# Hate Raids on Twitch: Echoes of the Past, New Modalities, and Implications for Platform Governance

CATHERINE HAN, Computer Science Department, Stanford University, USA
JOSEPH SEERING, Computer Science Department, Stanford University, USA
DEEPAK KUMAR, Computer Science Department, Stanford University, USA
JEFFREY T. HANCOCK, Communication Department, Stanford University, USA
ZAKIR DURUMERIC, Computer Science Department, Stanford University, USA

In the summer of 2021, users on the livestreaming platform Twitch were targeted by a wave of "hate raids," a form of attack that overwhelms a streamer's chatroom with hateful messages, often through the use of bots and automation. Using a mixed-methods approach, we combine a quantitative measurement of attacks across the platform with interviews of streamers and third-party bot developers. We present evidence that confirms that some hate raids were highly-targeted, hate-driven attacks, but we also observe another mode of hate raid similar to networked harassment and specific forms of subcultural trolling. We show that the streamers who self-identify as LGBTQ+ and/or Black were disproportionately targeted and that hate raid messages were most commonly rooted in anti-Black racism and antisemitism. We also document how these attacks elicited rapid community responses in both bolstering reactive moderation and developing proactive mitigations for future attacks. We conclude by discussing how platforms can better prepare for attacks and protect at-risk communities while considering the division of labor between community moderators, tool-builders, and platforms.

**133**

## 1 INTRODUCTION

**Content Warning**: This paper studies hateful online content. When necessary for clarity, this paper directly quotes user-generated content that contains offensive/hateful speech, profanity, and other potentially triggering content.

Authors' addresses: Catherine Han, cathan@stanford.edu, Computer Science Department, Stanford University, Stanford, USA; Joseph Seering, jseering@stanford.edu, Computer Science Department, Stanford University, Stanford, USA; Deepak Kumar, kumarde@stanford.edu, Computer Science Department, Stanford University, Stanford, USA; Jeffrey T. Hancock, hancockj@stanford.edu, Communication Department, Stanford University, Stanford, USA; Zakir Durumeric, zakird@stanford.edu, Computer Science Department, Stanford University, Stanford, USA.

Livestreaming platforms have boomed in popularity in recent years and become a major part of many users' Internet experience. The livestreaming industry saw a 45% uptick in viewership between March and April 2020 [54], likely in part due to the COVID-19 pandemic. Twitch is a popular livestreaming platform, and much like any other rapidly growing online platform, its communities have suffered from hate and harassment. July to October 2021 marked an intense period of harassment on Twitch with many streamers experiencing a surge of "hate raids." In the most common form of hate raid, a streamer's chatroom is overwhelmed by a rapid influx of Twitch accounts posting hateful messages. Because of their dissatisfaction with Twitch's handling of hate raids and poor treatment of marginalized-identity streamers, these streamers and their communities came together, gathering resources, developing tools and strategies to protect themselves, and organizing a major protest [18, 39]. This series of events reflected frustration within the community—particularly from minority streamers—toward Twitch and its perceived inaction on issues of trust, security, and safety on the platform.

In this paper, we investigate the nature of hate raids on Twitch, how they affected vulnerable communities, and how stakeholders reacted to hate raids. We combine an at-scale measurement of the hate raid phenomenon across 9,664 popular channels' chats on Twitch with interviews of seven LGBTQ+ and/or Black Twitch streamers. In addition, we interview two Twitch users that developed third-party moderation tools in response to hate raids. In our analysis, we explore the following three research questions:

**RQ1: What are hate raids: how are they orchestrated and who do they target?** We first seek to detail the fundamental characteristics of hate raids.[1] Our measurement of hate raids across 9,664 popular channels on Twitch reveals that 98% of hate raid messages consisted of identity-based attacks. However, while the *content* of these attacks was mostly anti-Black or antisemitic, the raids themselves selected targets indiscriminately with respect to streamer identity. These hate raids blurred the line between what prior work called "trolling" or disruptive behavior [41] and networked harassment [30]. To better understand how attackers selected their targets, we examined Twitch's streamer tags—a feature streamers use to categorize themselves and their community. Among streams that use tags, we find evidence that attackers may have leveraged these tags to discover and attack marginalized-identity streamers: particularly with Black, African American, and LGBTQ+ tags.

**RQ2: How do hate raids affect members of targeted groups?** Because a quantitative perspective on hate raids cannot fully depict the lived experiences of targeted community members, we interviewed seven Black and/or LGBTQ+ streamers on Twitch about the impact of these attacks. Through these interviews, we find that the perspectives of targeted streamers aligned with mainstream media portrayals of these attacks: hate raids are seen as highly-targeted attacks often persecuting Black and LGBTQ+ communities on Twitch. While identity-based attacks have always plagued these at-risk communities online, streamers found that this wave of hate raids was distinct in its highly-targeted nature and the persistence of its perpetrators. Furthermore, we find that the community saw hate raids as one piece of a larger campaign of harassment, often involving other platforms and in some cases extending into more extreme offline experiences (e.g., involving law enforcement, swatting[2]).

**RQ3: How did different groups of stakeholders respond?** To better understand the different ways community members and Twitch responded to hate raids, we further draw upon data from

---

[1]Some news reports stated that these raids began as an abuse of a built-in "raiding" feature originally intended to help grow a sense of community [38], but we did not find direct evidence of this in our dataset.

[2]A harassment tactic that involves calling emergency services or police to a target's residence

interviews with streamers and bot developers. We observed that streamers largely turned to their community and third-party bot developers for moderation, emotional, and technical support against hate raids. Volunteer bot developers created tools adopted by tens of thousands of streamers who felt that they might be targeted. These developers worked to constantly update their tools throughout the hate raid period, as the sophistication of hate raids evolved in response to developers' efforts to combat these attacks. In addition to an influx of support via resource aggregation, tool development, and volunteer moderation, the community rallied together for a social movement and virtual walkout to raise awareness for their longstanding frustrations with Twitch. While attitudes toward the degree of success of these movements varied among our interviewees, these community-driven movements gained attention and impacted overall platform engagement.

Our mixed-methods approach to understanding hate raids provides the following three primary research contributions:

(1) We characterize a novel form of long-term harassment campaigns on Twitch; not only do we observe that hate raids leverage the real-time nature of livestreaming platforms, but we also find that they exploit automation to select targets and amplify their attacks.
(2) We observe that the content and orchestration of hate raid messages indicate a dual motivation: first, hate-driven and second, attention-seeking, consistent with prior research into networked harassment and subcultural trolling.
(3) We find that members of these targeted communities, unhindered by the frictions platforms face when developing new features and policies, rapidly assembled high-quality resources and produced technical tools to address their needs and the limitations of Twitch's response.

Grounded in our data, we conclude by discussing the implications of our findings for livestreaming platform design and the broader community. We argue that platforms and researchers must proactively consider the unique experiences of targeted communities online, the dependency on and potential for community-based moderation and tool-building, and the range of motivations behind the actors coordinating hate-based attacks.

## 2 RELATED WORK

This paper builds on three key bodies of related work. First, we review literature on morally-motivated networked harassment [30] and subcultural trolling [41], and we identify characteristics of each that hate raids share. Second, we review online hate-based attacks documented in the literature, situating hate raids within taxonomies of their characteristics. Finally, we discuss ties to literature on volunteer moderation and coordinated action, identifying connections between hate raids and crisis informatics literature and highlighting how users' responses to hate raids parallel responses to natural disasters and other crises.

### 2.1 Harassment and "trolling" in online spaces

In this paper, we situate the Twitch hate raids within prior work that discusses online harassment and "trolling." Definitions for both of these terms have varied widely; for example, trolling has been defined as broadly as "behavior that falls outside acceptable bounds defined by [...] communities" [5, p. 1] and as specifically as in Phillips' description of "subcultural trolling" as a nuanced cultural phenomenon with historical and moral roots [40, 41]. Similarly, Marwick identified more than ten different types of behaviors listed under the umbrella term of "online harassment" in prior work [30, p. 2]. We operate under the definitions of the two terms provided by Marwick and Phillips, and we focus on the form of harassment that Marwick terms "networked harassment," where an individual is harassed by many people connected by social media.

Note that, subsequent to her original publications on subcultural trolling, Phillips wrote about the dangers of referring to something as "just" trolling [42, p. 2]. While in this paper, we compare aspects of hate raids to aspects of Phillips' characterization of subcultural trolling, this should not be construed to mean that hate raids are "just" trolling by any means; they cause real harm to targets that should not be taken lightly. Moreover, these attacks occurred in the context of a long history of racist, sexist, and transphobic behaviors in online spaces that have been especially prevalent in online gaming spaces [15, 17]. These behaviors have forced targeted users to hide their identities or even to withdraw from online spaces entirely [8, 15, 45, 60].

## 2.2 Characterizing hate-based attacks

Thomas et al. [56] identify three axes on which hate-based attacks can be classified:

(1) The *Audience* exposed to the attack, which can include the target and/or a different audience.
(2) The *Medium* through which the attacker reaches a target, which frequently includes media such as text, images, or video.
(3) The *Capabilities* that are required for the attack to succeed: whether the attack requires deception of an audience and/or a third-party authority, whether it requires amplification, and whether it requires privileged access to information, an account, or a device.

In the context of online hate and harassment behaviors, the most similar to hate raids is "brigading," where a single target (e.g., a YouTube video or Twitter account) is simultaneously attacked by a semi-coordinated set of antagonistic users. For example, 4chan users often coordinate to target YouTube videos that they are ideologically or otherwise opposed to [29]; Reddit users have previously, in large groups, entered other community spaces to harass and intimidate other subreddits [12]; Zoom users have leveraged legitimate insider access to join online meetings to disrupt and harass the other participants, otherwise known as "Zoombombing" [26].

Of the above criteria, the *medium* through which hate raids took place is primarily text, though in some cases other media on external platforms were involved. As we discuss later in this work, they required an *audience* that included both the target and a wider array of viewers. In some cases, the attacks included revealing personal information of targets ("doxxing"), and they benefited greatly from amplification.

However, as we discuss in Section 4.1, these attacks had a number of other attributes worth mentioning. For example, the *capabilities* required for this attack included that they were heavily automated and occurred over a significant period of time (several months), hearkening to more traditional cybersecurity attacks, such as Distributed Denial-of-Service (DDoS) [34] and for-profit spam and scam campaigns [23]. Though Zoombombing often operates under a notion of the infiltration of a private meeting, public Twitch streams share the capability of seeing the reactions and impact of the attack in Zoombombing attacks. Therefore, we draw upon prior work in the cybersecurity space to structure our understanding of abuse executed en masse via illegitimate accounts. Contextualizing the hate raids on Twitch through both a lens of subcultural trolling and morally-motivated networked harassment and a traditional cybersecurity lens better frames the underlying motivation and tactics of these activities.

## 2.3 Volunteer moderation, coordinated action, and crisis informatics

Prior work examining platform governance and volunteer labor in online social spaces has highlighted a variety of dynamics that inform our analysis of hate raids. While Twitch is a multi-modal platform incorporating text-based chat, video, and audio, the phenomenon of hate raids echoes the moderation challenges discussed by Jiang et al. for voice-based communities [22], as both Twitch and Discord share ephemeral and real-time components of user interactions. Additionally,

we discuss the experience of hate raids and the resulting mobilization of less visible streamers on Twitch and members of marginalized communities on the platform more broadly. Prior work details the obstacles that such communities in particular face with regards to platform visibility and accountability [55], further contextualizing the friction we observe between Twitch and its users. Several examples [2, 46, 48] in the literature emphasize the importance of volunteer labor in these communities, reporting that volunteer moderators on livestreaming platforms — both individually and in collaboration — have the capacity to effectively and quickly address norm-violating behaviors. In Section 4.2, we discuss the impact of community moderation and community-developed automated moderation tools, adding to conversation in prior work that has raised questions surrounding platform governance and the distribution of labor in content moderation [4, 24, 44, 49].

As we detail below, one of the core characteristics of users' responses was collective action to create tools and aggregate informational resources. A small number of examples of collective action to counter harassment at this scale have been documented in social computing literature. Blackwell et al. reported on "HeartMob," a platform where users can submit reports of being harassed and volunteers will provide support — supportive messages, help with reporting harassment, and/or help documenting abuse [1]. On a much smaller scale, Mahar, Zhang, and Karger's "Squadbox" allowed users to coordinate trusted friends to help shield them from harassment via email [28]. A small body of work from the early-mid 1990s [13, 27, 51] and early 2000s [20] also documented individual cases of harassment and communities' discussions about how to respond.

A broader related body of work, situated in part in CSCW literature, comes from the field of crisis informatics [36, 37]. Though this field has largely focused on responses to offline crises (e.g., natural disasters [32, 52, 53, 62], terrorist attacks and mass shootings [3, 6, 37], and in some cases ongoing violent conflict [33, 50]), many of the core principles are also mirrored in responses to hateful attacks based on social media. As we discuss in Section 4.2, we observe many of the same behaviors in our research on Twitch hate raids that occur during natural disaster response. Per this literature, we have organized our results to address questions about crisis response that parallel questions commonly asked in crisis informatics literature.

## 3 METHODS

We examine broad patterns in hate raids and common themes in individual messages, and we complement this analysis with insights from interviews with impacted individuals. In this section, we describe the methodologies of our (1) large-scale collection and analysis of Twitch chat messages, moderation actions, and channel attributes collected from 9,664 channels from September 2 to September 16, 2021, (2) interviews with seven Black and/or LGBTQ+ streamers, and (3) interviews with developers of two third-party Twitch moderation bots that were widely deployed in response to hate raids.

### 3.1 Twitch Chat Data Collection

To understand how hate raids impacted high-visibility streams on Twitch, we generated a corpus of channels to gather messages from. We used Twitch's API to pull information about online streamers ordered by their current number of viewers, from high to low. We pulled this data every hour for a week from May 4 to May 11, 2021 to compute an average number of viewers per stream when the channel was live. For our corpus, we only considered channels that had an average of at least 100 viewers each time they streamed and that also streamed at least three times over the course of a week.

We continuously gathered data from the channels on this list for two weeks in September, from September 2 to September 16, 2021, during which time many hate raid attacks occurred. Each channel on Twitch has an associated chatroom built on Internet Relay Chat (IRC) protocols. When

connecting to each channel's chat, we sent requests for information about the channel's chatroom modes—unique chat, subscribers-only mode, and slow mode.[3] We also sent requests for *command* and *membership* capabilities, which allow us to identify the usage of certain moderation and room state commands and to determine when users joined or left chat; the CLEARCHAT command indicates that all of a specific users' messages were purged from the chat, often as a result of a moderation action, like a timeout or a ban, while the membership capability reveals when specific users are joining and leaving the chat. In total, we collected 244,738,672 messages. For each message that was sent, we collected various pieces of metadata to contextualize it: what channel it was sent in, the account that sent the message, the text content of the message, the timestamp of when it was sent to the chat, the status of the chatroom (e.g., if it was in "slow mode"), and basic, publicly-visible information about the account that sent the message (e.g., if the account is a subscriber, follower, or moderator of the channel it is participating in). All data that was collected for this portion of this study was public to any user viewing the stream.

## 3.2 Detecting Hate Raids

We started with a collection of 1,319,890 likely malicious bot accounts curated by and shared among the Twitch community so that streamers could proactively ban and block these accounts from participating in their chats. We searched our Twitch chat dataset for messages sent by these accounts, creating a seed set of messages from 516 of these likely bot accounts. We then used approximate string matching computed using the Levenshtein distance with a threshold of 95% similarity to find messages with the same content despite some evasion techniques used by hate raid attackers, such as prepending randomness to the same message contents across different accounts. We continue this process of finding approximate message content until no new messages were discovered. Through this method, we found matching message contents found by an additional 1,067 discovered bot accounts for a total of 1,583 bots participating in hate raids (Figure 1). We then determined hate raid events to be windows of time where bot accounts in our dataset were seen sending messages within two minutes of prior messages sent by bots. We restricted this window to a short interval because raiding behavior (both benign and malicious) often involves an influx of similar messages sent across different accounts within a short period of time.

## 3.3 Streamer and Third-Party Bot Developer Interviews

We conducted semi-structured interviews with seven Twitch streamers who identified as Black and/or LGBTQ+ and with two Twitch users who created third-party moderation bots to combat hate raids. These interviews were conducted from early October through mid-November 2021, shortly after the major spike in hate raids in late September. Interviews lasted between 20 minutes and one hour, with length varying based on participants' exposure to hate raids, their roles within the community, and their knowledge of moderation tools. We recruited participants from lists of streamers who had previously participated in visible roles during LGBTQ+ focused events on Twitch, including featured streamers during Pride Month, streamers who were reported in news articles as having been heavily targeted by hate raids, and streamers who actively participated in hate raid-focused conversations in both public and semi-private spaces dedicated to hate raid responses. We recruited specifically from Black and LGBTQ+ streamers because these were the groups at the center of discourse surrounding hate raids and were the most visibly targeted. Interview questions focused on the same topics as the research questions, with a full list of primary questions presented in Appendix A. Due to the open-ended, semi-structured nature of these interviews, we asked additional follow-up questions when relevant.

---

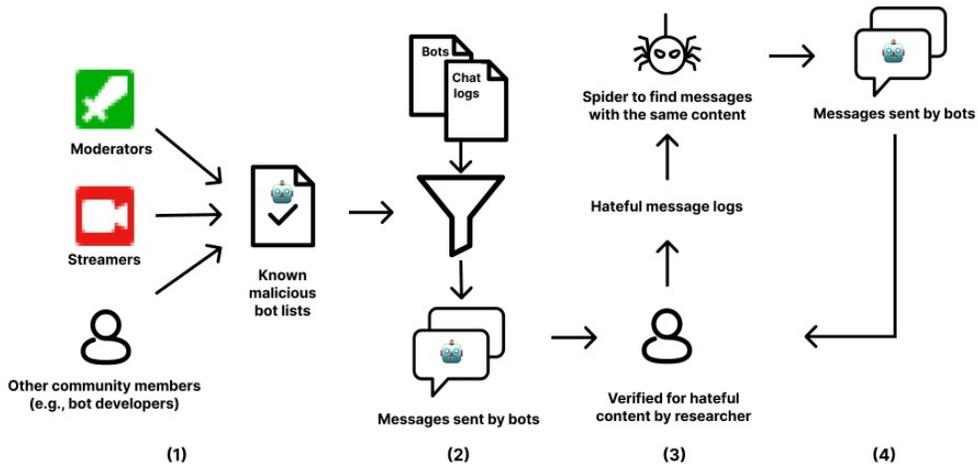[3]https://help.twitch.tv/s/article/chat-commands

Fig. 1. Compiling hate raid logs. We collected the hate raid logs in a series of four steps: (1) we began with a list of known malicious bot account names collected by moderators, streamers, and other community members on Twitch; (2) we used this to filter the chat logs for messages sent by malicious bots; (3) after confirming the contents of these messages are hateful, we used this seed set of messages to spider for similar content being sent by accounts not already in this list; (4) we iterated with this notion of message content similarity until no new bot accounts were detected.

We do not report identity characteristics for each streamer individually because doing so might identify them, but the following are aggregated, self-reported demographic categories: four streamers identified as Black, two as Hispanic, and one as white. Five identified as women, one identified as nonbinary, two identified as transgender, two identified as queer, and one also identified as aromantic and asexual. Some interviewees identified with more than one of these categories. The bot developer interviewees both identified as white, male, and heterosexual.

Following interview completion and transcription, interview text was separated into chunks. Each chunk contained a single idea, which ranged in length from several words to several sentences using a variant of the method described in Creswell [9, p. 86–89, 184–185]. These chunks were each given category labels, which included categories such as "Frequency of hate raids experienced," "Streamers' short-term responses to hate raids," and "Social support received by streamers." The full codebook is included in Appendix B. Initial category labels were defined by the research questions, but labels were iteratively added to the codebook when a chunk did not fit any existing labels. A Cohen's Kappa statistic was calculated to determine inter-rater reliability, with the final round of coding achieving a Cohen's Kappa of 0.91. The results of this analysis are summarized by category label in Section 4.2.

## 3.4 Ethical Considerations

We gathered data from 9,664 different Twitch channels, each with at least 100 viewers on average. Even though chat data from all channels on Twitch is publicly viewable, we elected to restrict the scope of our analysis to this set of larger channels to protect any assumption of privacy that smaller channels and their communities might have; channels with regular audiences of 100 or more viewers represent an exceedingly small proportion of Twitch channels overall—in May 2021, nearly 99% of streams had fewer than 50 average concurrent viewers [7]. This restriction applies a significant limitation to our quantitative analysis, as we cannot draw conclusions regarding hate

raid messages sent to smaller channels, but we believe that the ethical considerations in respecting privacy justify this limitation. The final list contains 9,664 active channels that match these criteria. In the chat data we collected, we took precautions to minimize the risk of inadvertently affecting communities: our script did not send any messages or interact with the chat, and we did not attempt to de-anonymize the involved accounts.

Furthermore, to gather our qualitative data, we interviewed members of Black and/or LGBTQ+ communities concerning their experiences with hate raids. Because of the sensitive nature of this research, participants were notified of the full purpose of the interview in advance, as well as what types of questions would be asked. Additionally, we reminded participants both on the consent form and at the beginning of the interview that they could decline to answer any questions or stop at any time. Interviews typically lasted between 20 to 60 minutes, and participants were compensated with an Amazon gift code for $15 or local currency equivalent. To protect participants' anonymity, we have removed any potentially personally identifiable information from their quotes. This work was approved by the Stanford University Institutional Review Board (IRB).

## 4  RESULTS

We measure hate raids across the platform and present our findings of their quantitative characteristics below (Section 4.1). We pair these measurements with a synthesis of the qualitative perspectives of streamers from at-risk communities and community bot developers on the responses of different stakeholders (4.2).

### 4.1  Characterizations of Hate Raids

Mainstream news outlets characterized the hate raids during late summer of 2021 as targeted, bot-mediated abuse often aimed toward marginalized streamers [11, 19, 38]. We find two forms hate raids: first, a broad, scattershot form of hate raids akin to classic subcultural trolling [41] that incorporates racist and antisemitic elements, and second, hate raids that targeted specific streamers based on their identities.

*4.1.1  Quantitative Perspectives.* We first sought to understand what hate raids looked like quantitatively across the platform.[4] To achieve this, we characterized hate raids observed in a corpus of 244M messages across 9,664 channels collected during a 14-day period from September 2 to September 16, 2021. Of these messages, 2,947 messages were identified as being part of hate raids through the methods discussed above. We observed 60 hate raid attacks in 57 unique channels—three of these channels were hate raided twice on separate occasions.

**Technical Characteristics**    We find that 50% of channels that were hate raided had at least 32 bot accounts involved in the attack (Figure 2). Some channels, however, experienced attacks with an acutely large number of bots. For example, one channel received hate raid messages from 222 unique bot accounts. We found that on average, there were 48 messages per raid. These messages were typically sent in close succession to one another. In the majority of raids, all of the messages were sent in less than 16 seconds (Figure 4), though a smaller proportion of raids lasted for minutes. Most messages were sent from unique bot accounts, with 302 bots (19.1%) sending more than one message in the same raid; even when these bots did send more than one message, the median number of messages sent by a single bot was two (Figure 5).

Overall, bots appear to have been largely throwaway, single-use accounts often created for the purpose of enacting these hate raids. The usernames of the bots we observe in our dataset

---

[4]Note that, as discussed above, we focus here on within-chat hate raids rather than on forms of follow-botting that were sometimes included under the umbrella term of hate raids.
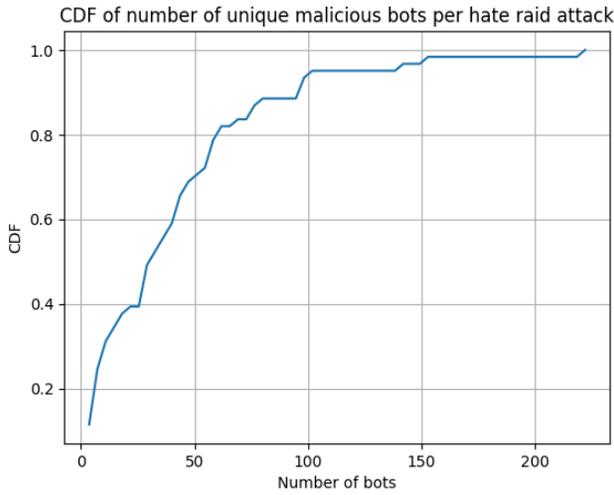
Fig. 2. CDF of the number of unique bots that participated in an instance of a hate raid. We find that the median bot count was 32, demonstrating the typical scale of these attacks.
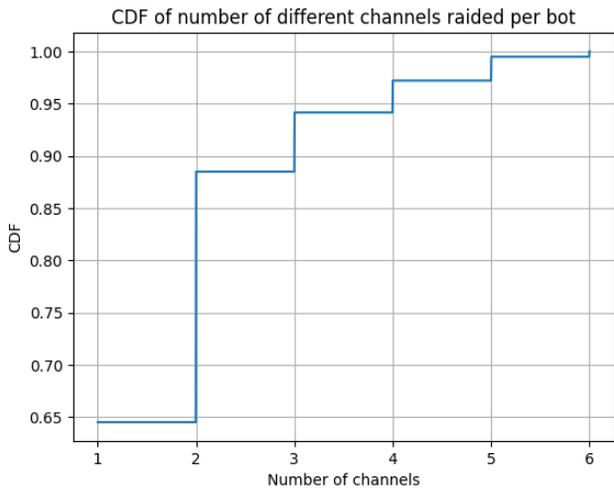


Fig. 3. CDF of the number of different channels a bot account hate raids. Almost 65% of the bot accounts we identified were found in only one channel, implying that the majority of these bots were created for a single use in our observation period, though it is possible that these bots were used additional times in channels that were not in our sample.

were predominantly (99.6%) strings of letters and numbers that appear to have been automatically generated (e.g., y9y7n18r0g6raem). Of the bots we detected sending messages ($N = 1,583$), many (25%) of the bots were created within a two-day window of their first use in our dataset. While these bots appear to have been created for use in a single hate raid, we find that 3% were made far in advance of their first use—these are accounts created at least several weeks before observed chatting in our corpus. This trend toward many single-use accounts likely controlled and created by
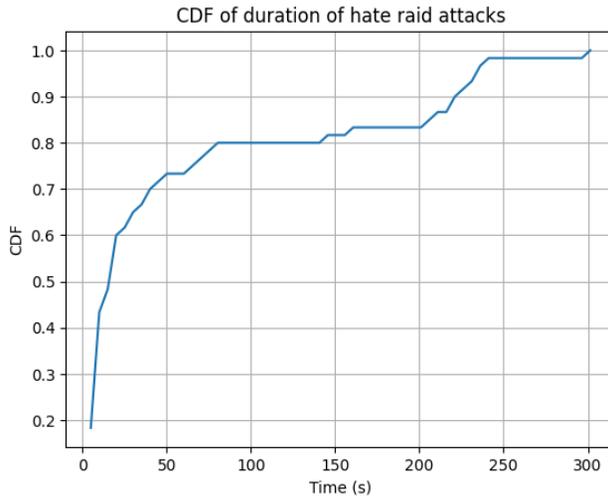
Fig. 4. CDF of the duration of hate raid instances in seconds. The hate raids we observed were largely a short burst of hateful messages sent in close succession, often within the span of seconds or minutes.



Fig. 5. CDF of the number of messages sent per hate raid instance. The median message count is 36, showing the usual scale of the observed hate raids. Note that this number reflects how many messages actually appeared in the stream chat; reactive moderation actions taken by the streamer and/or their volunteer moderators may have prevented bots from sending additional messages.

a single entity echoes the literature in "sockpuppetry," where accounts are created and controlled by a "puppetmaster" for engaging in deceptive behavior influencing the surrounding community [25]. Likewise, the small proportion of aged accounts may indicate what prior work defines as "zombie" accounts, which are ones that are created ahead of time but are dormant for a long period or indicate benign account compromise [14].

We found limited evidence of bot account reuse; the majority of bots were only observed participating in a single channel within our sample, and even then, most bots sent only one message (Figure 3). We also found that slightly more than 20% of the known bots sent hate raid messages in at least two different channels. As of the time of this analysis, there were two ways to sign up for a Twitch account: (1) through e-mail and (2) through phone number. The cost of creating Twitch accounts could thus be as low as the cost of creating e-mail accounts. During this period, Twitch did move quickly to disable accounts that participated in malicious activity. When querying Twitch's API to understand the ages of these accounts, we successfully fetched information for only 33 (4%) of the bot accounts in our dataset because the malicious accounts had been disabled by the time we queried the API for their information.

Only three of the 57 targeted channels experienced a hate raid attack more than once, meaning that our corpus contains 60 observed attacks. Two of these channels experienced two attacks in close succession—within 30 minutes of the first attack. Additionally, both channels experienced a similar pattern where, between the two attacks, the shared text of the messages spammed by different bot accounts in the same raid changed. For instance, in one of the channels, the message content spammed in the first attack was a violent, anti-Black racist statement mocking Twitch community efforts to organize against hate raids. Just 27 minutes later, a second attack began with accounts spamming a different anti-Black racist message. We note that, while both raids were anti-Black in nature, the target of this raid was white.

**Degree of Targeting**     Because of the tendency toward identity-based attacks we observed in the hate raid messages, we next more closely examined the contents of these messages for a semantic understanding of hate raid targeting. Although the bursty messaging pattern with identical message content that we observed in hate raids may appear similar to other, more benign behaviors on Twitch (e.g., chanting[5] and pro-social raiding), the content of hate raid messages distinguishes them as clearly malicious. We find that the hate raids in our dataset spanned several different kinds of hate—most often identity-based—and weaponized these hateful ideologies via graphic and threatening language.

To better characterize how hate was expressed in the raids, we categorized the content of the hate raid messages in our dataset. We evaluated each message along two axes: (1) what identities were attacked in the message, and (2) in what method this identity-based hate was operationalized. Following best practices for grounded theory coding, two researchers agreed upon a master codebook (Table 1) and independently coded 2,947 messages sent across 60 different attacks. The Kupper-Hafner agreement was computed to determine inter-rater agreement because some messages were assigned multiple labels, and the coding achieved an agreement of 0.85. Ultimately, the researchers met to agree on the final codes. We found that the most common category of hate expressed was anti-Black racism, which was present in nearly all of the hate raids in our dataset (59, 98.3%). We observed that anti-Black racism was most often operationalized through violent threats (43 of 59, 72.9%). We also noted that the content in hate raid messages was frequently an amalgam of hateful ideologies—for instance, while anti-Black racism is an explicit category of identity-based hate, messaging with anti-Black attacks often co-occured with QAnon propaganda (23 of 59, 39.0%). These ideologies often overlap with their hateful roots (e.g., white supremacy underlies anti-Black racism, antisemitism, and aspects of QAnon), but the way that these themes were presented together was typically disjointed, separated into different parts of a single message. For instance, the following message both expresses anti-Black racism and supports QAnon:

---

[5]An experimental feature introduced by Twitch in May 2021 that allows streamers and their moderators to "suggest" messages to be duplicated or "chanted" by other users.

| Hateful Ideology | Meaning |
| --- | --- |
| Antisemitic | Demonstrating prejudice toward Jewish people |
| Anti-Black | Demonstrating prejudice toward Black people |
| Anti-Trans | Demonstrating prejudice toward transgender people |
| Individual streamer | Harassing a particular streamer or individual (not necessarily the streamer whose channel the message is sent in) |

| Mode of Operation | Meaning |
| --- | --- |
| Violent threat | Threats of violence (describing explicit actions) |
| Known propaganda | Using known hate symbols or references (e.g., most commonly excerpts from the Great Replacement or 1488*) |
| Direct attack | Attacks that appear to directly address the streamer (e.g., attacks in the second person) or align with the streamer's identity |
| Fearmongering | Inspiring fear or resignation by emphasizing the futility of counter-hate raid efforts |
| Weaponized emote | Coopting Twitch emotes for harassment purposes (e.g., TriHard**) |
| Dehumanization | Implications that a group of people is not human (e.g., comparisons with animals) |

Table 1. **Codebook for hate raid message content**—Codebooks for hate raid message content separated into two axes: (1) hateful ideology and (2) the mode in which they were operationalized.

* 1488 is a pairing of two popular hate symbols, both regarding white supremacy and neo-Nazism.
** TriHard is a global Twitch emote depicting the face of a Black streamer, TriHex, and it has been used in the past to alienate Black streamers.

|  | Violent threat | Known propaganda | Direct attack | Fearmongering | Weaponized emote |
| --- | --- | --- | --- | --- | --- |
| Antisemitic | 10 | 12 | 2 | 0 | 3 |
| Anti-Black | 43 | 11 | 4 | 7 | 3 |
| Individual | 6 | 0 | 0 | 3 | 0 |

Table 2. **Number of raids containing categories of operationalized hate**—Anti-Black racism was the most common form of identity-based hate we identified, and it was most often operationalized through violent threats. We do not list anti-Trans attacks in this table because they were not present in our dataset of messages, though they may have been present in hate raids we were not able to capture.

> *"7i9nnde4k WayneLambright Legion | kiII -> bIacks | behead -> bIacks| <readacted> is a cloutchasing clown he isnt even hateraiding he is just following"*

There are four clear parts to this message, each with a separate meaning. This message's pattern is a common one throughout our dataset; in this case, pipes ("|") were used to delimit separate parts, as attackers presented several pieces of unrelated hateful content in one message, but other symbols (e.g., brackets, braces, "==>", etc.) were also used to separate or relate concepts. In this example, the first component promotes Wayne Lambright, who ran for president in the United States in 2020 on a campaign supporting QAnon, anti-Black racism, and pseudoscience. The second and third

|  | Stream Tags | | |
|  | LGBTQIA+ | Black+AfAm | No Tag |
| --- | --- | --- | --- |
| Hate raided streams (N=57) | 13.5% | 9.6% | 5.2% |
| Not hate raided streams (N=9,664) | 2 .4% | 0.2% | 15.6% |
| p-value | 0.01 | 0.01 | 0.11 |

Table 3. **Streamer tags**—The results of a two-sample proportion test of the self-assigned identity tags (e.g., LGBTQIA+, Black, African American) between the channels that were and were not hate raided. We find that both LGBTQIA+ and Black + African American tags were disproportionately represented among hate raided streams.

parts both express violent anti-Black threats. Finally, the last piece is an attack on an individual streamer. We note that in the messages attacking this streamer, they are neither present as a user in the chat nor are they the streamer of the channel itself; these attacks attempted to harass this streamer through fabricating negative associations between them and hate raid orchestration.

**Streamer Tags**    Because the contents of these hate raid messages were largely rooted in anti-Black racism and antisemitism, we next investigated the use of Twitch tags as potential vectors for targeting. When streaming on Twitch, streamers can choose to categorize their streams with "tags," which are ways to publicly describe the stream for viewers to better search for streams of interest. These tags are maintained by Twitch, but the list of available tags are updated based upon community feedback. In May 2021, Twitch introduced over 350 opt-in tags for streamers to better categorize their channels into a particular community. These tags were largely identity-based, including "gender, sexual orientation, race, nationality, ability, mental health, and more" [59]. However, some streamers feared that the same tags that were meant to increase visibility within a community could be abused to "single out minority streamers," [38] and other members of these communities discouraged use of these tags as a preventative measure [21] to hide themselves from potential attackers. We observe 54 of 57 channels (94.7%) tag themselves with at least one such category; however, we focus our attention on the tags that give insight into streamer identities (e.g., LGBTQ+, Black, etc.) because the messages consisted of identity-based attacks.

To better understand tags' potential usage as a mechanism for targeting marginalized communities, we performed two-sample proportion tests to compare the presence of these identity tags between channels that were and were not hate raided. We find quantitative evidence that suggests that tags may indeed have been used to find targets for harassment at scale: both the LGBTQ+ and Black/African American tags were disproportionately represented ($p < 0.05$) in the hate-raided streams (Table 3). We do *not* find, however, any usage of the Jewish identity tag despite the heavy usage of anti-Semitic language in hate raid messages, and we note that the disproportionate representation of LGBTQ+ streamers deviates from the identities attacked in the *contents* of the messages.

In order to more fully understand the disparity between the identities of the targeted streamers and the identities attacked in the content of the message, we categorized the racial identities of (1) the streamers who were raided and (2) a random sample ($N = 370$) of the broader corpus. Two researchers independently categorized these streamers by their perceived racial category in broad buckets: white, person of color (PoC), and unavailable. Two streamers in the hate-raided sample were unable to be categorized due to the streamer either not including a video feed of their face or using a racially ambiguous virtual avatar ("VTuber" model). We found that the majority of streamers were white in both the hate raided sample and the random sample. We found that the

mainstream sample consisted of 41% PoC streamers, which is higher than what we observed among hate raid victims (35%). We note that there is a very large discrepancy between the proportion of PoC streamers identified through manual coding versus the Black/African American tags. This may be because these tags are not assigned by default, and in order to apply them, streamers must explicitly select them. We performed a two-sample proportion test to evaluate whether the racial identities of the populations of (1) the victims of hate raids in our mainstream corpus and (2) our mainstream corpus differ. In this case, we failed to reject the null hypothesis ($p > 0.05$), meaning that we do not have evidence to say that the set of hate raids we quantified, which often missed the mark in the identities targeted in their content and the identities of the victim, disproportionately targeted PoC streamers as coded in this way. It is possible that this is in part due to an intersection of race and gender/sexual identity where the proportions of LGBTQ+-identifying streamers were unequal between racial groups, but we do not have the sample size to adequately test this within our sample. However, the above analysis does present evidence that hate raids occurred in different proportions across different identity tags, suggesting that tags may have been used as a targeting mechanism.

Our sample size of 57 hate raided channels is inadequate to perform rigorous statistical testing to determine whether the broader set of attacks disproportionately targeted Black streamers, but we note that the proportion of anti-Black content in hate raid messages (98.3%) is far greater than the proportion of Black streamers that we detected experiencing hate raids (10.5%); further, only 2 of these Black streamers received hate raid messages that specifically contained racist anti-Black language, though implicitly racist and/or antisemitic undertones and references were still present in some. This disparity between the identities of the streamers and the kinds of hate spewed in their chats indicates that many of these attacks were indiscriminate in their targets—in most but not all cases, they were not tailored to the specific streamer, but rather contained a consistent breadth of hate regardless of their target. Paired with our tag analysis, we find that the extent of targeting in the observed hate raids may have relied on the usage of tags due to the ease of automation. Through this large-scale, quantitative perspective, we find evidence of another mode of hate raids that included identity-based attacks but did not align the content of their messages with the targeted streamers' identities. Rather, we see recurring themes and shared message text across hate raids in different channels regardless of the streamers' racial identities; we describe such general, reused hateful content sent en masse as "canned hate."

*4.1.2 Perspectives of Streamers from Targeted Communities.* While the *content* of the hate raid messages primarily targeted Black and Jewish identities, analysis of tags revealed that *streamers* who were attacked were disproportionately likely to be those using Black and/or LGBTQ+ tags. To better understand these nuances, we consider the perspectives of streamers from these targeted communities. We conducted a series of interviews with seven Twitch streamers (labeled as TS in quote attributions) that identified as Black and/or LGBTQ+. In this section, we discuss their accounts of how members of these communities perceived the targeted nature of hate raids and the different channels through which they were executed.

**Degree of Targeting**     In our interviews with streamers, we found that streamers' experiences largely aligned with media descriptions of hate raids as a highly-targeted attack, often specifically targeted toward Black and LGBTQ+ creators on Twitch. Six of the seven streamers we interviewed explicitly described the primary targets of hate raids as Black, BIPOC, transgender, or LGBTQ+ communities. In addition to these commonly targeted demographics, two streamers noted that visibility also played a role in attackers' choice in targeted channels:

> *"Specifically like one of my friends who has a bigger viewership, they've been affected a lot more."* – TS01

> *"[My experience] was very mild in comparison to other streamers' who were either vocally and proudly trans or Black or both, as those were absolutely the target demographics."* – TS06

From these streamers' perspectives, viewership and reputation factored into which streamers were more likely to be targeted, in addition to their race and gender.

Additionally, while interviewees acknowledged that Black and LGBTQ+ communities have always been at-risk for hate and harassment on online platforms, three of seven participants stated that this wave of hate raids was drastically more severe than the attacks they had experienced before. For instance, one streamer described the hate raid they experienced in 2021 as "arguably the worst raid" and "most egregious iteration" they have seen to date (TS01). Per these accounts, we sought to examine what aspects of hate raids in 2021 distinguished them from previous attacks. Several streamers explained that this sharp peak in severity manifested in the persistence and scale of the attacks. The reported frequency of hate raids varied across the streamers; while one streamer stated that they were hate raided only once, two others observed a drastic increase in the duration and frequency of the hate raids they experienced firsthand and witnessed in other channels. One streamer recounted how they were hate raided for two weeks straight:

> *"The highlight was the first stream that they hit me in, I had a four and a half hour stream. They were in my stream for about three and a half of those hours, nonstop hate raiding me."* – TS01

Another streamer (TS03) contrasted their experience with hate raids before and during this particularly active period, where before, hate-driven attacks occurred as a single burst that "wasn't an all day, every day or an hours long thing" and would "die out for a while." However, in 2021, they witnessed hate raids that were far worse:

> *"They were raided for three hours straight, just three hours of just following and trying to put messages in chat, trying dox them, whether it was by putting an address in a message or making a username with the address and just following incessantly."* – TS03

While we did not find raids of this type within our dataset, we did not capture data from every targeted channel for reasons discussed above.

TS05 notes that the degree of automation played a key role in the impact of the threat; the usage of bots grew over time and later reached unprecedented scale—they would use a tool to block suspected bot accounts all night long, blocking 300,000 to 400,000 bots at a time. They described the churn of newly created and weaponized bot accounts as "incessant and overwhelming." In addition, TS05 commented on the sharp growth in attacks throughout the summer:

> *"It went from 0-100 in no time at all. But it got scary because they were finding personal information about me and throwing it into however many public internet locations as possible. I had 70+ people sending me screenshots of an address associated with me for weeks."* – TS05

Both TS03 and TS05's experiences of these raids raised another concern—the targeted nature of the content of these messages. The carefully-crafted contents of these messages, in addition to expressing identity attacks against their targets, sought to threaten even the physical well-being of their targets via doxxing. The impact of such targeted attacks and the violent threats underlying doxxing even pushed one streamer to escalate their mitigation strategies beyond their stream:

> *"Law enforcement got involved, I had to find a lawyer, [the attackers] were threatening violence against my children... It was a scary time for me."* – TS05

TS05's experience was not unique. TS01 also expressed that others also experienced swatting as a result of being doxxed. These experiences of online harassment have manifested in potential psychological, physical, and even financial harm for already marginalized groups.

Through both the incessant and bespoke nature of these attacks, we found that these streamers' perceptions and experiences of hate raids defined them as highly-motivated attacks on individuals based on their identities, targeting Black and LGBTQ+ communities in particular; while attacks on these marginalized communities have always existed in online spaces, the severity and persistence of hate raids distinguishes them from what many members of these communities had experienced before.

**Cross-Platform Attacks**    As explained by several streamers, the targeted nature of these attacks resulted in a varied set of vectors threatening their psychological and physical safety. To better define the range of threats hate raids posed to streamers, we asked each participant to describe their experiences with hate raids and what attack vectors were used. We find that four of the seven streamers envisioned hate raids on Twitch as one piece of a larger campaign of harassment, highlighting the multi-platform nature of these orchestrated attacks. For instance, TS01's address and phone number were released in public locations off Twitch, and attackers even made videos on other platforms to help disseminate their personal information. This was then leveraged to flood their phone with calls. Similarly, TS02 and TS04 noted that Discord was another platform of concern; Discord servers of targeted streamers were attacked, and some of the hate raids were organized in Discord servers. TS04 described the complexity of the multi-platform nature of these attacks:

> "Where I find that companies really fall flat is understanding the impact of things that happen on their platform, the things that are planned on their platform and committed on another one... I think this is part of the issue with some of these hate raids is that it is personal info being hit. It's people's personal stuff outside of hate raids, outside of Twitch being shared. It is also being called slurs in chat, and that's harmful absolutely to be called slurs in chat and stuff like that. But it's also the fear of, well, my full name just got shared or my address just got shared. For me, it was like my Discord got hit, which is a whole other platform." – TS04

TS05 echoes these concerns, acknowledging that while hate raids originated with Twitch, "when someone makes it their mission to harm you, they'll look for whatever they can to access you." As a result, the high motivation involved in these attacks has raised questions and frustration within the community regarding platform accountability.

In tandem, our quantitative and qualitative data on hate raids indicate that the experiences of the streamers from Black and LGBTQ+ communities align with the media's portrayal of hate raids—that is, as highly-targeted and motivated attacks. However, through our quantitative analysis (Section 4.1.1), we also identified a variation of hate raids that deviated from this depiction, a form of hate raids akin to subcultural trolling that did not target specific streamers according to their identities, instead using "canned hate" to spread hate against Black and Jewish identities en masse in popular channels with high visibility. That is, these hate raids contained identity-based attacks, but were spread across the platform indiscriminately, indicating that an eagerness to cause widely-visible, attention-grabbing chaos may have also motivated the attackers.

## 4.2  Community response to hate raids

As hate raids swept the platform, the community's need and urgency for tools and resources to mitigate the threat grew. We performed a series of interviews with both streamers and bot developers involved with marginalized communities on Twitch to understand the following: (1)

how streamers and their communities addressed the threat of hate raids in the short term, (2) what array of tools and resources were assembled to mitigate the impact of hate raids, (3) the efficacy of the grassroots organization for #TwitchDoBetter and #ADayOffTwitch, and (4) the longer-term effects of hate raids on streamers and their communities.

*4.2.1 Short term responses by streamers.* Four of the seven streamers we interviewed expressed that they employed both proactive and reactive mitigation techniques to protect themselves from hate raids in the short term. The kinds of techniques varied, often depending on the severity of the threat. On one end of the spectrum, one streamer explained that because their attackers had escalated to threatening violence against their children, they involved law enforcement and retained a lawyer. While these vectors of attack were impossible to address solely on-platform, the majority of streamers experienced attacks that manifested within the Twitch ecosystem of chat and engagement notifications (e.g., follows and raids). As such, these streamers were able to mitigate some of the impact of hate raids via modifications to their streams' moderation protocols. One streamer added more users to their moderation teams, recruiting them from longtime members of their community who were "constantly hanging out inside of the chat" and "offering up their services... so that they can keep an eye on the chat," a pattern previously identified in [49] and [47]. Several streamers detailed variations of informational-support seeking, resource aggregation, and development of new tools in ways similar to those previously detailed in crisis informatics literature. For instance, one streamer described Stream Deck presets that were helpful for an emergency response; a Stream Deck is a physical control pad with preset studio settings (e.g., switching media scenes, camera angles, executing chat commands). They described commands that they added for moderation purposes:

> "We added more commands to like basically put it in follower mode and to turn off the chat to where it's only emotes only so that they can't put in any hateful words. Shutting things down for 10 minutes, but with a push of a button." – TS02

Similarly, another streamer outlined channel lockdown protocols they followed for hate raids, incorporating the idea of a "panic button" into a human moderator pipeline to handle incidents post-facto:

> "I had a panic button that turned off alerts, locked down chat, my mods would record times of incidents such as follow botting, we added different terms to the banned words list, had the highest auto mod settings available." – TS05

One variation of the Twitch Panic Button was developed and publicly advertised by nutty, a Twitch streamer, to be a rapid response mechanism integrated with a Stream Deck so that a single push of a button (or in customized cases, a voice-activated trigger phrase) enabled subscribers-only mode and cleared existing chat from both the chat client and the stream display [10]. Furthermore, after performing damage control, nutty's tool attempted to reclaim the stream space; for instance, in nutty's stream, the button triggered changing background lights and snarky automatically generated messages. An official tool with functionality similar to a panic button, named "Shield Mode," was rolled out by Twitch in late November, 2022.

In addition to the automated tools like the panic button, we found that streamers were aware of bot developers that developed bespoke features or new bots altogether to help handle the wave of hate raids. In our interviews, a streamer mentioned one bot in particular, Sery_Bot:

> "Also there's an additional thing that has been added to a lot of... streamers chats called Sery_Bot. Someone who isn't working for Twitch created a bot where it kind of shuts down all the other bots. Like it blocks them from being able to say anything. Or once they can come into your chat, it blocks them out. So a lot of us have added that." – TS02

Sery_Bot was developed by Sery, a developer who also sometimes streamed on Twitch. On August 14, 2021, Sery publicly solicited the Twitch community via Twitter for examples of hate raid messages and other relevant information to begin developing his bot. In the span of just a couple of weeks, Sery developed a variety of features—for instance, text-based commands in IRC like !hateraidon that performed a similar function to the panic button. In addition to providing utility already provided by other tools, Sery_Bot also integrated community-based block lists of account usernames to automatically check new chat messages against. Subsequent months entailed list updates, more feature development (e.g., checking account age of chatters and profile picture scanning for repeat offensive images, like swastikas). With the rollout of so many features so rapidly, Sery_Bot became viral, and in just two months, it amassed over 55,000 integrations over different Twitch channels.

TS02 highlighted both the effectiveness and widespread adoption of such a third-party bot amongst streamers in the context of hate raids; however, they also indicated that discomfort with or distrust of technology—particularly third-party bots—may have inhibited the adoption of tools and resources meant to mitigate such threats within the community.

*4.2.2 Resources from tool developers and other community members.* The short-term responses from streamers alluded to the availability of community-sourced tools and streamers' reliance on them to combat hate raids. Our interviews with both streamers and bot developers illuminated the various kinds of tools and resources the community created and what the development process was like in response to real-time threats. Streamers' perspectives gave insights into what kinds of tools were visible and widespread throughout the community. Four streamers explicitly mentioned the use of third-party bots for hate raid mitigation or prevention. One streamer noted that in addition to guides for making the aforementioned panic buttons, they were well-informed of the various bots that were developed individually by different bot developers, all for the purpose of responding to hate raids:

> *"There was Smash Bot that was created. Mix It Up Bot, Sery Bot, StopHateBot, WiseBot, time out bot. All of those things were kind of developed."* – TS03

Another streamer contrasted the fast wave of tool development by the community members with the poor communication and delayed response from Twitch:

> *"And yet, for some reason we have six queers in a trench coat who have somehow made all these tools in the span of three days for us to use that no, they don't eradicate the issue, but they definitely helped kind of mitigate it immensely. We had [user] who was making full master lists along with [user] of all the bots that were being created which... Twitch I think should reach out to them and get that master list and pay them for their work and say, these are all bots that are out doing things that are worthwhile or valuable."* – TS01

TS01 underscored that these tools were developed by members of the communities targeted by hate raids in a rapid-response fashion that Twitch as a platform simply could not; furthermore, TS01 emphasized that these quickly-developed tools were also effective in mitigating specifically the bot-mediated harassment even just with fairly naive methods.

To supplement our understanding of moderation bots' roles in the hate raid ecosystem, we interviewed two bot developers (referred to as BD in quote attributions) to understand (1) their perspectives and experiences with the community's needs, (2) Twitch as a platform for development, and (3) technical challenges they encountered while developing features for hate raids. One bot developer noted a sort of cat-and-mouse game between the community members and attackers,

resulting in fast-paced changes in the sophistication of hate raids and applicable mitigation techniques. This bot developer remarked on the primitive, manual nature of how hate raids started, only to quickly become more coordinated and varied:

> *"So they weren't that organized back then, started with a couple guys who just came into [a streamer]'s voice chat while he was streaming Phasmophobia, and they were shouting the N word, and then he gave them a really strong reaction by immediately ending the stream, deleting the VOD and so on. So they came back and spammed all nasty stuff in chat."* – BD02

This bot developer also illuminated some of the motivation for hate raid participants: to elicit a strong, disruptive reaction from the streamer. Similarly, there were approaches that the larger community initially employed, such as joining the attackers' Discord servers where hate raids are organized, that the attackers responded to:

> *"And it's also a little bit of a double-edged sword because we originally used to enter the offensive Discord servers where they would gather up and organize these hate raids with a second account, and then report the Discord server and the respective messages. And now of course, they learned that we do that, and they also set up a similar set of security measures."* – BD02

As such, both the attackers and the defenders in the hate raid ecosystem were made aware of what strategies the other side was employing, and they adapted their methods in response. From one bot developer's perspective, they analyzed the threat of hate raid messages and identified that attackers used "automated tools to just spam the chat," and in response began developing a bot as a counter measure. This bot developer noted that the hate spam sent by an army of bots was difficult to mitigate with manual moderation, so identified a bot-mediated moderation approach as an effective way to respond to the attack.

*4.2.3 Collective action for #TwitchDoBetter/#ADayOffTwitch.* In addition to the different ways the community attempted to better protect themselves from hate raids, there were also attempts by the community to raise awareness through collective action. Two hashtags, #TwitchDoBetter and #ADayOffTwitch, rallied support from Twitch users via Twitter. #TwitchDoBetter was started in an attempt to raise awareness of the harassment of targeted creators on Twitch. Subsequently, #ADayOffTwitch was a boycott of Twitch that took place on September 1, 2021, meaning participants would not stream, watch streams, or participate in any chats. This day took place with hopes that reduced engagement on the platform would highlight the urgency of better safety for creators on Twitch. We found through our interviews with streamers that while the community was able to raise awareness through these movements, there was some disagreement about their long-term impact within the community. One streamer we interviewed who was a co-organizer of the walkout felt that the movements were successful in meeting what they perceived to be the goal, which was to raise awareness:

> *"I think it worked in the way that I had hoped. We raised awareness. We actively called on a MAJOR streaming platform to make changes and we're seeing the fruits of our labor. It's not always about money like so many bigger streamers commented. Sometimes we have to understand that reputation is a currency."* – TS05

However, setting an open-ended goal of raising awareness for these community-organized movements did not satisfy another participant. In contrast, TS01 felt that ultimately, these movements failed to narrow in on concrete demands of Twitch and therefore failed to be as effective and impactful as they could have been:

*"I feel like they should have done a better job of sitting down and figuring out an actualized list of demands. Why are we taking a day off of Twitch? Because at face value, the reason we're taking the day off of Twitch is because hate raids suck and that's a true assessment. But what about after that? Why do those hate raids suck? What do we want to see to address those hate raids? How do we want Twitch to address it? How do we want Twitch to actually get engaged more about it? How do we want Twitch to respond to this? How does that carry over into future endeavors? What does that look like for a conversation around how they need to update their security? There was just little to no demand to be had whatsoever. And so, what could have been an actual movement or an actual kind of protest type of thing, just wound up being plainly speaking, a bunch of people just not logging in."* – TS01

While these movements were organized within the community, perceptions of their goals and effectiveness varied throughout. Still, despite conflict around the goal and organization of the movement, #ADayOffTwitch did indeed significantly impact the number of viewers and streamers engaging with the platform; per one external estimate, this movement led to up to 15% less engagement on Twitch overall during the walkout [39].

*4.2.4   Longer-term impacts of hate raids on streamers.* Even with all of the mitigation attempts and movements to raise awareness, hate raids undoubtedly caused distress and skepticism throughout the community. Our interviews with streamers indicate that the visceral nature of these attacks paired with Twitch's response has largely shaped their views of the platform as unprepared and detached from the community's suffering. All seven of the participants expressed disappointment with Twitch's failure to consider abuse protections proactively, the slow rollout of features and tools to mitigate the harm of hate raids, and poor communication with between Twitch and stakeholders. One streamer expressed resigned frustration that this experience had been consistent with Twitch's attitudes toward protecting its at-risk communities in the past:

*"I think Twitch's response has been absolutely abysmal. I think that very frankly speaking, it's pretty pathetic. Twitch has a longstanding historical track record of not knowing how to communicate ever at all. So, while I do think that their communication for this was abysmal, I would be remiss to omit the part where it is exactly what I expected them to do. And Twitch is going to continue falling on their face over topics of this nature and conversations of this type every single time so long as they insist that silence is the best solution. And they need to do better than put out a simple tweet saying, 'We hear you, we see you and we want you to know that we care, we promise.'"* – TS01

Streamers also felt that poor communication even around existing mitigation tools led to unnecessary chaos during hate raids. One streamer noted that enabling two-factor authentication was a common suggestion by the community and Twitch for streamers to protect themselves. However, these suggestions conflate the verification of user identity with that of accounts chatting in that user's channel. Therefore, even if there are tools that may better protect individuals, poor communication may lead to the misuse of the tool or misinterpretation of the protections it actually offers.

    In addition to disappointment with Twitch's communication response, several streamers lamented the lack of tooling Twitch had prepared for such attacks, even for features that had been requested in the past or existed on other platforms:

*"The chat verification tools [released at the tail end of the hate raids] are really nice, and I think that that's what a lot of people have wanted for so long. I'm not sure exactly why it took them that long to implement. I feel like it should have been implemented."* – TS03

> *"I went to Twitch HQ for [a Black History summit in the past], and that was one of the things that all of us echoed and said and was like, 'If I banned someone, they should not be able to continue consuming my content.' It needs to be like some of these other sites. Twitter and Facebook are perfect examples. When I block somebody, it's scorched earth. As far as they're concerned, I no longer exist to them. That's what it needs to be."* – TS01

These embody some of the frustrations that streamers have had for baseline protection features that the community had been wanting for years. The overall emotional harm of hate raids was even enough to dissuade some members to leave Twitch or even streaming altogether; one streamer explains that the hate raids gave them a lot of anxiety, and that they know streamers who "walked away entirely because of that anxiety and distress." Another streamer expressed concern that their protective measures to prevent anomalous viewership might even affect their stream's long-term growth:

> *"We have things that we're trying to do at all times and if we blockade the people who want to watch us, they are going to inherently want to move on elsewhere."* – TS01

More broadly, the threat of hate raids and their impact on streamers in the future still looms over several of the participants. One streamer noted that, while there may be a sense of fatigue among the community concerning hate raids, the lack of recent publicity over hate raids may not be an indication that the larger threat has passed:

> *"I really think that it's either happen[ing] less or because we've been dealing with it now for so long, people are just... There's only so many times that you can post and be like, 'Yep, got hate raided again today. Yep, got hate raided again today.' So that could also be a factor as to why I'm not seeing it as much on Twitter."* – TS03

## 5 DISCUSSION

In this paper, we make three primary contributions: (1) the descriptive characterization of a novel form of long-term harassment campaigns on livestreaming platforms; (2) the definition of hate raids as a dually-motivated phenomenon: first as a hate-driven attack, and second as an act of seeking attention; (3) the observation that members of targeted communities rapidly responded to the threat of hate raids to address the shortcomings of protections provided by Twitch. In each of the following paragraphs, we elaborate on these contributions.

### 5.1 Characterization of Hate Raids

Although hate raids on Twitch caused significant disruption and emotional harm to streamers, these attacks were relatively technically unsophisticated. Accounts were created en masse (likely in an automated fashion) to serve a single purpose, hateful comments were largely identical across channels, and user-specified identity tags were operationalized to attack marginalized groups. Many of these tactics might have been prevented if Twitch had followed established trust and safety practices like rate-limiting account creation [57], adding a delay between account creation and platform participation, deploying additional identity verification requirements (e.g., SMS or phone) [58],[6] and protecting at-risk streamers by safeguarding automated access to sensitive data, such as identity-based channel tags.

Although there are trade-offs between adding friction in joining communities and protecting users from abuse, the security practices employed by Twitch at the time of this wave of hate raids did not deter these relatively unsophisticated attacks. The broader history of Trust and Safety is often characterized by *reactive* feature development as attack vectors become apparent on each

---

[6]Phone verification was added as a feature during the later phases of hate raids, indicating that it may have been under development but not yet released when this wave of hate raids began.

specific platform, but a common set of forms of attacks have appeared many times throughout the history of social platforms and, as in this case, they do not become substantially more sophisticated as they are ported from platform to platform [16, 29, 31]. The appearance of these attacks and the form that they take on any new platform is often predictable, and it is much easier to build safeguards during earlier development phases than to be forced to reactively add them under time pressure when crises arise. Future community-driven platforms should prioritize the allocation of resources to teams developing defensive tactics a necessary first step for curbing online abuse of this nature before it causes significant harm.

The qualitative accounts of streamers' experiences that we examine affirm that highly-targeted hate raids can lead to long-term emotional distress and can even threaten streamers' physical safety. While Twitch has begun to take steps to combat hate raids via automated tooling (e.g., AutoMod), optional account verification methods, and the aforementioned Shield Mode, the threat of highly-motivated hate raids coordinated off-platform continues to loom over its streamers. In March 2022, a wave of hate raids orchestrated by streamers on Cozy.tv, a livestreaming platform founded by far-right white nationalist Nick Fuentes, hit Twitch, this time targeting women and LGBTQ+ streamers with homophobic, transphobic, and misogynistic messages in their Twitch channels, direct messages, and Discord servers [43].[7] As hate raids continue to threaten streamers with varying degrees of off-platform coordination, legitimate user participation, and bot account manipulation, the need for platforms to consider both increasingly sophisticated threat models and historically common patterns of attacks has only grown. Platforms must consider preemptively what their policies, protection, and communication processes will be, and by designing these mechanisms for the needs of their at-risk communities, they can better protect all of their users. The design of proactive prevention measures that do not disproportionately burden or disadvantage marginalized communities—with respect to their online engagement and technical overhead—remains an important question for future research and development.

### 5.2 Dual Motivation of Hate Raids

We draw several primary characteristics from Marwick [30] and Phillips [41] as a baseline to compare hate raids with: first, per Phillips, subcultural trolling benefits from (and to some extent relies on) amplification [41, pp. 3–6, 56–61], and fits within existing media narratives, often referencing mainstream concepts and/or publicized events [41, pp. 115–118] in absurd or repurposed ways. Second, per Marwick, morally-motivated networked harassment also benefits from amplification, but it also relies heavily on identity and identity conflicts to justify harassment campaigns that have none of the underlying absurd logic that characterize subcultural trolling [30, pp. 5–8]. Moreover, where subcultural trolling originates from specific communities, often in planned, targeted attacks, morally-motivated networked harassment often originates more organically and is partially self-amplifying through the properties of networks such as those on Twitter. As we discussed in Section 4.1, the hate raids on Twitch share variants of each of these characteristics, and we therefore argue that they lie in a space between subcultural trolling and morally-motivated networked harassment.

### 5.3 Stakeholder Rapid Response

We observed many similar behaviors in Twitch hate raids that occur during natural disaster response as documented in crisis informatics literature — informational-support seeking, aggregation of

---

[7]These hate raids were performed manually, and as such would likely not have been deterred by security measures designed to prevent automated attacks from bots; however, their occurrence represents a continued threat to streamers from marginalized groups on the platform.

resources, and development of new tools and technologies to address specialized needs arising from the crisis. We also observed the use of social media for social support-seeking and solidarity, even leading to the organization of a significant protest.

During the hate raids, there was no formal organization (e.g., a state or federal government) coordinating public response, as is often the case in the aftermath of natural disasters; while Twitch did respond to the hate raids in several ways, these did not involve coordinating with its users at any scale. As such, responses to hate raids more closely resembled those documented in literature on longer-term conflicts where public institutions play less of a role because they have been weakened as a result of the conflict [33, pp. 2–3].

Users were able to respond to rapidly evolving situations during the hate raids in ways that brought relief to their communities far more quickly than Twitch was able to. They developed and rapidly iterated on tools to counter the attacks and improved those tools as attacks changed. The first of these tools appeared within days of when the hate raids started to gain public attention. Users also created guides on how to use Twitch's moderation features and Discord's moderation features (for streamers who had servers affiliated with their streams), and on how streamers could better protect their personal information. Guides for all of these already existed in forms created by the respective platforms, but the community-created guides gained significantly more traction in this case because of their applicability to the specific circumstances of hate raids and because of the shared trust between community members.

With this work we do not mean to suggest that, because users were effective in rapidly responding to these issues, Twitch should cede their authority to users on Trusty & Safety issues. Instead, we note that Twitch and its users each have different strengths in how they are able to respond, and that Twitch and other platforms with similar moderation structures could gain much value from better communication and collaboration with users on moderation problems that arise. Volunteer moderators' domain-specific knowledge and reputational trust paired with the findings from prior work showing that experienced moderators can successfully onboard volunteers into new moderation contexts [46] suggests that Twitch as a platform can gain insight and trust from their users by building connections with power users (e.g., prominent community tool developers like Sery). By consulting with such users, Twitch can also improve the dissemination of resource guides and the visibility of community-built tools. This access to information may be particularly effective in enhancing coordinated action because users are far more agile than the platform in organizing and producing tools to respond to imminent threats. As both Seering et al. and Roberts suggest [44, 46], the use of volunteer moderation for commercial platforms brings to question the ethics in the division of labor between volunteers and platforms. We argue that platforms should consult with power users to improve communication and tooling, and that these platforms should consider paying such power users for their valuable, contextual knowledge to compensate them for their large contributions to their communities.

Finally, we reiterate the recommendations of prior work [1] rooted in intersectional feminist theory: that platforms must center the needs of their most marginalized, vulnerable users in their design. Platforms designed around existing structural inequalities recreate and further disseminate these systems of oppression [35]. We argue that addressing the needs of the oppressed more effectively encompasses the needs of all users, allowing platforms to be better prepared to mitigate inevitable attempts of abuse.

## 5.4 Limitations

We acknowledge that our analysis is not based on a comprehensive view of the platform. Because smaller communities may have a tacit expectation of privacy, we intentionally did not collect chat data from channels with less than an average of 100 viewers. However, many communities targeted

by hate raids were not necessarily large, mainstream channels. According to Twitch, as of 2018, 81.5% of its creators and viewers were male [61], and user surveys have shown that a majority are white. As such, we expect that some of the highly-targeted hate raiding behavior was not captured in our large-scale data collection methodology. Furthermore, our hate raid detection mechanism was based on community-aggregated lists of known malicious bots. Because of this, we may not have detected categories of hate raids that were not actively documented by community members. This likely narrows the variance in attack structure and message content flagged in our dataset. Even with these limitations, however, we argue that our quantitative perspective still provides insights into various technical characteristics and attacker motivations of hate raids. Particularly, when paired with our qualitative results that *specifically* seek the perspectives of targeted community members, we believe that we are able to capture multiple facets of a nuanced and dynamic threat model.

## 6 CONCLUSION

Our large-scale quantitative measurement of hate raids across mainstream channels on Twitch and interviews with community members from targeted groups confirm that hate raids are indeed highly-targeted and hate-driven attacks. Our quantitative analysis reveals an additional mode of hate raid, however, that is similar to subcultural trolling and networked harassment. We find that the technical characteristics of these attacks mirror many of the naïve methods of other forms of online abuse, such as spam. The content of these hate raid messages are deeply entrenched in two main hateful ideologies: anti-Black racism and antisemitism. Our interviews demonstrate the various approaches—both proactive and reactive—to defense that the community took in response to hate raids. Our analysis furthers our understanding of the complexities in the ecosystem surrounding hate raids, highlights lessons to be learned in designing proactive harassment mitigation into a platform from the start, and brings attention to the interplay between platform and community governance in the face of a collective crisis.

## 7 ACKNOWLEDGEMENTS

## REFERENCES

[1] Lindsay Blackwell, Jill Dimond, Sarita Schoenebeck, and Cliff Lampe. 2017. Classification and Its Consequences for Online Harassment: Design Insights from HeartMob. *Proc. ACM Hum.-Comput. Interact.* 1, CSCW, Article 24 (Dec. 2017), 19 pages. https://doi.org/10.1145/3134659

[2] Jie Cai and Donghee Yvette Wohn. 2022. Coordination and Collaboration: How do Volunteer Moderators Work as a Team in Live Streaming Communities? *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems* (2022).

[3] Carlos Castillo. 2016. *Big crisis data: social media in disasters and time-critical situations.* Cambridge University Press, Cambridge, UK.

[4] Eshwar Chandrasekharan, Chaitrali Gandhi, Matthew Wortley Mustelier, and Eric Gilbert. 2019. Crossmod: A Cross-Community Learning-Based System to Assist Reddit Moderators. *Proc. ACM Hum.-Comput. Interact.* 3, CSCW, Article 174 (Nov. 2019), 30 pages. https://doi.org/10.1145/3359276

[5] Justin Cheng, Michael Bernstein, Cristian Danescu-Niculescu-Mizil, and Jure Leskovec. 2017. Anyone Can Become a Troll: Causes of Trolling Behavior in Online Discussions. In *Proceedings of the 2017 ACM Conference on Computer

*Supported Cooperative Work and Social Computing* (Portland, Oregon, USA) *(CSCW '17)*. ACM, New York, NY, USA, 1217–1230. https://doi.org/10.1145/2998181.2998213

[6] Marc Cheong and Vincent C S Lee. 2011. A microblogging-based approach to terrorism informatics: Exploration and chronicling civilian sentiment and response to terrorism events via Twitter. *Information Systems Frontiers* 13, 1 (2011), 45–59. https://doi.org/10.1007/s10796-010-9273-x

[7] Bill Cooney. [n.d.]. New Twitch stats reveal how few viewers are needed to be a "top" streamer. *Dexerto* ([n. d.]). https://www.dexerto.com/entertainment/new-twitch-stats-reveal-how-few-viewers-are-needed-to-be-a-top-streamer-1527638/

[8] Amanda C Cote. 2017. "I Can Defend Myself": Women's Strategies for Coping With Harassment While Gaming Online. *Games and Culture* 12, 2 (2017), 136–155. https://doi.org/10.1177/1555412015587603

[9] John W Creswell. 2013. *Qualitative Inquiry and Research Design: Choosing Among Five Traditions*. SAGE, Thousand Oaks, CA.

[10] Ralston Dacanay. [n.d.]. Twitch Streamer Creates Third-Party 'Panic Button' to Counter Hate Raids. *DBLTAP* ([n. d.]). https://www.dbltap.com/posts/twitch-streamer-creates-third-party-panic-button-to-counter-hate-raids-01fekqmgacvw

[11] Cecilia D'Anastasio. [n.d.]. Twitch Sues Users Over Alleged 'Hate Raids' Against Streamers. *Wired* ([n. d.]). https://www.wired.com/story/twitch-sues-users-over-alleged-hate-raids/

[12] Srayan Datta and Eytan Adar. 2019. Extracting inter-community conflicts in reddit. In *Proceedings of the international AAAI conference on Web and Social Media*.

[13] Julian Dibbell. 1993. A Rape in Cyberspace: How an Evil Clown, a Haitian Trickster Spirit, Two Wizards, and a Cast of Dozens Turned a Database Into a Society. *The Village Voice* December 23 (1993), 36–42. https://www.villagevoice.com/2005/10/18/a-rape-in-cyberspace/

[14] Tuğrulcan Elmas, Rebekah Overdorf, Ahmed Furkan Özkalay, and Karl Aberer. 2022. Ephemeral Astroturfing Attacks: The Case of Fake Twitter Trends. In *EuroS&P '22*.

[15] Jesse Fox and Wai Yen Tang. 2017. Women's experiences with general and sexual harassment in online video games: Rumination, organizational responsiveness, withdrawal, and coping strategies. *New Media & Society* 19, 8 (2017), 1290–1307. https://doi.org/10.1177/1461444816635778

[16] Sara Gordon. 1994. IRC and Security — Can the two co-exist? (1994).

[17] Kishonna L Gray. 2017. They're just too urban": Black gamers streaming on Twitch. In *Digital Sociologies*, Jessie Daniels, Karen Gregory, and Tressie McMillan Cottom (Eds.). Vol. 1. Policy Press, Bristol, England, Chapter 22, 355–368.

[18] Nathan Grayson. [n.d.]. Marginalized streamers beg Twitch to 'do better' in wake of hate raids, poor pay. *The Washington Post* ([n. d.]). https://www.washingtonpost.com/video-games/2021/08/11/twitch-do-better-hate-raids/

[19] Nathan Grayson. [n.d.]. Twitch hate raids are more than just a Twitch problem, and they're only getting worse. *The Washington Post* ([n. d.]). https://www.washingtonpost.com/video-games/2021/08/25/twitch-hate-raids-streamers-discord-cybersecurity/

[20] Susan Herring, Kirk Job-Sluder, Rebecca Scheckler, and Sasha Barab. 2002. Searching for Safety Online: Managing "Trolling" in a Feminist Forum. *The Information Society* 18, 5 (2002), 371–384. https://doi.org/10.1080/01972240290108186

[21] Dylan Horetski. [n.d.]. Twitch hate raids return in massive wave of attacks on LGBTQIA+ streamers. *Dexerto* ([n. d.]). https://www.dexerto.com/entertainment/twitch-hate-raids-return-in-massive-wave-of-attacks-on-lgbtqia-streamers-1781574/

[22] Jialun Aaron Jiang, Charles Kiene, Skyler Middler, Jed R. Brubaker, and Casey Fiesler. 2019. Moderation Challenges in Voice-based Online Communities on Discord. *Proc. ACM Hum.-Comput. Interact.* 3, CSCW, Article 55 (Nov. 2019), 23 pages. https://doi.org/10.1145/3359157

[23] Chris Kanich, Christian Kreibich, Kirill Levchenko, Brandon Enright, Geoffrey M Voelker, Vern Paxson, and Stefan Savage. 2008. Spamalytics: An empirical analysis of spam marketing conversion. In *15th ACM conference on Computer and communications security*.

[24] Charles Kiene, Andrés Monroy-Hernández, and Benjamin Mako Hill. 2016. Surviving an "Eternal September": How an Online Community Managed a Surge of Newcomers. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems* (San Jose, California, USA) *(CHI '16)*. ACM, New York, NY, USA, 1152–1156. https://doi.org/10.1145/2858036.2858356

[25] Srijan Kumar, Justin Cheng, Jure Leskovec, and V.S. Subrahmanian. 2017. An Army of Me: Sockpuppets in Online Discussion Communities. In *WWW '17*.

[26] Chen Ling, Utkucan Balci, Jeremy Blackburn, and Gianluca Stringhini. 2021. A First Look at Zoombombing. In *IEEESP*.

[27] Richard MacKinnon. 1997. Virtual Rape. *Journal of Computer-Mediated Communication* 2, 4 (1997), 1–2. https://doi.org/10.1111/j.1083-6101.1997.tb00200.x

[28] Kaitlin Mahar, Amy X. Zhang, and David Karger. 2018. Squadbox: A Tool to Combat Email Harassment Using Friendsourced Moderation. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems* (Montreal

QC, Canada) *(CHI '18)*. ACM, New York, NY, USA, Article 586, 13 pages. https://doi.org/10.1145/3173574.3174160

[29] Enrico Mariconti, Guillermo Suarez-Tangil, Jeremy Blackburn, Emiliano De Cristofaro, Nicolas Kourtellis, Ilias Leontiadis, Jordi Luque Serrano, and Gianluca Stringhini. [n.d.]. "You Know What to Do" Proactive Detection of YouTube Videos Targeted by Coordinated Hate Attacks. *Proceedings of the ACM on Human-Computer Interaction* CSCW ([n. d.]).

[30] Alice E Marwick. 2021. Morally Motivated Networked Harassment as Normative Reinforcement. *Social Media + Society* 7, 2 (2021), 20563051211021378. https://doi.org/10.1177/20563051211021378

[31] Adrienne Massanari. 2017. #Gamergate and The Fappening: How Reddit's algorithm, governance, and culture support toxic technocultures. *New Media & Society* 19, 3 (2017), 329–346.

[32] Danaë Metaxa-Kakavouli, Paige Maas, and Daniel P. Aldrich. 2018. How Social Ties Influence Hurricane Evacuation Behavior. *Proc. ACM Hum.-Comput. Interact.* 2, CSCW, Article 122 (nov 2018), 16 pages. https://doi.org/10.1145/3274391

[33] Andrés Monroy-Hernández, danah boyd, Emre Kiciman, Munmun De Choudhury, and Scott Counts. 2013. The New War Correspondents: The Rise of Civic Media Curation in Urban Warfare. In *Proceedings of the 2013 Conference on Computer Supported Cooperative Work* (San Antonio, Texas, USA) *(CSCW '13)*. Association for Computing Machinery, New York, NY, USA, 1443–1452. https://doi.org/10.1145/2441776.2441938

[34] David Moore, Colleen Shannon, Douglas J Brown, Geoffrey M Voelker, and Stefan Savage. 2006. Inferring internet denial-of-service activity. *ACM Transactions on Computer Systems (TOCS)* (2006).

[35] Safiya Umoja Noble. 2018. *Algorithms of Oppression.* NYU Press, New York, NY, USA. https://doi.org/10.2307/j.ctt1pwt9w5

[36] Leysia Palen and Kenneth M Anderson. 2016. Crisis informatics–New data for extraordinary times. *Science* 353, 6296 (2016), 224–225. https://doi.org/10.1126/science.aag2579

[37] Leysia Palen, Sarah Vieweg, Jeannette Sutton, Sophia B Liu, and Amanda Hughes. 2007. Crisis informatics: Studying crisis in a networked world. In *Proceedings of the Third International Conference on E-Social Science.* 7–9.

[38] Manish Pandey. [n.d.]. Twitch announces new tools to fight hate raids. *BBC* ([n. d.]). https://www.bbc.com/news/newsbeat-58594732

[39] Ash Parrish. [n.d.]. Twitch viewership noticeably dropped when streamers took a day off in protest. *The Verge* ([n. d.]). https://www.theverge.com/2021/9/2/22654534/streamers-twitch-walkout-viewership-drop

[40] Whitney Phillips. 2011. LOLing at tragedy: Facebook trolls, memorial pages and resistance to grief online. *First Monday* 16, 12 (2011), 12 pages. https://doi.org/10.5210/fm.v16i12.3168

[41] Whitney Phillips. 2015. *This is why we can't have nice things: Mapping the relationship between online trolling and mainstream culture.* MIT Press, Cambridge, MA, USA.

[42] Whitney Phillips. 2019. It Wasn't Just the Trolls: Early Internet Culture, "Fun," and the Fires of Exclusionary Laughter. *Social Media + Society* 5, 3 (2019), 2056305119849493. https://doi.org/10.1177/2056305119849493

[43] Blaine Polhamus. [n.d.]. Hate raids return to Twitch, another wave of attacks target LGBTQIA+ streamers. *DOT ESPORTS* ([n. d.]). https://dotesports.com/streaming/news/hate-raids-return-to-twitch-another-wave-of-attacks-target-lgbtqia-streamers

[44] Sarah T. Roberts. 2016. Commercial Content Moderation: Digital Laborers' Dirty Work. In *The Intersectional Internet: Race, Sex, Class and Culture Online*, Safiya Umoja Noble and Brendesha M. Tynes (Eds.). Peter Lang Digital Formations series, New York, NY, USA, 147–160.

[45] Morgan Klaus Scheuerman, Stacy M. Branham, and Foad Hamidi. 2018. Safe Spaces and Safe Places: Unpacking Technology-Mediated Experiences of Safety and Harm with Transgender People. *Proc. ACM Hum.-Comput. Interact.* 2, CSCW, Article 155 (Nov. 2018), 27 pages. https://doi.org/10.1145/3274424

[46] Joseph Seering, Brianna Dym, Geoff Kaufman, and Michael Bernstein. 2022. Pride and Professionalization in Volunteer Moderation: Lessons for Effective in Platform-User Collaboration. *Journal of Online Trust and Safety* (2022).

[47] Joseph Seering and Sanjay R. Kairam. 2022. Who Moderates on Twitch and What Do They Do? Quantifying Practices in Community Moderation on Twitch. *Proc. ACM Hum.-Comput. Interact.* 7, GROUP, Article 18 (dec 2022), 18 pages. https://doi.org/10.1145/3567568

[48] Joseph Seering, Robert Kraut, and Laura Dabbish. 2017. Shaping Pro and Anti-Social Behavior on Twitch Through Moderation and Example-Setting. In *Proceedings of the 2017 ACM Conference on Computer Supported Cooperative Work and Social Computing* (Portland, Oregon, USA) *(CSCW '17)*. ACM, New York, NY, USA, 111–125. https://doi.org/10.1145/2998181.2998277

[49] Joseph Seering, Tony Wang, Jina Yoon, and Geoff Kaufman. 2019. Moderator engagement and community development in the age of algorithms. *New Media & Society* 21, 7 (2019), 1417–1443. https://doi.org/10.1177/1461444818821316

[50] Bryan Semaan and Gloria Mark. 2011. Technology-Mediated Social Arrangements to Resolve Breakdowns in Infrastructure during Ongoing Disruption. *ACM Trans. Comput.-Hum. Interact.* 18, 4, Article 21 (12 2011), 21 pages. https://doi.org/10.1145/2063231.2063235

[51] Anna DuVal Smith. 1999. Problems of Conflict Management in Virtual Communities. In *Communities in Cyberspace* (1st ed.), Marc A Smith and P Kollock (Eds.). Routledge, New York, NY, USA, 135–166.

[52] Robert Soden and Leysia Palen. 2014. From Crowdsourced Mapping to Community Mapping: The Post-earthquake Work of OpenStreetMap Haiti. In *COOP 2014 - Proceedings of the 11th International Conference on the Design of Cooperative Systems, 27-30 May 2014, Nice (France)*, Chiara Rossitto, Luigina Ciolfi, David Martin, and Bernard Conein (Eds.). Springer International Publishing, Cham, 311–326.

[53] Robert Soden and Leysia Palen. 2016. Infrastructure in the Wild: What Mapping in Post-Earthquake Nepal Reveals about Infrastructural Emergence. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems* (San Jose, California, USA) *(CHI '16)*. Association for Computing Machinery, New York, NY, USA, 2796–2807. https://doi.org/10.1145/2858036.2858545

[54] Bijan Stephen. [n.d.]. The lockdown live-streaming numbers are out, and they're huge. *The Verge* ([n. d.]). https://www.theverge.com/2020/5/13/21257227/coronavirus-streamelements-arsenalgg-twitch-youtube-livestream-numbers

[55] Hibby Thach, Samuel Mayworm, Daniel Delmonaco, and Oliver Haimson. 2022. (In)visible moderation: A digital ethnography of marginalized users and content moderation on Twitch and Reddit. *New Media & Society* (2022).

[56] Kurt Thomas, Devdatta Akhawe, Michael Bailey, Dan Boneh, Sunny Consolvo, Nicola Dell, Zakir Durumeric, Patrick Gage Kelley, Deepak Kumar, Damon McCoy, Sarah Meiklejohn, Thomas Ristenpart, and Gianluca Stringhini. 2021. SoK: Hate, Harassment, and the Changing Landscape of Online Abuse. In *IEEESP*.

[57] Kurt Thomas, Chris Grier, Dawn Song, and Vern Paxson. 2011. Suspended accounts in retrospect: an analysis of twitter spam. In *ACM SIGCOMM conference on Internet measurement conference*.

[58] Kurt Thomas, Dmytro Iatskiv, Elie Bursztein, Tadek Pietraszek, Chris Grier, and Damon McCoy. 2014. Dialing back abuse on phone verified accounts. In *ACM SIGSAC Conference on Computer and Communications Security*.

[59] Inc. Twitch Interactive. [n.d.]. Celebrate Yourself and Your Community with 350+ New Tags. *Twitch Blog* ([n. d.]). https://blog.twitch.tv/en/2021/05/26/celebrate-yourself-and-your-community-with-350-new-tags/

[60] Jessica Vitak, Kalyani Chadha, Linda Steiner, and Zahra Ashktorab. 2017. Identifying Women's Experiences With and Strategies for Mitigating Negative Effects of Online Harassment. In *Proceedings of the 2017 ACM Conference on Computer Supported Cooperative Work and Social Computing* (Portland, Oregon, USA) *(CSCW '17)*. ACM, New York, NY, USA, 1231–1245. https://doi.org/10.1145/2998181.2998337

[61] Adam Yosilewitz. [n.d.]. StreamElements Analysis on Twitch Bullying. *StreamElements* ([n. d.]). https://blog.streamelements.com/streamelements-analysis-on-twitch-bullying-c3f2b2240318

[62] Himanshu Zade, Kushal Shah, Vaibhavi Rangarajan, Priyanka Kshirsagar, Muhammad Imran, and Kate Starbird. 2018. From Situational Awareness to Actionability: Towards Improving the Utility of Social Media Data for Crisis Response. *Proc. ACM Hum.-Comput. Interact.* 2, CSCW, Article 195 (nov 2018), 18 pages. https://doi.org/10.1145/3274464

## A PRIMARY INTERVIEW QUESTIONS

(1) Over the course of the past few months, have you been impacted either directly or indirectly by hate raids?
   (a) If so, how?
   (b) (If they were hate raided or observed hate raids) Can you describe what the hate raid(s) were like?
(2) What did you do to protect yourself from hate raids if anything?
   (a) (If not mentioned) Did you add any new moderation tools?
   (b) (If not mentioned) Did you add new moderators or request additional moderator support?
(3) How effective were each of these strategies for protecting yourself?
(4) How did you learn about new ways to protect yourself?
   (a) (If not mentioned) How did you learn about how to use new tools?
   (b) (If not mentioned) How did you learn about any new strategies for protecting your personal information?
   (c) (If not mentioned) Were you involved in any spaces were these topics were discussed?
(5) Were you involved in any forms of collective action like #TwitchDoBetter or #ADayOffTwitch?
(6) How do you feel about Twitch's response to the Hate Raids?
   (a) How do you feel about the lawsuit that Twitch has announced against the perpetrators?
   (b) Do you think that the new moderation features Twitch released have helped?
   (c) How do you feel about the way Twitch communicated about the Hate Raids in August and September?

(7) What do you think Twitch could do better in the future in handling cases like this one?

## B INTERVIEW CODEBOOKS

| Streamer Category Label | Description |
| --- | --- |
| Effectiveness of Twitch's responses | Chunks about specific aspects of Twitch's responses, including communication, lawsuit, tools they added, etc. |
| Instrumental community support | Chunks about community resource sharing and instrumental/informational support |
| Social community support | Chunks about social/interpersonal support they received |
| Community organization | Chunks about collective action or group-organized things, e.g., #TwitchDoBetter |
| Degree of raid targeting | Degree of attack personalization (how targeted?) |
| Frequency of hate raids experienced | Frequency (how often?) |
| Raid vectors | Different attack vectors (how many different ways) |
| Raid responses (short-term) | Things the streamer did in-the-moment or during the weeks while the hate raids were going on to protect themselves/others |
| Raid impact (long-term) | Longer-term impact on streamers' careers, health, well-being |

| Bot Developer Category Label | Description |
| --- | --- |
| Community need | Chunks about how a need for specific third-party resources for the community revealed, if streamers ask specifically for features, developers' own observations/pro-social motivations, and gaps/shortcomings in Twitch-provided tools |
| Developer dependence | Chunks about the degree to which streamers (large vs. small might have different experiences) depend on third-party bot developers to better protect themselves from hate raids, how many channels (how was the adoption), how effective (numbers of raids/bots/messages intercepted or moderated) |
| Hate raid arms race | Chunks about the kind of arms race or "cat and mouse game" bot developers experienced while rolling out features to combat hate raids |
| Effectiveness of Twitch tools | Chunks about bot developers' perspectives on the efficacy of Twitch's technical tools before/after the hate raids |
| Twitch development obstacles | Chunks about Twitch obstacles/hurdles that made effective bot development difficult |
| Twitch communication | Chunks about Twitch communications with the community (and how it fueled their dissatisfaction with the platform) |
| Developer coordination | Chunks about how the community of developers organized |