlikely

that an admin would abuse access for personal reasons. P31, an admin herself, noted that although it was true admins could look at all data in their back-end, doing so was time-consuming and required manual searches: "*The majority of us just don't care about that enough.*" Given that other aspects of being an admin could require a great deal of volunteered time, admins emphasized a lack of interest and incentives to do extra work.

*5.3.3 Having direct access to the admin.* Participants appreciated that they could identify a specific admin for an instance and reach out to them. P84 explained why the ability to connect to specific admins mattered: "*It's not like we are some anonymous entities on this big website. It's more like, I can somehow relate to you. And, it's more like a friendship and not like a business communication.*" These positive interactions with an admin as an individual and fellow user gave people the sense that the admin would operate the instance in good faith, reliably, and responsively. For example, P50 believed that their admin would reach out to them if there was a situation that might require their information to be shared somewhere, such as responding to a moderation report. P50 also emphasized that this trust had been cultivated over time, having seen the specific admin be proactive in communicating with the community. In a similar vein, participants also felt that they could hold admins accountable if needed. P22 believed that the community on his instance — focused on technology issues — would make sure the admin would not put the community at risk:

> The guy who runs the server, we make sure that he's on top of it. I mean, he's on top of it anyway. But like, if he weren't, we would notice, and it would be — you know, we deal with it. [...] The right people are paying attention, I suppose.

## 6 GOVERNANCE FRICTIONS BETWEEN COMMUNITIES

As communities made governance decisions, they naturally differed in how they sought to shape information flows. Differences could lead to *governance frictions* between communities: incompatibilities between communities that, on one hand, could enforce the goals of one community but on the other hand, potentially undermine the goals of the other. In the context of privacy concerns, governance frictions could pose new risks to privacy in a community or limit the effectiveness of a community's strategies to protect the privacy of its users. We identified three main incompatibilities between instances that had consequences for privacy: value, security, and software.

### 6.1 Value incompatibilities

The majority of the incompatibilities mentioned by the participants were due to differences in "*communication cultures*", as referred to by P08, that impacted what would be acceptable content and behavior. These differences arose because communities had "*different interpretations of [rules] and different levels of granularity of how far they'll go in enforcing different bits of [them]*" (P73).

Participants generally viewed normative differences as a positive feature, not a bug, of the Fediverse. As P73 put it: "*I feel like it really makes it so everyone can find their place [...] on the Fediverse [...] that's one of the things that's really beautiful about the Fediverse to me.*" Similarly, P54 felt that the diversity of norms across the Fediverse made it much more interesting and vibrant: "*The great thing about the Fediverse is the combination of these local communities at the instance level within this broader conversation. Instances have different norms, instances have different vibes.*"

However, at times the values behind these norms differed enough that people did not want to interact with parts of the Fediverse. For example, P73 noted that some instances did not want to see nudity or other NSFW content because it clashed with their goals for being on the Fediverse (such as professional branding and posting). Because values were social and political, the stakes could be high. P45 observed that the queer members of their community had concerns about being outed (as part of the LGBTQIA+ community) or "*being harassed or doxxed or whatever [for their*

*views],*" as queer people are disproportionate targets of online harassment [9]. Meanwhile, some instances embraced "*free speech absoluti[sm]*" (P73) and allowed abusive content in the name of free speech rights. These differences raised questions of how federation could potentially put one's instance members' reputations at risk as well as put them in harm's way of abusive behavior. As P45 stressed: "*The main thing that I think about when it comes to privacy concerns [on social media] is that you don't want to share stuff that makes you feel vulnerable.*"

Defederating (and blocking more generally) enabled communities to address feelings of vulnerability by cutting information flows between themselves and value-incompatible instances. Participants perceived blocking access to one's information (profile information, posts, online activity) or blocking "*people who are harassing you [...] or attempting to surveil you*" (P52) to be "*intimately linked with privacy*" (P52). Several participants who were admins described making decisions to defederate to protect their members from discriminatory behavior and content, as well as morally reprehensible content such as child sexual abuse materials (which is also illegal, putting instances at legal risk as well). P33 explained that defederating was a swift decision when it came to value incompatibilities:

> *More often than not, I'm really referring to abusive behavior. You know, it's not, it's not illegal for somebody to use curse words at me and call me slurs online. But, you know, I don't want to see it. And I don't want our community to be involved in that.*

However, echoing other participants, P16 noted that blocking individuals or defederating from an instance has limited effectiveness in governing the flow of information online:

> *If I block them, you know, maybe they can't follow me with their regular account, but they can make another account where they can do whatever. If something's on my profile, I know it's not going to be something that I can hide from anybody on the internet — like, it's out there. So, you know, I prefer that if I can block somebody, they just suddenly couldn't see my stuff at all. But I know, that's not reality.*

In addition, in a variant of "doxxing," people could maliciously expose a person to another instance or person one had blocked. P73 and P47 recalled seeing instances that *intentionally* operated bots that would find out who blocked whom, and tag both users in the same post to force them back into contact with each other. In other cases, this exposure could be accidental, like if a post got shared to a different instance that didn't block the same instances as the original poster's instance. In short, being federated to many communities meant having many potential points of leakage.

To further reduce the risk of these kinds of privacy violations, some instances opted to operate in "allow-list" mode in which the server would only federate with approved instances (instead of the more widely-used practice, where servers defederate from instances over time) so that "*random people couldn't see their [content]*" (P52). However, participants noted that the "allow-list" mode is a drastic and somewhat controversial measure that could be at odds with the information flows people *did* expect on social media.

Ultimately, the interview participants did not demonstrate consensus about how to mitigate privacy risks posed to users due to the value incompatibilities. However, across our interviews, participants stressed the importance of keeping up to date with meta-discussion about evolving safety issues across the Fediverse. As P12 put it: if one wanted to properly prevent information from spreading to certain spaces, there needed to be "*a lot of intentional chatting.*"

## 6.2 Security incompatibilities

Participants described potential issues around instance security, which spanned all aspects of how an instance might choose to govern access to the instance (e.g., registration, federation) and its

back-end (e.g., admin access to logs and databases, encryption of stored data, credit card info for donations to support server costs).

Participants expressed concerns that an instance's poorly-devised security processes could undermine the security of an instance's operations. For P67, the fact that some instances could simply use "*informal and social arrangements*" could be "*anxiety-inducing.*" For example, having open registrations could let bad actors and bots into an instance, inadvertently impacting people on instances federated with the compromised instances. Most admins with whom we spoke had semi-closed registrations on their instances and often saw open registrations as a "*red flag*" (P17). As P51 suggested, poor security configurations of an instance created a weak point that could affect federated instances even with relatively strong security settings:

> Even if the admin on my instance — AKA, me — tries to be very security conscious and all, if the fine admin of some other instance were to make a security oopsie, a gaffe, a failure... Then that is going to affect anyone [...] who has had dealings with them. But beyond that, I don't think there's really much of a threat model that other servers pose [...] The only legitimate added threat is if private user information on my instance ends up on another instance through any number of reasons, and if it's not secure over there, that poses a risk to the users of my instance, not to that instance alone.

As with value incompatibilities, participants shared that they mainly use the "blocking" feature on compromised instances to avoid potential security breaches. While crude in the granularity at which it could be applied, blocking enabled participants to draw "*boundary lines with the rest of the Fediverse*" (P70). However, participants also noted that instances varied in what they considered a security threat and what configurations would be sufficient to preserve user privacy, even if they shared the same overarching community values. For example, participants pointed to the aforementioned debate about federating with corporate-run instances such as Threads. While some like P90 felt federating with Threads would be a major threat by giving data to corporate actors, others did not. They argued that one's posts "*could still get sucked up [...] if your posts end up on a server that is friendly with Threads*" (P17). As P70 observed, dealing with security differences was important but required a great deal of work as an ongoing negotiation, "*even [between] adjacent administrators who are maybe on the same wavelength.*"

## 6.3 Software incompatibilities

Instances used software for running services on the Fediverse with different technical affordances for privacy to their users. P12, for example, recalled how he and a fellow Fediverse user had realized that the technical implementation of marking messages as "private" differed between their instances because their instances used different software (Mastodon and Friendica):

> I was following [a Mastodon user] and I boosted a post of theirs. And then they sort of screamed at me because I was boosting a private post. And I was like: Well, then mark it private! And they were like: But I did! And then we figured it out. [...] That's what I mean when I talk about it being impossible to control at that point. Like, if I reply to something [marked private on Mastodon], that post would be shown to other people who use Friendica, because that post isn't marked private [on Friendica]

By itself, the protocol underpinning the Fediverse does not define the notion of messages marked "private." Instead, developers of softwares like Mastodon or Friendica can build additional functionality. As P73 noted, marking a message private was a soft "*hack*" (P73) in the Mastodon software alone that had no meaning to the Friendica software. Moreover, these softwares were usually updated over time, which meant instances needed to keep posted on software updates; in a large federated environment, it is difficult to guarantee that all instances are using the up-to-date version of the

software. Importantly, an instance's software wasn't compelled to respect privacy options used by another instance: "*[The software] isn't really a wall. It's a stop sign, you know?*" (P51). Differences in implementations and versions of softwares revealed, as P51 continued, "*an illusion of privacy*" with "*no actual privacy at the heart of [the protocol].*"

The technical differences across instance software described above could generate conflict between users who believed the other had carelessly violated their privacy. Nevertheless, in our interviews, relatively few participants raised this issue because most people were on an instance using the Mastodon software and primarily interacted with people on instances using the same software. Some, like P84, even stated that they wanted to see more softwares on the Fediverse because they liked "*that [they] can interact with all kinds of different software [...] implementations that rely on ActivityPub, and that software works together instead of being always its own silo.*"

The rise of Mastodon as the dominant software of the Fediverse puts it under larger scrutiny as to whether it offers adequate privacy options to users across the many instances running it. P84 felt that moderation tools on Mastodon were insufficient, while P51 observed a sense of *ad hoc* standardization on the ActivityPub protocol driven primarily by wanting to federate with Mastodon instances (in particular, the flagship instance with nearly a million users, `mastodon.social`) rather than by better security or privacy practices. P90 highlighted a recent bug found on Mastodon, noting that instances using different software (including his own) had avoided the problem:

> It was basically something like, you could send a fake post and it would overwrite the real one or someone's profile. That's a pretty serious vulnerability. You could impersonate anybody, and apparently it was really easy to do. But Sharkey was not vulnerable. So we didn't have that vulnerability even to begin with.

Where technical differences in softwares could pose a challenge to privacy in earlier examples, P90's story suggested that technical differences could also preserve user privacy by limiting the impact of vulnerabilities from one particular software. In sum, software incompatibilities could lead to significant privacy-related harms but participants noted that there was ongoing debate among Fediverse users about how to resolve these issues. Participants suggested that improving the core privacy-functionalities of the ActivityPub protocol instead of through software like Mastodon or Friendica could provide a consistent way forward for privacy-preserving federation.

## 7 DISCUSSION

### 7.1 The unique challenges of decentralization in social media communities

Echoing work in other social media contexts [e.g., 32, 35, 88], our interviews emphasized safety as an important consideration when governing privacy at the community level, particularly in discussing community rules (§5.1). However, the varying goals of communities meant that they had different ideas of what constituted safety. In traditional online communities, normative difference is expected [see 19, 33] and can even be an advantage [101]. On community-based platforms like Reddit, prior work notes that despite such difference, inter-community conflict appears to be relatively rare [59]. However, on the Fediverse, communities are *federated*: creating and sharing content with one another, inducing information flows that may or may not be appropriate for each of their social contexts. As evidenced by the governance frictions in §6, the increase in potential cross-community information flows produces a unique challenge for decentralized social media in defining and enforcing safety in privacy decisions.

One admin noted they were highly selective about other communities they federated with, and a few suggested that this strategy could ensure a sense of safety. However, participants also noted that this approach was in tension with the general purpose of the Fediverse: to be a social media network, possibly one that could be as connected and influential as centralized social media platforms. Our

participants represented communities with a diverse range of goals. As such, some communities saw benefits from broader networks (e.g., professional, research, hobbies) while others wanted to limit reach and access (e.g., subcultures, queer communities, personal groups). Hence, future work will need to identify ways to implement safety in accordance with the different breadths of networked reach that communities want, whether by tools, policies, or new protocols.

Users' past experiences on centralized social media also exacerbated another problem on the Fediverse: newcomer onboarding. As noted in §5.2, many communities found that new members were coming from Twitter with a strong but misaligned mental model of how social media "works." Admin efforts were thus drawn to correcting misunderstandings about the technical inner workings of the Fediverse. As noted in our description of security and software incompatibilities, correcting misunderstandings was especially important because technical differences between communities could produce disjunctions in features intended to protect privacy. Newcomer onboarding is a classic problem in the development of any social computing system [17, 57, 73]. Focused on privacy as the guiding heuristic, our work underscores that onboarding practices in decentralized social media must pay further attention to ensuring newcomers understand the technical and infrastructural aspects of each community: where the server is hosted, who has backend access to data, how information is shared between servers, and so on. Future work might draw on the rich body of prior research building interactive tutorials [76], testing new social spaces for newcomers [73, 113], and strategically placing rules [68] to strengthen and improve the onboarding process.

## 7.2 Designing to mitigate governance frictions

We identified three kinds of incompatibilities that produce governance frictions for privacy in decentralized social media. Designing tools and interfaces that help communities better respond to and anticipate governance frictions may help improve users' experiences in the Fediverse, especially given how these frictions are how problems emerge. Participants noted that the main tool available to them was a variety of blocklists, which enabled them to defederate from communities *en masse*. Blocklists, while an important option to have, are primarily about value incompatibilities that become salient in content moderation practices and rules. As a result, security and software incompatibilities remain mostly unaddressed. As each incompatibility will require different tooling and design to resolve, future work should avoid conflating the three.

Our interviews emphasized a lack of consensus about how to deal with incompatibilities alongside a pressing need to resolve them. In §5, our findings show how communities handled potential confusion about privacy (including issues that emerge due to incompatibilities) *within* the community through consistent communication, such as an admin account posting updates. We suggest that, likewise, designing opportunities for signals, communication, and negotiation *between* communities can help address the kinds of incompatibilities leading to governance frictions shown in our work. Below we briefly provide an example of a potential approach per incompatibility:

- **Value incompatibilities**: While participants look to written rules to get a sense of a community's values, differences in rules do not always equate value incompatibilities. However, recent work shows that communities often share or have similar rules that they copy from reliable sources to signal legitimacy [55]. Thus when differences are significant, warning flags during moments of federation might prompt admins or members to review or open dialogue with another community. This might be done through natural language processing (NLP) methods to evaluate similarity or contradiction between rule sets.
- **Security incompatibilities**: Communities might develop "standards" for levels of security an instance implements. This could be communicated clearly with badges that indicates a certain set of criteria related to the instance's security are met. Badges can enable other

instances to make quicker and more effective decisions about whether an instance's standards are up to par. How these standards are developed is an open question and suggests it may be good to have a mechanism for admins or representatives from different instances to convene, discuss, and later change Fediverse-wide decisions as a collective.

- **Software incompatibilities**: Our interviews noted that the software used to run instances on the Fediverse were inconsistent in how they treated posts as "private" (i.e., which audience the post will be shown to). We suggest developing concise abstractions that convey such incompatibilities. For example, widgets might automatically show, based on an instance's software, where a certain post will travel to. Given that this should be consistent to be useful, focusing on new features or capabilities at the protocol-level will likely be a helpful first step. As with security practices, how this idea might be developed as a standard in the protocol is an open question and will require coordination and input from diverse communities.

These suggestions emphasize encouraging interaction between communities. In many ways, this contrasts with defederation, which is the main mode of dealing with conflicts or incompatibilities on the Fediverse. While defederation is a crucial feature for addressing online harm [24], our interviews indicated that there were many simple workarounds that could be leveraged by malicious actors. Meanwhile, security and software incompatibilities less frequently warrant defederation but still must be resolved. Strategies that complement defederation as an option available to communities are crucial in enabling communities to strengthen privacy on decentralized social media systems like the Fediverse. Our suggestions above are just a few potential paths moving forward.

To this end, because incompatibilities are fundamentally about information flows, returning to the framework of contextual integrity can prompt normative evaluation to generate recommendations and design ideas. Choksi et al. [21] identifies social roles in online groups and places them in information flows as actors to analyze privacy challenges within groups. A similar exercise evaluating information flows across groups could be useful to guide design in surfacing incompatibilities or reasoning why incompatibilities raise privacy risks. For example, developing badges for security incompatibilities helps communities better understand the transmission principles of information flows between two actors (communities), one of which may need much stricter security practices that make another's weaker practices norm-violating.

## 7.3   Toward more participatory decision-making in decentralized social media

Whether or not a user could trust their community admin was a major theme in our findings, as admins held executive power over much of how a community (i.e., instance) worked. This arrangement echoes the "feudalistic" forms of governance commonly found in many online communities [89], where only a few individuals have the power to make decisions and function as "benevolent dictators." In *Governable Spaces*, Schneider [90] critiques this pattern as cultivating undemocratic digital spaces. Recent work in CSCW has attempted to enable more participatory forms of decision-making [60, 91]. Our findings suggested that people trusted their admins would operate in good faith, while also recognizing the inherent risk in trusting a person they might not know. In §5.3, we describe how trust was driven by how accessible an admin was and how proactive they were. This suggests that even in the absence of formal or explicit participatory structures, communities can avoid *de facto* feudalism. For example, participants of one of the communities we spoke to noted that the average tech expertise of community members meant that community members were quick to call out privacy discrepancies. In other communities, admins were beholden to social or reputation costs (e.g., a server for personal friends), which motivated them to engage in participatory decision-making informally. To move toward participatory decision-making, future work must systematically evaluate the conditions under which having an admin provides an

advantage for community outcomes — and what mechanisms help safeguard a community from a "dictator." In the context of privacy, articulating these conditions can help us understand when and where more formal participatory structures of governance would strengthen the maintenance of privacy in community-based approaches to governing social media. Future work investigating how participatory decision-making subsequently shapes how issues like governance frictions are resolved will also be crucial in improving the design of decentralized systems overall.

### 7.4 Limitations

We highlight study limitations that would benefit from further consideration in future work. Participants emphasized how membership in marginalized identity groups shifted the privacy risks people faced, especially due to racism and queerphobia. While many of our participants were part of the LGBTQIA+ community, we did not systematically collect information about race or ethnicity. Given the importance of safety in our findings, future work should explicitly center the perspectives of marginalized communities toward developing "equitable and privacy-protective technologies" [86]. We also note that all but one of our participants were located in the Global North, but the Fediverse has seen growing popularity in countries like Brazil. Different legal paradigms and political contexts impact privacy concerns as well as shape community practices. Future work on non-United States and non-European contexts in particular will offer valuable insights.

Our interview protocol involved discussing scenarios with CI-based vignettes (as a survey) based on a privacy policy commonly seen in a sample of Fediverse instances. However, privacy policies often omit contextual parameters that impact how people judge information flows to be appropriate [96], which may lead people to rely on preconceived (and inaccurate) mental models of privacy [65]. We used the survey as an artifact for open discussion so the interviewer could clarify interpretations and ask follow-up questions. As we found the CI-based vignettes useful for evoking conversation about information flows on social media, future work might design and leverage CI vignettes using other guidelines that might make other privacy issues (concerns of marginalized communities, different legal paradigms or geopolitical contexts, etc.) more salient in analysis.

## 8 CONCLUSION

Grounded in the theory of contextual integrity, this study describes how community governance enabled participants to engage in an active negotiation of privacy expectations by shaping the bounds, trustworthiness, and technical understandings of information flows. Our findings suggest that communities develop shared expectations around information flows through concrete rules and proactive leadership. Together, these practices seem to foster trust among users that privacy expectations will be met. As such, community governance can offer valuable opportunities to re-imagine privacy on social media, particularly toward prioritizing values like safety. These opportunities do not come without hurdles. *Governance frictions* could generate tensions due to incompatibilities between communities' values, security concerns, and software, which can give rise to novel risks and privacy challenges. Future work must take stock of the unique challenges decentralization brings and work toward stronger community-based governance mechanisms that support processes of negotiation around how information flows move through communities.

## REFERENCES

[1] Tawfiq Ammari, Sarita Schoenebeck, and Daniel Romero. 2019. Self-Declared Throwaway Accounts on Reddit: How Platform Affordances and Shared Norms Enable Parenting Disclosure and Support. *Proceedings of the ACM on Human-Computer Interaction* 3, CSCW (Nov. 2019), 135:1–135:30. https://doi.org/10.1145/3359237

[2] Anonymous. 1998. To Reveal or Not to Reveal: A Theoretical Model of Anonymous Communication. *Communication Theory* 8, 4 (1998), 381–407. https://doi.org/10.1111/j.1468-2885.1998.tb00226.x

[3] Hans Asenbaum. 2018. Anonymity and Democracy: Absence as Presence in the Public Sphere. *American Political Science Review* 112, 03 (Aug. 2018), 459–472. https://doi.org/10.1017/S0003055418000163

[4] Karla Badillo-Urquiola, Xinru Page, and Pamela J. Wisniewski. 2019. Risk vs. Restriction: The Tension between Providing a Sense of Normalcy and Keeping Foster Teens Safe Online. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems (CHI '19)*. Association for Computing Machinery, New York, NY, USA, 1–14. https://doi.org/10.1145/3290605.3300497

[5] Louise Barkhuus. 2012. The Mismeasurement of Privacy: Using Contextual Integrity to Reconsider Privacy in HCI. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems (CHI '12)*. Association for Computing Machinery, New York, NY, USA, 367–376. https://doi.org/10.1145/2207676.2207727

[6] Solon Barocas and Helen Nissenbaum. 2014. Big Data's End Run around Anonymity and Consent. In *Privacy, Big Data, and the Public Good*, Julia Lane, Victoria Stodden, Stefan Bender, and Helen Nissenbaum (Eds.). Cambridge University Press, New York NY.

[7] Michael S. Bernstein, Eytan Bakshy, Moira Burke, and Brian Karrer. 2013. Quantifying the Invisible Audience in Social Networks. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems (CHI '13)*. Association for Computing Machinery, New York, NY, USA, 21–30. https://doi.org/10.1145/2470654.2470658

[8] Jeremy Birnholtz, Nicholas Aaron Ross Merola, and Arindam Paul. 2015. "Is It Weird to Still Be a Virgin": Anonymous, Locally Targeted Questions on Facebook Confession Boards. In *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems*. ACM, Seoul Republic of Korea, 2613–2622. https://doi.org/10.1145/2702123.2702410

[9] Lindsay Blackwell, Jill Dimond, Sarita Schoenebeck, and Cliff Lampe. 2017. Classification and Its Consequences for Online Harassment: Design Insights from Heartmob. *Proceedings of the ACM on Human-Computer Interaction* 1, CSCW (Dec. 2017), 1–19. https://doi.org/10.1145/3134659

[10] Lindsay Blackwell, Jean Hardy, Tawfiq Ammari, Tiffany Veinot, Cliff Lampe, and Sarita Schoenebeck. 2016. LGBT Parents and Social Media: Advocacy, Privacy, and Disclosure during Shifting Social Movements. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems (CHI '16)*. Association for Computing Machinery, New York, NY, USA, 610–622. https://doi.org/10.1145/2858036.2858342

[11] Anita L. Blanchard and M. Lynne Markus. 2004. The Experienced "Sense" of a Virtual Community: Characteristics and Processes. *ACM SIGMIS Database: the DATABASE for Advances in Information Systems* 35, 1 (Feb. 2004), 64–79. https://doi.org/10.1145/968464.968470

[12] Hannah Bloch-Wehba. 2021. Content Moderation as Surveillance.

[13] Christoph Bösch, Benjamin Erb, Frank Kargl, Henning Kopp, and Stefan Pfattheicher. 2016. Tales from the Dark Side: Privacy Dark Strategies and Privacy Dark Patterns. *Proceedings on Privacy Enhancing Technologies* 2016, 4 (2016), 237–254. https://doi.org/10.1515/popets-2016-0038

[14] Anne Bowser, Katie Shilton, Jenny Preece, and Elizabeth Warrick. 2017. Accounting for Privacy in Citizen Science: Ethical Research in a Context of Openness. In *Proceedings of the 2017 ACM Conference on Computer Supported Cooperative Work and Social Computing (CSCW '17)*. Association for Computing Machinery, New York, NY, USA, 2124–2136. https://doi.org/10.1145/2998181.2998305

[15] Virginia Braun and Victoria Clarke. 2006. Using Thematic Analysis in Psychology. *Qualitative Research in Psychology* 3, 2 (Jan. 2006), 77–101. https://doi.org/10.1191/1478088706qp063oa

[16] Amy Bruckman. 2006. A New Perspective on "Community" and Its Implications for Computer-Mediated Communication Systems. In *CHI '06 Extended Abstracts on Human Factors in Computing Systems*. ACM, Montréal, Québec, Canada, 616–621. https://doi.org/10.1145/1125451.1125579

[17] Moira Burke, Cameron Marlow, and Thomas Lento. 2009. Feed Me: Motivating Newcomer Contribution in Social Network Sites. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems (CHI '09)*. ACM, New York, NY, USA, 945–954. https://doi.org/10.1145/1518701.1518847

[18] Fred H. Cate. 2010. The Limits of Notice and Choice. *IEEE Security & Privacy* 8, 2 (March 2010), 59–62. https://doi.org/10.1109/MSP.2010.84

[19] 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. *Proc. ACM Hum.-Comput. Interact.* 2, CSCW (2018), 32:1–32:25. https://doi.org/10.1145/3274301

[20] Kathy Charmaz. 2015. *Constructing Grounded Theory: A Practical Guide through Qualitative Analysis* (2nd ed.). SAGE, Thousand Oaks, California.

[21] Madiha Zahrah Choksi, Ero Balso, Frauke Kreuter, and Helen Nissenbaum. 2024. Privacy for Groups Online: Context Matters. *Proc. ACM Hum.-Comput. Interact.* 8, CSCW2 (Nov. 2024), 406:1–406:23. https://doi.org/10.1145/3686945

[22] Camille Cobb. 2019. *User-to-User Privacy in Social and Communications Applications.* Thesis. University of Washington.

[23] Cindy Cohn and Rory Mir. 2022. The Fediverse Could Be Awesome (If We Don't Screw It up).

[24] Carl Colglazier, Nathan TeBlunthuis, and Aaron Shaw. 2024. The Effects of Group Sanctions on Participation and Toxicity: Quasi-experimental Evidence from the Fediverse. *Proceedings of the International AAAI Conference on Web and Social Media* 18 (May 2024), 315–328. https://doi.org/10.1609/icwsm.v18i1.31316

[25] Juliet Corbin and Anselm Strauss. 2015. *Basics of Qualitative Research: Techniques and Procedures for Developing Grounded Theory* (fourth edition ed.). SAGE Publications, Inc, Los Angeles, CA.

[26] Sauvik Das, W. Keith Edwards, DeBrae Kennedy-Mayo, Peter Swire, and Yuxi Wu. 2021. Privacy for the People? Exploring Collective Action as a Mechanism to Shift Power to Consumers in End-User Privacy. *IEEE Security & Privacy* 19, 5 (Sept. 2021), 66–70. https://doi.org/10.1109/MSEC.2021.3093135

[27] Munmun De Choudhury and Sushovan De. 2014. Mental Health Discourse on Reddit: Self-Disclosure, Social Support, and Anonymity. *Proceedings of the International AAAI Conference on Web and Social Media* 8, 1 (May 2014), 71–80.

[28] Ralf De Wolf, Koen Willaert, and Jo Pierson. 2014. Managing Privacy Boundaries Together: Exploring Individual and Group Privacy Management Strategies in Facebook. *Computers in Human Behavior* 35 (June 2014), 444–454. https://doi.org/10.1016/j.chb.2014.03.010

[29] Michael A. DeVito, Ashley Marie Walker, and Jeremy Birnholtz. 2018. 'Too Gay for Facebook': Presenting Lgbtq+ Identity throughout the Personal Social Media Ecosystem. *Proceedings of the ACM on Human-Computer Interaction* 2, CSCW (Nov. 2018), 44:1–44:23. https://doi.org/10.1145/3274313

[30] Bryan Dosono and Bryan Semaan. 2020. Decolonizing Tactics as Collective Resilience: Identity Work of AAPI Communities on Reddit. *Proceedings of the ACM on Human-Computer Interaction* 4, CSCW1 (May 2020), 1–20. https://doi.org/10.1145/3392881

[31] Evelyn Douek. 2022. Content Moderation as Systems Thinking. https://doi.org/10.2139/ssrn.4005326

[32] Brianna Dym and Casey Fiesler. 2018. Vulnerable and Online: Fandom's Case for Stronger Privacy Norms and Tools. In *Companion of the 2018 ACM Conference on Computer Supported Cooperative Work and Social Computing (CSCW '18 Companion).* Association for Computing Machinery, New York, NY, USA, 329–332. https://doi.org/10.1145/3272973.3274089

[33] Casey Fiesler, Jialun" Aaron" Jiang, Joshua McCann, Kyle Frye, and Jed R. Brubaker. 2018. Reddit Rules! Characterizing an Ecosystem of Governance.. In *Proceedings of the International AAAI Conference on Web and Social Media*, Vol. 12. AAAI, Stanford, CA, 72–81. https://doi.org/10.1609/icwsm.v12i1.15033

[34] Jeremy Foote and Sohyeon Hwang. 2023. Online Communities and Big Data. In *Group Communication: An Advanced Introduction*, Torsten Reimer, Ernest S Park, and Joseph A. Bonito (Eds.). Routledge, New York, NY.

[35] Andrea Forte, Nazanin Andalibi, and Rachel Greenstadt. 2017. Privacy, Anonymity, and Perceived Risk in Open Collaboration: A Study of Tor Users and Wikipedians. In *Proceedings of the 2017 ACM Conference on Computer Supported Cooperative Work and Social Computing.* ACM, New York, NY, 1800–1811. https://doi.org/10.1145/2998181.2998273

[36] Seth Frey, Maarten W. Bos, and Robert W. Sumner. 2017. Can You Moderate an Unreadable Message? 'Blind' Content Moderation via Human Computation. *Human Computation* 4, 1 (July 2017), 78–106. https://doi.org/10.15346/hc.v4i1.5

[37] Patricia I. Fusch and Lawrence R. Ness. 2015. Are We There yet? Data Saturation in Qualitative Research. *The Qualitative Report* 20, 9 (Feb. 2015), 1408–1416. https://doi.org/10.46743/2160-3715/2015.2281

[38] Robert W. Gehl. 2023. ActivityPub, the Non-Standard Standard.

[39] Robert W. Gehl and Diana Zulli. 2023. The Digital Covenant: Non-Centralized Platform Governance on the Mastodon Social Network. *Information, Communication & Society* 26, 16 (Dec. 2023), 3275–3291. https://doi.org/10.1080/1369118X.2022.2147400

[40] R. Stuart Geiger. 2016. Bot-Based Collective Blocklists in Twitter: The Counterpublic Moderation of Harassment in a Networked Public Space. *Information, Communication & Society* 19, 6 (June 2016), 787–803. https://doi.org/10.1080/1369118X.2016.1153700

[41] Sebastiaan Gorissen and Robert W. Gehl. 2023. When Wikipedia Met Tor: Trials of Legitimacy at a Key Moment in Internet History. *Internet Histories* 7, 2 (April 2023), 105–121. https://doi.org/10.1080/24701475.2021.2015967

[42] Kishonna L. Gray. 2012. Intersecting Oppressions and Online Communities. *Information, Communication & Society* 15, 3 (April 2012), 411–428. https://doi.org/10.1080/1369118X.2011.642401

[43] Tamy Guberek, Allison McDonald, Sylvia Simioni, Abraham H. Mhaidli, Kentaro Toyama, and Florian Schaub. 2018. Keeping a Low Profile? Technology, Risk and Privacy among Undocumented Immigrants. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems (CHI '18)*. Association for Computing Machinery, New York, NY, USA, 1–15. https://doi.org/10.1145/3173574.3173688

[44] Oliver L. Haimson, Justin Buss, Zu Weinger, Denny L. Starks, Dykee Gorrell, and Briar Sweetbriar Baron. 2020. Trans Time: Safety, Privacy, and Content Warnings on a Transgender-Specific Social Media Site. *Proceedings of the ACM on Human-Computer Interaction* 4, CSCW2 (Oct. 2020), 1–27. https://doi.org/10.1145/3415195

[45] Joanne Hinds, Emma J. Williams, and Adam N. Joinson. 2020. "It Wouldn't Happen to Me": Privacy Concerns and Perspectives Following the Cambridge Analytica Scandal. *International Journal of Human-Computer Studies* 143 (Nov. 2020), 102498. https://doi.org/10.1016/j.ijhcs.2020.102498

[46] Evey Jiaxin Huang, Abhraneel Sarma, Sohyeon Hwang, Eshwar Chandrasekharan, and Stevie Chancellor. 2024. Opportunities, Tensions, and Challenges in Computational Approaches to Addressing Online Harassment. In *Proceedings of the 2024 ACM Designing Interactive Systems Conference (DIS '24)*. Association for Computing Machinery, New York, NY, USA, 1483–1498. https://doi.org/10.1145/3643834.3661623

[47] Hsiao Ying Huang and Masooda Bashir. 2016. The Onion Router: Understanding a Privacy Enhancing Technology Community. In *Proceedings of the 79th ASIS&T Annual Meeting: Creating Knowledge, Enhancing Lives Through Information & Technology (ASIST '16)*. American Society for Information Science, Silver Springs, MD, USA, 34:1–34:5.

[48] Sohyeon Hwang, Charles Kiene, Serene Ong, and Aaron Shaw. 2024. Adopting Third-Party Bots for Managing Online Communities. *Proceedings of the ACM on Human-Computer Interaction: Computer Supported Cooperative Work* 8, CSCW1 (April 2024), 216:1–216:26. https://doi.org/10.1145/3653707

[49] Jane Im, Ruiyi Wang, Weikun Lyu, Nick Cook, Hana Habib, Lorrie Faith Cranor, Nikola Banovic, and Florian Schaub. 2023. Less Is Not More: Improving Findability and Actionability of Privacy Controls for Online Behavioral Advertising. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems (CHI '23)*. Association for Computing Machinery, New York, NY, USA, 1–33. https://doi.org/10.1145/3544548.3580773

[50] L. Irani. 2023. Algorithms of Suspicion: Authentication and Distrust on the Amazon Mechanical Turk Platform. https://doi.org/10.2139/ssrn.4482508

[51] Shagun Jhaver, Seth Frey, and Amy X. Zhang. 2023. Decentralizing Platform Power: A Design Space of Multi-Level Governance in Online Social Platforms. *Social Media + Society* 9, 4 (Oct. 2023), 20563051231207857. https://doi.org/10.1177/20563051231207857

[52] Ruogu Kang, Stephanie Brown, and Sara Kiesler. 2013. Why Do People Seek Anonymity on the Internet? Informing Policy and Design. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems (CHI '13)*. Association for Computing Machinery, New York, NY, USA, 2657–2666. https://doi.org/10.1145/2470654.2481368

[53] Tero Karppi, Britt Paris, Robert W. Gehl, Corinne Cath, and Sarah Myers West. 2023. IF NOT, ELSE: STANDARDS, PROTOCOLS, NETWORKS AND HOW THEY MAKE A DIFFERENCE. *AoIR Selected Papers of Internet Research* (Dec. 2023). https://doi.org/10.5210/spir.v2023i0.13525

[54] Ra'il I'Nasah Kiam. 2023. Blackness in the Fediverse: A Conversation with Marcia X. *Logic(s) Magazine* 20 (Dec. 2023).

[55] Charles Kiene. 2024. *Identity, Legitimacy, and Voice: Understanding Rule Adoption, Compliance, and Evolution in Online Communities*. Ph. D. Dissertation. University of Washington, Seattle WA USA.

[56] Sara E. Kiesler, Robert E. Kraut, Paul Resnick, and Aniket Kittur. 2012. Regulating Behavior in Online Communities. In *Building Successful Online Communities: Evidence-Based Social Design*, Robert E. Kraut and Paul Resnick (Eds.). MIT Press, Cambridge, MA.

[57] Robert E. Kraut and Paul Resnick. 2012. *Building Successful Online Communities: Evidence-based Social Design*. MIT Press, Cambridge, MA.

[58] Priya C. Kumar, Michael Zimmer, and Jessica Vitak. 2024. A Roadmap for Applying the Contextual Integrity Framework in Qualitative Privacy Research. *Proceedings of the ACM on Human-Computer Interaction* 8, CSCW1 (April 2024), 219:1–219:29. https://doi.org/10.1145/3653710

[59] Srijan Kumar, William L. Hamilton, Jure Leskovec, and Dan Jurafsky. 2018. Community Interaction and Conflict on the Web. In *Proceedings of the 2018 World Wide Web Conference (WWW '18)*. International World Wide Web Conferences Steering Committee, Lyon, France, 933–943. https://doi.org/10.1145/3178876.3186141

[60] Michelle S. Lam, Mitchell L. Gordon, Danaë Metaxa, Jeffrey T. Hancock, James A. Landay, and Michael S. Bernstein. 2022. End-User Audits: A System Empowering Communities to Lead Large-Scale Investigations of Harmful Algorithmic Behavior. *Proceedings of the ACM on Human-Computer Interaction* 6, CSCW2 (Nov. 2022), 512:1–512:34. https://doi.org/10.1145/3555625

[61] Hanlin Li, Brent Hecht, and Stevie Chancellor. 2022. All That's Happening behind the Scenes: Putting the Spotlight on Volunteer Moderator Labor in Reddit. In *Proceedings of the International AAAI Conference on Web and Social Media*, Vol. 16. AAAI, Atlanta, Georgia, 584–595.

[62] Hanlin Li, Nicholas Vincent, Stevie Chancellor, and Brent Hecht. 2023. The Dimensions of Data Labor: A Road Map for Researchers, Activists, and Policymakers to Empower Data Producers. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency (FAccT '23)*. Association for Computing Machinery, New York, NY, USA, 1151–1161. https://doi.org/10.1145/3593013.3594070

[63] Eden Litt and Eszter Hargittai. 2016. The Imagined Audience on Social Network Sites. *Social Media + Society* 2, 1 (Jan. 2016), 2056305116633482. https://doi.org/10.1177/2056305116633482

[64] Aymeric Mansoux and Roel Rocsam Abbing. 2020. Seven Theses on the Fediverse and the Becoming of Floss. In *The Eternal Network: The Ends and Becomings of Network Culture*, Kristoffer Gansing and Inga Luchs (Eds.). Institute of Network Cultures and transmediale e.V.

[65] Kirsten Martin. 2015. Privacy Notices as Tabula Rasa: An Empirical Investigation into How Complying with a Privacy Notice Is Related to Meeting Privacy Expectations Online. *Journal of Public Policy & Marketing* 34, 2 (Sept. 2015), 210–227. https://doi.org/10.1509/jppm.14.139

[66] A. E. Marwick and danah boyd. 2011. I Tweet Honestly, I Tweet Passionately: Twitter Users, Context Collapse, and the Imagined Audience. *New Media & Society* 13, 1 (Feb. 2011), 114–133. https://doi.org/10.1177/1461444810365313

[67] Alice E Marwick and danah boyd. 2014. Networked Privacy: How Teenagers Negotiate Context in Social Media. *New Media & Society* 16, 7 (Nov. 2014), 1051–1067. https://doi.org/10.1177/1461444814543995

[68] J. Nathan Matias. 2019. Preventing Harassment and Increasing Group Participation through Social Norms in 2,190 Online Science Discussions. *Proceedings of the National Academy of Sciences* 116, 20 (May 2019), 9785–9789. https://doi.org/10.1073/pnas.1813486116

[69] Nora McDonald and Nazanin Andalibi. 2023. "I Did Watch 'The Handmaid's Tale'": Threat Modeling Privacy Post-roe in the United States. *ACM Transactions on Computer-Human Interaction* 30, 4 (Sept. 2023), 63:1–63:34. https://doi.org/10.1145/3589960

[70] Nora McDonald and Andrea Forte. 2020. The Politics of Privacy Theories: Moving from Norms to Vulnerabilities. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*. ACM, Honolulu HI USA, 1–14. https://doi.org/10.1145/3313831.3376167

[71] Nora McDonald and Andrea Forte. 2021. Powerful Privacy Norms in Social Network Discourse. *Proceedings of the ACM on Human-Computer Interaction* 5, CSCW2 (Oct. 2021), 421:1–421:27. https://doi.org/10.1145/3479565

[72] Nora McDonald, Benjamin Mako Hill, Rachel Greenstadt, and Andrea Forte. 2019. Privacy, Anonymity, and Perceived Risk in Open Collaboration: A Study of Service Providers. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems (CHI '19) (CHI '19)*. ACM, New York, NY, 671:1–671:12. https://doi.org/10.1145/3290605.3300901

[73] Jonathan T. Morgan, Siko Bouterse, Heather Walls, and Sarah Stierch. 2013. Tea and Sympathy: Crafting Positive New User Experiences on Wikipedia. In *Proceedings of the 2013 Conference on Computer Supported Cooperative Work (CSCW '13)*. ACM, New York, NY, USA, 839–848. https://doi.org/10.1145/2441776.2441871

[74] Kimberlee Morrison. 2014. Domestic Abuse Survivors Go 'Underground' With the Tor Network. *AdWeek* (May 2014).

[75] Deirdre K. Mulligan and Jennifer King. 2011. Bridging the Gap between Privacy and Design. *University of Pennsylvania Journal of Constitutional Law* 14 (2011), 989.

[76] Sneha Narayan, Jake Orlowitz, Jonathan Morgan, Benjamin Mako Hill, and Aaron Shaw. 2017. The Wikipedia Adventure: Field Evaluation of an Interactive Tutorial for New Users. In *Proceedings of the 2017 ACM Conference on Computer Supported Cooperative Work and Social Computing (CSCW '17)*. ACM, New York, NY, USA, 1785–1799. https://doi.org/10.1145/2998181.2998307

[77] Helen Nissenbaum. 2009. *Privacy in Context: Technology, Policy, and the Integrity of Social Life*. Stanford University Press, Stanford, CA. https://doi.org/10.1515/9780804772891

[78] Helen Nissenbaum. 2015. Respect for Context as a Benchmark for Privacy Online: What It Is and Isn't. In *Social Dimensions of Privacy: Interdisciplinary Perspectives*, Beate Roessler and Dorota Mokrosinska (Eds.). Cambridge University Press, Cambridge, UK.

[79] Jonathan A. Obar and Anne Oeldorf-Hirsch. 2018. The Biggest Lie on the Internet: Ignoring the Privacy Policies and Terms of Service Policies of Social Networking Services. *Information, Communication & Society* 23, 1 (July 2018), 1–20. https://doi.org/10.1080/1369118X.2018.1486870

[80] Nicholas O'Donnell. 2021. Have We No Decency? Section 230 and the Liability of Social Media Companies for Deepfake Videos Notes. *University of Illinois Law Review* 2021, 2 (2021), i–740.

[81] Pamela Oliver, Gerald Marwell, and Ruy Teixeira. 1985. A Theory of the Critical Mass: Interdependence, Group Heterogeneity, and the Production of Collective Action. *Amer. J. Sociology* 91, 3 (Nov. 1985), 522–556. https://doi.org/10.1086/228313

[82]   Leysia Palen and Paul Dourish. 2003. Unpacking "Privacy" for a Networked World. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems (CHI '03)*. Association for Computing Machinery, New York, NY, USA, 129–136. https://doi.org/10.1145/642611.642635

[83]   Sandra Petronio. 2002. *Boundaries of Privacy: Dialectics of Disclosure.* SUNY Press, Albany, NY.

[84]   Pincus. 2024. Ah, Lemmy: Racism and Denial in the Threadiverse (DRAFT). https://privacy.thenexus.today/racism-and-denial-in-the-threadiverse/.

[85]   Shruti Sannon, Natalya N. Bazarova, and Dan Cosley. 2018. Privacy Lies: Understanding How, When, and Why People Lie to Protect Their Privacy in Multiple Online Contexts. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*. ACM, Montreal QC Canada, 1–13. https://doi.org/10.1145/3173574.3173626

[86]   Shruti Sannon and Andrea Forte. 2022. Privacy Research with Marginalized Groups: What We Know, What's Needed, and What's Next. *Proceedings of the ACM on Human-Computer Interaction* 6, CSCW2 (Nov. 2022), 455:1–455:33. https://doi.org/10.1145/3555556

[87]   Sarah Scheffler and Jonathan Mayer. 2023. SoK: Content Moderation for End-to-End Encryption. https://doi.org/10.48550/arXiv.2303.03979 arXiv:2303.03979 [cs]

[88]   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. *Proceedings of the ACM on Human-Computer Interaction* 2, CSCW (Nov. 2018), 155:1–155:27. https://doi.org/10.1145/3274424

[89]   Nathan Schneider. 2022. Admins, Mods, and Benevolent Dictators for Life: The Implicit Feudalism of Online Communities. *New Media & Society* 24, 9 (Sept. 2022), 1965–1985. https://doi.org/10.1177/1461444820986553

[90]   Nathan Schneider. 2024. *Governable Spaces: Democratic Design for Online Life.* University of California Press, Oakland, CA.

[91]   Nathan Schneider, Primavera De Filippi, Seth Frey, Joshua Z. Tan, and Amy X. Zhang. 2021. Modular Politics: Toward a Governance Layer for Online Communities. *Proceedings of the ACM on Human-Computer Interaction* 5, CSCW1 (April 2021), 1–26. https://doi.org/10.1145/3449090

[92]   Sarita Yardi Schoenebeck. 2013. The Secret Life of Online Moms: Anonymity and Disinhibition on Youbemom.Com. *Proceedings of the 7th International Conference on Weblogs and Social Media, ICWSM 2013* 7, 1 (2013), 555–562.

[93]   Angela M. Schöpke-Gonzalez, Shubham Atreja, Han Na Shin, Najmin Ahmed, and Libby Hemphill. 2022. Why Do Volunteer Content Moderators Quit? Burnout, Conflict, and Harmful Behaviors - Angela M. Schöpke-Gonzalez, Shubham Atreja, Han Na Shin, Najmin Ahmed, Libby Hemphill, 2022. *New Media & Society* 0(0) (Dec. 2022), 25.

[94]   Joseph Seering. 2020. Reconsidering Self-Moderation: The Role of Research in Supporting Community-Based Models for Online Content Moderation. *Proceedings of the ACM on Human-Computer Interaction* 4, CSCW2 (Oct. 2020), 107:1–107:28. https://doi.org/10.1145/3415178

[95]   Erica Shusas, Patrick Skeba, Eric P. S. Baumer, and Andrea Forte. 2023. Accounting for Privacy Pluralism: Lessons and Strategies from Community-Based Privacy Groups. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems (CHI '23)*. Association for Computing Machinery, New York, NY, USA, 1–12. https://doi.org/10.1145/3544548.3581331

[96]   Yan Shvartzshnaider, Noah Apthorpe, Nick Feamster, and Helen Nissenbaum. 2018. Analyzing Privacy Policies Using Contextual Integrity Annotations. https://doi.org/10.48550/arXiv.1809.02236 arXiv:1809.02236 [cs]

[97]   Yan Shvartzshnaider, Madelyn Rose Sanfilippo, and Noah Apthorpe. 2022. GKC-CI: A Unifying Framework for Contextual Norms and Information Governance. *Journal of the Association for Information Science and Technology* 73, 9 (2022), 1297–1313. https://doi.org/10.1002/asi.24633

[98]   Daniel Solove, J. 2015. The Meaning and Value of Privacy. In *Social Dimensions of Privacy: Interdisciplinary Perspectives*, Beate Roessler and Dorota Mokrosinska (Eds.). Cambridge University Press, Cambridge, UK.

[99]   Gerry Stoker. 1998. Governance as Theory: Five Propositions. *International social science journal* 50, 155 (1998), 17–28.

[100]  John Suler. 2004. The Online Disinhibition Effect. *CyberPsychology & Behavior* 7, 3 (June 2004), 321–326. https://doi.org/10.1089/1094931041291295

[101]  Nathan TeBlunthuis, Charles Kiene, Isabella Brown, Laura (Alia) Levi, Nicole McGinnis, and Benjamin Mako Hill. 2022. No Community Can Do Everything: Why People Participate in Similar Online Communities. *Proceedings of the ACM on Human-Computer Interaction: Computer Supported Cooperative Work* 6 (April 2022), 1–25. https://doi.org/10.1145/3512908

[102]  Emma Tosch, Luis Garcia, Cynthia Li, and Chris Martens. 2024. Privacy Policies on the Fediverse: A Case Study of Mastodon Instances. *Proceedings on Privacy Enhancing Technologies* 2024, 4 (Oct. 2024), 700–733. https://doi.org/10.56553/popets-2024-0138

[103]  Chau Tran, Kaylea Champion, Andrea Forte, Benjamin Mako Hill, and Rachel Greenstadt. 2020. Are Anonymity-Seekers Just like Everybody Else? An Analysis of Contributions to Wikipedia from Tor. In *2020 IEEE Symposium on Security and Privacy (SP)*, Vol. 1. IEEE Computer Society, San Francisco, California, 974–990. https://doi.org/10.1109/SP40000.2020.00053

[104] Sabine Trepte, Michael Scharkow, and Tobias Dienlin. 2020. The Privacy Calculus Contextualized: The Influence of Affordances. *Computers in Human Behavior* 104 (March 2020), 106115. https://doi.org/10.1016/j.chb.2019.08.022

[105] Anthony Henry Triggs, Kristian Møller, and Christina Neumayer. 2019. Context Collapse and Anonymity among Queer Reddit Users. *New Media & Society* 23, 1 (Nov. 2019), 5–21. https://doi.org/10.1177/1461444819890353

[106] Jan E. Trost. 1986. Statistically Nonrepresentative Stratified Sampling: A Sampling Technique for Qualitative Studies. *Qualitative Sociology* 9, 1 (March 1986), 54–57. https://doi.org/10.1007/BF00988249

[107] Salome Viljoen. 2021. A Relational Theory of Data Governance. *Yale Law Journal* 131 (2021), 573.

[108] Nicholas Vincent, Hanlin Li, Nicole Tilly, Stevie Chancellor, and Brent Hecht. 2021. Data Leverage: A Framework for Empowering the Public in Its Relationship with Technology Companies. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (FAccT '21)*. Association for Computing Machinery, New York, NY, USA, 215–227. https://doi.org/10.1145/3442188.3445885

[109] Jessica Vitak. 2012. The Impact of Context Collapse and Privacy on Social Network Site Disclosures. *Journal of Broadcasting & Electronic Media* 56, 4 (Oct. 2012), 451–470. https://doi.org/10.1080/08838151.2012.732140

[110] Jessica Vitak, Stacy Blasiola, Sameer Patil, and Eden Litt. 2015. Balancing Audience and Privacy Tensions on Social Network Sites: Strategies of Highly Engaged Users. *International Journal of Communication* 9, 0 (May 2015), 20.

[111] Xinyu Wang, Sai Koneru, and Sarah Rajtmajer. 2024. The Failed Migration of Academic Twitter. https://doi.org/10.48550/arXiv.2406.04005 arXiv:2406.04005 [cs]

[112] Yi-Chia Wang, Moira Burke, and Robert Kraut. 2016. Modeling Self-Disclosure in Social Networking Sites. In *Proceedings of the 19th ACM Conference on Computer-Supported Cooperative Work & Social Computing (CSCW '16)*. Association for Computing Machinery, New York, NY, USA, 74–85. https://doi.org/10.1145/2818048.2820010

[113] Morten Warncke-Wang, Rita Ho, Marshall Miller, and Isaac Johnson. 2023. Increasing Participation in Peer Production Communities with the Newcomer Homepage. *Proceedings of the ACM on Human-Computer Interaction* 7, CSCW2 (Sept. 2023), 1–26. https://doi.org/10.1145/3610071

[114] Galen Weld, Amy X. Zhang, and Tim Althoff. 2021. What Makes Online Communities 'Better'? Measuring Values, Consensus, and Conflict across Thousands of Subreddits. *arXiv:2111.05835 [cs]* 16 (Nov. 2021), 1121–1132. arXiv:2111.05835 [cs]

[115] Pamela Wisniewski, A.K.M. Najmul Islam, Bart P. Knijnenburg, and Sameer Patil. 2015. Give Social Network Users the Privacy They Want. In *Proceedings of the 18th ACM Conference on Computer Supported Cooperative Work & Social Computing (CSCW '15)*. Association for Computing Machinery, New York, NY, USA, 1427–1441. https://doi.org/10.1145/2675133.2675256

[116] Yuxi Wu, W. Keith Edwards, and Sauvik Das. 2022. "A Reasonable Thing to Ask For": Towards a Unified Voice in Privacy Collective Action. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems (CHI '22)*. Association for Computing Machinery, New York, NY, USA, 1–17. https://doi.org/10.1145/3491102.3517467

[117] Sijia Xiao, Danaë Metaxa, Joon Sung Park, Karrie Karahalios, and Niloufar Salehi. 2020. Random, Messy, Funny, Raw: Finstas as Intimate Reconfigurations of Social Media. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems (CHI '20)*. Association for Computing Machinery, New York, NY, USA, 1–13. https://doi.org/10.1145/3313831.3376424

[118] Amy X. Zhang, Grant Hugh, and Michael S. Bernstein. 2020. PolicyKit: Building Governance in Online Communities. In *Proceedings of the 33rd Annual ACM Symposium on User Interface Software and Technology*. ACM, Virtual Event USA, 365–378. https://doi.org/10.1145/3379337.3415858 arXiv:2008.04236

[119] Zhilin Zhang, Jun Zhao, Ge Wang, Samantha-Kaye Johnston, George Chalhoub, Tala Ross, Diyi Liu, Claudine Tinsman, Rui Zhao, Max Van Kleek, and Nigel Shadbolt. 2024. Trouble in Paradise? Understanding Mastodon Admin's Motivations, Experiences, and Challenges Running Decentralised Social Media. *Proc. ACM Hum.-Comput. Interact.* 8, CSCW2 (Nov. 2024), 520:1–520:24. https://doi.org/10.1145/3687059

[120] Shoshana Zuboff. 2019. *The Age of Surveillance Capitalism: The Fight for a Human Future at the New Frontier of Power*. Public Affairs, New York.

# A  RESEARCH METHODS

## A.1  Interview Protocol - Survey Exercise

As noted in §4, the first part of our interviews involved using a survey as an artifact for structured discussion. The interviewer shared their screen showing the survey and walked through survey items with participants, who said their choices aloud and then explained their reasons for their choice (at times followed by further prompting or clarification questions by the interviewer). Each item (bullet dash) on the survey had the following options: "Not at all acceptable," "Somewhat

acceptable," "Acceptable," "Completely acceptable," and "I do not wish to answer this question." Items were arranged under two blocks, as follows:

- Rate the extent to which you find it acceptable that an instance collects and uses data (e.g., personal information someone has entered, posts, DMs and comments someone toots on the server) for each of the following reasons.
  - To share your posts and activity with other communities or users
  - To send you information, notifications about other people interacting with your content or sending you messages
  - To automatically save your preferences for future visits with cookies
  - To aid moderation of the community, like comparing your IP address with known banned IP addresses
  - To retain server logs, like all IP addresses
  - To share with trusted third parties who assist the server in operating the site, conducting business, and servicing users, so long as those parties agree to keep this information confidential
  - To comply with the law, enforce site policies, or protect the server's or others rights, property, or safety
- Rate the extent to which you find it acceptable for each potential type of person to come across your public information, posts, and comments originally shared on your instance.
  - Someone you DM or @
  - Your followers and/or people you follow
  - Someone you've blocked on the Fediverse or elsewhere
  - A person on the Fediverse you have no connections to
  - Admins, moderators, and other operators of instances (yours and/or others)
  - An app collecting online data generally