

# The Economics of Content Moderation: Theory and Experimental Evidence from Hate Speech on Twitter

Rafael Jiménez-Durán\*

This draft: December 1, 2022

## Abstract

Social media platforms ban users and remove posts to moderate their content. This “speech policing” remains controversial because little is known about its consequences and the costs and benefits for different individuals. I conduct two pre-registered field experiments on Twitter to examine the effect of moderating hate speech on user behavior and welfare. Randomly reporting posts for violating the rules against hateful conduct increases the likelihood that Twitter removes them. Reporting does not affect the activity on the platform of the posts’ authors or their likelihood of reposting hate, but it does increase the activity of those attacked by the posts. These results are consistent with a model in which content moderation is a quality decision for platforms that increases user engagement and hence advertising revenue. The second experiment shows that changing users’ perceived content removal does not change their willingness to pause using social media, a measure of consumer surplus. My results imply that content moderation does not necessarily moderate users, but it can marginally increase advertising revenue. It can be consistent with both profit and welfare maximization as long as out-of-platform externalities are small.

JEL codes: C93, D12, D85, D90, I31, J15, L82, L86, Z13

Keywords: social media, moderation, report, hate speech, experiment, welfare

---

\*Social Science Research Council and Chicago Booth Stigler Center. Email: [rafaeljjd@uchicago.edu](mailto:rafaeljjd@uchicago.edu). I’m deeply grateful to Leonardo Bursztyn, Pietro Tebaldi, Ali Hortaçsu, and Kevin Murphy for their guidance and support. I also thank Hunt Allcott, Michael Dinerstein, Matt Gentzkow, Justin Holz, Diego Jiménez, John List, Karsten Müller, Casey Mulligan, David Novgorodsky, Javier Pérez, Marta Prato, Carlo Schwarz, David Yanagizawa-Drott, Pinar Yildirim, Luigi Zingales, seminar participants at the 2022 NBER Summer Institute on IT and Digitization, Bocconi, Microsoft Research, Télécom Paris, Georgetown, Oxford, Ofcom, Purdue, ITAM, Banco de México, Zurich ETH, Northwestern Law, Queen’s, Stigler Center, and many, many more for helpful comments and suggestions. The Political Economics Initiative from the Becker Friedman Institute at the University of Chicago and the Stigler Center at Booth funded this research. The University of Chicago IRB approved the experiments (IRB18-1586, IRB21-1067), and I pre-registered them at the AEA RCT Registry (0005131, 0008189). Any errors are mine.

# 1 Introduction

Social media is the “modern public square,” according to the U.S. Supreme Court<sup>1</sup>—a place where speech happens among individuals with different backgrounds and ideologies. Yet, the biggest strengths of platforms—their size and diversity—also represent their greatest challenges. Forty percent of people have experienced online harassment (Anti-Defamation League, 2021), and studies document real-world consequences of online speech on hate crimes (Bursztyn et al., 2019; Müller and Schwarz, 2020a), and election outcomes (Fujiwara et al., 2021). Despite these consequences, few governments have crafted laws or regulation of online content (Carlson, 2021).

As a result, social media companies *self-regulate* and issue community guidelines that forbid not only illegal content but also some combination of hate speech, misinformation, harassment, spam, sexual content, and graphical content (Gillespie, 2018). Platforms moderate content by enforcing these guidelines with sanctions such as removing posts or accounts. Still, even if it is widespread, “policing speech” remains controversial (Kaye, 2019), in part due to scarce data and studies about this practice. The debate oscillates between arguments about freedom of expression (Strossen, 2018) and the harms that can be caused by online content (Waldron, 2012).

This paper contributes to the discussion by providing theory and experimental evidence of how moderation works, how it changes online behavior, and how to weigh its welfare gains and losses to different users. Guided by a model, I run two large-scale field experiments to document the consequences of content moderation on user behavior and welfare. The focus is on hate speech on Twitter as a prominent example—one-quarter of U.S. adults use this platform, and hate speech is its most sanctioned violation (Pew Research Center, 2021; Twitter, 2020b).

I begin by modeling a platform on which users spend time and interact. The platform maximizes profits by choosing its prices—the frequency at which it displays ads—and its content removal rate, which reduces spillovers between users. As in Weyl (2010), the pricing policy is what allows the platform to effectively choose the amount of time that users spend on it. The moderation policy is a quality decision that maximizes the users’ willingness to engage with ads. When setting its moderation rate, the platform trades off the change in censored and non-censored users’ engagement with ads. Because moderation is costly, it is only profitable if it increases the activity of at least some users. In other words, it makes sense for profit-maximizing platforms to restrict the content of some of their users if this increases the overall engagement with ads. Thus, a key parameter of the platform’s decisions is how users change their time on the platform in response to moderation.

The first experiment provides information about this parameter by leveraging the report-

---

<sup>1</sup>Packingham v. North Carolina, 137 S. Ct. 1730, 1737 (2017).

ing tool of the platform that allows flagging content that violates the rules. Twitter combines these reports with algorithms to detect violations and enforce its guidelines. It then chooses from a wide range of sanctions at the Tweet or user level, such as reducing Tweet or user visibility (also called shadowbanning), temporarily locking accounts, removing Tweets, or suspending (banning) users. Because reports increase the likelihood that Twitter detects content and, plausibly, do not affect users directly, they instrument moderation as long as sanctions are perfectly observed. Thus, reporting overcomes the challenge of moderation not being randomly assigned.

Over the course of two months, I sampled over 6,000 Tweets containing slurs about disability, which constitute 98% of the sample, or that deny the Holocaust. These slurs are covered by Twitter’s hateful-conduct policy.<sup>2</sup> The sample included different spellings of the slurs to capture attempts to evade detection algorithms, excluded bots, inactive accounts, and other quality filters. Users enter the sample once, so there is one Tweet per user.

The day after they were posted, I randomly reported half of the Tweets for violating the rules against hateful conduct. I then collected daily server-level data of users’ sanctions, their behavior and their followers’ behavior, and the behavior of the users that the Tweets replied to, if any. The data comes from Twitter’s Application Programming Interface (API) and other sources such as Google’s Perspective API (Wulczyn et al., 2017), Botometer’s API (Yang et al., 2020), and shadowban.eu’s API (Merrer et al., 2020).

The first set of results show that reporting has a first-stage impact on sanctions. Reported Tweets are 66% (1.4 percentage points or 0.08 standard deviations) more likely to be removed within three weeks by Twitter than non-reported Tweets, with an F-statistic of 11. The treatment does not significantly change user suspensions and shadowbans, the other observable sanctions.<sup>3</sup> However, I also find evidence of “unobservable” sanctions, such as temporarily locking users’ accounts, which I obtained from the updates that Twitter sent me after I reported users and from the timing gap between user posts.<sup>4</sup> Hence, reports remain a valid instrument for all sanctions, even if they violate the exclusion restriction for observed sanctions. I also found evidence that Twitter does not treat all user reports equally: the accounts that I used for reporting were not equally successful at taking down Tweets.

The second set of results concern the reduced-form impact of reports on user behavior on

---

<sup>2</sup>The policy mentions the Holocaust and slurs that reinforce negative stereotypes about a protected category, including disability. These slurs are only a subset of hate speech, but most other slurs are appropriated by minorities (Bianchi, 2014) and led to high false positives in a pilot study. Another option was to sample Tweets with a detection algorithm from the computer science literature, but even state-of-the-art methods suffer from low external validity (Arango et al., 2019).

<sup>3</sup>There is no evidence of users self-censoring (deleting their posts or accounts or locking their accounts from public view) in response to reports. There is also no evidence that reports induce their other Tweets to disappear.

<sup>4</sup>No observable sanctions were implemented in 4% of the accounts I reported, but Twitter provided an update that it had found rule violations. This number is likely biased downward because Twitter does not always send updates, even for reports for which a sanction is observed.

the platform. This estimation is possible because accounts do not disappear after reporting. I find that reports do not reduce the users’ Twitter activity or their likelihood of reposting hate. A proxy of the hours spent on Twitter, constructed with the daily number of Tweets and likes, increases by 7.5% (0.042 standard deviations) in the three weeks after reports, but it is not statistically significant.<sup>5</sup> The fraction of hateful Tweets that users post in the three weeks after the treatment, measured using Google’s toxicity score, decreases by an insignificant 1.8%.

The third set of results show that reporting has significant spillovers on other users. The main measure of spillover is the activity of the users to whom the Tweets in my sample are replying, which I call “replied users”—86% of Tweets reply or quote a post from another user. Reports increase the time the replied users spend on the platform over the course of three weeks by 10%, or 10 minutes per week. Furthermore, the estimate is stronger among Tweets that attack the other user rather than, for example, those that are just replies among friends. The effect is 13.4% among those Tweets that were labeled as attacks by human annotators. Beyond this effect, there were no other evident spillovers (i.e., on the reported users’ followers).

There is no evidence of differential attrition, and the results are robust to alternative measures of user activity and hatefulness, winsorizing variables, specifications with different sets of controls, different inference methods, and adjusting for multiple-hypothesis testing. Together, these findings imply that sanctions induced by reporting do not change the behavior of those who posted the Tweets; content moderation does not seem to moderate users. Reports, however, increase the activity of those attacked by the hateful posts. Hence, the evidence supports the model’s prediction that content moderation in a profit-maximizing platform marginally increases the advertising revenue from some users.

Does this evidence mean that moderation increases welfare? Not necessarily. Following Spence (1975), another result from my model is that a platform can, in principle, remove too much or too little content relative to a surplus-maximizing planner. Intuitively, the monopolist caters to the marginal consumer, whereas the planner caters to the average consumer. From the consumers’ point of view, a utilitarian test of whether the platform underprovides or overprovides moderation is whether a small change in censorship, all else equal, increases or decreases consumer surplus.<sup>6</sup> Even if this test ignores offline externalities, many costs and benefits associated with moderation, such as free speech and direct harms from hateful expressions, occur inside platforms.

I conducted the test in a survey of 3,000 U.S. Twitter users sampled through Luc.id, a

---

<sup>5</sup>Moreover, the impact on time spent might be biased downward, because Twitter restricts some accounts temporarily (Twitter, 2021e). In these cases, the number of Tweets and likes will be mechanically lower, even if users do not change their behavior.

<sup>6</sup>This test can be generalized to a model of multiple platforms by measuring the change in users’ social-media valuation, not just their valuation of a given platform.

widely used online survey panel provider.<sup>7</sup> I shift users’ beliefs about the likelihood that Twitter moderates hate speech, and I elicit their willingness to accept (WTA) to stop using social media. I vary the perceived likelihood of moderation using an information-provision design with an active control group (Haaland et al., 2020). I randomize survey participants into two treatment arms that receive different information about the likelihood of moderation among hateful Tweets.

The information provided comes from a random sample of 10,000 Tweets that I collected in August 2020 and classified as hate speech with the help of human annotators. I vary the likelihood of moderation without deception by using different rules to classify hate speech. Under a majority decision rule, in which a post is hateful if most annotators label it as such, Twitter removes 3.6% of hateful Tweets or suspends their authors within one month. Under a consensus decision rule, in which a post is hateful if all annotators label it as such, the likelihood of moderation is 9.1%. Under both rules, the prevalence of hate—that is, the fraction of hateful Tweets—is less than 1%. Both treatment arms receive the same information about hate prevalence, which allows isolating the effect of moderation.

The survey first elicits beliefs about the prevalence of hate speech and the likelihood of moderation with incentives for accuracy and then provides participants with the randomized information. Respondents are told that some of them will be randomly selected for a small follow-up study that compensates participants to stop using social media (Twitter, Facebook, Instagram, Snapchat, YouTube, Reddit, and TikTok) for one week. I then elicit the WTA to participate in this follow-up, using an incentive-compatible procedure in the form of an iterative multiple price list (iMPL).<sup>8</sup>

I find large misperceptions about hate speech and moderation. Most users overestimate the prevalence of hate speech on Twitter and the likelihood of sanctions. Ninety-six percent of users believe the prevalence of hate speech is above the observed value of less than 1%, with a median of 33%. Eighty-four percent of respondents believe the likelihood of moderation is above 9.1%, with a median of 36%.<sup>9</sup>

Informing participants of a higher likelihood of moderation does not change their valuation of social media. The WTA falls by 15 cents (0.5% or 0.004 standard deviations), from \$33.7 to \$33.6. This result is robust to different specifications and measures of WTA, and I find

---

<sup>7</sup>I reweight observations to match a representative sample of Twitter users based on gender, age, race or ethnicity, region, and political orientation. I also present unweighted results.

<sup>8</sup>In this procedure, participants have to choose if they are willing to participate in the follow-up for different compensation offers. The sequence of offers starts at \$50, and subsequent amounts decrease or increase depending on whether participants accept or reject. The sequence stops until the WTA is classified in 11 intervals, which I then convert into a continuous measure following Allcott and Kessler (2019).

<sup>9</sup>Platforms’ lack of transparency could be driving these misperceptions. The likelihood of moderation on Facebook remained unknown until a whistleblower revealed internal documents some weeks after my survey (Giansiracusa, 2021).

no evidence that it is explained by inattention or experimenter demand effects.<sup>10</sup> At the end of the survey, I asked participants to repeat the information about the percentage of Tweets that get sanctioned. The treatment effect on this recollection was 5.6 percentage points, significantly different from zero ( $F$ -statistic = 36) and not statistically different from 5.5, which is the gap between the information provided in both arms, 9.1% and 3.6%. The treatment also significantly shifted the posterior beliefs about the likelihood of moderation on Facebook and there is suggestive evidence that it increased the time that minorities spent on Twitter one week after the survey.

Overall, my results suggest that moderation on Twitter is consistent with profit maximization, and they rule out large moderation distortions from the consumers' point of view, holding constant the prevalence of hate speech. These findings have two policy implications. First, cost-benefit analyses of online moderation can emphasize its offline consequences, such as hate crimes (Jiménez Durán et al., 2022). Second, authorities might want to deal with hate speech on social media not by directly regulating moderation, but by supervising platforms' pricing (advertising) policies; for example, Twitter could still be setting its advertising loads suboptimally, leading to inefficient amounts of hate speech.

The paper makes four contributions to a multi-disciplinary literature. First, a growing body of work in economics focuses on the offline consequences of online content and social-media penetration, including voting and political participation (Fergusson and Molina, 2019; Enikolopov et al., 2020; Fujiwara et al., 2021; Zhuravskaya et al., 2020; Beknazar-Yuzbashev and Stalinski, 2022), polarization (Allcott and Gentzkow, 2017; Boxell et al., 2019; Levy, 2021; Melnikov, 2021), hate crimes (Müller and Schwarz, 2020a,b; Bursztyn et al., 2019; Jiménez Durán et al., 2022), and mental health and well-being (Allcott et al., 2020; Mosquera et al., 2020; Allcott et al., 2022; Braghieri et al., 2021). This paper contributes to this work by investigating online outcomes.

Experimentally varying moderation is challenging due to limited cooperation with platforms. Thus, a second contribution is using social media's infrastructure experimentally, as Levy (2021) did on Facebook, which is useful for independent research. The reporting treatment is similar to other exercises by academics (Carlson and Rousselle, 2020), Governmental organizations (Jourová, 2016; Reynders, 2020), and non-profits (Matias et al., 2015; Center for Countering Digital Hate, 2021) who report content to monitor platforms' responsiveness. However, these exercises are non-experimental (they contain no control group) and do not analyze the impact on other outcomes beyond the platform's response. Most empirical evidence of the effects of platform-initiated moderation is non-causal,<sup>11</sup> with the exception of

<sup>10</sup>The experiment was ex-ante powered to detect effects of 0.1 standard deviations, and the sample size is more than double what Haaland et al. (2020) recommend for information-provision designs.

<sup>11</sup>The computer science literature provides correlational evidence (Ali et al., 2021; Rauchfleisch and Kaiser, 2021; Jhaver et al., 2021; Zannettou, 2021). A challenge with observational studies is isolating the effect of moderation from confounders, e.g., Chandrasekharan et al. (2017) study hateful subReddit bans, but this

more recent work by Katsaros et al. (2022); Ribeiro et al. (2022a); Wojcik et al. (2022). To the best of my knowledge, this paper is one of the first to provide evidence of the behavioral and welfare effects of moderation with a field experiment.

A third contribution is combining an information-provision design with a welfare elicitation of social media. Haaland et al. (2020) and Bursztyn and Yang (2021) provide overviews of information-provision designs, and Bottan and Perez-Truglia (2017) and Bursztyn et al. (2020) are some applications. The WTA that I elicit is in the ballpark of other social-media welfare studies such as Mosquera et al. (2020), and Allcott et al. (2020); the median and mean WTA per week were \$15 and \$34, respectively.<sup>12</sup> Providing information required computing other basic statistics, surprisingly scarce in the literature, such as the prevalence of hate speech in a random sample of posts (0.1%-5.6% depending on the measure) and the occurrence of Tweet deletions and user suspensions (2.5%-9.1% among hateful Tweets, within one month).<sup>13</sup> As other surveys find (Anti-Defamation League, 2021), minorities experience more harassment online, but I also find that they are more likely to be sanctioned by Twitter.

The fourth contribution is to develop a simple model of user behavior and platform moderation decisions that captures spillovers between users, using the two-sided market framework of Weyl (2010). Prices—advertising loads—allow the platform to determine its amount of hateful and non-hateful content, which clarifies the separation between pricing distortions and moderation distortions as in Spence (1975).<sup>14</sup> Liu et al. (2021) are among the first to model moderation decisions and discuss the implications of different revenue models on platform incentives. One difference with their model is that, in my framework, users respond to the pricing policy (advertising frequency) of the platform (see also Madio and Quinn (2021)).<sup>15</sup> Acemoglu et al. (2021) model online misinformation and show that engagement-maximizing platforms have incentives to create filter bubbles and propagate extremist content.

---

includes the effect of banning non-hateful posts and users. Experimental interventions include counter-speech treatments (Munger, 2017, 2021; Siegel and Badaan, 2020; Hangartner et al., 2021), reminders of Twitter suspensions (Yildirim et al., 2021), and lab moderation (Han and Brazeal, 2015; Rösner et al., 2016; Cheng et al., 2017; Álvarez-Benjumea and Winter, 2018; Kim et al., 2021). There is also field evidence of community-based moderation (Seering et al., 2017; Jhaver et al., 2019; Srinivasan et al., 2019; Ribeiro et al., 2022b).

<sup>12</sup>This is after reweighting my sample to match representative U.S. Twitter users. Allcott et al. (2020) find a median and average WTA of \$25 and \$45 per week, respectively, for deactivating Facebook for four weeks.

<sup>13</sup>Relia et al. (2019) find that 0.5% of Tweets in a sample of 73.42 million posts contained hate speech keywords. Founta et al. (2018) found a 4% prevalence in a random sample of 10,000 Tweets. Facebook (2021) reports a prevalence of 0.05% of hate speech among all views. Few studies report the occurrence of sanctions. An exception is Merrer et al. (2020), who document that 2.34% of accounts are shadowbanned. Seyler et al. (2021) find that 5.1% of accounts from a 2009 sample are suspended.

<sup>14</sup>There is evidence that consumers respond to platforms' advertising policies; Huang et al. (2018) traced out a downward-sloping demand curve for a music platform by randomizing ad-loads across consumers. See Aridor (2022) for more references of how ad-loads act as prices on social media.

<sup>15</sup>In both frameworks (under an advertising business model), platforms use moderation as a tool to increase revenue. In Liu et al. (2021) and Madio and Quinn (2021), moderation increases revenue through increases in the consumer base (quantity decision). However, in my model, prices determine the customer base and moderation increases revenue through the willingness of users to engage with ads, for a fixed customer base (quality decision).

The next section develops the model. Section 3 gives background information about hate speech, moderation, Twitter, and the data sources in this study. Sections 4 and 5 describe the experimental designs and results. Section 6 concludes.

## 2 Model

**Users and platform.** Users can be one of two types,  $\theta \in \{A, H\}$ . “Acceptable” users ( $\theta = A$ ) post content that is not subject to content moderation. “Hateful” users ( $\theta = H$ ) post content that is censored by the platform with probability  $c$ , the censorship or moderation rate. Users who join the platform experience utility that increases on the time that they spend consuming or posting content.<sup>16</sup>

The utility of spending  $t$  minutes on the platform for user  $i$  of type  $\theta_i = \theta$  is

$$\underbrace{U_i^\theta(t; \mathbf{T}, c)}_{\text{Utility from consuming content}} - \underbrace{t \times w(1 + p^\theta)}_{\text{Time cost}}, \quad (1)$$

where  $U_i^\theta(0; \mathbf{T}, c) = 0$  for all  $i$ .  $\mathbf{T} = (T^A, T^H)$  is the aggregate content of the platform and captures spillovers and network effects, and  $T^\theta$  is the total content posted by  $\theta$  users. The sign of spillovers is flexible; users can be positively or negatively affected by each type of content. For instance,  $A$  users could dislike encountering hate speech, but haters might like trolling  $A$  users. The unrestricted heterogeneity in the utility function (which can be interpreted as a value function maximized over users’ following and blocking decisions), allows users to be diversely exposed to the different types of content (e.g., by adjusting their following and blocking). The censorship rate  $c$  enters the utility function because it reduces spillovers from hateful content (since users see less of it), but users can also obtain direct utility or disutility from  $c$ ; for example, haters might dislike having their account locked.

The time-cost of  $t$  minutes spent enjoying content is proportional to the value of time  $w > 0$ . Moreover, the “price” that users pay is the advertising load  $p^\theta$ ; the time they have to spend watching ads per minute of content consumed. Following Weyl (2010), the platform can set a different price for each type of user.<sup>17</sup>

The time that user  $i$  spends enjoying content is  $t_i^*$ , which maximizes (1) with respect to

<sup>16</sup>There is no difference between consuming or producing content (e.g., as in Filippas and Horton (2021)). In practice, however, users differ substantially in the amount of content they post. On Twitter, few users post the majority of Tweets (Wojcik and Hughes, 2019). Yet, it is not obvious whether users who like posts are less responsible for their diffusion than those who write them. For instance, sometimes Twitter alerts the followers of a user when she likes a post.

<sup>17</sup>Indeed, anecdotal evidence from Levy (2021) suggests that users might experience heterogeneous ad loads. In his data, the interquartile range of ad loads goes from 1 ad per 6 posts to 1 ad per 9 posts. There is similar evidence of heterogeneity of Twitter ad loads in Beknazar-Yuzbashev et al. (2022).



$t \geq 0$ . The aggregate content demand  $\mathbf{T}$  is then computed setting

$$T^\theta \equiv \int_{\{i:\theta_i=\theta\}} t_i^* di, \quad \text{for each } \theta. \quad (2)$$

Since the time spent on the platform by any user is decreasing in the advertising load, one can define the inverse demand functions  $P^A(\mathbf{T}, c)$  and  $P^H(\mathbf{T}, c)$ , where  $P^A(\mathbf{T}, c)$  is equal to the  $p^A$  inducing  $T^A$  given  $\mathbf{T}$  and  $c$ ; similarly for  $P^H(\mathbf{T}, c)$ .<sup>18</sup> In other words, the pricing policy—not moderation—is what allows the platform to choose the amount of content of both types of user.

The platform maximizes profits by solving

$$\max_{\mathbf{T}, c} \underbrace{a \times \left( \underbrace{P^A(\mathbf{T}, c) T^A + P^H(\mathbf{T}, c) T^H}_{\text{Time spent watching ads}} \right)}_{\text{Advertising revenue}} - \underbrace{\phi(\mathbf{T}, c)}_{\text{Cost of moderation}}, \quad (3)$$

where the platform takes as given the price per unit of advertising— $a > 0$ , although this assumption can be relaxed—and  $\phi$  is a function describing the costs of censorship. For instance, Gillespie (2018) documents that moderation is a labor-intensive task, and Kaye (2019) argues that regulatory fines push platforms to remove borderline content.<sup>19</sup>

The first-order condition (FOC) with respect to quantity  $T^\theta$  is similar to a standard monopoly problem. The FOC with respect to  $c$  requires that

$$a \times \underbrace{\left( \frac{\partial P^A(\mathbf{T}, c)}{\partial c} T^A + \frac{\partial P^H(\mathbf{T}, c)}{\partial c} T^H \right)}_{\text{Change in time spent watching ads}} = \underbrace{\frac{\partial \phi(\mathbf{T}, c)}{\partial c}}_{\text{Marginal cost of moderation}}. \quad (4)$$

This condition is analogous to the quality decision in Spence (1975); the platform moderates such that the marginal benefit—the value of the marginal increase in the willingness to watch ads—equals the marginal cost. The left-hand side of equation (4) clarifies the main trade-off faced by the platform when choosing its moderation policy. Consistently with the observations in Kaye (2019) regarding controversial pages,

These kinds of pages seem to put Facebook in a no-win position: If they leave up the page, they anger opponents who see hateful content or disinformation; if they take it down, they offend free-expression advocates who do not think the rules very clearly

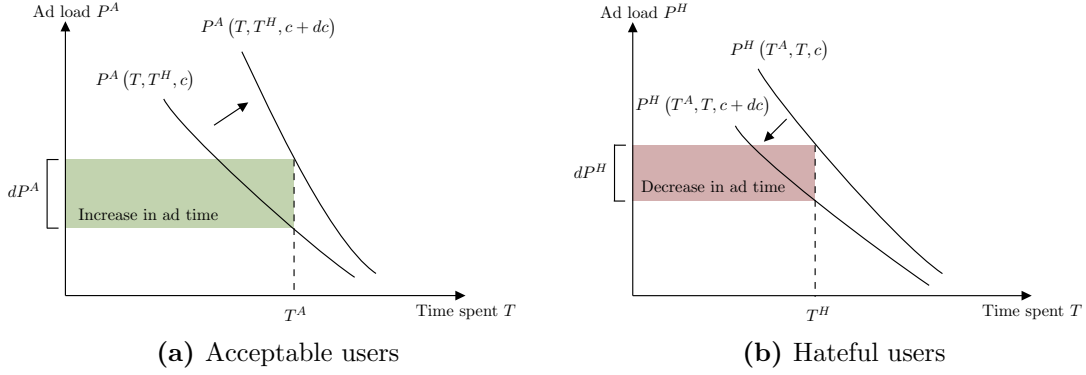
---

<sup>18</sup>Formally, imposing rational expectations  $\widetilde{T}^\theta = T^\theta(p^\theta, c, \widetilde{T}^A, \widetilde{T}^H)$ , one can invert  $T^\theta$  in an interior equilibrium point, where  $\widetilde{T}^\theta > 0$ . This procedure requires demands  $T^\theta$  to be strictly decreasing in  $p^\theta$ , which results from imposing Inada conditions on utilities or full-support assumptions as in Weyl (2010).

<sup>19</sup>Platforms might worry about future fines, even if current ones are small; e.g., in 2019, Germany fined Facebook for 2 million euros for violating the NetzDG law (Bundesamt für Justiz, 2019).

articulate hate speech standards.

Figure 1 illustrates this trade-off for the case in which  $\partial p^A/\partial c > 0$ , and  $\partial p^H/\partial c < 0$ . For a fixed amount of time that users spend on the platform, moderation changes the number of ads they are willing to watch. The platform increases revenue from A users, who dislike hateful content, while it loses revenue from H users, who do not like to be censored. The optimal level of content moderation balances the net change in revenue with the marginal cost of increasing the censorship rate.



**Figure 1:** Graphical intuition of the platform's moderation decision

Notes: These figures plot the change in the inverse demands of acceptable and hateful users in response to an increase in moderation,  $dc$ , holding quantities fixed, and assuming that the moderation elasticity is positive for  $A$  and negative for  $H$ . The colored areas are the change in time spent watching ads, which equals the change in ad load  $dP^\theta$  times the time spent  $T^\theta$ . The net revenue gains equal the green minus the red area, multiplied by ad prices.

The FOC is, equivalently,<sup>20</sup>

$$-\frac{\partial T^A/\partial c}{\partial T^A/\partial p^A} aT^A - \frac{\partial T^H/\partial c}{\partial T^H/\partial p^H} aT^H = \frac{\partial \phi(\mathbf{T}, c)}{\partial c}. \quad (5)$$

We know that the right-hand side is strictly positive (by assumption), and that demand decreases in prices,  $\partial T^\theta/\partial p^\theta < 0$ . Therefore, at the optimal level of  $c$  for the platform it must be that either  $\partial T^A/\partial c > 0$ , or  $\partial T^H/\partial c > 0$ , or both. In words, the first implication of the model is that, for at least one type of user, small increases to the (interior) level of moderation must increase their platform activity, holding constant the aggregate quantities. The derivatives of  $T^\theta$  with respect to  $c$ , one for each type, are the main parameters of interest of my first experiment.

The implication that small increases in the equilibrium (interior) level of moderation should increase the engagement of at least some users does not rely on the assumption of a

<sup>20</sup>This equation uses the implicit function theorem. For example, letting  $\widetilde{T}^A = T^A(p^A, c, \widetilde{T}^A, \widetilde{T}^H)$ , taking the total derivative implies  $0 = \frac{\partial T^A}{\partial p^A} dp^A + \frac{\partial T^A}{\partial c} dc$ , so that  $\frac{dp^A}{dc} = -\frac{\partial T^A/\partial c}{\partial T^A/\partial p^A}$ .

fixed ad price  $a$ . Note that the implication remains the same if we allow a different ad price for haters ( $a^H$ ) and non haters ( $a^A$ ) in Equation (5). The implication is also unchanged if we allow for price-setting behavior in the ads market or assume that advertisers have brand-safety concerns, as long as advertisers care only about the (surviving) amount of hate speech and not about moderation directly.<sup>21</sup>

**Welfare.** In principle, the platform-optimal level of censorship could differ from the socially-optimal level. Similarly to Spence (1975), the platform in my model optimizes moderation with respect to the marginal users. The social planner, however, chooses the level of censorship that maximizes total welfare, which includes the impact of moderation on inframarginal consumers. I formalize this argument in Appendix A; the platform can moderate more or less than a surplus-maximizing social planner, holding quantities  $\mathbf{T}$  fixed.<sup>22</sup> Hence, two distortions exist: the usual monopolist pricing distortion that leads to inefficient quantities and an additional quality distortion.

The goal of my second experiment is to test the presence of the second distortion. The test consists in evaluating whether Twitter provides too little or too much content moderation from the perspective of the user; that is, whether consumer surplus increases or decreases, respectively, in response to increases in the perceived moderation rate. I follow the approach of Mosquera et al. (2020) and Allcott et al. (2020) to measure consumer surplus. In practice, I quantify the impact of different levels of censorship on the willingness to accept a monetary reward to pause the use of social media. I ask users to pause the use of social media, not just a single platform, to allow for substitution between platforms as argued in Appendix A.

## 3 Background and Data Sources

### 3.1 Twitter and Moderation of Hate Speech

Twitter is a microblogging social media platform. Users of this platform create profiles that display self-reported information such as their name, a short biography, and a profile picture. They also post messages to their profiles called Tweets, which contain a combination of text of

---

<sup>21</sup>To see why, let  $T^{H,s} = (1 - c)T^H$  be the surviving amount of hate speech and redefine  $\mathbf{T} \equiv (T^A, T^{H,s})$ . The platform can be a price setter or have brand-safety concerns ( $a(\mathbf{T})$ ) and its first-order conditions with respect to  $c$  remain unchanged—as platforms can determine user behavior through their advertising policies. Redefining the problem in terms of the surviving hate speech illuminates another implication from the model: at least some users need to derive direct utility from moderation. If users cared about moderation only through its effect on hate speech, the equilibrium amount of moderation would be zero. The reason is that platforms would only need to use their advertising policy to determine their quantities.

<sup>22</sup>Liu et al. (2021) argue that platforms undermoderate in an ad-based business model and overmoderate in a subscription-based business model. However, in their ad-based business model there are no prices, so the platform has to use moderation as the only tool to adjust its quantity. In their subscription-based case there are prices, but the sign of the Spence distortion is determined by their extremeness aversion assumption, which in practice implies that the demand curve becomes steeper in response to a small increase in moderation.

up to 280 characters, photos, and videos. Users can follow other accounts to see their Tweets more readily, but they can interact with others without following them. They interact with others’ Tweets by giving them a like (or favorite), replying to them, Retweeting (reposting) them, or quoting them.

Like all social media platforms in the Surface Web, Twitter has rules that delimit the content that users are allowed to post. Besides illegal activity, the rules tend to cover hate speech (as well as misinformation, harassment, spam, sexual content, and graphic content). There is no single legal definition of hate speech (Waldron, 2012; Strossen, 2018). Still, most platforms define it using common elements such as the concept of protected categories from U.S. anti-discrimination law (Gillespie, 2018). Twitter’s hateful-conduct policy (Twitter, 2021c) says, “You may not promote violence against or directly attack or threaten other people on the basis of race, ethnicity, national origin, caste, sexual orientation, gender, gender identity, religious affiliation, age, disability, or serious disease.”

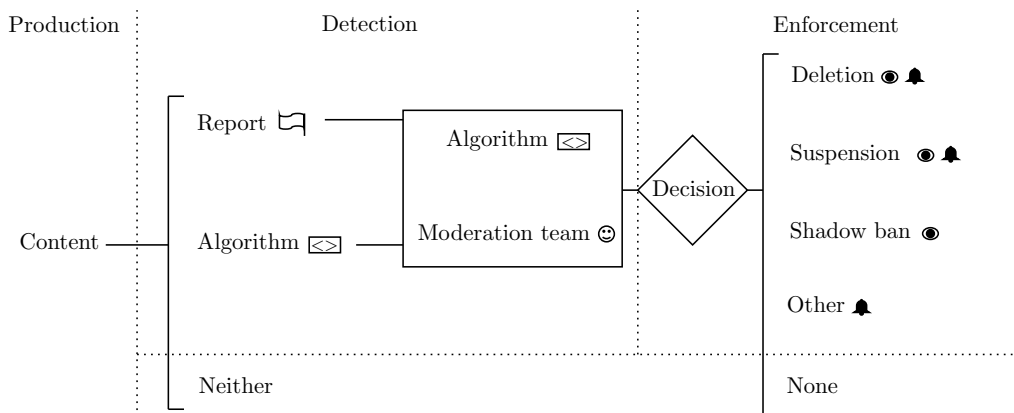
Twitter enforces these rules and moderates content by sanctioning users. Figure 2 illustrates the process of content moderation. Twitter can detect content by algorithms, or by the “flagging” mechanism that allows users to report Tweets or accounts for violating the rules. After the content is detected, a team of human moderators or an algorithm decide whether to enforce the rules by imposing post-level or account-level sanctions. The range of sanctions include a combination of removing Tweets from the platform, shadowbanning (reducing the visibility of) users or Tweets,<sup>23</sup> and suspending or banning users (that is, deleting their accounts). Other sanctions, such as locking accounts, prevent users from posting or liking content and can last from 12 hours to seven days. More recently, Twitter also started moderating content at the production stage (Katsaros et al., 2022). See Twitter (2021e) for the full list of sanctions.

There is no public information about the correspondence between the severity or frequency of rule violations and sanctions, with the exception of the COVID-19 misinformation policy. In that case, the first violation does not grant account-level action (but might include Tweet-level actions), the second and third violation lead to 12-hour account locks, the fourth violation leads to a 7-day account lock and 5 or more violations lead to a permanent suspension (Twitter, 2021a).

Sanctions differ by their observability, that is, whether they are privately observed by the user or publicly observable, and whether the user is notified about them. When Twitter removes a Tweet, users are notified that they violated the rules, and must remove the Tweet to be able to use the platform again. Twitter replaces the Tweet with a notice that indicates it violated the rules. Anyone with access to the Tweets can see the notice, starting from

---

<sup>23</sup>Twitter has stated it does not shadowban users (Gadde and Beykpour, 2018), even if it ranks Tweets “to create a more relevant experience” (Twitter, 2021b). However, Merrer et al. (2020) document evidence of shadowbanning.



**Figure 2:** Content moderation process

Notes: Eye icons indicate that the sanction is observable to others besides the user. Bell icons indicate that the user receives a notification. This diagram omits interventions at the production stage, such as recent tests in which Twitter asked users if they wanted to review offensive Tweets before posting. It also omits the appealing process, in which users can contest a sanction.

the moment that Twitter sanctions the Tweet and up to 14 days after the user agrees to remove the post. When Twitter suspends a user, she can no longer log in to the account, and her profile and Tweets are replaced by suspension notices, which seem to last indefinitely. In principle, suspended users cannot create new accounts, but in practice they do. Users are not notified when they are shadowbanned, but Twitter sometimes hides their Tweets behind a notice—especially those that reply to another user. Twitter notifies users when their accounts are locked, but whether others can observe this sanction is unknown.<sup>24</sup> Figure B.1 shows examples of public notices and the notifications that users receive.

### 3.2 Measuring Hate Speech

Platforms rely to some extent on algorithms to detect hate speech and enforce their rules. Most of the detection algorithms in the computer science literature share the following procedure (see Fortuna and Nunes (2018) for a review of the literature). The first step is to obtain a training dataset, consisting of a sample of texts—usually social media posts—paired with labels, for example, hate speech or not hate speech. Often, these labels or “ground truth” result from aggregating the opinions of multiple humans or “annotators” into a single category. For example, Davidson et al. (2017) ask three or more crowd workers to annotate each Tweet as “hate,” “offensive,” or “neither.” Then, they aggregate these annotations into a single label with the majority decision rule, that is, the category chosen by most annotators. The second step is to convert the text into vectors of features with text analysis, reviewed in

<sup>24</sup>Twitter (2021d) shows examples of notices of locked accounts, but anecdotal evidence suggests accounts are locked without any notice. For instance, in 2020 Twitter locked actor James Woods and his account did not show any notice (Whalen, 2020).

Gentzkow et al. (2019). The final step is to use machine learning to predict the labels with the features.

One challenge in the hate-speech-detection literature is the algorithms’ low external validity; see Arango et al. (2019) and Fortuna et al. (2021). For this reason, this study uses three approaches to classify hate speech and limit measurement error. First, for large-scale tasks, I use the Perspective toxicity score developed by Google. This score is widely used in the industry and as a benchmark in academic articles. It is a number between 0 and 1 that reflects the likelihood that a text is an attack or harassment.<sup>25</sup> Many studies classify posts as hate speech if their toxicity is higher than a 0.8 cutoff (ElSherief et al., 2018; Han and Tsvetkov, 2020; Vidgen et al., 2020). Second, I sample hate speech using keywords, instead of an algorithm, to minimize false positives in the reporting experiment. Third, I use human annotation by MTurk workers to account for measurement error in the first experiment and to increase the interpretability of the information treatment of the second experiment.

### 3.3 Data Sources

Most of the variables analyzed in this paper come from Twitter’s API. This data source provides publicly available information about all Twitter users, such as their number of followers, number of accounts they follow, date of account creation, total number of Tweets and likes, and biography. The API provides additional information about users who do not restrict their profile visibility, such as their list of followers and accounts followed, and a collection of up to 3,200 of their most recent Tweets. The API returns detailed information for these Tweets, such as their timestamp, text and media, likes, and Retweets. This source also allows me to sample Tweets by searching for specific keywords or sampling at random 1% of all Tweets. Lastly, I also use this API to detect whether Twitter removes specific Tweets or suspends users, following the procedure outlined in Appendix C.1. To the best of my knowledge, only these two sanctions are observable using the response codes from Twitter’s API.

Besides the API, I also collect some information manually from the website. Twitter occasionally notifies users when it sanctions an account they previously reported, even if the sanction might not correspond to the reported content.<sup>26</sup> Figure B.3 has a screenshot of some of these updates. I collect this information for the reporting experiment, because it provides a signal of “unobservable” sanctions.

I also use other APIs. I retrieve the toxicity score of posts from Google’s Perspective API. I also obtain a measure of the likelihood that users are bots from the Botometer API

---

<sup>25</sup>The algorithm is a convolutional neural network trained on Wikipedia Talk Pages; see Wulczyn et al. (2017) and Dixon et al. (2018).

<sup>26</sup>Twitter says: “You will receive an in-product notification if an action is taken on an account that you recently reported. This action may or may not be related to your report” (Twitter, 2021f).

(see Yang et al. (2020)). Finally, I retrieve a (noisy) measure of shadowbans from the API of `shadowban.eu`,<sup>27</sup> because Twitter does not give an official shadowban measure. This API measures different forms of shadowbanning, for example, whether Twitter hides accounts, Tweets, or replies from search results (see Merrer et al. (2020) for more details). I combine the different measures into a single indicator of whether users are shadowbanned.

Another data source is human annotation; I ask MTurk workers to annotate posts. For example, I follow the approach in Davidson et al. (2017) and ask workers to classify posts as “Hate speech,” “Offensive but not hate speech,” and “Neither offensive nor hate speech.” I assign three workers to annotate each post. I give them Twitter’s definition of hate speech for reference, offer a \$20 bonus to the five most accurate workers (measured by the inter-annotator agreement), and include attention checks to improve the quality of annotations. Figure D.1 in the Appendix includes screen shots of the instructions. Then, I aggregate workers’ annotations into a single label using either the majority decision rule, in which a post is hate speech if two or three workers label it a such, or the consensus decision rule, in which all three workers have to agree.

Lastly, I obtain demographics of representative Twitter users from the American Trends Panel (ATP) of September 2020. The Pew Research Center conducts this nationally representative panel of randomly selected U.S. adults.

### 3.4 Summary Statistics

**Accounts and Tweets.** According to the ATP, 25% of adults in the U.S. use Twitter. Table 1 displays selected summary statistics of Twitter users and their accounts. Twitter users are younger, more educated, and more likely to be Democrats than the general population. Thirty-one percent of them are between 18 and 29 years old, 40% are at least college graduates, and 35% are Democrats, compared to 18%, 33%, and 30%, respectively, in the overall ATP respondents. The table also shows statistics from a sample of 200,000 Tweets that I collected in August 2020 from the 1% random sample of Twitter’s API. On average, the accounts in this sample were five years old, posted 12 Tweets per day, gave 13 likes per day, followed 1,000 users, and had 4,800 followers. Ten percent of these accounts are bots; that is, they have a Botometer score of 0.5 or more.

**Prevalence of hate speech.** The random sample of Tweets allows me to quantify the percent of Tweets that are hate according to different measures. Using the 0.8 toxicity cutoff, I find 5.6% of Tweets are hate. To compare this number with human annotation, I annotated a subsample of 10,000 Tweets from the random sample. As Table 1 shows,

---

<sup>27</sup>This website went down on multiple occasions during my data collection, which is why this measure is noisy.

**Table 1:** Summary statistics

	Mean	Std. Dev	p10	Median	p90	Obs.	Sample
<i>Accounts</i>							
Account years	5.24	3.82	1.17	4.04	11.21	191,835	Random
Tweets per day	12.02	39.26	0.23	4.35	29.62	191,835	Random
Likes per day	13.03	24.08	0.07	3.98	36.16	191,835	Random
Followers	4,804	169,343	15	340	3,167	191,835	Random
Followed	1,071	6,755	45	381	2,078	191,835	Random
Is bot (%)	9.90	29.88	0.00	0.00	0.00	1,000	Random
Age 18-29 (%)	30.99	46.25	0.00	0.00	100.00	2,463	ATP
Male (%)	53.20	49.91	0.00	100.00	100.00	2,464	ATP
White (%)	57.92	49.38	0.00	100.00	100.00	2,464	ATP
College graduate (%)	40.47	49.09	0.00	0.00	100.00	2,464	ATP
Republican (%)	20.72	40.54	0.00	0.00	100.00	2,464	ATP
Democrat (%)	35.16	47.76	0.00	0.00	100.00	2,464	ATP
<i>Tweets</i>							
Is reply or quote (%)	62.53	48.40	0.00	100.00	100.00	201,038	Random
Is toxic (%)	5.61	23.01	0.00	0.00	0.00	201,038	Random
Is hate (% ,majority)	0.56	7.47	0.00	0.00	0.00	9,991	MTurk
Is hate (% ,consensus)	0.11	3.32	0.00	0.00	0.00	9,991	MTurk

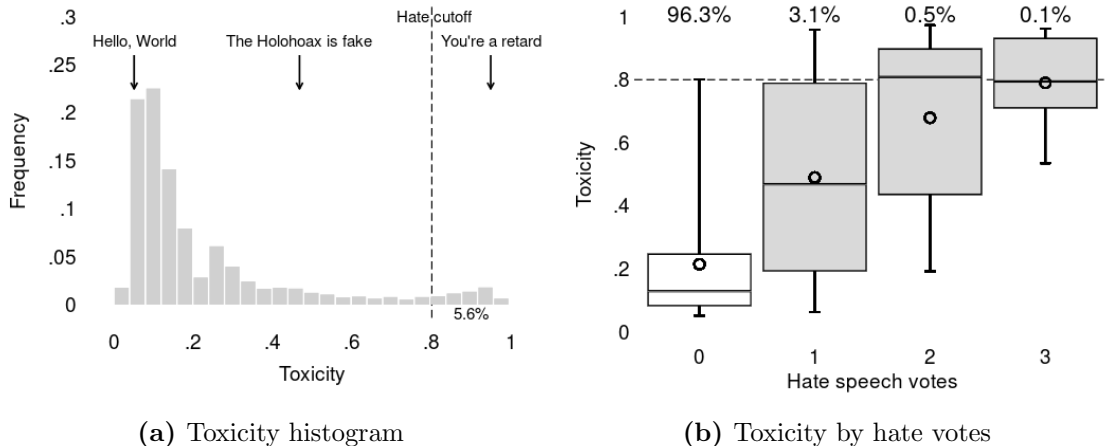
Notes: The random sample indicates a random extraction of 201,308 Tweets from Twitter’s API on August 2020. The bot score was computed on a subsample of 1,000 accounts from the random sample of Tweets, due to rate limits from the Botometer API. The ATP sample is a subsample of Twitter users from the Pew Research Center’s ATP. The MTurk sample is a random subsample of Tweets that I annotated on MTurk following Davidson et al. (2017).

less than 1% of Tweets are considered hate speech using human annotation, under both the majority decision rule and the consensus rule. Thus, hate speech is a low-probability event.

The long-tailed nature of hate is more evident in Figure 3a, which plots a histogram of the toxicity score in the random sample of Tweets. The figure also includes the toxicity scores of three example texts: the neutral phrase “Hello, World” (toxicity = 0.05), one phrase related to disability (toxicity = 0.95), and one that denies the Holocaust (toxicity = 0.47). These examples, which are relevant for the reporting experiment, illustrate that the toxicity cutoff adequately separates some slurs from neutral expressions, but it fails to identify more subtle hate. Still, toxicity is closely correlated with human annotation. Figure 3b shows the distribution of toxicity scores shifts to the right as more workers label Tweets as hate speech.

**Occurrence of sanctions and reports.** Table 2 presents the fraction of removals and suspensions in the random sample of Tweets and the different subsamples of hate speech. Depending on the measure of hate, the fraction of Tweets that Twitter removed or suspended within one month is 2.6% to 9.1%—higher than the 2% in a random sample. These numbers match the statistics recently revealed in Facebook’s whistleblowing event, that the platform removes 3% to 5% of hateful content (Giansiracusa, 2021). In this table we can also see that removals are a rare event. I did not measure shadowbans in this sample, but Merrer et al.





**Figure 3:** Toxicity scores and annotation in a random sample of Tweets

Notes: Panel (a) displays a histogram of toxicity scores based on a random sample of 201,038 Tweets from August 2020. The dashed line is the 0.8 toxicity cutoff to classify hate speech; 5.6% of Tweets have a toxicity above that cutoff. The phrases “Hello, World,” “The Holohoax is fake,” and “You’re a retard” have toxicities of 0.05, 0.47, and 0.95, respectively. Panel (b) has toxicity box plots by the number of workers who voted that a Tweet is hateful. The data is from a subsample of 10,000 Tweets annotated by MTurk workers. The boxes indicate percentiles 25, 50, and 75; the circles indicate the means; and the lines indicate percentiles 5 and 95. The percentages at the top indicate the fraction of Tweets by number of votes.

(2020) document that 2.3% of accounts are shadowbanned. Figure C.1 in the Appendix plots the fraction of sanctions by the type of rule violation; hateful conduct and harassment are the most commonly sanctioned violations in the platform.

**Table 2:** Likelihood of sanctions by subsample

	Random	Hate speech		
		Toxicity $\geq 0.8$	MTurk annotation	
			Majority	Consensus
Removal	0.01	0.1	0	0
Suspension	1.9	2.5	3.6	9.1

Notes: This table shows the fraction of Tweets or accounts that get removed from the platform within 1 month of posting hate speech by each subsample. The random sample of posts is based on 201,038 Tweets and the MTurk annotation is based on a subsample of 9,991 annotated Tweets.

In the second half of 2020, 11% of active accounts were reported according to official Twitter data,<sup>28</sup> and 1% of accounts concentrate the majority of reports (Twitter, 2018). Recently, Twitter’s (former) CEO reported that algorithms detect 51% of the content that

<sup>28</sup>This number results from dividing the total number of accounts reported, from the Rules Enforcement Report (Twitter, 2020b), by the monetizable daily active usage published on the letter to shareholders from Q4 2020 (Twitter, 2020a).

the platform finds in violation of the rules and that the company’s goal is to increase this percentage to 90% (Melendez, 2020). Users can report content even if they are not its targets; in a small study by a nonprofit, 57% of reports were filed on behalf of someone else (Matias et al., 2015).

## 4 Reporting Experiment

### 4.1 Experimental Design

**Sample.** I sampled 6,148 Tweets containing hateful keywords during July and August 2020. I collected the Tweets every day with an algorithm that uses the search function of Twitter’s API, which queries a subset of recent English-language Tweets excluding Retweets.<sup>29</sup> I searched posts containing two slurs: one that denies the Holocaust (Holohoax), and a disability slur (retard), the latter constituting 98% of the sample. Both terms are prevalent on social media and considered by many to be hate speech; see Guhl and Davey (2020) and Sherry (2019). Even if some people use the disability slur frequently (Albert et al., 2016), it is precisely the removal of this type of slurs that is controversial and policy-relevant. Moreover, Twitter’s hateful-conduct policy covers the Holocaust and slurs that reinforce negative stereotypes about a protected category, which includes disability (Twitter, 2021c).<sup>30</sup>

Because the disability slur has alternative meanings, for example, to retard the progress of something, I refine the search with sentence structures such as, “You are a retard.” This refinement captures directed hate speech (ElSherief et al., 2018) and facilitates identifying the targets of hate speech. I also consider multiple misspellings and word distortions to sample Tweets that attempt to bypass detection algorithms. Table B.1 in the Appendix contains the full list of queries used to search Tweets.

After my search algorithm detects a Tweet, it filters users to increase the quality of the sample and to reduce false positives, that is, Tweets that are not hate speech even if they contain the slurs. The filter drops users who self-report being under 18 in their profile biographies, those with new accounts (opened less than 2 weeks before the Tweet), inactive users or non-English speakers (with less than 10 posts in English and more than 50% of posts in another language), and bots (those with a Botometer score higher than 0.5). I also exclude users who display their preferred pronouns on their profile biographies,<sup>31</sup> Tweets that enclose

---

<sup>29</sup>The algorithm conducted the search every 20 minutes. This timing allowed the data processing to be spread throughout the day to comply with the API’s rate limits.

<sup>30</sup>In a pilot study, I included a broader list of slurs about race, ethnicity, religion, gender and sexual orientation. However, the sample contained many false positives, because most slurs are used by the members of the group that they target. Bianchi (2014) refers to this practice as appropriated or reclaimed uses of slurs. The two keywords that I use seem to have lower false positives: only 1.6% of the Tweets were not considered hate speech or offensive by human annotators.

<sup>31</sup>Arguably, these users might be more empathetic and more likely to refer to the slurs rather than use them to attack.

the slurs in quotation marks (to capture users who are only referring to the slur), and those in the Holocaust sample who self-report being Jewish in their biographies. Users enter the sample once, so the filter also drops Tweets from duplicate users. This way, every observation in the sample is a user-Tweet pair, and I report users at most once.

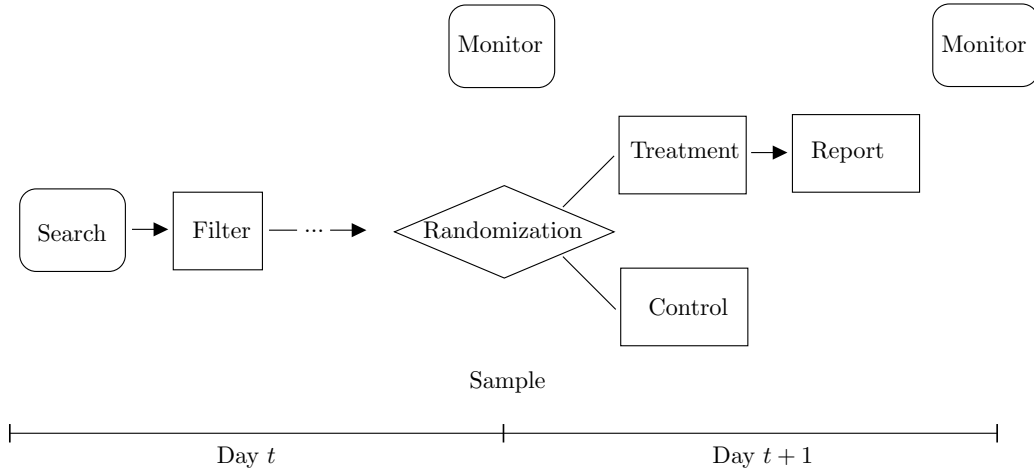
At midnight every day, immediately before randomization into treatment, my algorithm checks whether the users or Tweets collected the previous day were removed from the platform and only those that have not been removed at this point enter the final sample. Table 3 compares descriptive statistics between the experimental sample and the random sample of Tweets from section 3. These samples are quite different. Experimental subjects have more recent accounts, give more likes per day, and are more likely to have posted toxic Tweets in the past. Tweets in the experimental sample are more toxic, as expected. The Holocaust and disability samples are also different; for example, the Tweets and timelines of users from the Holocaust sample have a lower toxicity. Figure B.5 in the Appendix plots the most common topics in each subsample, which I obtained by annotating the Tweets on MTurk. Some common topics include politics, religion, sports, and COVID-19.

**Table 3:** Characteristics of the reporting experiment sample

	Means			Difference <i>t</i> -statistic	
	Full Sample	Holocaust	Disability	Random-Full	Hol.-Disab.
<i>Observations</i>	6,148	123	6,025		
<i>Accounts</i>					
Account years	3.22	3.29	3.22	40.2	0.2
Tweets per day	11.62	19.69	11.46	2.2	3.7
Likes per day	24.17	33.64	23.98	-32.3	2.1
Followers	634.85	1,436.41	618.49	2.1	1.7
Followed	433.75	554.98	431.27	7.6	1.6
Initial shadow ban	0.71	0.71	0.71		0.1
<i>Tweets</i>					
Word count	15.98	23.98	15.81	-14.1	6.8
Is toxic	0.80	0.06	0.82	-244.3	-22.0
Is hate (MTurk)	0.30	0.43	0.30	-63.6	3.1
Is reply	0.84	0.56	0.84	-48.4	-8.4
Is attack (MTurk)	0.78	0.24	0.79		-14.8
Is quote	0.07	0.02	0.07	7.9	-2.0
Is mention	0.85	0.67	0.85	-42.5	-5.5
Tweet from phone	0.79	0.49	0.80	-9.0	-8.3
<i>Timelines</i>					
Previous toxicity	0.93	0.69	0.94	-28.2	-11.1
Previous disability	0.39	0.15	0.40	-179.1	-5.6
Previous Holocaust	0.10	0.66	0.09	-6.0	21.9

Notes: This table presents means of characteristics in the reporting experiment sample and subsamples. It also presents *t*-statistics from tests of difference in means between the random and the experimental samples and between the Holocaust and disability subsamples.

**Treatment.** Figure 4 summarizes the experimental design and timing of the algorithm. Every day at midnight, my algorithm randomly splits users or Tweets sampled in the previous 24 hours, who have not been removed from the platform, into a control or a treatment arm. The assignment is stratified by sampling date and slur; every day, half of the Tweets using each slur enter each experimental arm. Users in the control arm do not receive any intervention. The treatment consists of reporting Tweets for violating Twitter’s rules against hateful conduct on the next day after they enter the sample, so Tweets can be reported between five and 48 hours after they are posted. The algorithm assigns the Tweets in the reporting arm evenly to one out of the 11 accounts that I use for reporting. Table B.4 in the Appendix displays summary statistics of the accounts that I used for reporting and Figure 5 includes screenshots of the reporting process.<sup>32</sup>

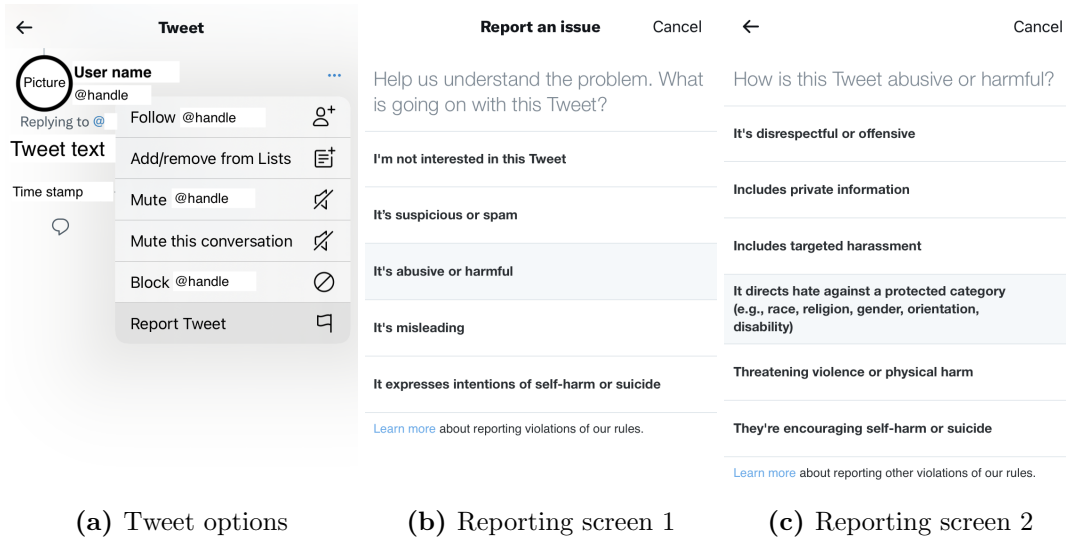


**Figure 4:** Design of the reporting experiment

Notes: Two main programs collect the sample and outcomes. The search program looks for hateful posts every 20 minutes and the monitor program keeps track of user activity every day. Randomization takes place at the beginning of every day with the sample of users collected the previous day.

Table B.3 in the Appendix shows that the two experimental arms are balanced in pre-treatment characteristics. Normalized differences—for each characteristic and all of them jointly—are well below the 0.25 value that Imbens and Rubin (2015) suggest, confirming that randomization was successful. I did not report 3.26% of the Tweets that were assigned to the reporting arm; that is, there is one-sided non-compliance. Three percent of Tweets in the reporting arm disappeared after treatment assignment and before I could report them, because users deleted them or deleted or protected their account, or because Twitter deleted the Tweets or suspended the users. Additionally, I did not report eight Tweets (0.26% of the

<sup>32</sup>When reporting Tweets, I click “It’s abusive or harmful,” then “It directs hate against a protected category (e.g., race, religion, gender, orientation, disability).” Due to logistics, 1% of the reported subjects were reported using a different account than the one that was assigned at the moment of randomization.



**Figure 5:** Procedure to report Tweets

Tweets in the reporting arm) that were clearly not hate speech.<sup>33</sup> Because of this one-sided non-compliance, the estimates can be interpreted as an intention-to-treat (ITT).

Reports are an instrument for content moderation, that is, for receiving any sanction from Twitter. First, Twitter uses reports to detect content and enforce its rules, which implies the relevance and monotonicity conditions of instrumental variables hold. Second, reports only affect user behavior through their effect on sanctions, so the exclusion restriction holds if sanctions are perfectly observed. To the best of my knowledge, Twitter does not notify users that they have been reported.<sup>34</sup>

**Outcomes.** I measure two types of outcomes: first-stage outcomes are the sanctions that Twitter enforces on users and second-stage outcomes are the users' activity on Twitter, their hatefulness, and spillovers to the activity of others. These outcomes allow testing whether reports influence moderation, whether Twitter's sanctions moderate users, and whether sanctions affect other users. I construct these outcomes with data that my algorithm collects every day. I gather users' cumulative number of Tweets, likes, accounts followed, and followers. I also collect the 100 most recent Tweets of each user (posted within 24 hours), and select 20 Tweets at random per user to compute their toxicity score by calling Perspective's API.

I measure sanctions as an absorbing state; that is, once users receive a sanction, they

<sup>33</sup>For example, one user quoted some people using the disability slur to refer to him or her. Other users posted the Holocaust-denial term quoting a study that was published around those dates (Center for Countering Digital Hate, 2021).

<sup>34</sup>Some users have received notifications from Twitter saying their posts were reported. According to the survey of section 5, 9% of users have received a notification that someone reported their Tweets. However, users seem to receive these notifications only when an account from Germany reports content, due to the Network Enforcement Act. Figure B.2 in the Appendix has a screenshot of one of these notifications.

remain sanctioned. By construction, at the time of entering the sample, none of the Tweets have been removed by Twitter and none of the users are suspended. However, 71% of users are initially shadowbanned.

I measure activity on Twitter as the time that users spend posting or liking Tweets, which corresponds to  $t$  in the model of Section 2. I do not directly observe time spent, but I construct a proxy using the number of Tweets that users post (that is, the statuses count object from the API) and the number of likes that they give (that is, the favorites count of the API). I then approximate the total number of words that users wrote and read during the period, by multiplying the Tweets and likes times the average number of words per Tweet in the random sample of Tweets, which is 13.81. Then, I convert words into time by using the average reading and typing speeds that have been documented in the literature.<sup>35</sup>

The main measure of hatefulness is the fraction of Tweets with a toxicity score higher than 0.8, but I consider alternative measures for robustness. Spillovers focus on the time spent by the users to whom the Tweets in the sample are replying (“replied users”); 86% of Tweets in the sample are replies to others. I focus on replied users because, arguably, users mentioned in a Tweet are more likely to notice sanctions related to the Tweet than others. Figure B.4 illustrates a reply to another Tweet.

**Empirical strategy.** This paper reports cross-sectional estimates of the effect of reporting users on different outcomes, three weeks after treatment assignment. I focus on first-stage and ITT estimates because, as the next subsection shows, I find evidence of unobservable sanctions which means that reports violate the exclusion restriction. In other words, reports affect outcomes not only through their impact on observable sanctions, but also through unobservable sanctions. Thus, I estimate regressions of the form:

$$Y_i = \alpha + \beta Z_i + \delta X_i + \varepsilon_i, \quad (6)$$

where  $i$  indexes user-Tweet pairs,  $Y_i$  denotes first-stage or second-stage outcomes,  $Z_i$  denotes treatment assignment (reports), and  $X_i$  is a vector of controls. I estimate specifications without controls, controlling for stratum—sampling date and slur—fixed effects, and adding controls from the rich set of pre-treatment characteristics of Table B.3. I select controls with a two-step method using lasso as suggested in Urminsky et al. (2016) with the methodology of Belloni et al. (2014). Regressions use robust standard errors unless noted otherwise.

---

<sup>35</sup>I use the words per Tweet from the random sample, as opposed to the value from the experimental sample, because this is the value that I pre-registered, before the experimental sample existed. The average typing speed on a desktop computer is 51.56 words per minute (WPM) according to Dhakal et al. (2018). The average typing speed on a mobile device is 36.2 WPM (Palin et al., 2019). Elliott et al. (2019) estimate a reading speed of 179 WPM that is constant across different devices and screen sizes. I obtain the device of a user from the source object of Tweets; I consider the device to be a desktop when the source is “Twitter Web App” and mobile for all the other sources.

I also estimate dynamic treatment effects, which is possible because my algorithm collects outcomes every day after users enter the sample. As preregistered, the main treatment effects are measured over a three-week period, but I also present estimates of a longer time horizon (close to 5 months). I use the efficient estimator proposed by Roth and Sant’Anna (2021), which is robust to heterogeneous treatment effects. Because reports are randomized every day, the design satisfies their assumptions of random treatment timing and no anticipation.<sup>36</sup> I use their method to obtain event-study estimates, where the event date is the number of days since a report. I construct the estimates on balanced panels but also report the treatment effect on attrition. The results use their Neyman-style pointwise confidence intervals and the sup- $t$  confidence bands of Montiel Olea and Plagborg-Møller (2019).

## 4.2 Results

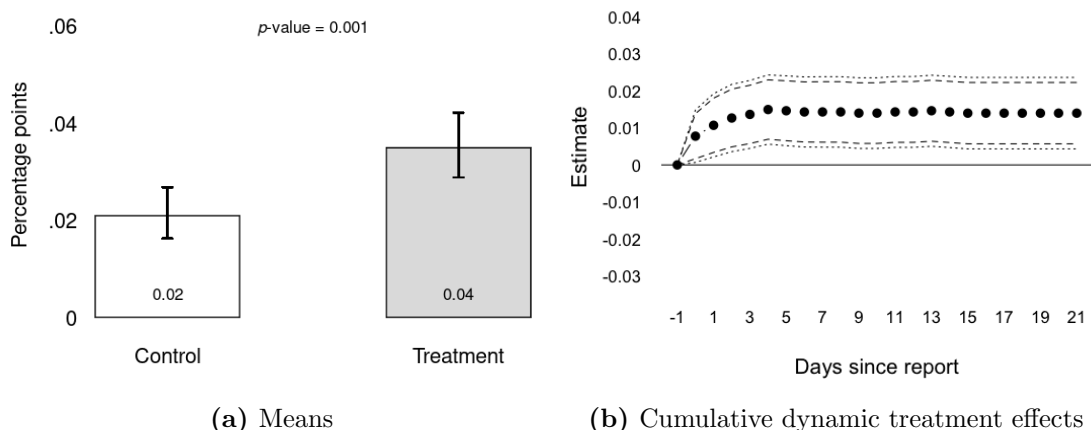
**Sanctions.** Reporting Tweets increases the likelihood that Twitter deletes them. Figure 6a shows the impact of assignment to treatment on the likelihood that Twitter removes the Tweets in the sample. Twitter removed 2.1% of the Tweets in the control arm within three weeks (21 days) after they entered the sample, and it removed 3.5% of them in the treatment arm. The treatment effect is 1.4 percentage points, which is 0.08 standard deviations, or a 66% increase. The  $p$ -value of the difference in proportions is 0.001, and the  $F$  statistic from a regression of deletions on treatment assignment is 11.01.<sup>37</sup> Figure 6b displays cumulative dynamic treatment effects on the likelihood of having a Tweet removed. Most of the effect on deletions occurs within the next few days after reporting (Figure C.5 confirms that the effect remains flat in a longer time horizon). Moreover, Figure C.2 shows that the effect is stable across reporting dates.

Table C.2 in the Appendix shows estimates of the effect of reporting on other Twitter sanctions and user self-censorship. Reporting does not significantly influence the other observable sanctions; that is, suspensions or shadowbans (see also Figure C.5 for dynamic effects). The table also displays insignificant effects on the likelihood of users deleting their own posts or accounts, or protecting their accounts (making them private) within three weeks after reporting. Moreover, it shows that reporting does not change the likelihood that other Tweets in the users’ profiles go missing, which includes self-removals and Twitter removals.<sup>38</sup>

<sup>36</sup>These assumptions hold within each stratum (sampling date and slur). However, since the Holocaust denial slur has few observations per day, I compute the estimators within each sampling date, pooling observations from both slurs.

<sup>37</sup>These numbers include cases in which Twitter required the removal of a Tweet but the user did not remove it within the three weeks. Eleven percent of Tweets were not removed by users in the control arm, and 5% were not removed in the treatment arm. These estimates keep all users, even those whose accounts were deleted. Results are unchanged if we drop them. Results from a two-stage least-squares regression that uses treatment assignment as an instrument for reports are the same.

<sup>38</sup>These numbers include the Tweets that users post after the sampling date and up to three weeks after the end of the sampling period. For these Tweets, distinguishing user deletions from Twitter deletions was



**Figure 6:** Likelihood that Twitter removes a post

Notes: Panel (a) displays the mean and 95% confidence intervals of the likelihood that Twitter removes a Tweet in the three weeks after reporting by treatment arm. The  $p$ -value is from a test of proportion differences. Panel (b) presents cumulative dynamic treatment effects of the likelihood of deletions, pointwise confidence intervals (dashed), and sup- $t$  simultaneous confidence bands (dotted). Dynamic effects use the estimator from Roth and Sant’Anna (2021).

The null effects persist after adding strata fixed effects and other controls, and the size of all estimates is below 0.033 standard deviations.

However, there is evidence of “unobservable” sanctions, such as accounts being temporarily locked (see Figure B.1f for an example), coming from two sources. First, Twitter sent updates informing that it found that 270 (8.8%) out of the 3,074 accounts on the reporting arm violated the rules, within three weeks of the reports (see Figure B.3 for a screenshot of an update and Table C.3 for updates and sanctions in the reported sample). Roughly four percent (117, 3.8%) of reports received an update but were not accompanied by Tweet deletions, user suspensions, or new shadowbans. This proportion—available only for reported users—provides one possible measure of unobservable sanctions. However, it likely understates the true fraction of unobservable sanctions because Twitter does not always send updates whenever it imposes a sanction. For instance, Twitter sent updates only for 13.4% of the 967 accounts in the reporting arm that received an observable sanction, and Figure C.6 in the Appendix shows that the likelihood of receiving updates as a fraction of reported accounts or of sanctioned reported accounts decreases over time.<sup>39</sup>

Indeed, as Appendix A.3 shows, under the (weak) assumption that Twitter does not send

not possible due to the API’s rate limits.

<sup>39</sup>The probability of receiving an update conditional on reporting was 8.8%. As a benchmark, an exercise conducted by the European Commission (Reynders, 2020) observed that Twitter sent an update on 26% of reports filed by the monitoring authorities. As will be clear in the heterogeneity section below, it is likely that Twitter treats reporting accounts differently. See also Gillespie (2018) for a discussion of platforms having “trusted flaggers.”



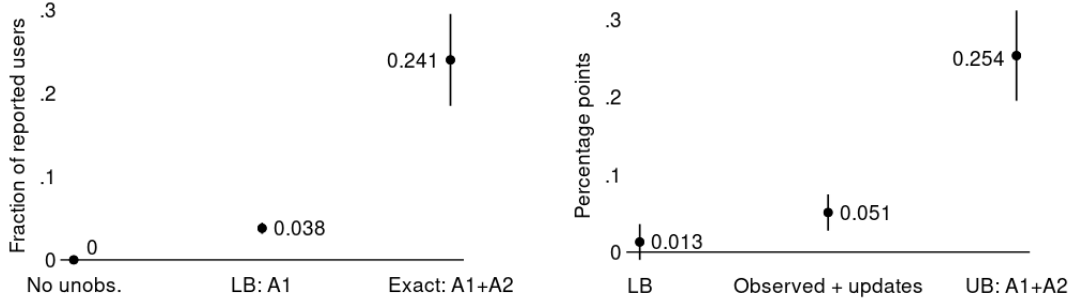
updates if it does not impose a sanction (A1), the proportion of reports with updates but no observable sanction actually gives a lower bound on the fraction of unobservable sanctions among reported accounts. As Figure 7a shows, if we assume A1 we can reject that there are no unobservable sanctions in the reporting treatment arm. Moreover, if we additionally impose that Twitter sends updates for observable and unobservable sanctions at the same rate (A2), we can identify the fraction of unobserved sanctions among reported accounts by rescaling the lower bound—dividing it by the probability of receiving an update among sanctioned, reported accounts.

On Figure 7b I report the overall first-stage effect on the “true” sanctions under these different assumptions about the unobservables. The lower bound on this effect is 0.013pp (a 4.8% increase relative to the control), assuming away the unobservable sanctions. The upper bound is 0.254pp (an 84.2% increase), imposing A1 and A2 and assuming no unobservables in the control group. An intermediate estimate is 0.051pp (a 16.9% increase), which adds the observable sanctions and the updates; that is, it also assumes that there are no unobservables in the control group, but uses the lower bound for the unobservables in the treatment group.

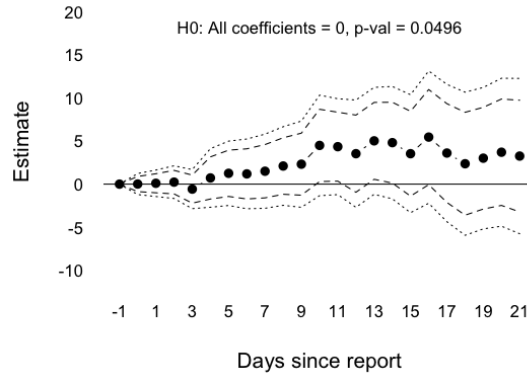
Besides using Twitter updates as a signal of unobservable sanctions, I also use information about the amount of time between users’ social media posts. Presumably, if their accounts are being locked, they would take longer hours in between posting. Figure 7c provides additional evidence in this vein; it plots daily treatment effects on the number of hours since the users’ last post, computed at midnight. The treatment effect is positive and pointwise significant around day 10 after reporting, although not significant with the simultaneous bands. However, we can reject the null hypothesis (with  $p$ -value = 0.0496) that the treatment effect is zero for all days jointly using a Wald test, which means that reporting increases the gap in between posting activity for at least one day. This test is conservative because the daily number of hours since last post might not reflect locking periods of less than 24 hours.

In short, there is evidence that reports induced Twitter to delete posts and to impose unmeasurable sanctions such as locking user accounts for some time. It is unlikely that the treatment induced other observable sanctions such as shadowbanning or suspending user accounts. While there is no similar information in the context of hateful conduct, the observed pattern is consistent with Twitter’s COVID-19 misinformation policy, which is to remove Tweets and lock accounts for users with a low number of violations and to suspend users with a high number of violations (Twitter, 2021a).

**Activity.** Reporting does not significantly decrease user activity on Twitter. Figure 8a displays the treatment effect on the number of hours that users spend posting and liking Tweets in the three weeks after reporting. Both treatment and control spent around three and a half hours, and the treatment effect is 0.25 hours (5 minutes per week), which is a 7.5% increase or .042 standard deviations. This effect, however, is not significant at conventional



(a) Unobservable sanctions, by assumption      (b) Bounds for the effect on true sanctions



(c) Effect on hours since last post

**Figure 7:** Evidence of “unobservable” sanctions

Notes: LB and UB denote lower and upper bounds, respectively. Panel (a) presents estimates of the fraction of unobservables under different assumptions. A1 assumes that Twitter does not send updates if it does not impose a sanction. A2 assumes that Twitter notifies observable and unobservable sanctions equally. The exact fraction was estimated with the Delta method. Panel (b) presents bounds of the treatment effects. LB assumes away unobservable sanctions in the control group. “Observed + updates” and UB assume away unobservable sanctions in the control group. The former also assumes that Twitter updates without an observable sanction perfectly indicate an unobservable sanction and the latter assumes A1 and A2. Panel (c) presents dynamic treatment effects (estimated as in Roth and Sant’Anna (2021)) on the number of hours since the last post, recorded at midnight every day after users entered the sample. It also shows results from a Wald test of joint equality of the coefficients. Pointwise 95% confidence intervals are dashed (or vertical, solid) lines and sup- $t$  simultaneous bands are dotted lines.

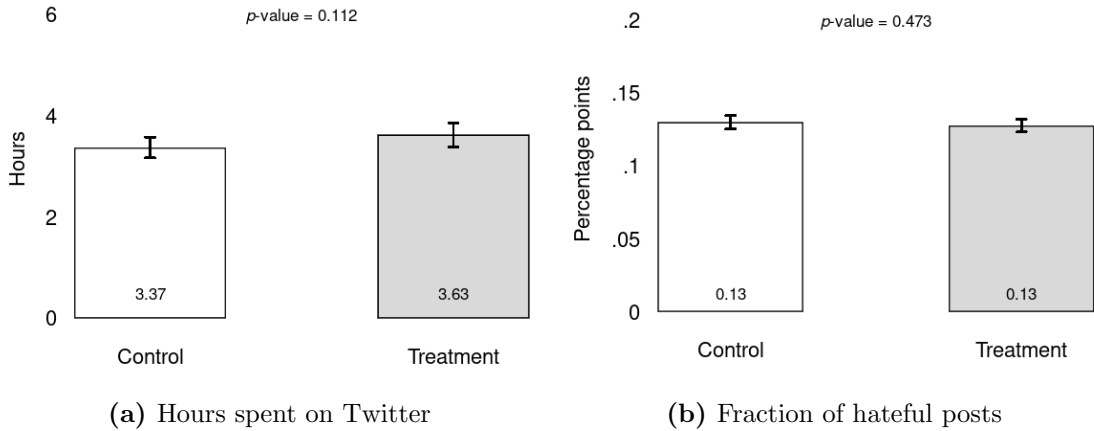
levels, because the  $p$ -value of the difference in means is 0.11. Figures C.7 in the Appendix shows that treatment effects remain flat throughout the period.

Table C.4 shows regression estimates using alternative measures of activity: Tweets and likes separately, a winsorized measure of time spent online removing the top and bottom percentiles, and an extensive-margin measure (the fraction of days that users post, like, or follow someone). The results remain unchanged using these alternative measures; if anything, the effect on Tweets is positive and significant at the 10% level under some specifications.

Moreover, these estimates are mechanically biased downward since Twitter might temporarily lock user accounts. These findings are consistent with more recent experimental and quasi-experimental work that finds that moderated users do not decrease their engagement beyond some transitory effect in the very short run (Katsaros et al., 2022; Ribeiro et al., 2022a).

**Hatefulness.** Reporting does not significantly decrease the likelihood of posting hate on Twitter. Figure 8b shows that the fraction of hateful Tweets (toxicity bigger than 0.8) that users post in the three weeks after the treatment is the same for both experimental arms. The treatment effect is -0.02 percentage points of hateful Tweets, which is a 1.7% decrease (-0.02 standard deviations). Figure C.7 shows a decrease in hatefulness in the first three to five days after reporting, but the effect returns to zero afterward. Table C.5 considers other measures of hatefulness; two extensive-margin measures (whether users post any Tweet with toxicity  $\geq 0.8$  or they repeat the slur), the average toxicity, and the average severe toxicity (another measure developed by Google). None of these measures yield significant effects, and the treatment effect is less than 0.011 standard deviations across all variables using different specifications.

These results are in line with more recent existing quasi-experimental evidence: Ribeiro et al. (2022a) find a null effect of hiding Facebook posts on users’ subsequent rule-breaking, although they find a negative effect of post deletions. In contrast, Katsaros et al. (2022) find that nudging users who are about to post toxic replies reduces their subsequent toxicity.



**Figure 8:** Hours spent on Twitter and fraction of hateful posts

Notes: This figure displays means and 95% confidence intervals of outcomes in the three weeks after reporting by treatment arm. Hours spent is calculated using statuses and favorites. Hateful posts are those with toxicity higher than 0.8. The  $p$ -value is from a test of difference in means.

**Spillovers.** Even if reporting does not seem to moderate the authors of the Tweets, that is, decrease their activity or hatefulness, it impacts other users. Figure 9a shows that reporting increases the time the replied users spend Tweeting and liking by 0.51 hours, which is 10 minutes per week, 10%, or 0.064 standard deviations ( $p$ -value = 0.028). The treatment effect seems somewhat persistent; Figure C.8 shows the cumulative effect increases continuously after the reporting day, and starts slowing down after roughly 2 months.

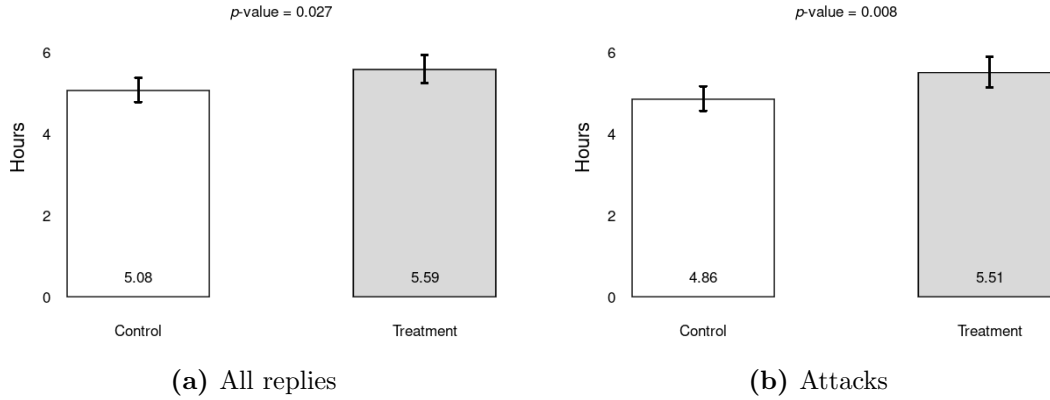
Although many of these replies are attacks, some are replies between social media friends. As specified in the pre-analysis plan, I asked MTurk workers to read the context of both posts and classify whether the replies in my sample were attacks on the replied user. Under the majority decision rule, in which Tweets are attacks if the majority of workers agree, 87% of replies were attacks on others. Figure 9b shows that the effect of reporting is stronger among attacks; it is 0.65 hours (13 minutes per week, a 13.4% increase or 0.08 standard deviations,  $p$ -value = 0.008). Again, Figure C.8 shows that the cumulative effect only becomes statistically insignificant after roughly 2 months following the reports.

As Table C.7 shows, the coefficient remains significant at the 5% level across specifications considering Tweets and likes separately or winsorizing time spent.<sup>40</sup> The effect is not, however, driven by the extensive margin, as there is an insignificant effect on the number of days that attacked users are active on the platform. Moreover, Table C.8 shows that there is an insignificant effect on the time spent of the followers of the posts' authors, and on their follower and following count. Hence, reporting seems to increase the activity of those users that are attacked by the Tweets in the sample without other effects on the network of users.

**Heterogeneity.** While this dimension of heterogeneity was not pre-registered—it was not anticipated, Figure 10 compares the first-stage success (likelihood of Tweet deletion) of each reporting account relative to the control group. As the figure shows, there was substantial heterogeneity in terms of how effective each of the 11 reporting accounts was at taking down content. A statistical test at conventional significance levels rules out not only that the effect of all accounts is jointly zero ( $p = 0.013$ ), but it also rules out that these accounts have the same effect ( $p = 0.028$ ). Only accounts 1, 2, and 4 managed to remove a statistically significant higher fraction of content than the rest of accounts, while the rest had an indistinguishable effect relative to the control group.

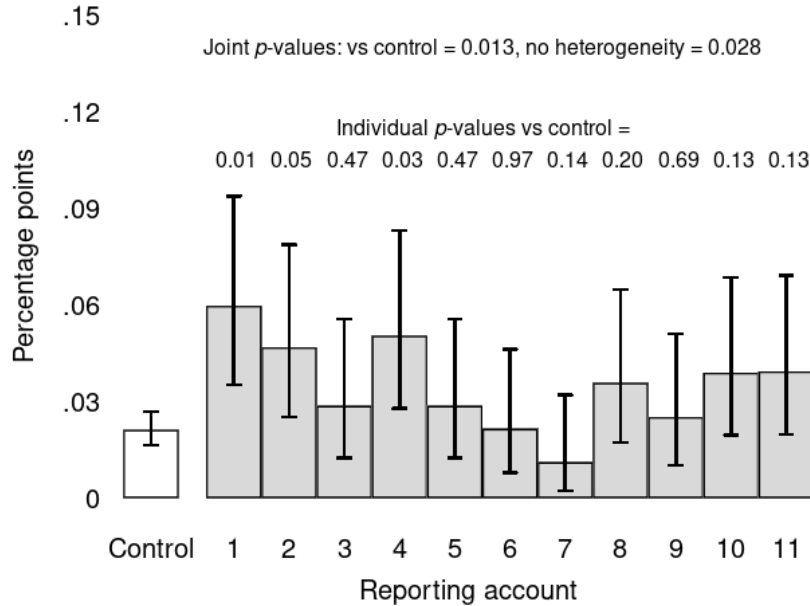
Why were some accounts more successful than others at taking down Tweets? Table B.4 in the Appendix displays summary characteristics of the reporting accounts and Table C.10 presents correlations between first-stage success and reporting account characteristics, from an OLS regression on the subsample of treated users. There is suggestive evidence that

<sup>40</sup>The findings are also robust to dropping those users who were replies to more than one Tweet in the hateful sample. More than 93% of the replied users corresponded to a single user in the hateful sample, so concerns about SUTVA violations are minimal.



**Figure 9:** Spillover on the time spent of users replied by the posts

Notes: This figure displays means and 95% confidence intervals of the time spent on Twitter in the three weeks after reporting by treatment arm. Panel (a) includes all users replied by the Tweets. Panel (b) includes users that were attacked by the Tweets, according to MTurk annotators. The  $p$ -value is from a test of difference in means.



**Figure 10:** Likelihood that Twitter removes a post, by reporting account

Notes: This figure displays means and 95% confidence intervals of the likelihood that Twitter removes a Tweet in the three weeks after reporting for the control group and by reporting account.  $p$ -values come from a regression of an indicator variable of Tweet deletion on dummies for each of the reporting accounts. Joint  $p$ -values are from  $F$  tests of 1) all coefficients being equal to zero, and 2) all coefficients being equal to each other. Individual  $p$ -values are from tests of each coefficient being equal to zero.

characteristics such as whether reports were made from a mobile device and whether the account's profile had a location were significantly correlated with a higher Tweet takedown.

While the usual caveat of omitted variables applies, these results are consistent with Gillespie (2018), who documents that some platforms assign internal reputation scores to flaggers, and hence do not treat equally all user reports.

As for other dimensions of heterogeneity of the first stage, Table C.1 presents a characterization of the compliers (1.4% of the sample; whose Tweets are deleted thanks to the reports), never-takers (96.5%; whose Tweets are never deleted), and always-takers (2.1%; whose Tweets are deleted regardless of reports). While there is not enough power to detect differences between compliers and others, the always-takers seem to have fewer followers, fewer accounts followed, and are more likely to have previous posts about the Holocaust. I also estimate conditional average treatment effects for Tweet deletions following the “generalized random forest” approach of Athey et al. (2019). Figure C.3 in the Appendix displays estimates of the “importance” each variable used to estimate the heterogeneous treatment effects. The word count of the sampled posts, the number of accounts that the users followed, their account age, the number of likes they give per day, and whether the posts were quoting another user seem to be important predictors of treatment effect heterogeneity.

The pre-analysis plan specified two other dimensions for the heterogeneity analysis: by slur (Holocaust vs disability) and by human annotation (among the hate sample). I report these results in Figure C.10 in the Appendix due to their low informational content.<sup>41</sup>

**Attrition, Multiple-Hypothesis Testing, and Inference Robustness** In this experiment, attrition occurs because accounts go missing after treatment assignment; 7% of them were missing after three weeks. Attrition happens when users delete their own accounts or Twitter suspends them. Given that the previous results showed no treatment effect on account suspensions or on the likelihood that users delete their accounts, finding no differential attrition by treatment arm is not surprising. Table C.9 shows insignificant treatment effects on the likelihood that users leave the sample at the end or on any day of the three weeks after users enter the sample. Figure C.9 shows dynamic treatment effects on attrition; the effect is not significant pointwise or with the simultaneous bands (within the 3 weeks after reporting or within the almost 5 months of full data collection).

Table C.11 presents Lee bounds of the main estimates. It also presents  $p$ -values based on different inference methods (robust standard errors, wild bootstraps, and permutation tests). Due to the multiple outcomes analyzed in this experiment, I also report the adjusted

---

<sup>41</sup>The experiment is not powered to detect the effect on the small Holocaust sample. The heterogeneity analysis by human annotation was intended to capture measurement error (false positives) in the sampling of hate speech. More than false positives, these labels seem to capture heterogeneity due to the subjective nature of hate speech. Thirty percent of Tweets in the sample were labeled as hate speech by the majority of annotators, 61% were considered offensive, 1.6% were not considered offensive or hate, and the remaining did not have a majority label. Hence, splitting the sample between “hate” and “not hate,” as preregistered, captures the difference between hateful and offensive Tweets.

$p$ -values for multiple-hypothesis testing developed by List et al. (2019). Overall, the first-stage effect on Tweet deletions, the spillover effect on the attacked users, the null effect on other sanctions, and the null effect on the authors’ activity are robust to these adjustments.

### 4.3 Discussion

The previous results indicate that reports instrument for sanctions, particularly Tweet removals and, potentially, unobservable sanctions. Moreover, the treatment did not decrease user activity or the likelihood of posting hate within three weeks; reports did not moderate users. The effect on the users’ activity is insignificant, which I interpret as a low elasticity of time spent with respect to moderation among the users in my sample;  $\partial T^H / \partial c \approx 0$  in the notation of the model of Section 2. This finding is consistent with more recent experimental and quasi-experimental studies (Katsaros et al., 2022; Ribeiro et al., 2022a).

Yet, reports spill over to other users; they increased the amount of time that the attacked users spent posting and liking. I interpret these findings as evidence of a positive elasticity of the time spent of some users in my sample with respect to moderation;  $\partial T^A / \partial c > 0$  in the notation of the model of Section 2. Thus, this result is consistent with the implication from the model that, in an internal equilibrium, small increases in moderation should increase the engagement of at least some users. This finding is in line with previous work from Matias (2019), who finds that community-driven moderation can result in an increase in engagement. Moreover, these results are consistent with a model in which platforms decide which user reports they ignore and which ones they act upon. Conditional on incurring in the cost of reviewing a report, platforms should act upon reports that generate a small cost for rule violators and a large benefit for others, and ignore those that generate the opposite results.

Three main mechanisms may explain why reports impacted the replied users. First, the reported users could have changed their behavior or their interactions with the replied users. Second, if Twitter removed the Tweets, the replied users could have noticed the legends that Twitter placed on the Tweets, as in Figure B.1a. Third, if the replied users also reported these Tweets, Twitter could have sent them an update on their reports, as in Figure B.3.

Regarding the first mechanism, the results in subsection 4.2 rule out that the reported users substantially changed their behavior. Additionally, Figure C.11 shows an insignificant effect on the likelihood that the users in the sample mention the replied users again within three weeks. Hence, the evidence in favor of this mechanism is weak. The same is true for the second mechanism. Table C.12 shows that the treatment effect on deletions is smaller in the sample of replies relative to the full sample, and insignificant in the sample of Tweets that attack others.

As for the third mechanism, the percentage of reports for which I received an update and

found no observable sanction is similar in the full sample, among replies, and among attacks (3.8%, 3.9%, and 4.0%, respectively). Hence, Twitter may have imposed an unobservable sanction (e.g., locking accounts) on the users who attacked others, and the attacked users who reported these Tweets may have received an update about the sanction. This hypothesis is of course difficult to test without access to internal data because user reports are unobservable, and it highlights some of the limitations of conducting an independent field experiment.<sup>42</sup>

How does reporting affect monetization? I obtain a back-of-the-envelope estimate as follows. The treatment increased by 10-15 minutes per week the time that reported users and replied users spent liking and posting. The advertising load on a small sample of 50 Tweets was one ad per four regular Tweets. Assume this number translates into an ad load of 0.25 minutes per minute of content consumed. Twitter’s Ad website has a default bid of \$0.21 per six-second video advertisement.<sup>43</sup> Ignoring effects on others, the treatment amounts to a \$5.25-\$7.88 increase in ad revenue per week per report.<sup>44</sup>

## 5 A Test of Overprovision or Underprovision

### 5.1 Experimental Design

**Sample.** I recruited 3,027 respondents in September 2021 through Luc.id, a widely used online marketplace that matches researchers with survey providers (Coppock and McClellan, 2019; Bursztyn et al., 2020). I pre-screened participants to select English speakers who live in the U.S., are over 18 years old, are willing to provide their email, self-report using Twitter, and pass a basic attention check. After the pre-screen, participants entered the online survey and had to answer demographic questions. The survey also asked them for their Twitter handle (optionally), which I used to get their account creation date, Tweet counts, and like counts. Sixty-four percent of participants provided a handle, and 74% of the handles were valid. This results in a sample size of 1,427 respondents, which satisfies the recommendation of Haaland et al. (2020) of 700 respondents per treatment arm.

Table 4 compares the characteristics of the sample with representative adult Twitter users from the ATP survey and with accounts from the random sample of Tweets. My survey undersamples users in the 18-29 age range, college graduates, and politically Independents, and oversamples white respondents and Democrats.<sup>45</sup> Users who provided their Twitter

---

<sup>42</sup>It is unclear whether users would reveal that they reported a particular Tweet in a survey, even if reporting is common (indeed, the next section shows that one-third of users have reported content).

<sup>43</sup>This price is for the general audience of U.S. adults. The ad price did not change when I tried targeting an ad to the list of users in the sample.

<sup>44</sup>Besides being a rough estimate, this calculation is based on a selected sample and ignores equilibrium effects, so it does not imply that Twitter would like to increase reports. Moreover, these numbers do not consider the marginal costs of moderating.

<sup>45</sup>I pre-registered introducing quotas to match representative Twitter users on gender, age, race or ethnicity, region, and political orientation, but relax the quotas to obtain the desired sample size was necessary.



handle have an older account, and fewer Tweets and likes per day relative to accounts in a random sample of Tweets.

**Table 4:** Characteristics of the welfare experiment sample

<i>Panel A: Demographics, N= 3,027</i>		
	Means (Survey)	ATP-Survey <i>t</i> -stat.
Age 18-29 (%)	24.48	3.01
Male (%)	53.88	-0.34
White (%)	68.19	-5.04
College graduate (%)	31.68	4.89
Republican (%)	22.76	-1.34
Democrat (%)	52.89	-9.41
<i>Panel B: Twitter accounts, N= 1,427</i>		
	Means (Survey)	Random-Survey <i>t</i> -stat.
Account years	7.93	-23.65
Likes per day	2.34	29.37
Tweets per day	1.54	41.76

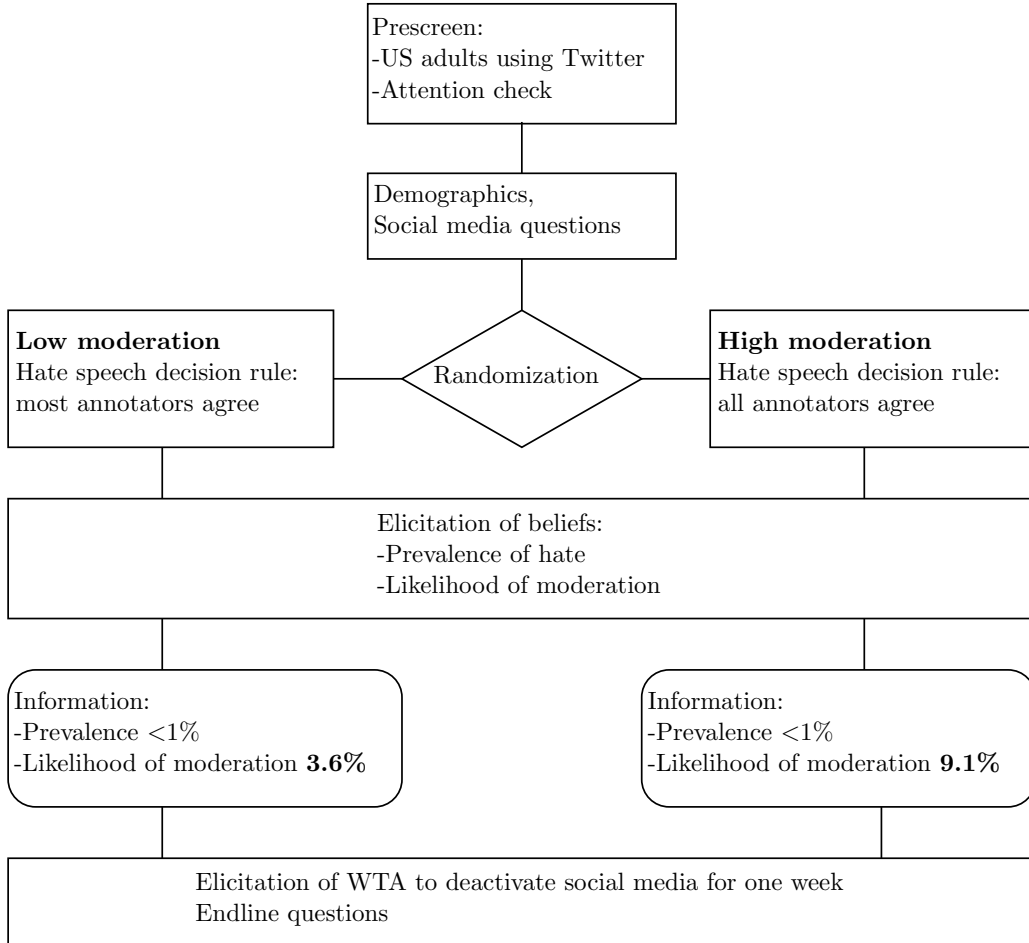
Notes: This tables presents means of characteristics in the welfare experiment sample. It also presents *t*-statistics from tests of difference in means between the ATP or the random sample of Tweets, and the experimental samples.

Afterward, the survey asked questions about social media use, online harassment, hate speech, and Twitter sanctions. These questions provide further insights about the previous experiment. Figure C.12 shows that the API-based measure of time spent on Twitter correlates closely with users' self-reported hours, so it is a good proxy measure. Table C.13 includes additional statistics. For instance, 32% of users have reported content for violating the rules, 10% have had a Tweet removed, and 5% have been suspended. Moreover, the experience on the platform differs by minority status, which I define based on religion (Jewish, Muslim, Buddhist, Hindu, or other), sexual preference (not heterosexual), gender (other than man or woman), and race (other than white). Consistent with other surveys (Anti-Defamation League, 2021), minorities are more likely to experience harassment online, to self-report seeing more hate speech in their feed, and to report content. However, they also receive more sanctions and reports, which, to the best of my knowledge, has not been documented before.<sup>46</sup>

**Treatment.** Figure 11 summarizes the experimental design. I use an information-provision treatment with an active control group (Haaland et al., 2020). After the baseline questions, I randomize survey participants into two treatment arms that receive different information about the likelihood of moderation among hateful Tweets. The information provided comes

<sup>46</sup>This finding is related to the literature on racial biases in detection algorithms; see, for example, Cowgill and Tucker (2019).

from the annotated random sample of 10,000 Tweets. To vary the likelihood of moderation without deception, I use different decision rules to classify hate speech. As Table 2 shows, 3.6% of hateful Tweets are removed or their authors are suspended within one month of the post under the majority decision rule, that is, if most annotators agree. That percentage changes to 9.1% under the consensus rule, that is, if all annotators agree. Half of participants are randomized into the low-moderation arm (3.6%) and half into the high-moderation arm (9.1%). The treatment is stratified by whether respondents are male, minorities, and have been sanctioned by Twitter, and whether they provided a Twitter handle.




**Figure 11:** Design of the welfare experiment

After randomizing participants, I inform them, for transparency, of the rule that I use to classify hate. As pictured in Figure B.6, I tell them that a crowd-sourced team of annotators identified hate speech using 10,000 Tweets, and that a Tweet is hate speech if [most/all] annotators label it as hateful. I then elicit their beliefs about (1) the prevalence of hate speech in this sample and (2) the fraction of Tweets that are removed or suspended within

one month. These elicitations are incentivized, because they know that one participant with the closest guess will get a \$50 Amazon gift card.


After the elicitation, I provide information about the likelihood of moderation, as displayed in Figure 12. I also hold constant the prevalence of hate speech in both arms, by telling respondents that less than 1% of Tweets are classified as hate (recall that 0.56% Tweets are hate under the majority rule and 0.11% are hate under the consensus rule). The message also shows that other popular platforms, such as YouTube, Facebook, and Reddit, have a similar prevalence of hate, according to different sources (Kennedy et al., 2020; Vidgen et al., 2020; Facebook, 2021). They can consult the sources by clicking a button on this screen.

Twitter **removed (de-platformed)**  
**3.6%** of hate speech Tweets or the  
accounts that posted them, within 1  
month

**Less than 1%** of Tweets in our sample  
were classified as hate speech. Other  
popular platforms (Youtube,  
Facebook, and Reddit) have a similar  
prevalence of hate 

(a) Low moderation

Twitter **removed (de-platformed)**  
**9.1%** of hate speech Tweets or the  
accounts that posted them, within 1  
month

**Less than 1%** of Tweets in our sample  
were classified as hate speech. Other  
popular platforms (Youtube,  
Facebook, and Reddit) have a similar  
prevalence of hate 

(b) High moderation

**Figure 12:** Information provision by treatment arm

Table B.5 shows that both experimental arms are balanced on pre-treatment characteristics. The table also rules out that changing the decision rule to classify hate influences the participants' concept of hate; the treatment has no effect on the belief about the prevalence of hate or the likelihood of moderation.<sup>47</sup>

**Outcomes.** There are two outcomes of interest. Based on the results of Section 2, the main outcome is the willingness to accept (WTA) to stop using social media, that is, Twitter, Facebook, Instagram, YouTube, Snapchat, TikTok, and Reddit, for one week. I first tell participants that the research team will conduct a small follow-up study that compensates some participants to deactivate their social media for one week. I inform them that similar studies have been conducted in the past (Hunt et al., 2018; Mosquera et al., 2020; Allcott et al., 2020). I then elicit their WTA with an iterative multiple price list (iMPL, see Harrison

<sup>47</sup>This finding is similar to what other studies obtain, such as Bottan and Perez-Truglia (2017), who argue that changing the source of information does not have an impact on participants who do not have expertise on the data.

et al. (2005); Andersen et al. (2006)).<sup>48</sup> Subjects have to decide whether they are willing to stop using social media for different Amazon gift card offers. The first offer is for \$50, and subsequent amounts increase or decrease until the WTA is placed in intervals that go from  $(-\infty, 0]$  to  $[100, \infty)$  and increase by \$10, as Figure B.7 illustrates. I transform these intervals into a continuous measure using the triangular distribution procedure from Allcott and Kessler (2019).

This elicitation is incentivized. I inform respondents that a computer will randomly choose some eligible participants whom the research team will contact for the follow-up.<sup>49</sup> If the participant is selected, the computer will also choose one of her answers at random. If the answer is “yes,” the research team will ask her to stop using social media for one week and pay the offered amount. If the answer is “no,” the participant will not be asked to stop using social media. This information is truthful; I recontacted 50 participants at random in October 2021 and implemented the follow-up study.<sup>50</sup>

The second outcome of interest is the API-based time spent on Twitter one week after the survey, which I compute for the participants who provided valid Twitter handles following the procedure outlined in section 4. At the end of the survey, I ask questions to measure attention, experimenter demand effects, and posterior beliefs.

**Empirical strategy.** The empirical strategy consists of OLS regressions of outcomes on an indicator of treatment status. All estimates use robust standard errors. I run regressions without controls, controlling for stratum fixed effects, and a specification adding controls as in Urminsky et al. (2016). As pre-registered, I report estimates of the main outcomes reweighting observations to match the ATP on first moments of gender, age, race or ethnicity, region, and political orientation, but I also report unweighted estimates. I obtain the weights using the maximum entropy approach of Hainmueller (2012).

## 5.2 Results

**Misperceptions about hate speech and moderation.** Most users overestimate the prevalence of hate speech on Twitter and the likelihood that Twitter sanctions hateful content. Figure C.13 displays histograms of beliefs among respondents. Ninety-six percent of

<sup>48</sup>The iMPL has two advantages over a regular multiple price list. First, it induces monotonicity on responses by construction. Second, it saves time by omitting redundant questions.

<sup>49</sup>Following Allcott et al. (2020), I did not tell participants the likelihood of being selected into the follow-up; previous research has shown that, at least on Becker-DeGroot-Marschak elicitations, informing participants can bias WTA estimates.

<sup>50</sup>Thirteen participants replied to the recontact email. Seven of them had been randomized into the deactivation treatment, and six to the control group. I asked participants in the deactivation arm to upload screenshots of the time-tracking app of their phones as proof of deactivation, as in Hunt et al. (2018). Five out of seven participants self-reported that they had stopped using social media, and four submitted the screenshots.

Twitter users overestimate the prevalence of hate speech, that is, their belief is above 1%, and 84% guess a moderation rate above the higher 9.1% value. These results add another example to the literature on misperceptions about others (Bursztyn and Yang, 2021).

There are several explanations for these facts. An “echo-chamber” argument is that users might not notice what happens outside their curated feeds, which they personalize with the help of Twitter’s algorithms. Consistent with this argument, I find that 74% of users believe that the prevalence of hate in the random sample of Tweets is higher than what they see in their feed. Platforms’ lack of transparency might also contribute to misperceptions. Even Facebook, which publishes a substantial amount of information (Facebook, 2021), informs only about the prevalence of hate speech but not about the likelihood of moderation (Bradford et al., 2019). The only information about the likelihood of moderation, between 3 to 5% of hateful content, was revealed thanks to the recent whistleblowing incident (Giansiracusa, 2021).

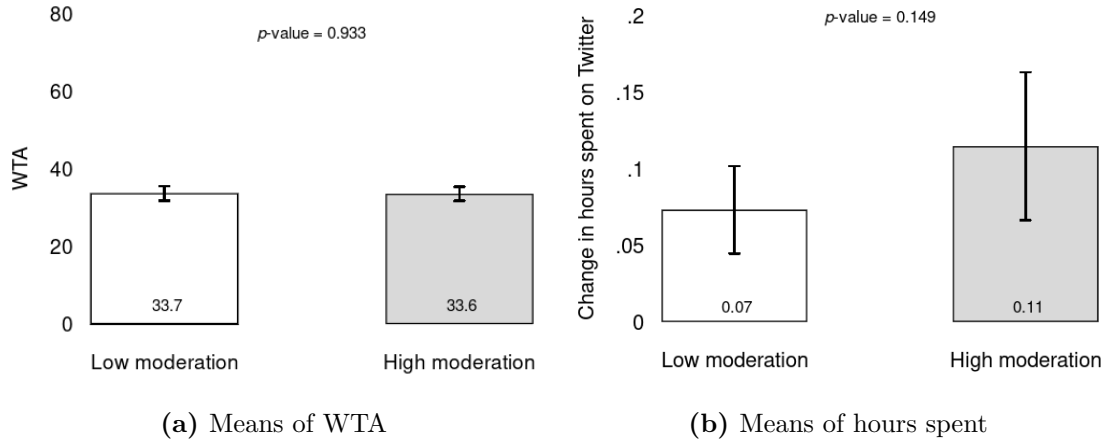
**WTA to stop using social media.** Providing information about a higher likelihood of moderation has little effect on the users’ social-media valuation. Figure 13a displays the treatment effect on the WTA to stop using social media during one week. The average WTA was \$33.6 in the low moderation arm, and \$33.7 in the high moderation arm. The treatment effect is -15 cents per week, which is 0.004 standard deviations, or a 0.5% decrease. The null effect is not just on average; Figure C.14 in the Appendix shows that the cumulative distribution function of WTA is the same for both arms. Table C.14 presents regression estimates with alternative measures of social-media valuation. As in Allcott and Kessler (2019), I assume a uniform distribution of WTA beyond the endpoints instead of the triangular distribution. I also use -\$50 and \$150 for the endpoints as benchmarks, or a take-it-or-leave-it dummy for the first \$50 offer. The results remain unchanged using these alternative measures.

**Activity.** The information provision treatment has an positive but insignificant effect on the time that users spent on Twitter one week after the survey. Figure 13a plots the effect on the number of hours spent by users who provided their Twitter handle. The effect is 0.04 hours, which is 2.4 minutes (57% increase relative to the low-moderation arm, or 0.077 standard deviations).<sup>51</sup>

**Posterior beliefs, attention, attrition, and experimenter demand.** Respondent inattention cannot explain the previous null results; providing information significantly shifts

---

<sup>51</sup>Table C.15 shows estimates using the same alternative measures of activity as in section 4; Tweets and likes separately, winsorized hours, and an extensive-margin measure of the fraction of days in which users post or like. The effect remains insignificant with these measures across specifications. Figure C.15 confirms that dynamic treatment effects remain flat throughout the week post-survey.



**Figure 13:** WTA to stop using social media and hours spent on Twitter

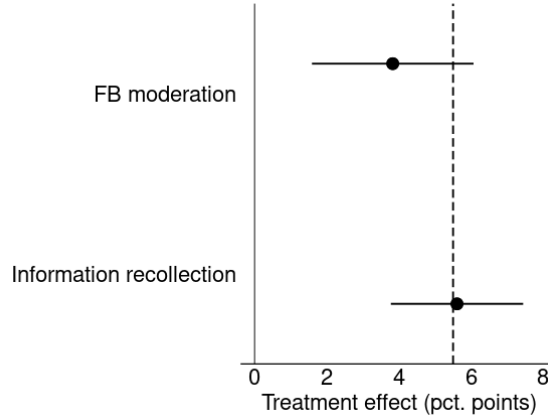
Notes: Panel (a) displays the mean and 95% confidence intervals of the WTA to stop using social media one week by treatment arm. Panel (b) presents the mean and 95% confidence intervals of the hours spent on Twitter one week after the survey. The  $p$ -values are from a test of difference in means and observations are reweighted to match Twitter users from the ATP on observables.

participant’s recollection of the information provided and their posterior beliefs about moderation. At the end of the survey, I asked participants to repeat the moderation rate that I gave them, and I incentivized the closest answer with a \$50 gift card. Figure 14 plots the effect on the respondents’ recollection of the moderation information. Sixty percent of participants recalled a number within one percentage point of the true value.<sup>52</sup> The treatment effect on this recollection is 5.6 percentage points (53% or 0.425 standard deviations, with an  $F$ -statistic of 36), not statistically different from 5.5 ( $p$ -value = 0.907), which is the gap between the high moderation rate (9.1%) and the low moderation rate (3.6%).

Figure 14 also plots the treatment effect of information on users’ beliefs about the likelihood of moderation on Facebook. This follows the recommendation of Haaland et al. (2020), of measuring posteriors by asking post-treatment beliefs about a related but different variable. The average belief of the moderation rate on Facebook was 19% for users in the low-moderation arm and 22.9% in the high-moderation arm. The treatment effect was 3.8 percentage points (20% or 0.16 standard deviations, with an  $F$ -statistic of 11.2).

Additionally, there is no evidence of differential inattention, attrition, or experimenter demand effects by treatment arm. Table C.16 presents insignificant treatment effects on inattention, measured as the absolute difference between participants’ recollection and the information provided. Thirty-four participants (1.1% of the sample) completed the prescreening questions but did not finish the survey, and Table C.16 shows null treatment effects on attrition under different specifications. Following Allcott et al. (2020), the last part of the

<sup>52</sup>Because of left-digit bias, many participants in the low-moderation arm remembered 3%.



**Figure 14:** Posterior beliefs about moderation on Facebook and attention check

Notes: This figure presents coefficients and 95% confidence intervals of OLS regressions on an indicator of the high-moderation arm. FB moderation is the users’ beliefs about the fraction of posts or users that Facebook moderates. The attention check is the participants’ recollection, at the end of the survey, of the information provided about the moderation rate on Twitter. The dashed line is at 5.5 percentage points, the difference between the high moderation rate (9.1%) and the low moderation rate (3.6%). Observations are reweighted to match Twitter users from the ATP on observables.

survey asked a question to test for experimenter demand effects: “Do you think the researchers in this study had an agenda?” Similar to that study, 57% of respondents in both arms thought I had no particular agenda or were not sure. Figure C.16 shows insignificant treatment effects on the responses to that question.

**Heterogeneity.** I do not find substantial heterogeneity of the effect on the WTA across most of the pre-registered covariates, including minority status (as defined above), whether participants have experienced a sanction on Twitter, and whether their beliefs are above or below the median moderation belief of 33% of hateful Tweets. The exception is the time spent on Twitter after the survey. Figure C.17 in the Appendix shows suggestive evidence that minorities spend more time on Twitter when they receive the high moderation information. The treatment effect in this subsample is 0.054 hours (three minutes, 100% increase relative to the control group, 0.17 standard deviations,  $p$ -value = 0.03).<sup>53</sup>

### 5.3 Discussion

The previous results indicate that providing information about a higher moderation rate shifted users’ beliefs, but had little impact on their social-media valuation. Taken at face value, these results mean that Twitter does not moderate too much or too little from the

<sup>53</sup>The treatment effect among minorities is significant at the 10% without reweighting observations. Figure C.17 also shows large point estimates on the subsample of users who have been sanctioned and those with high prior beliefs, although these are noisily estimated.

consumers’ point of view, for a fixed prevalence of hate speech. One explanation for this finding is that Twitter internalizes the impact of moderation on users’ willingness to pay for the platform, which requires that marginal and inframarginal users respond similarly to sanctions.

Another option is that users do not directly care about moderation, holding constant the hate they encounter. Indeed, it is possible that the experiment did not change users’ perceptions about hate in their own feed. In that case, users could have differentially updated their beliefs about how effective the algorithms are at hiding content without moderating. This is consistent with platforms providing a wide range of tools that allow users to customize their experience. For instance, Twitter allows users to mute and block accounts and words, and to hide sensitive content from their feeds.

One challenge to the interpretation of these findings comes from the welfare discussion in Allcott et al. (2020). They argue that users might misperceive Facebook’s value, and thus the WTA might overstate consumer surplus. These value misperceptions could explain why increasing perceived moderation did not impact users’ WTA. Another challenge is that the treatment not only shifted users’ beliefs about moderation on Twitter; it also impacted beliefs about moderation on other platforms (at least Facebook). Based on Appendix A, the correct measure of the change in consumer surplus is to consider current social media users, not just current Twitter users. Table C.17 in the Appendix shows that results are unchanged after reweighting observations to match representative social media users, or without reweighting.

There is also suggestive evidence that the treatment increased minorities’ time spent on Twitter. Given that these individuals are more likely to experience harassment online (Table C.13), this is consistent with the finding from the previous experiment that reporting increases the activity of the targets of hate speech.

## 6 Conclusions

Simple economics explain why it makes sense for profit-maximizing social media companies to ban some of their customers or restrict their content: because this increases the willingness to pay of marginal users. In an advertising-driven business model, platforms remove content only if this increases the time that some users spend consuming content, and hence interacting with ads. I find evidence consistent with this implication, by running a natural field experiment in which I report content that violates Twitter’s rules against hateful conduct. Reports increase Tweet removals and, potentially, unobservable sanctions, and they do not decrease user activity or hatefulness. Yet, the targets of hateful posts increase their activity after the reports. While this treatment provides some evidence of the behavioral effects of moderation, further work is needed to understand repeated sanctions, different classes of platform interventions, or the effects of moderation on other types of content.



In terms of policy, both sides in the discussion of how to regulate platforms often mention a tension between profit maximization and optimality of content moderation. While platforms can, in theory, remove too little or too much content relative to a surplus-maximizing planner, this study finds no evidence of distortions from the consumers' point of view. There are, however, two caveats to these findings. First, consumer surplus ignores the costs that hate speech imposes outside platforms. Hence, an avenue for future research is to examine the costs and benefits of the real-world consequences of content moderation. Second, even without moderation distortions, imperfect competition between platforms likely leads to pricing distortions, so they might be setting the ad loads of haters or non-haters suboptimally. These distortions can be empirically confirmed by future work.

## References

- Acemoglu, D., A. Ozdaglar, and J. Siderius (2021). Misinformation: Strategic sharing, homophily, and endogenous echo chambers. Technical report, National Bureau of Economic Research.
- Albert, A. B., H. E. Jacobs, and G. N. Siperstein (2016). Sticks, stones, and stigma: Student bystander behavior in response to hearing the word “retard”. *Intellectual and developmental disabilities* 54(6), 391–401.
- Ali, S., M. H. Saeed, E. Aldreabi, J. Blackburn, E. De Cristofaro, S. Zannettou, and G. Stringhini (2021). Understanding the effect of deplatforming on social networks. In *13th ACM Web Science Conference 2021*, pp. 187–195.
- Allcott, H., L. Braghieri, S. Eichmeyer, and M. Gentzkow (2020). The welfare effects of social media. *American Economic Review* 110(3), 629–76.
- Allcott, H. and M. Gentzkow (2017). Social media and fake news in the 2016 election. *Journal of Economic Perspectives* 31(2), 211–36.
- Allcott, H., M. Gentzkow, and L. Song (2022). Digital addiction. *American Economic Review* 112(7), 2424–63.
- Allcott, H. and J. B. Kessler (2019). The welfare effects of nudges: A case study of energy use social comparisons. *American Economic Journal: Applied Economics* 11(1), 236–76.
- Álvarez-Benjumea, A. and F. Winter (2018). Normative change and culture of hate: An experiment in online environments. *European Sociological Review* 34(3), 223–237.
- Andersen, S., G. W. Harrison, M. I. Lau, and E. E. Rutström (2006). Elicitation using multiple price list formats. *Experimental Economics* 9(4), 383–405.
- Anti-Defamation League (2021). Online hate and harassment. the american experience 2021. *Center for Technology and Society*. Accessed: 2021-10-23.
- Arango, A., J. Pérez, and B. Poblete (2019). Hate speech detection is not as easy as you may think: A closer look at model validation. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 45–54.
- Aridor, G. (2022). Drivers of digital attention: Evidence from a social media experiment. *Available at SSRN 4069567*.
- Athey, S., J. Tibshirani, S. Wager, et al. (2019). Generalized random forests. *Annals of Statistics* 47(2), 1148–1178.
- Becker, G. S. and K. M. Murphy (1993). A simple theory of advertising as a good or bad. *The Quarterly Journal of Economics* 108(4), 941–964.
- Beknazar-Yuzbashev, G., R. Jiménez-Durán, M. Stalinski, and J. McCrosky (2022). Toxic content and user engagement on social media: Evidence from a field experiment. *Working paper*.

- Beknazar-Yuzbashev, G. and M. Stalinski (2022). Do social media ads matter for political behavior? A field experiment. *Journal of Public Economics* 214, 104735.
- Belloni, A., V. Chernozhukov, and C. Hansen (2014). Inference on treatment effects after selection among high-dimensional controls. *The Review of Economic Studies* 81(2), 608–650.
- Berry, S., A. Gandhi, and P. Haile (2013). Connected substitutes and invertibility of demand. *Econometrica* 81(5), 2087–2111.
- Bianchi, C. (2014). Slurs and appropriation: An echoic account. *Journal of Pragmatics* 66, 35–44.
- Bottan, N. L. and R. Perez-Truglia (2017). Choosing your pond: Revealed-preference estimates of relative income concerns. *Available at SSRN 2944427*.
- Boxell, L., M. Gentzkow, and J. M. Shapiro (2019). Cross-Country trends in affective polarization. *Working Paper*.
- Bradford, B., F. Grisel, T. L. Meares, E. Owens, B. L. Pineda, J. N. Shapiro, T. R. Tyler, and D. E. Peterman (2019). Report of the Facebook Data Transparency Advisory Group. *Yale Justice Collaboratory*.
- Braghieri, L., R. Levy, and A. Makarin (2021). Social media and mental health. *Available at SSRN*.
- Bundesamt für Justiz (2019). Federal Office of Justice issues fine against Facebook. [https://www.bundesjustizamt.de/DE/Presse/Archiv/2019/20190702\\_EN.html](https://www.bundesjustizamt.de/DE/Presse/Archiv/2019/20190702_EN.html). Accessed: 2021-09-30.
- Bursztyn, L., G. Egorov, R. Enikolopov, and M. Petrova (2019). Social media and xenophobia: evidence from Russia. Technical report, National Bureau of Economic Research.
- Bursztyn, L., A. L. González, and D. Yanagizawa-Drott (2020). Misperceived social norms: Women working outside the home in Saudi Arabia. *American Economic Review* 110(10), 2997–3029.
- Bursztyn, L., I. K. Haaland, A. Rao, and C. P. Roth (2020). Disguising prejudice: Popular rationales as excuses for intolerant expression. Technical report, National Bureau of Economic Research.
- Bursztyn, L. and D. Y. Yang (2021). Misperceptions about others. Technical report, National Bureau of Economic Research.
- Carlson, C. R. (2021). *Hate Speech*. MIT Press.
- Carlson, C. R. and H. Rousselle (2020). Report and repeat: Investigating Facebook’s hate speech removal process. *First Monday*.
- Center for Countering Digital Hate (2021). Failure to protect: How tech giants fail to act on user reports of antisemitism. Technical report.
- Chandrasekharan, E., U. Pavalanathan, A. Srinivasan, A. Glynn, J. Eisenstein, and E. Gilbert

- (2017). You can’t stay here: The efficacy of reddit’s 2015 ban examined through hate speech. *Proceedings of the ACM on Human-Computer Interaction 1* (CSCW), 1–22.
- Cheng, J., M. Bernstein, C. Danescu-Niculescu-Mizil, and J. 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*, pp. 1217–1230.
- Chowdhury, F. A., L. Allen, M. Yousuf, and A. Mueen (2020). On Twitter purge: A retrospective analysis of suspended users. In *Companion Proceedings of the Web Conference 2020*, pp. 371–378.
- Coppock, A. and O. A. McClellan (2019). Validating the demographic, political, psychological, and experimental results obtained from a new source of online survey respondents. *Research & Politics 6*(1), 2053168018822174.
- Correia-da Silva, J., B. Jullien, Y. Lefouili, and J. Pinho (2019). Horizontal mergers between multisided platforms: Insights from cournot competition. *Journal of Economics & Management Strategy 28*(1), 109–124.
- Cowgill, B. and C. E. Tucker (2019). Economics, fairness and algorithmic bias. *preparation for: Journal of Economic Perspectives*.
- Davidson, T., D. Warmley, M. Macy, and I. Weber (2017). Automated hate speech detection and the problem of offensive language. In *Proceedings of the International AAAI Conference on Web and Social Media*, Volume 11.
- Dhakal, V., A. M. Feit, P. O. Kristensson, and A. Oulasvirta (2018). Observations on typing from 136 million keystrokes. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*, pp. 1–12.
- Dixon, L., J. Li, J. Sorensen, N. Thain, and L. Vasserman (2018). Measuring and mitigating unintended bias in text classification. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*, pp. 67–73.
- Elliott, L. J., M. Ljubijanac, and D. Wiczorek (2019). The effect of screen size on reading speed: A comparison of three screens to print. In *International Conference on Applied Human Factors and Ergonomics*, pp. 103–109. Springer.
- ElSherief, M., V. Kulkarni, D. Nguyen, W. Y. Wang, and E. Belding (2018). Hate lingo: A target-based linguistic analysis of hate speech in social media. In *Proceedings of the International AAAI Conference on Web and Social Media*, Volume 12.
- Enikolopov, R., A. Makarin, and M. Petrova (2020). Social media and protest participation: Evidence from Russia. *Econometrica 88*(4), 1479–1514.
- Facebook (2021). Community standards enforcement report, second quarter 2021. <https://about.fb.com/news/2021/08/community-standards-enforcement-report-q2-2021/>. Accessed: 2021-10-23.
- Fergusson, L. and C. Molina (2019). Facebook causes protests. *Documento CEDE* (41).

- Filippas, A. and J. J. Horton (2021). The production and consumption of social media. Technical report, National Bureau of Economic Research.
- Fortuna, P. and S. Nunes (2018). A survey on automatic detection of hate speech in text. *ACM Computing Surveys (CSUR)* 51(4), 1–30.
- Fortuna, P., J. Soler-Company, and L. Wanner (2021). How well do hate speech, toxicity, abusive and offensive language classification models generalize across datasets? *Information Processing & Management* 58(3), 102524.
- Founta, A. M., C. Djouvas, D. Chatzakou, I. Leontiadis, J. Blackburn, G. Stringhini, A. Vakali, M. Sirivianos, and N. Kourtellis (2018). Large scale crowdsourcing and characterization of twitter abusive behavior. In *Twelfth International AAAI Conference on Web and Social Media*.
- Fujiwara, T., K. Müller, and C. Schwarz (2021). The effect of social media on elections: Evidence from the United States. Technical report, National Bureau of Economic Research.
- Gadde, V. and K. Beykpour (2018). Setting the record straight on shadow banning. [https://blog.twitter.com/en\\_us/topics/company/2018/Setting-the-record-straight-on-shadow-banning](https://blog.twitter.com/en_us/topics/company/2018/Setting-the-record-straight-on-shadow-banning). Accessed: 2021-10-13.
- Gentzkow, M. (2007). Valuing new goods in a model with complementarity: Online newspapers. *American Economic Review* 97(3), 713–744.
- Gentzkow, M., B. Kelly, and M. Taddy (2019). Text as data. *Journal of Economic Literature* 57(3), 535–74.
- Giansiracusa, N. (2021). Facebook uses deceptive math to hide its hate speech problem. *Wired*.
- Gillespie, T. (2018). *Custodians of the Internet*. Yale University Press.
- Guhl, J. and J. Davey (2020). Hosting the ‘Holohoax’: A snapshot of holocaust denial across social media. *The Institute for Strategic Dialogue*.
- Haaland, I., C. Roth, and J. Wohlfart (2020). Designing information provision experiments.
- Hainmueller, J. (2012). Entropy balancing for causal effects: A multivariate reweighting method to produce balanced samples in observational studies. *Political analysis* 20(1), 25–46.
- Han, S.-H. and L. M. Brazeal (2015). Playing nice: Modeling civility in online political discussions. *Communication Research Reports* 32(1), 20–28.
- Han, X. and Y. Tsvetkov (2020). Fortifying toxic speech detectors against veiled toxicity. *arXiv preprint arXiv:2010.03154*.
- Hangartner, D., G. Gennaro, S. Alasiri, N. Bahrach, A. Bornhoft, J. Boucher, B. B. Demirci, L. Derksen, A. Hall, M. Jochum, et al. (2021). Empathy-based counterspeech can reduce racist hate speech in a social media field experiment. *Proceedings of the National Academy of Sciences* 118(50), e2116310118.

- Harrison, G. W., M. I. Lau, E. E. Rutström, and M. B. Sullivan (2005). Eliciting risk and time preferences using field experiments: Some methodological issues. In *Field experiments in economics*. Emerald Group Publishing Limited.
- Huang, J., D. Reiley, and N. Riabov (2018). Measuring consumer sensitivity to audio advertising: A field experiment on Pandora internet radio. *Available at SSRN 3166676*.
- Hunt, M. G., R. Marx, C. Lipson, and J. Young (2018). No more FOMO: Limiting social media decreases loneliness and depression. *Journal of Social and Clinical Psychology* 37(10), 751–768.
- Imbens, G. W. and D. B. Rubin (2015). *Causal inference in statistics, social, and biomedical sciences*. Cambridge University Press.
- Jhaver, S., C. Boylston, D. Yang, and A. Bruckman (2021). Evaluating the effectiveness of deplatforming as a moderation strategy on Twitter.
- Jhaver, S., A. Bruckman, and E. Gilbert (2019). Does transparency in moderation really matter? User behavior after content removal explanations on Reddit. *Proceedings of the ACM on Human-Computer Interaction* 3(CSCW), 1–27.
- Jiménez Durán, R., K. Müller, and C. Schwarz (2022). The effect of content moderation on online and offline hate: Evidence from Germany’s NetzDG. *Available at SSRN 4230296*.
- Jourová, V. (2016). Code of conduct on countering illegal hate speech online: First results on implementation. Technical report, European Commission, Directorate-General for Justice and Consumers.
- Katsaros, M., K. Yang, and L. Fratamico (2022). Reconsidering Tweets: Intervening during Tweet creation decreases offensive content. In *Proceedings of the International AAAI Conference on Web and Social Media*, Volume 16, pp. 477–487.
- Kaye, D. (2019). *Speech Police: The Global Struggle to Govern the Internet*. Columbia Global Reports.
- Kennedy, C. J., G. Bacon, A. Sahn, and C. von Vacano (2020). Constructing interval variables via faceted rasch measurement and multitask deep learning: a hate speech application. *arXiv preprint arXiv:2009.10277*.
- Kim, J. W., A. Guess, B. Nyhan, and J. Reifler (2021). The distorting prism of social media: How self-selection and exposure to incivility fuel online comment toxicity. *Journal of Communication* 71(6), 922–946.
- Levy, R. (2021). Social media, news consumption, and polarization: Evidence from a field experiment. *American Economic Review* 111(3), 831–70.
- List, J. A., A. M. Shaikh, and Y. Xu (2019). Multiple hypothesis testing in experimental economics. *Experimental Economics* 22(4), 773–793.
- Liu, Y., P. Yildirim, and Z. J. Zhang (2021). Implications of revenue models and technology for content moderation strategies. *Implications of Revenue Models and Technology for Content Moderation Strategies (November 23, 2021)*.

- Madio, L. and M. Quinn (2021). Content moderation and advertising in social media platforms. *Available at SSRN 3551103*.
- Marbach, M. and D. Hangartner (2020). Profiling compliers and noncompliers for instrumental-variable analysis. *Political Analysis* 28(3), 435–444.
- Matias, J., A. Johnson, W. E. Boesel, B. Keegan, J. Friedman, and C. DeTar (2015). Reporting, reviewing, and responding to harassment on Twitter. *Available at SSRN 2602018*.
- Matias, J. N. (2019). Preventing harassment and increasing group participation through social norms in 2,190 online science discussions. *Proceedings of the National Academy of Sciences* 116(20), 9785–9789.
- Melendez, S. (2020). Twitter automatically flags more than half of all Tweets that violate its rules. *Fast Company*. Accessed: 2021-10-13.
- Melnikov, N. (2021). Mobile internet and political polarization. *Available at SSRN 3937760*.
- Merrer, E. L., B. Morgan, and G. Trédan (2020). Setting the record straighter on shadow banning. *arXiv preprint arXiv:2012.05101*.
- Montiel Olea, J. L. and M. Plagborg-Møller (2019). Simultaneous confidence bands: Theory, implementation, and an application to svars. *Journal of Applied Econometrics* 34(1), 1–17.
- Mosquera, R., M. Odunowo, T. McNamara, X. Guo, and R. Petrie (2020). The economic effects of facebook. *Experimental Economics* 23(2), 575–602.
- Müller, K. and C. Schwarz (2020a). Fanning the flames of hate: Social media and hate crime. *Journal of the European Economic Association*.
- Müller, K. and C. Schwarz (2020b). From hashtag to hate crime: Twitter and anti-minority sentiment. *Available at SSRN 3149103*.
- Munger, K. (2017). Tweetment effects on the tweeted: Experimentally reducing racist harassment. *Political Behavior* 39(3), 629–649.
- Munger, K. (2021). Don’t@ me: Experimentally reducing partisan incivility on Twitter. *Journal of Experimental Political Science* 8(2), 102–116.
- Palin, K., A. M. Feit, S. Kim, P. O. Kristensson, and A. Oulasvirta (2019). How do people type on mobile devices? observations from a study with 37,000 volunteers. In *Proceedings of the 21st International Conference on Human-Computer Interaction with Mobile Devices and Services*, pp. 1–12.
- Pew Research Center (2021). Social media use in 2021. <https://www.pewresearch.org/internet/2021/04/07/social-media-use-in-2021/>. Accessed: 2021-10-24.
- Rauchfleisch, A. and J. Kaiser (2021). Deplatforming the far-right: An analysis of YouTube and BitChute. *Available at SSRN*.
- Relia, K., Z. Li, S. H. Cook, and R. Chunara (2019). Race, ethnicity and national origin-based discrimination in social media and hate crimes across 100 us cities. In *Proceedings of the International AAAI Conference on Web and Social Media*, Volume 13, pp. 417–427.

- Reynders, D. (2020). Countering illegal hate speech online: 5th evaluation of the code of conduct. Technical report, European Commission, Directorate-General for Justice and Consumers.
- Ribeiro, M. H., J. Cheng, and R. West (2022a). Automated content moderation increases adherence to community guidelines. *arXiv preprint arXiv:2210.10454*.
- Ribeiro, M. H., J. Cheng, and R. West (2022b). Post approvals in online communities. In *Proceedings of the International AAAI Conference on Web and Social Media*, Volume 16, pp. 335–346.
- Rösner, L., S. Winter, and N. C. Krämer (2016). Dangerous minds? effects of uncivil online comments on aggressive cognitions, emotions, and behavior. *Computers in Human Behavior* 58, 461–470.
- Roth, J. and P. H. Sant’Anna (2021). Efficient estimation for staggered rollout designs. *arXiv preprint arXiv:2102.01291*.
- Seering, J., R. Kraut, and L. 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*, pp. 111–125.
- Seyler, D., S. Tan, D. Li, J. Zhang, and P. Li (2021). Textual analysis and timely detection of suspended social media accounts.
- Sherry, M. (2019). *Disability Hate Speech: Social, Cultural and Political Contexts*, Chapter Disablist hate speech online. Routledge.
- Siegel, A. A. and V. Badaan (2020). #no2sectarianism: Experimental approaches to reducing sectarian hate speech online. *American Political Science Review* 114(3), 837–855.
- Spence, A. M. (1975). Monopoly, quality, and regulation. *The Bell Journal of Economics*, 417–429.
- Srinivasan, K. B., C. Danescu-Niculescu-Mizil, L. Lee, and C. Tan (2019). Content removal as a moderation strategy: Compliance and other outcomes in the changemyview community. *Proceedings of the ACM on Human-Computer Interaction* 3(CSCW), 1–21.
- Strossen, N. (2018). *Hate: Why we should resist it with free speech, not censorship*. Oxford University Press.
- Tan, G. and J. Zhou (2021). The effects of competition and entry in multi-sided markets. *The Review of Economic Studies* 88(2), 1002–1030.
- Tibshirani, J., S. Athey, S. Wager, R. Friedberg, L. Miner, and M. Wright (2020). grf: Generalized random forests. *R package version 1*(0), 7–3.
- Twitter (2018). Serving healthy conversation. [https://blog.twitter.com/official/en\\_us/topics/product/2018/Serving\\_Healthy\\_Conversation.html](https://blog.twitter.com/official/en_us/topics/product/2018/Serving_Healthy_Conversation.html). Accessed: 2021-10-19.
- Twitter (2020a). Q4 2020 letter to shareholders. <https://s22.q4cdn.com/826641620/>



- files/doc\_financials/2020/q4/FINAL-Q4'20-TWTR-Shareholder-Letter.pdf. Accessed: 2021-10-13.
- Twitter (2020b). Rules enforcement report. <https://transparency.twitter.com/en/reports/rules-enforcement.html#2020-jul-dec>. Accessed: 2021-10-13.
- Twitter (2021a). Covid-19 misleading information policy. <https://help.twitter.com/en/rules-and-policies/medical-misinformation-policy>. Accessed: 2022-11-28.
- Twitter (2021b). Debunking twitter myths. <https://help.twitter.com/en/using-twitter/debunking-twitter-myths>. Accessed: 2021-10-13.
- Twitter (2021c). Hateful conduct policy. <https://help.twitter.com/en/rules-and-policies/hateful-conduct-policy>. Accessed: 2021-06-20.
- Twitter (2021d). Notices on Twitter and what they mean. <https://help.twitter.com/en/rules-and-policies/notices-on-twitter>. Accessed: 2021-10-13.
- Twitter (2021e). Our range of enforcement options. <https://help.twitter.com/en/rules-and-policies/enforcement-options>. Accessed: 2021-06-21.
- Twitter (2021f). Report abusive behavior. <https://help.twitter.com/en/safety-and-security/report-abusive-behavior>. Accessed: 2021-06-21.
- Twitter (2021g). Response codes. <https://developer.twitter.com/ja/docs/basics/response-codes>. Accessed: 2021-06-21.
- Urminsky, O., C. Hansen, and V. Chernozhukov (2016). Using double-lasso regression for principled variable selection. *Available at SSRN 2733374*.
- Vidgen, B., S. Hale, S. Staton, T. Melham, H. Margetts, O. Kammar, and M. Szymczak (2020). Recalibrating classifiers for interpretable abusive content detection. In *Proceedings of the Fourth Workshop on Natural Language Processing and Computational Social Science*, pp. 132–138.
- Waldron, J. (2012). *The harm in hate speech*. Harvard University Press.
- Weyl, E. G. (2010). A price theory of multi-sided platforms. *American Economic Review* 100(4), 1642–72.
- Whalen, A. (2020). What Did Twitter Do to James Woods? The Story Behind the Trend. Accessed: 2021-10-13.
- White, A. and E. G. Weyl (2016). Insulated platform competition. *Available at SSRN 1694317*.
- Wojcik, S., S. Hilgard, N. Judd, D. Mocanu, S. Ragain, M. Hunzaker, K. Coleman, and J. Baxter (2022). Birdwatch: Crowd wisdom and bridging algorithms can inform understanding and reduce the spread of misinformation. *arXiv preprint arXiv:2210.15723*.
- Wojcik, S. and A. Hughes (2019). Sizing up Twitter users. *Pew Research Center* 24.
- Wulczyn, E., N. Thain, and L. Dixon (2017). Ex machina: Personal attacks seen at scale. In *Proceedings of the 26th international conference on world wide web*, pp. 1391–1399.

- Yang, K.-C., O. Varol, P.-M. Hui, and F. Menczer (2020). Scalable and generalizable social bot detection through data selection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Volume 34, pp. 1096–1103.
- Yildirim, M. M., J. Nagler, R. Bonneau, and J. A. Tucker (2021). Short of suspension: How suspension warnings can reduce hate speech on Twitter. *Perspectives on Politics*, 1–13.
- Zannettou, S. (2021). “I won the election!”: An empirical analysis of soft moderation interventions on Twitter. *arXiv preprint arXiv:2101.07183*.
- Zhuravskaya, E., M. Petrova, and R. Enikolopov (2020). Political effects of the internet and social media. *Annual Review of Economics* 12, 415–438.

## A Formal Propositions and Model Extensions

### A.1 Overprovision or underprovision of moderation

The following proposition reproduces Spence’s result for the model of Section 2.<sup>54</sup>

**Proposition 1.** *For fixed quantities  $T^\theta$  and assuming that second-order conditions hold, the platform can overprovide or underprovide moderation relative to a surplus-maximizing planner. A sufficient condition for underprovision is that  $P_{T^\theta}^\theta < 0$ , for overprovision is that  $P_{T^\theta}^\theta > 0$ , and for efficient provision is that  $P_{T^\theta}^\theta = 0$ .*

*Proof.* A social planner chooses  $\mathbf{T}$  and  $c$  to maximize total surplus  $W$ , which equals:

$$\begin{aligned}
 W(\mathbf{T}, c) &= w \underbrace{\left( \int_0^{T^A} p^A(t, T^H, c) dt + \int_0^{T^H} p^H(T^A, t, c) dt - p^A T^A - p^H T^H \right)}_{\text{Consumer surplus}} \\
 &\quad + \underbrace{a(p^A(\mathbf{T}, c) T^A + p^H(\mathbf{T}, c) T^H) - \phi(\mathbf{T}, c)}_{\text{Producer surplus}} \\
 &= w \left( \int_0^{T^A} p^A(t, T^H, c) dt + \int_0^{T^H} p^H(T^A, t, c) dt \right) \\
 &\quad + (a - w)(p^A(\mathbf{T}, c) T^A + p^H(\mathbf{T}, c) T^H) - \phi(\mathbf{T}, c).
 \end{aligned}$$

The first two terms in the second equality are the areas below the inverse demand curves and the last term is the cost function, but the third term is new. This new term appears because the platform collects time with an opportunity cost  $w$  and sells it to advertisers for a price  $a$ . To the best of my knowledge there are no analyses that compare the price of advertisements of social media to the opportunity cost of time, so the magnitude of  $a - w$  is unknown.<sup>55</sup>

The first-order condition with respect to  $c$  from this problem is:

$$w \left( \int_0^{T^A} \frac{\partial p^A}{\partial c} dt + \int_0^{T^H} \frac{\partial p^H}{\partial c} dt \right) + (a - w) \left( \frac{\partial p^A}{\partial c} T^A + \frac{\partial p^H}{\partial c} T^H \right) = \frac{\partial \phi}{\partial c} \quad (7)$$

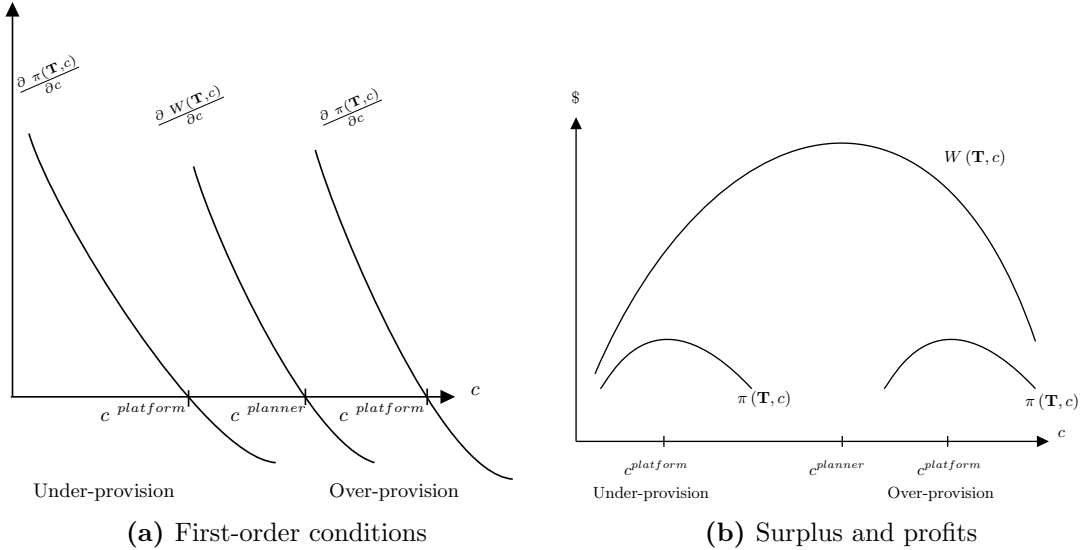
<sup>54</sup>The model differs from Spence’s framework in two ways. First, the monopolist sells two “products” instead of one. Second, there is a gap between the opportunity cost of time ( $w$ ) and the value of time spent watching ads ( $a$ ).

<sup>55</sup>A no-arbitrage argument suggests that  $a \approx w$ . Suppose that ad prices were higher than the opportunity cost of time. This creates incentives for companies to pay users to watch advertisements. While Becker and Murphy (1993) argue that this might not be profitable, since consumers would “buy” a large number of ads and ignore as many as possible, current technology might facilitate this. Indeed, websites like [adwallet.com](http://adwallet.com) reward consumers for watching ads. On the other hand, if  $w > a$ , platforms would find it more profitable to have consumers complete tasks rather than show them ads; e.g., “Fill out this survey in order to proceed to your feed”.

Suppose that  $p_{t^\theta}^\theta < 0$ .<sup>56</sup> Then  $\partial p^A(t, T^H, c)/\partial c > \partial p^A(\mathbf{T}, c)/\partial c$  for all  $t < T^A$  and likewise for  $H$ . Then, the left-hand side of equation (7) satisfies:

$$\begin{aligned}
& w \left( \int_0^{T^A} \frac{\partial p^A(t, T^H, c)}{\partial c} dt + \int_0^{T^H} \frac{\partial p^H(T^A, t, c)}{\partial c} dt \right) \\
& + (a - w) \left( \frac{\partial p^A(\mathbf{T}, c)}{\partial c} T^A + \frac{\partial p^H(\mathbf{T}, c)}{\partial c} T^H \right) \\
& > w \left( \frac{\partial p^A(\mathbf{T}, c)}{\partial c} T^A + \frac{\partial p^H(\mathbf{T}, c)}{\partial c} T^H \right) \\
& + (a - w) \left( \frac{\partial p^A(\mathbf{T}, c)}{\partial c} T^A + \frac{\partial p^H(\mathbf{T}, c)}{\partial c} T^H \right) \\
& = a \left( \frac{\partial p^A(\mathbf{T}, c)}{\partial c} T^A + \frac{\partial p^H(\mathbf{T}, c)}{\partial c} T^H \right),
\end{aligned}$$

which is identical to the left-hand side of equation (4). Since equations (4) and (7) both have  $\partial \phi / \partial c$  on the right-hand side, this means that the planner's first-order condition is above the monopolist's one for fixed  $t^\theta$  and all  $c$ :  $\partial W(\mathbf{T}, c) / \partial c < \partial \pi(\mathbf{T}, c) / \partial c$ . Assuming that second-order conditions hold, this means that the root of the planner's first-order condition,  $c^{planner}$ , is higher than the root of the monopolist's condition,  $c^{platform}$ , so there is under-provision of moderation. Figure A.1 illustrates the proof.



**Figure A.1:** Illustration of moderation overprovision and underprovision

□

<sup>56</sup>The proof is analogous for the opposite case.

## A.2 Generalization to Multiple Platforms

Assume without loss of generality that there are two platforms  $j \in \{1, 2\}$ . The solution concept of the model is a Cournot equilibrium as in Correia-da Silva et al. (2019).<sup>57</sup> First, platforms simultaneously set the amount of content  $T_j^\theta$  on each side of the market and the moderation rates  $c_j$ . Then, given the quantities and moderation rates, prices adjust to equate demand and supply on each platform.

A fraction  $\mu^\theta$  of users are of type  $\theta \in \{A, H\}$ . Consumers are now characterized by the parameter vectors  $\gamma^\theta = (\gamma_1^\theta, \gamma_2^\theta, \delta_1^\theta, \delta_2^\theta)$ . The  $\gamma$ 's govern how utility responds to spillovers and the  $\delta$ 's govern membership benefits. As in (Weyl, 2010), the conditional density of membership benefits has full support. Users decide whether to join one of the platforms or neither. Below I discuss an extension to a multi-homing case. Once consumers join a platform, they decide how much time to spend on it. If they join platform  $j$ , they obtain membership benefits  $\delta_j^\theta$  and indirect utility:

$$v_j^\theta(\mathbf{T}_j, c_j, p_j^\theta, \gamma_j^\theta) = \max_{t \in [0, T]} u^\theta(t, \mathbf{T}_j, c_j; \gamma_j^\theta) - t \times w(1 + p_j^\theta),$$

where  $\mathbf{T}_j = (T_j^A, T_j^H)$

Define the vectors  $\mathbf{T} = (\mathbf{T}_1, \mathbf{T}_2)$ ,  $\mathbf{T}^\theta = (T_1^\theta, T_2^\theta)$ ,  $\mathbf{p}^\theta = (p_1^\theta, p_2^\theta)$ , and  $\mathbf{c} = (c_1, c_2)$ , and the set of types that decide to use platform  $j$  as:

$$\bar{\gamma}_j^\theta(\mathbf{T}, \mathbf{c}, \mathbf{p}^\theta) \equiv \left\{ \gamma^\theta : v_j^\theta(\mathbf{T}_j, c_j, p_j^\theta, \gamma_j^\theta) + \delta_j^\theta \geq \max\{v_{-j}^\theta(\mathbf{T}_{-j}, c_{-j}, p_{-j}^\theta, \gamma_{-j}^\theta) + \delta_{-j}^\theta, 0\} \right\},$$

where  $-j$  denotes the other platform. Let  $t_j(\mathbf{T}_j, c_j, p_j^\theta, \gamma_j^\theta)$  be the optimal time spent on platform  $j$ . Aggregate demands are:

$$T_j^\theta(\mathbf{T}, \mathbf{c}, \mathbf{p}^\theta) = \mu^\theta \int_{\bar{\gamma}_j^\theta(\mathbf{T}, \mathbf{c}, \mathbf{p}^\theta)} t_j(\mathbf{T}_j, c_j, p_j^\theta, \gamma_j^\theta) f^\theta(\gamma^\theta) d\gamma^\theta$$

The consumer equilibrium constraints are, for all  $j$  and  $\theta$ :

$$t_j^\theta = T_j^\theta(\mathbf{T}, \mathbf{c}, \mathbf{p}^\theta)$$

Inverting the demand curves is not as straightforward as in Weyl (2010), since demands now depend on the other platform's prices. We can, however, use the global inverse function theorem from Berry et al. (2013) to obtain the twice-continuously differentiable inverse demands  $P_j^\theta(\mathbf{T}, \mathbf{c})$ .<sup>58</sup>

<sup>57</sup> Alternative solution concepts are flat pricing (Tan and Zhou, 2021) and insulating equilibrium (White and Weyl, 2016). See Correia-da Silva et al. (2019) for more discussion.

<sup>58</sup> Note that demands  $T_j^\theta$  are twice-continuously differentiable and strictly decreasing in prices  $p_{-j}^\theta$  and weakly decreasing in prices  $p_{-j}^{-\mathcal{I}}$  and  $p_j^{-\mathcal{I}}$ , where  $-\mathcal{I}$  denotes the other side. Moreover, the demand of the

The problem of platform  $j$  is now:

$$\max_{T_j^A, T_j^H, c_j} \pi_j(\mathbf{T}, \mathbf{c}) \equiv a(P_j^A(\mathbf{T}, \mathbf{c}) T_j^A + P_j^H(\mathbf{T}, \mathbf{c}) T_j^H) - \phi_j(\mathbf{T}_j, c_j).$$

The first-order condition with respect to the moderation rate is identical to equation (4), but using residual inverse demands instead of the market inverse demand curve:

$$a \left( \frac{\partial P_j^A}{\partial c_j} T_j^A + \frac{\partial P_j^H}{\partial c_j} T_j^H \right) = \frac{\partial C_j}{\partial c_j}$$

Hence, the same intuition of the platform's moderation decision holds in a model with two platforms. Moderation is a quality decision that allows platforms to increase their advertising revenue. The increase in ad revenue is the weighted change in willingness to pay of both types of users.

The following proposition shows that it is sufficient to measure the change in surplus on a sample of existing consumers; one can ignore the change in marginal users since they get zero surplus by definition.

**Proposition 2.** *The derivative of consumer surplus with respect to the moderation rate of platform  $j$  equals the average derivative of consumer surplus among users of that platform:*

$$\frac{\partial CS(\mathbf{T}, \mathbf{c}, \mathbf{p})}{\partial c_j} = \sum_{\theta} \mu^{\theta} \int_{\tilde{\gamma}_j^{\theta}(\mathbf{T}, \mathbf{c}, \mathbf{p}^{\theta})} \frac{\partial v_j^{\theta}(\mathbf{T}_j, c_j, p_j^{\theta}, \gamma_j^{\theta})}{\partial c_j} f^{\theta}(\gamma^{\theta}) d\gamma^{\theta}.$$

*Proof.* Define the membership benefit from joining platform  $-j$  relative to the membership benefit from  $j$  as  $\tilde{\delta}_{-j}^{\theta} \equiv \delta_{-j}^{\theta} - \delta_j^{\theta}$ . Define also the vector of network parameters of both platforms  $\gamma^{\theta} \equiv (\gamma_1^{\theta}, \gamma_2^{\theta})$ , the vector of parameters  $\tilde{\gamma}^{\theta} \equiv (\gamma^{\theta}, \delta_j^{\theta}, \tilde{\delta}_{-j}^{\theta})$  and the distribution of types  $\tilde{f}^{\theta}(\tilde{\gamma}^{\theta}) \equiv f^{\theta}(\gamma^{\theta}, \delta_j^{\theta}, \delta_{-j}^{\theta} + \delta_j^{\theta})$ . The membership benefits of those users who join platform  $j$  are bounded as follows:

$$\begin{aligned} \delta_j^{\theta} &\geq -v_j^{\theta}(\mathbf{T}_j, c_j, p_j^{\theta}, \gamma_j^{\theta}), \\ \tilde{\delta}_{-j}^{\theta} &\leq v_j^{\theta}(\mathbf{T}_j, c_j, p_j^{\theta}, \gamma_j^{\theta}) - v_{-j}^{\theta}(\mathbf{T}_{-j}, c_{-j}, p_{-j}^{\theta}, \gamma_j^{\theta}) \end{aligned}$$

Likewise, the bounds of the membership benefits of those users who join platform  $-j$  are:

$$\begin{aligned} \delta_{-j}^{\theta} &\geq -v_{-j}^{\theta}(\mathbf{T}_{-j}, c_{-j}, p_{-j}^{\theta}, \gamma_j^{\theta}) - \tilde{\delta}_{-j}^{\theta}, \\ \tilde{\delta}_{-j}^{\theta} &\geq v_j^{\theta}(\mathbf{T}_j, c_j, p_j^{\theta}, \gamma_j^{\theta}) - v_{-j}^{\theta}(\mathbf{T}_{-j}, c_{-j}, p_{-j}^{\theta}, \gamma_j^{\theta}) \end{aligned}$$

---

outside option is strictly increasing in all prices. Hence, this model satisfies all the conditions of Corollary 2 from Berry et al. (2013).

Omitting the arguments of  $v_j^\theta$  and  $v_{-j}^\theta$  for brevity, the consumer surplus is:

$$\begin{aligned}
CS(\mathbf{T}, \mathbf{c}, \mathbf{p}) = \sum_{\theta} \mu^{\theta} & \left( \underbrace{\int \int_{-\infty}^{v_j^\theta - v_{-j}^\theta} \int_{-v_j^\theta}^{\infty} (v_j^\theta + \delta_j^\theta) \tilde{f}^\theta(\tilde{\gamma}^\theta) d\delta_j^\theta d\tilde{\delta}_{-j}^\theta d\gamma^\theta}_{\text{Surplus of } j\text{'s consumers}} \right. \\
& \left. + \underbrace{\int \int_{v_j^\theta - v_{-j}^\theta}^{\infty} \int_{-v_{-j}^\theta - \tilde{\delta}_{-j}^\theta}^{\infty} (v_{-j}^\theta + \tilde{\delta}_{-j}^\theta + \delta_j^\theta) \tilde{f}^\theta(\tilde{\gamma}^\theta) d\delta_j^\theta d\tilde{\delta}_{-j}^\theta d\gamma^\theta}_{\text{Surplus of } -j\text{'s consumers}} \right) \quad (8)
\end{aligned}$$

Use the Leibniz integral rule to differentiate the first row after the equality sign from the previous expression with respect to  $c_j$ :

$$\begin{aligned}
& \int \frac{\partial v_j^\theta}{\partial c_j} \int_{-v_j^\theta}^{\infty} (v_j^\theta + \delta_j^\theta) \tilde{f}^\theta(\gamma^\theta, \delta_j^\theta, v_j^\theta - v_{-j}^\theta) d\delta_j^\theta d\gamma^\theta \\
& + \int \int_{-\infty}^{v_j^\theta - v_{-j}^\theta} \frac{\partial v_j^\theta}{\partial c_j} \left( \underbrace{v_j^\theta - v_{-j}^\theta}_{=0} \right) \tilde{f}^\theta(\gamma^\theta, -v_j^\theta, \tilde{\delta}_{-j}^\theta) d\tilde{\delta}_{-j}^\theta d\gamma^\theta \\
& + \int \int_{-\infty}^{v_j^\theta - v_{-j}^\theta} \int_{-v_j^\theta}^{\infty} \frac{\partial v_j^\theta}{\partial c_j} \tilde{f}^\theta(\tilde{\gamma}^\theta) d\delta_j^\theta d\tilde{\delta}_{-j}^\theta d\gamma^\theta \quad (9)
\end{aligned}$$

Likewise, differentiating the second row of equation (8):

$$\int -\frac{\partial v_j^\theta}{\partial c_j} \int_{-v_j^\theta}^{\infty} (v_j^\theta + \delta_j^\theta) \tilde{f}^\theta(\gamma^\theta, \delta_j^\theta, v_j^\theta - v_{-j}^\theta) d\delta_j^\theta d\gamma^\theta \quad (10)$$

The first row from equation (9) cancels with equation (10), so adding these two expressions together gives

$$\int \int_{-\infty}^{v_j^\theta - v_{-j}^\theta} \int_{-v_j^\theta}^{\infty} \frac{\partial v_j^\theta}{\partial c_j} \tilde{f}^\theta(\tilde{\gamma}^\theta) d\delta_j^\theta d\tilde{\delta}_{-j}^\theta d\gamma^\theta,$$

which equals the average derivative of consumer surplus among users of  $j$ ,  $\gamma^\theta \in \bar{\gamma}_j^\theta(\mathbf{T}, \mathbf{c}, \mathbf{p}^\theta)$ .  $\square$

The previous single-homing model can be extended to allow multi-homing and flexible substitution or complementarity patterns by considering bundles of platforms as different choices, by using the approach of Gentzkow (2007) as Berry et al. (2013) suggest. Proposition 2 still holds in such a model, but the indirect utility  $v_j^\theta$  would represent the utility from using any bundle of platforms that includes  $j$ , and  $v_{-j}^\theta$  would correspond to bundles without  $j$ .

Hence, in that setting,  $\partial v_j^\theta / \partial c_j$  would be the derivative of the surplus of using social media, not just platform  $j$ , for users who join any bundle of platforms that includes  $j$ .

In a multi-platform model, the steps to test whether a platform under-provides or over-provides moderation, for fixed quantities, are as follows. First, introduce a small change in the moderation rate of the desired platform. Second, compute the average change in consumer surplus from using social media—not just the surplus from using the platform—among representative consumers of the platform.

### A.3 Identifying unobservable sanctions

Let  $U$ ,  $S$ ,  $N$ , and  $Z$  be indicators of receiving unobservable sanctions, observable sanctions, Twitter notifications (updates), and reports, respectively.

**Assumption 1.** *Twitter does not send updates if it does not impose a sanction (observable or not):  $\Pr(N = 1, S = 0, U = 0 | Z = 1) = 0$ .*

**Assumption 2.** *Twitter sends updates for observable and unobservable sanctions at the same rate:  $\Pr(N = 1 | S = 1, Z = 1) = \Pr(N = 1 | U = 1, Z = 1)$ .*

**Proposition 3.** *Under Assumption 1, the fraction of reports for which a notification is not accompanied by an observable sanction gives a lower bound on the probability of receiving unobservable sanctions in the reporting arm:*

$$\underbrace{\Pr(U = 1 | Z = 1)}_{\text{Not observable in the data}} \geq \underbrace{\Pr(N = 1, S = 0 | Z = 1)}_{\text{Observable in the data}}.$$

*Proof.* We have  $\Pr(U = 1 | Z = 1) \geq \Pr(U = 1, N = 1, S = 0 | Z = 1)$ . By the law of total probability (LTP),  $\Pr(N = 1, S = 0 | Z = 1) = \Pr(U = 1, N = 1, S = 0 | Z = 1) + \Pr(U = 0, N = 1, S = 0 | Z = 1)$ . Assumption 1 says that this last term is zero, so  $\Pr(U = 1 | Z = 1) \geq \Pr(N = 1, S = 0 | Z = 1)$ .  $\square$

**Proposition 4.** *Under Assumptions 1 and 2, the probability of receiving unobservable sanctions in the reporting arm is equal to the fraction of reports for which a notification is not accompanied by an observable sanction divided by the probability of receiving an update among sanctioned, reported accounts:*

$$\underbrace{\Pr(U = 1 | Z = 1)}_{\text{Not observable in the data}} = \frac{\Pr(N = 1, S = 0 | Z = 1)}{\underbrace{\Pr(N = 1 | S = 1, Z = 1)}_{\text{Observable in the data}}}.$$

*Proof.* By the LTP and Assumption 1 we have:  $\Pr(N = 1, S = 0 | Z = 1) = \Pr(N = 1, S = 0, U = 1 | Z = 1)$ . By definition of unobservable sanctions,  $\Pr(N = 1, S = 1, U = 1 | Z = 1) = 0$ ,



so by the LTP we have  $\Pr(N = 1, S = 0, U = 1|Z = 1) = \Pr(N = 1, U = 1|Z = 1)$ . Hence:

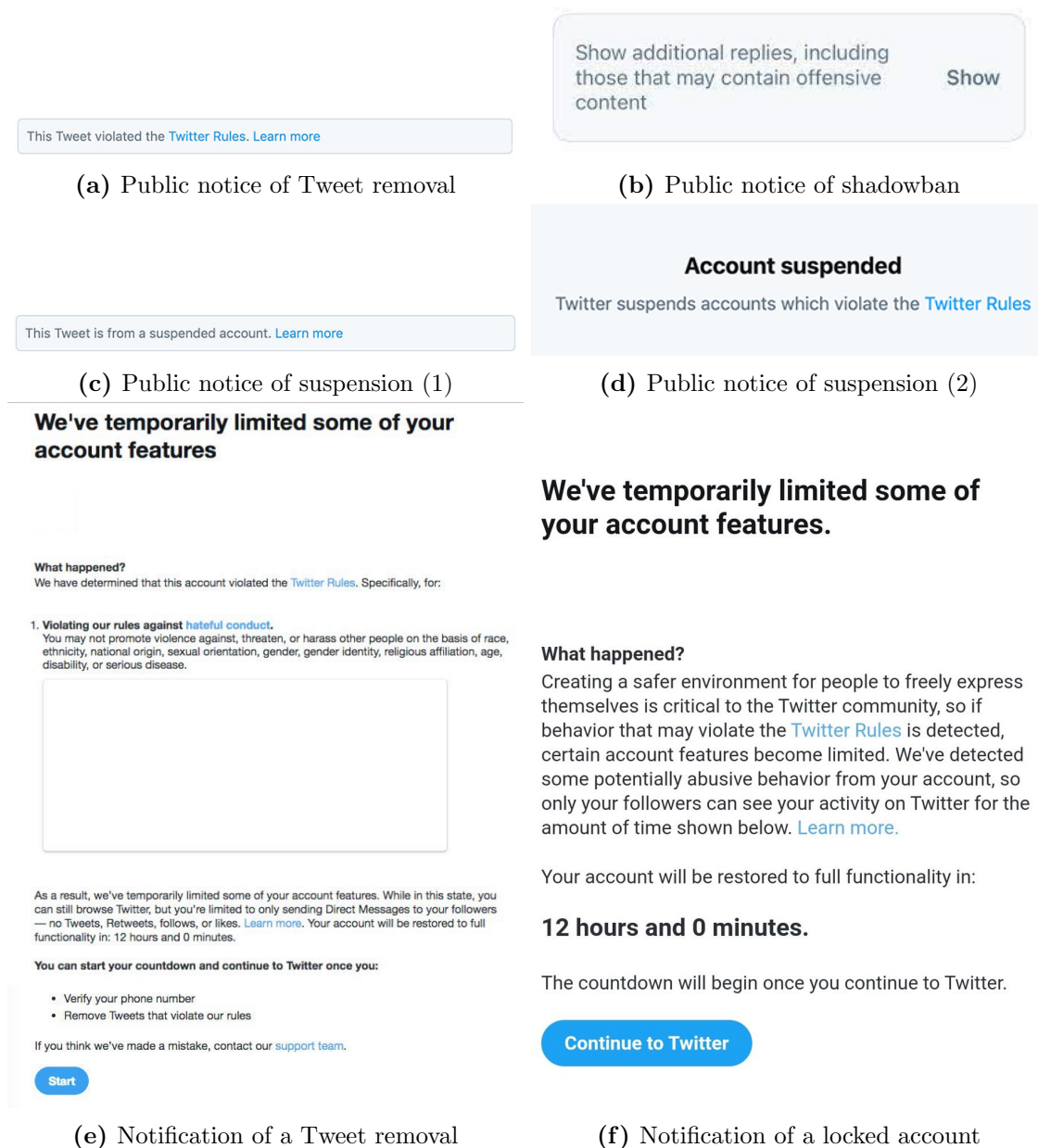
$$\frac{\Pr(N = 1, S = 0|Z = 1)}{\Pr(N = 1|S = 1, Z = 1)} = \frac{\Pr(N = 1, U = 1|Z = 1)}{\Pr(N = 1|S = 1, Z = 1)}.$$

We plug in Assumption 2 in the denominator of the right-hand side of the previous equation and we use the definition of conditional probability:

$$\frac{\Pr(N = 1, U = 1|Z = 1)}{\Pr(N = 1|S = 1, Z = 1)} = \frac{\Pr(N = 1, U = 1|Z = 1)}{\Pr(N = 1|U = 1, Z = 1)} = \Pr(U = 1|Z = 1).$$

□

## B Experimental design



**Figure B.1:** Public notices and notifications

Notes: This figure includes images of public notices of different sanctions and an example of a notification that a users receive when their account is locked.

**Table B.1:** Query list

Holocaust	holohoax, h0lohoax, hol0hoax, holoh0ax, holoho4x, holoho@x, holoho*x, h0l0h0ax, h0l0h04x, holo hoax, h0lo hoax, hol0 hoax, holo h0ax, holo ho4x, holo ho@x, holo ho*x, h0l0 h0ax, h0l0 h04x, holo-hoax, h0lo-hoax, hol0-hoax, holo-h0ax, holo-ho4x, holo-ho@x, holo-ho*x, h0l0-h0ax, h0l0-h04x
Disability	<p>youre a retard, youre a r3tard, youre a r*tard, youre a ret*rd, youre a returd, youre a ret@rd, youre a ret4rd, youre a r3t4rd, youre a r3t@rd, your a retard, your a r3tard, your a r*tard, your a ret*rd, your a returd, your a ret@rd, your a ret4rd, your a r3t4rd, your a r3t@rd, you're a retard, you're a r3tard, you're a r*tard, you're a ret*rd, you're a returd, you're a ret@rd, you're a ret4rd, you're a r3t4rd, you're a r3t@rd, ure a retard, ure a r3tard, ure a r*tard, ure a ret*rd, ure a returd, ure a ret@rd, ure a ret4rd, ure a r3t4rd, ure a r3t@rd, ur a retard, ur a r3tard, ur a r*tard, ur a ret*rd, ur a returd, ur a ret@rd, ur a ret4rd, ur a r3t4rd, ur a r3t@rd, u're a retard, u're a r3tard, u're a r*tard, u're a ret*rd, u're a returd, u're a ret@rd, u're a ret4rd, u're a r3t4rd, u're a r3t@rd</p> <p>youre retarded, youre r3tarded, youre r*tarded, youre ret*rded, youre returded, youre ret@rded, youre ret4rded, youre r3t4rded, youre r3t@rded, you're retarded, you're r3tarded, you're r*tarded, you're ret*rded, you're returded, you're ret@rded, you're ret4rded, you're r3t4rded, you're r3t@rded, ure retarded, ure r3tarded, ure r*tarded, ure ret*rded, ure returded, ure ret@rded, ure ret4rded, ure r3t4rded, ure r3t@rded, u're retarded, u're r3tarded, u're r*tarded, u're ret*rded, u're returded, u're ret@rded, u're ret4rded, u're r3t4rded, u're r3t@rded</p> <p>youre a retarded, youre a r3tarded, youre a r*tarded, youre a ret*rded, youre a returded, youre a ret@rded, youre a ret4rded, youre a r3t4rded, youre a r3t@rded, your a retarded, your a r3tarded, your a r*tarded, your a ret*rded, your a returded, your a ret@rded, your a ret4rded, your a r3t4rded, your a r3t@rded, you're a retarded, you're a r3tarded, you're a r*tarded, you're a ret*rded, you're a returded, you're a ret@rded, you're a ret4rded, you're a r3t4rded, you're a r3t@rded, ure a retarded, ure a r3tarded, ure a r*tarded, ure a ret*rded, ure a returded, ure a ret@rded, ure a ret4rded, ure a r3t4rded, ure a r3t@rded, ur a retarded, ur a r3tarded, ur a r*tarded, ur a ret*rded, ur a returded, ur a ret@rded, ur a ret4rded, ur a r3t4rded, ur a r3t@rded, u're a retarded, u're a r3tarded, u're a r*tarded, u're a ret*rded, u're a returded, u're a ret@rded, u're a ret4rded, u're a r3t4rded, u're a r3t@rded</p>

Hello,

Twitter is required by German law to provide notice to users who are reported by people from Germany via the Network Enforcement Act reporting flow.

We have received a complaint regarding your account, @handle , for the following content:

Tweet ID:

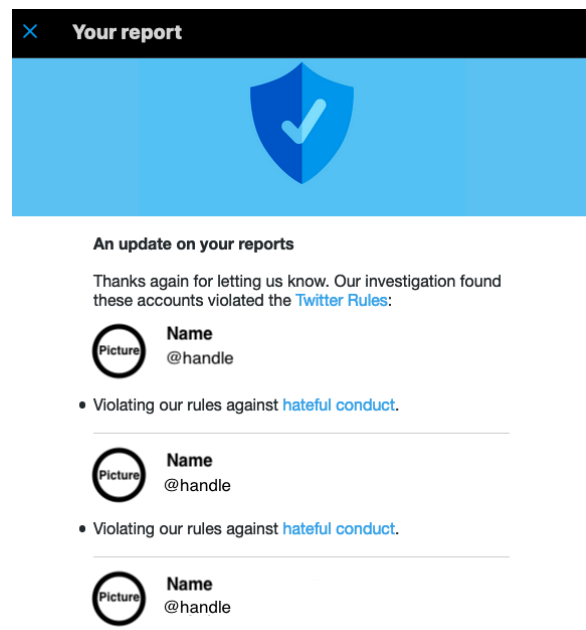
Tweet Text:

We have investigated the reported content and have found that it is not subject to removal under the Twitter Rules (<https://support.twitter.com/articles/18311>) or German law. Accordingly, we have not taken any action as a result of this specific report.

Sincerely,

Twitter

**Figure B.2:** Screenshot of a notification of a user report



**Figure B.3:** Screenshot of an update on reports

**Table B.2:** Variable definition

Variable	Definition
Account years	Years from account creation date until measurement date
Tweets per day	Statuses count divided by days since account creation
Likes per day	Likes count divided by days since account creation
Followers	Number of accounts that follow a user
Followed	Number of accounts that the user follows
Bot score	Probability of being a bot, from Botometer API
Is bot	Indicates whether bot score $\geq 0.5$
Initial shadow ban	Whether an account is shadow banned at the time of sampling
Word count	Number of words in a tweet
Is toxic	Indicates whether toxicity $\geq 0.8$
Is hate (MTurk)	Indicates whether a majority of MTurkers label the post as hate
Is reply	Indicates whether the tweet is a reply to another user
Is attack (MTurk)	Indicates whether the majority of MTurkers consider the post to be an attack on another user
Is quote	Indicates whether the tweet is a quote to another user
Is mention	Indicates whether the tweet mentions another user
Has media	Indicates whether the tweet contains a video or picture
Disability key word	Indicates whether the tweet contains the expression “r*t*rd”
Holocaust key word	Indicates whether the tweet contains the expressions “h*l*h**x”, “h*l*c**st”, “jew”
Tweet from phone	Indicates whether the source of the tweet is Twitter for iPhone or Twitter for Android
Has description	Indicates whether a profile has a description
Has location	Indicates whether a profile has a location
Default picture	Indicates whether a profile has a default profile picture
Is verified	Indicates whether an account is verified
Has Instagram	Indicates whether a profile description, location or URL contains an Instagram handle
Has backup	Indicates whether a profile description, location or URL contains an alternative or backup Twitter handle
Has pronouns	Indicates whether a profile description or location contains pronouns or a <b>carrd.co</b> link
Under 18	Indicates whether a profile description or location contains numbers 13 to 17 (in number or word), years 2003 to 2008 or words like “minor” or “teen”
Previous toxicity	Indicates whether any of a user’s most recent 50 tweets has toxicity $\geq 0.8$
Previous disability	Indicates whether any of a user’s most recent 50 tweets has the expression “r*t*rd”
Previous Holocaust	Indicates whether any of a user’s most recent 50 tweets has the expressions “h*l*h**x”, “h*l*c**st”, “jew”

**Table B.3:** Balance in the reporting experiment

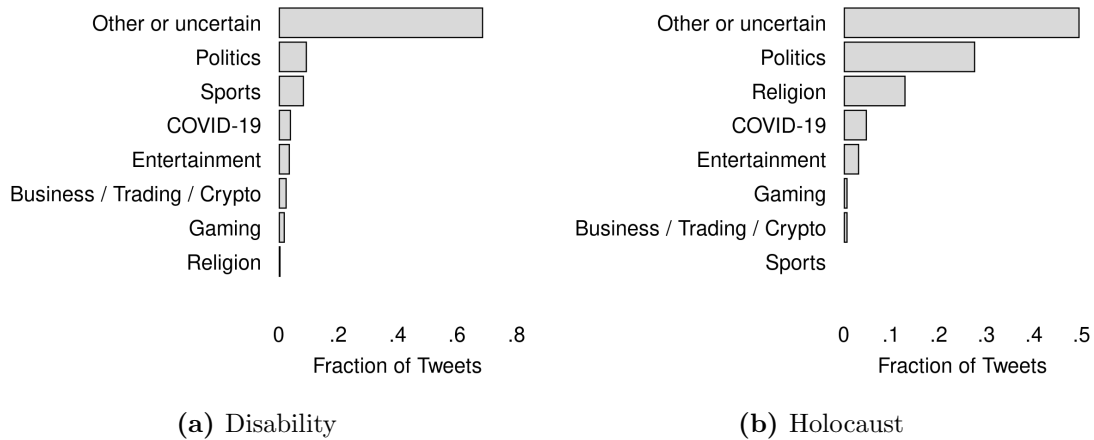
Characteristic	Control		Treatment		Difference	
	Mean	SD	Mean	SD	Normalized	<i>p</i> -value
<i>Observations</i>	3,074		3,074			
<i>Accounts</i>						
Account years	3.21	3.4	3.23	3.5	-0.01	0.77
Tweets per day	11.19	24.2	12.05	25.2	-0.03	0.20
Likes per day	23.39	50.2	24.95	51.6	-0.03	0.22
Followers	517.07	3,439.5	752.64	6,476.9	-0.05	0.08
Followed	426.84	751.4	440.65	946.0	-0.02	0.55
Bot score	0.24	0.1	0.24	0.1	0.01	0.54
Initial shadow ban	0.71	0.5	0.71	0.5	0.01	0.78
<i>Tweets</i>						
Word count	16.02	13.2	15.93	13.4	0.01	0.80
Is toxic	0.81	0.4	0.80	0.4	0.04	0.19
Is hate (MTurk)	0.31	0.5	0.30	0.5	0.00	0.81
Is reply	0.84	0.4	0.84	0.4	0.00	0.95
Is attack (MTurk)	0.78	0.4	0.78	0.4	-0.02	0.42
Is quote	0.07	0.3	0.07	0.3	0.02	0.53
Is mention	0.85	0.4	0.85	0.4	0.01	0.76
Has media	0.04	0.2	0.04	0.2	-0.04	0.13
Tweet from phone	0.80	0.4	0.78	0.4	0.03	0.25
<i>Profiles</i>						
Has description	0.82	0.4	0.82	0.4	-0.01	0.70
Has location	0.51	0.5	0.51	0.5	-0.01	0.53
Default picture	0.04	0.2	0.03	0.2	0.02	0.42
Is verified	0.00	0.0	0.00	0.0	-0.02	0.51
Has Instagram	0.01	0.1	0.02	0.1	-0.02	0.46
Has backup	0.01	0.1	0.01	0.1	0.04	0.15
<i>Timelines</i>						
Previous toxicity	0.94	0.2	0.93	0.2	0.01	0.83
Previous disability	0.39	0.5	0.39	0.5	0.01	0.74
Previous Holocaust	0.10	0.3	0.10	0.3	-0.01	0.98
<i>Joint tests/differences</i>						
<i>F</i> -test ( <i>p</i> -value)						0.70
Multivariate normalized difference					0.12	

Notes: Columns 2 to 5 display means and standard deviations (SD). Column 6 displays normalized differences (Imbens and Rubin, 2015); all variables have differences below the recommended 0.25. Column 7 has  $p$ -values from regressions of characteristics on a treatment dummy and strata fixed-effects.  $F$ -tests are from regressions of a treatment indicator on pre-treatment variables.



**Figure B.4:** Screenshot of a reply

Notes: Some Tweets in my sample are replies or comments to other users' Tweets.



**Figure B.5:** Topic classification by slur

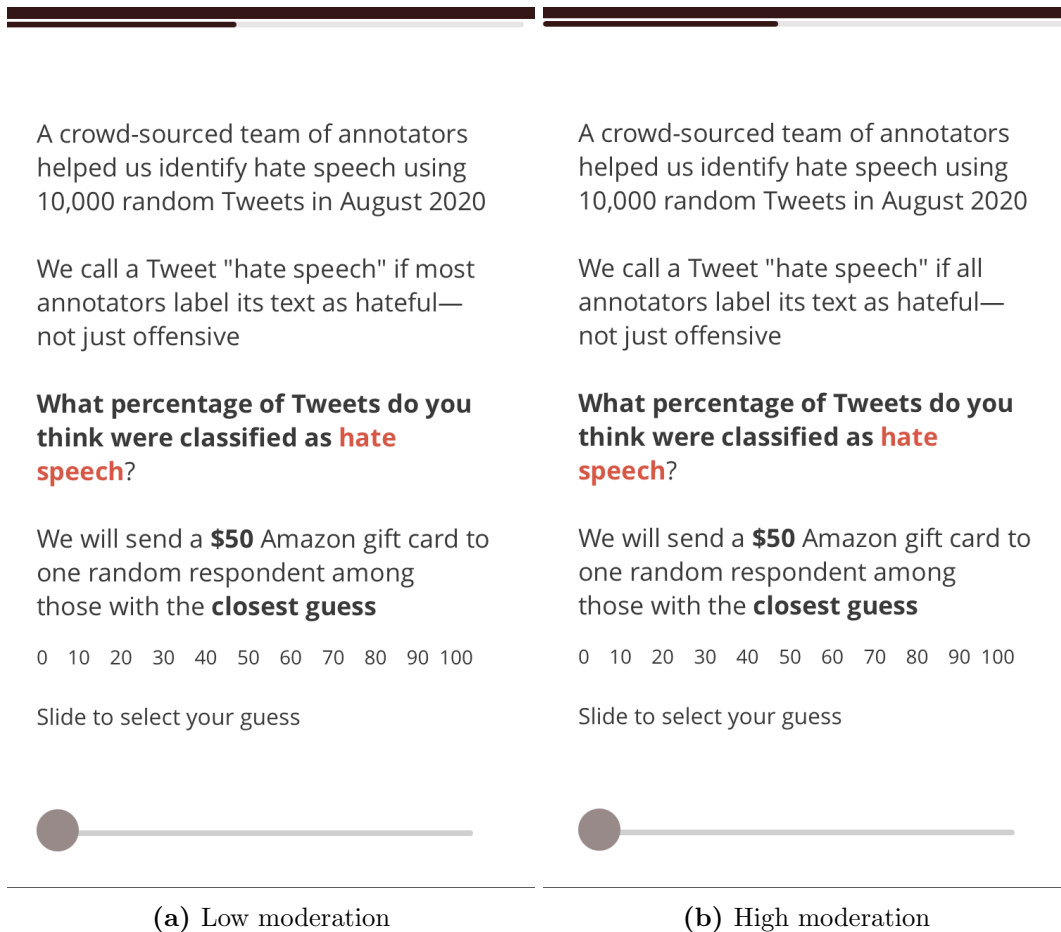
Notes: Each figure presents the distribution of Tweets by their main topic. Three MTurk workers read each Tweet and decided its most relevant topic among the eight options in the figures. The main topic is the one that two or three workers agreed upon. If there was no agreement, the topic of the Tweet is set to “Other or uncertain”.

**Table B.4:** Reporting accounts summary statistics

	Account										
	1	2	3	4	5	6	7	8	9	10	11
<i>Accounts</i>											
City	CHI	CHI	NYC	MIA	LA	LA	DAL	SF	ATL	CHI	DC
Email	Yes	Yes	Yes	No	Yes	No	Yes	No	No	Yes	No
Phone	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Mobile	Yes	No	No	No	No	Yes	No	No	Yes	No	Yes
App	Yes	No	No	No	No	No	Yes	No	No	Yes	No
Shadow ban	No	No	No	No	Yes	No	No	No	Yes	No	No
Ads	Yes	No	No	No	No	No	No	No	No	No	No
Account years	2.7	2.5	2.5	2.3	2.3	2.3	2.3	2.2	0.3	9.1	0.1
Tweets/month	0	1.1	0.3	2.4	2.2	0.1	0.2	0.3	3.5	0.5	4.1
Likes/month	0	1.0	0.6	1.8	1.8	0.8	1.1	0.5	4.7	1.2	8.2
Followers	0	18	3	2	1	0	0	0	0	168	1
Followed	6	22	28	48	43	25	19	16	14	134	17
Bot score	.	0.4	0.5	0.2	0.3	0.4	0.4	0.5	0.4	0.2	0.5
<i>Profiles</i>											
Description	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Location	Yes	Yes	No	Yes	Yes	No	No	Yes	No	Yes	Yes
Default pic	No	No	No	No	No	No	No	No	No	No	No
Verified	No	No	No	No	No	No	No	No	No	No	No
Protected	No	No	No	No	No	No	No	No	No	No	No

Notes: Each column corresponds to one of the 11 Twitter accounts used for the reporting treatment. City is the location of the virtual private network used for reporting. Email and Phone indicate whether the account had an associated email and phone number, respectively. Mobile indicates whether reporting was done using a phone or a computer. App indicates whether the account was accessed using the official Twitter app or a browser. Ads denotes whether the account had been used in the past to run ads. Data gathered in August, 2021.



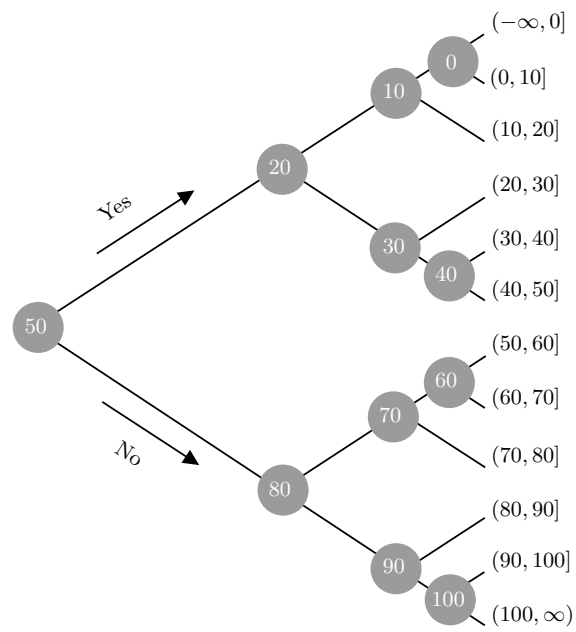


**Figure B.6:** Instructions and elicitation of beliefs about prevalence

**Table B.5:** Balance in the welfare experiment

Characteristic	Control		Treatment		N. Dif.	<i>p</i> -value
	Mean	SD	Mean	SD		
<i>Observations</i>	1,515		1,512			
<i>Demographics</i>						
Age	38.05	12.8	38.10	12.3	-0.00	0.92
Female	0.45	0.5	0.45	0.5	0.01	0.44
College graduate +	0.32	0.5	0.31	0.5	0.03	0.35
Some college	0.33	0.5	0.33	0.5	-0.01	0.73
White non-Hispanic	0.67	0.5	0.69	0.5	-0.05	0.06
Black non-Hispanic	0.15	0.4	0.14	0.4	0.01	0.82
Hispanic	0.09	0.3	0.08	0.3	0.04	0.27
Asian non-Hispanic	0.03	0.2	0.02	0.2	0.06	0.10
Northeast	0.22	0.4	0.25	0.4	-0.07	0.07
Midwest	0.18	0.4	0.18	0.4	0.01	0.84
South	0.39	0.5	0.36	0.5	0.06	0.12
Republican	0.23	0.4	0.22	0.4	0.03	0.47
Democrat	0.52	0.5	0.54	0.5	-0.04	0.33
Christian	0.62	0.5	0.61	0.5	0.02	0.54
Jewish	0.02	0.1	0.02	0.1	-0.00	0.91
Muslim	0.04	0.2	0.04	0.2	-0.01	0.78
Buddhist or Hindu	0.01	0.1	0.01	0.1	0.02	0.65
Income	64.09	32.9	64.22	33.4	-0.00	0.90
Minority	0.48	0.5	0.48	0.5	-0.00	1.00
<i>Twitter / Social media</i>						
Daily hours on Twitter	1.52	2.3	1.52	2.2	-0.00	0.99
Provided handle	0.64	0.5	0.64	0.5	-0.00	1.00
User exists	0.47	0.5	0.47	0.5	-0.01	0.69
Tweets per day	1.29	6.4	1.79	10.7	-0.06	0.29
Likes per day	1.82	6.3	2.86	18.1	-0.08	0.14
Account years	8.06	4.3	7.80	4.2	0.06	0.25
Platforms other than Twitter	5.12	2.1	5.10	2.0	0.01	0.74
Has been harassed online	0.28	0.5	0.29	0.5	-0.01	0.85
Prevalence of hate in feed	20.06	23.1	20.85	24.5	-0.03	0.31
<i>Moderation</i>						
Has been sanctioned	0.23	0.4	0.23	0.4	0.00	1.00
Has reported	0.36	0.5	0.37	0.5	-0.00	0.88
Has been reported	0.12	0.3	0.12	0.3	0.03	0.36
<i>Beliefs</i>						
Prevalence of hate	36.73	25.7	36.12	26.2	0.02	0.51
Likelihood of moderation	39.73	28.6	40.91	28.7	-0.04	0.25
<i>Joint tests/differences</i>						
<i>F</i> -test ( <i>p</i> -value)						0.33
Multivariate normalized difference						0.25

Notes: Columns 2 to 5 display means and standard deviations (SD). Column 6 displays normalized differences (Imbens and Rubin, 2015). Column 7 has *p*-values from a regression of characteristics on treatment and strata fixed-effects. *F*-tests are from regressions of a treatment indicator on characteristics.



**Figure B.7:** Iterative multiple price list

Notes: The circles denote compensation (Amazon gift card) offers to deactivate social media. The intervals correspond to the willingness to accept.

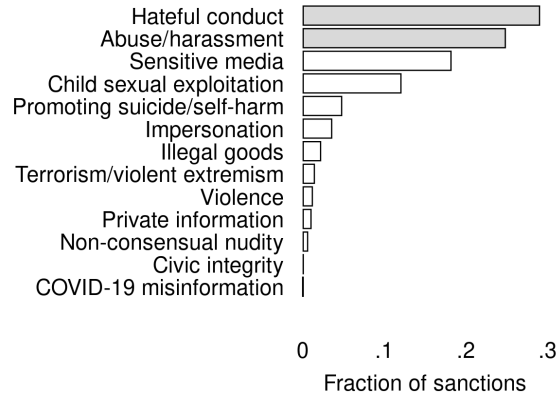
## C Data appendix

### C.1 Measurement of sanctions

The “lookup statuses” or “lookup users” endpoints of the Twitter API indicate when a tweet or account go missing. Among missing accounts and statuses, the “show users” or “show statuses” endpoints of the API return an error code that details why they were missing (see Twitter (2021g) for a full list of error codes). With the error code information one can measure the following events:

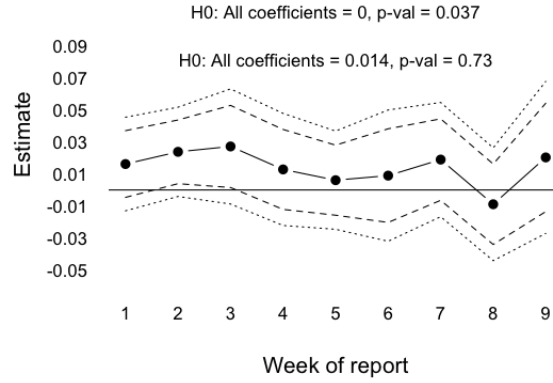
- Twitter required the removal of a post, but it has not been removed by the user. This is reflected in a missing status with error code 421.
- Twitter required the removal of a post, and it has been removed by the user. This is reflected in a missing status with error code 422. After some days, the status transitions to code 144 (deleted status). Twitter claims that the notice will be available 14 days after the tweet is removed (Twitter, 2021e) but empirically it seems like this period varies.
- A post is missing because the user deleted it. This is reflected in a missing status with error code 144.
- A post is missing because the user protected their account or because the user blocked my developer account. This is reflected in a missing status with error code 179 or 136, respectively. Protected accounts are also detected with the lookup users endpoint. It is rare to encounter a user that blocks my developer account. Most likely is due to users mass-blocking all the followers of some famous account.
- A post and the account are missing because the user is suspended. This is reflected in a missing user and potentially missing status with code 63 (Chowdhury et al., 2020).
- A post and the account are missing because the user deleted their account. This is reflected in a missing user with code 50.

## C.2 Figures and Tables



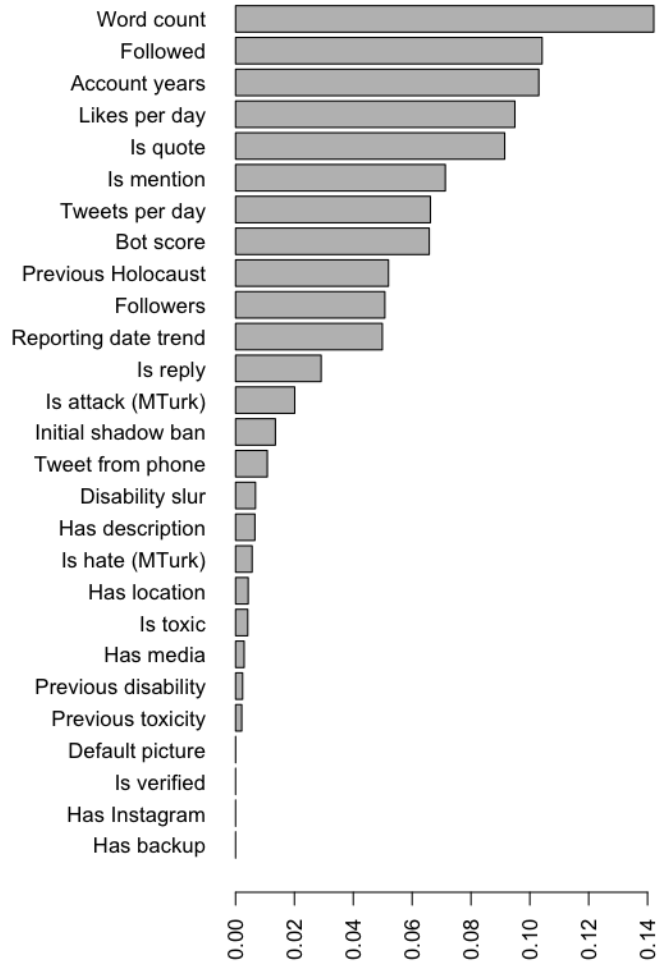
**Figure C.1:** Histogram of sanctions by rule violation

Notes: This figure plots the fraction of sanctions (actioned accounts) by the type of rule violation. It uses data from the second half of 2020 from Twitter’s Transparency Rules Enforcement Report (Twitter, 2020b).



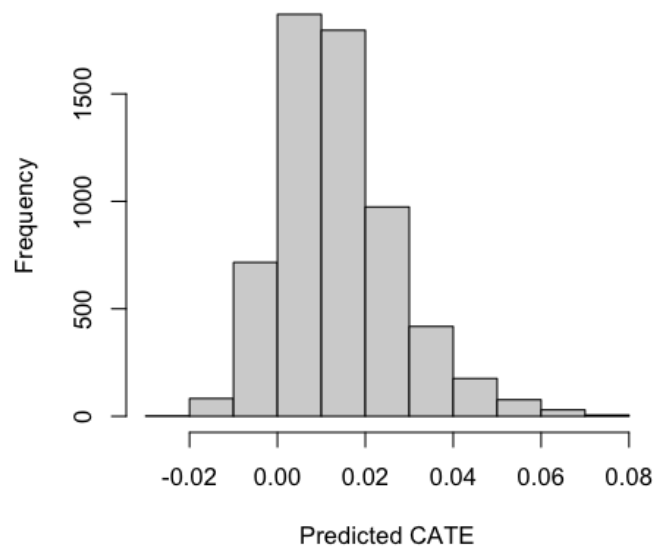
**Figure C.2:** Treatment effect on Tweet deletions by week of the report

Notes: This figure presents the treatment effects on the cumulative likelihood of deletions at day 21, by week of the report. Pointwise confidence intervals are dashed and sup- $t$  simultaneous confidence bands are dotted. Dynamic effects use the estimator from Roth and Sant’Anna (2021). The  $p$ -values are from joint Wald tests. Note that there is not enough power to measure the weekly treatment effects precisely (since there are around 680 observations per week).



**Figure C.3:** Importance for heterogeneous treatment effects in the first stage

Notes: This figure plots the measure of “importance” for each variable used to estimate the heterogeneous treatment effects (on the full sample). This measure is a weighted sum of how many times each feature was split on at each depth in the forest (Tibshirani et al., 2020).



**Figure C.4:** Histogram of predicted CATEs

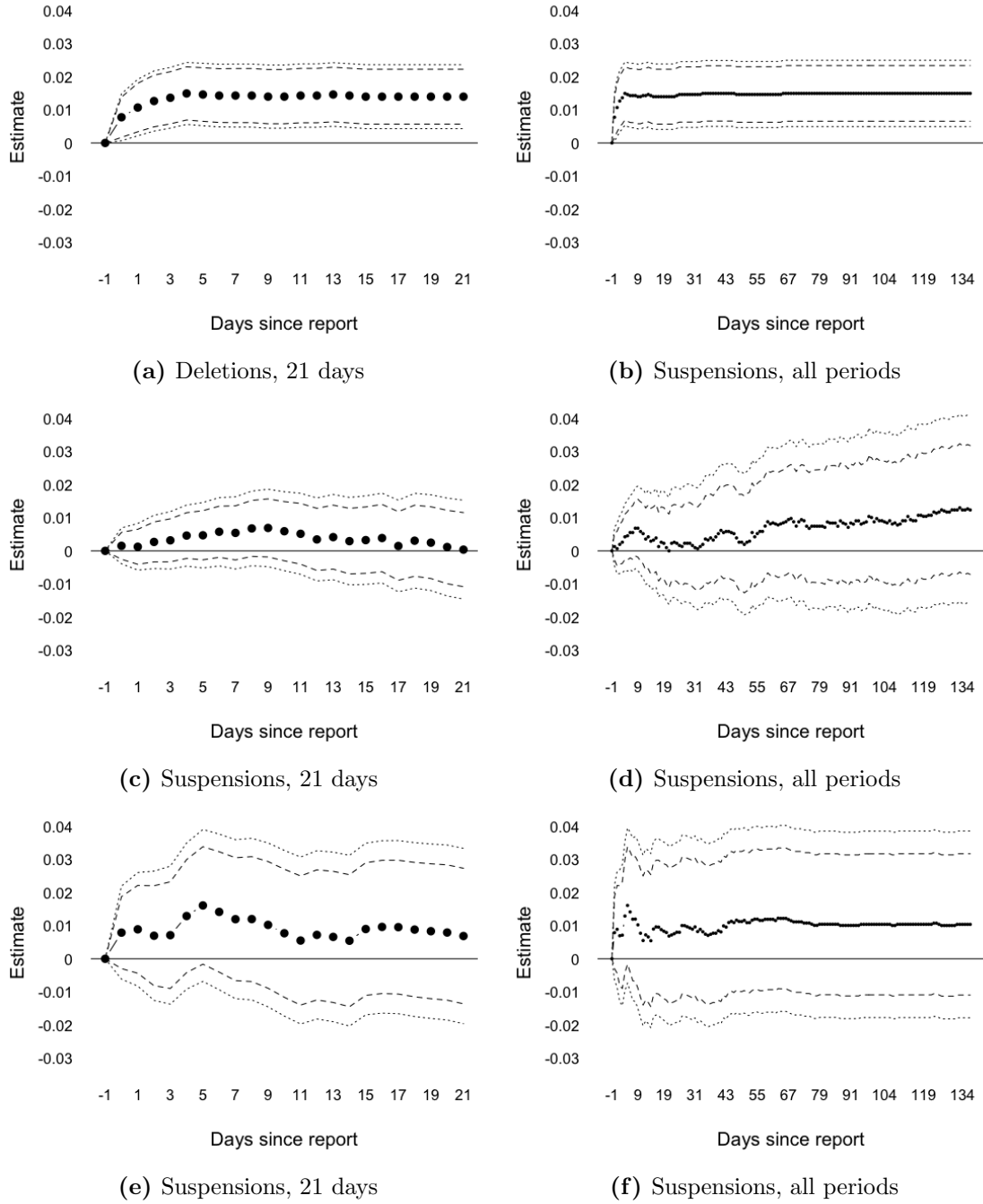
Notes: This figure plots the estimated CATEs for the full sample (Tibshirani et al., 2020).

**Table C.1:** Profiling compliers and non-compliers

Characteristic	Sample	Complier	Never-taker	Always-taker	C-A	C-N	A-N
<i>Proportion</i>	1.000 (0.000)	0.014 (0.004)	0.965 (0.003)	0.021 (0.003)			
<i>Accounts</i>							
Account years	3.22 (0.04)	3.47 (3.16)	3.22 (0.06)	3.06 (0.47)	0.90	0.94	0.73
Tweets per day	11.62 (0.32)	-16.56 (24.07)	12.05 (0.46)	10.51 (2.77)	0.27	0.24	0.58
Likes per day	24.17 (0.65)	-64.14 (53.67)	25.16 (0.96)	37.49 (12.06)	0.07	0.10	0.31
Followers	634.85 (66.15)	-7,858.02 (5353.87)	767.85 (121.02)	184.42 (64.32)	0.14	0.11	0.00
Followed	433.75 (10.89)	222.29 (780.69)	440.73 (17.58)	254.88 (37.47)	0.97	0.78	0.00
Bot score	0.24 (0.00)	0.28 (0.12)	0.24 (0.00)	0.26 (0.01)	0.86	0.72	0.11
Initial shadow ban	0.71 (0.01)	0.75 (0.50)	0.71 (0.01)	0.76 (0.05)	0.97	0.94	0.32
<i>Tweets</i>							
Word count	15.98 (0.17)	27.50 (12.56)	15.83 (0.25)	15.15 (1.45)	0.33	0.36	0.65
Is toxic	0.80 (0.01)	1.22 (0.39)	0.80 (0.01)	0.75 (0.05)	0.23	0.28	0.40
Is hate (MTurk)	0.30 (0.01)	0.36 (0.42)	0.30 (0.01)	0.42 (0.06)	0.90	0.89	0.07
Is reply	0.84 (0.00)	0.44 (0.35)	0.84 (0.01)	0.88 (0.04)	0.23	0.26	0.40
Is attack (MTurk)	0.78 (0.01)	0.09 (0.43)	0.79 (0.01)	0.83 (0.05)	0.09	0.11	0.39
Is quote	0.07 (0.00)	0.50 (0.26)	0.07 (0.00)	0.05 (0.03)	0.09	0.10	0.46
Is mention	0.85 (0.00)	0.55 (0.34)	0.85 (0.01)	0.92 (0.03)	0.27	0.36	0.04
Has media	0.04 (0.00)	-0.16 (0.19)	0.04 (0.00)	0.05 (0.03)	0.28	0.28	0.89
Tweet from phone	0.79 (0.01)	1.03 (0.38)	0.79 (0.01)	0.78 (0.05)	0.52	0.51	0.99
<i>Profiles</i>							
Has description	0.82 (0.00)	0.60 (0.36)	0.82 (0.01)	0.82 (0.05)	0.56	0.55	0.92
Has location	0.51 (0.01)	0.26 (0.46)	0.51 (0.01)	0.46 (0.06)	0.66	0.58	0.41
Default picture	0.03 (0.00)	0.16 (0.17)	0.03 (0.00)	0.02 (0.02)	0.40	0.45	0.25
Is verified	0.00 (0.00)	-0.02 (0.03)	0.00 (0.00)	0.00 (0.00)	0.49	0.46	0.03
Has Instagram	0.01 (0.00)	-0.06 (0.11)	0.02 (0.00)	0.00 (0.00)	0.60	0.51	0.00
Has backup	0.01 (0.00)	0.13 (0.10)	0.01 (0.00)	0.00 (0.00)	0.19	0.22	0.00
<i>Timelines</i>							
Previous toxicity	0.93 (0.00)	0.98 (0.23)	0.93 (0.00)	0.95 (0.03)	0.92	0.85	0.45
Previous disability	0.39 (0.01)	0.49 (0.45)	0.39 (0.01)	0.35 (0.06)	0.77	0.83	0.55
Previous Holocaust	0.10 (0.00)	0.24 (0.27)	0.09 (0.01)	0.20 (0.05)	0.87	0.58	0.04

Notes: Columns 2-5 present means and standard errors (in parenthesis) for all users in the sample, compliers, never-takers and always-takers, following Marbach and Hangartner (2020) (assuming Tweet deletions as first stage). Columns 6-8 present p-values from *t*-tests of difference in means with unequal variances.





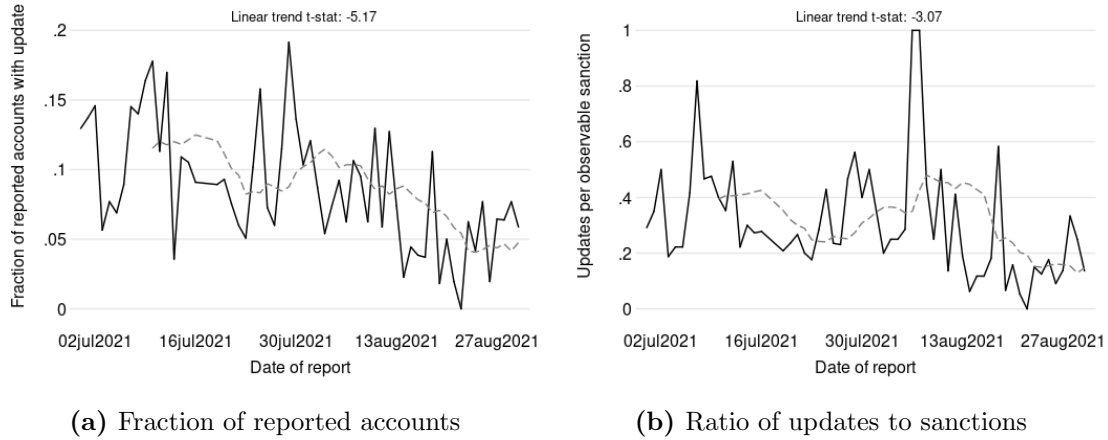
**Figure C.5:** Cumulative dynamic treatment effects on observable sanctions

Notes: This figure presents cumulative dynamic treatment effects, pointwise confidence intervals (dashed), and sup- $t$  simultaneous confidence bands (dotted).

**Table C.2:** Effects of reporting on other observable sanctions and self-censorship

<i>Panel A: other Twitter sanctions</i>									
	Suspensions			Shadow-bans			Missing Other Tweets		
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
Treatment	-0.000 (0.006)	-0.000 (0.006)	-0.000 (0.006)	0.001 (0.012)	0.001 (0.012)	-0.002 (0.011)	0.004 (0.005)	0.004 (0.005)	0.004 (0.005)
<i>y</i> Mean	0.05	0.05	0.05	0.26	0.26	0.26	0.05	0.05	0.05
<i>y</i> SD	0.22	0.22	0.22	0.44	0.44	0.44	0.18	0.18	0.18
$R^2$	0.00	0.03	0.03	0.00	0.02	0.10	0.00	0.02	0.03
Obs.	6,148	6,134	6,134	5,692	5,675	5,675	5,381	5,360	5,360
<i>Panel B: self-censorship</i>									
	Tweet deletion			Account deletion			Protecting account		
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
Treatment	0.006 (0.005)	0.005 (0.005)	0.005 (0.005)	0.001 (0.003)	0.001 (0.003)	0.001 (0.003)	0.000 (0.004)	0.000 (0.004)	-0.000 (0.004)
<i>y</i> Mean	0.03	0.03	0.03	0.02	0.02	0.02	0.03	0.03	0.03
<i>y</i> SD	0.18	0.18	0.18	0.13	0.13	0.13	0.17	0.17	0.17
$R^2$	0.00	0.01	0.01	0.00	0.02	0.02	0.00	0.02	0.04
Obs.	6,148	6,134	6,134	6,148	6,134	6,134	6,148	6,134	6,134
Strata FE	No	Yes	Yes	No	Yes	Yes	No	Yes	Yes
Controls	No	No	Yes	No	No	Yes	No	No	Yes

Notes: This table reports estimates from OLS regressions on treatment assignment. Robust standard errors are parenthesized. Controls are selected using the double-lasso method of Belloni et al. (2014) recommended by Urmitsky et al. (2016). \*, \*\*, and \*\*\* denote significance at the 10%, 5%, and 1% levels, respectively.

**Figure C.6:** Fraction of accounts that received an update and ratio of updates to observable sanctions, by reporting date

Notes: Panel a) presents the fraction of reported accounts that received an update within 21 days of the report, by reporting date (black line). Panel b) presents the ratio of updates to number of observable sanctions (black line). Gray dashed lines denote 10-day moving averages.

**Table C.3:** Reported accounts with updates and sanctions

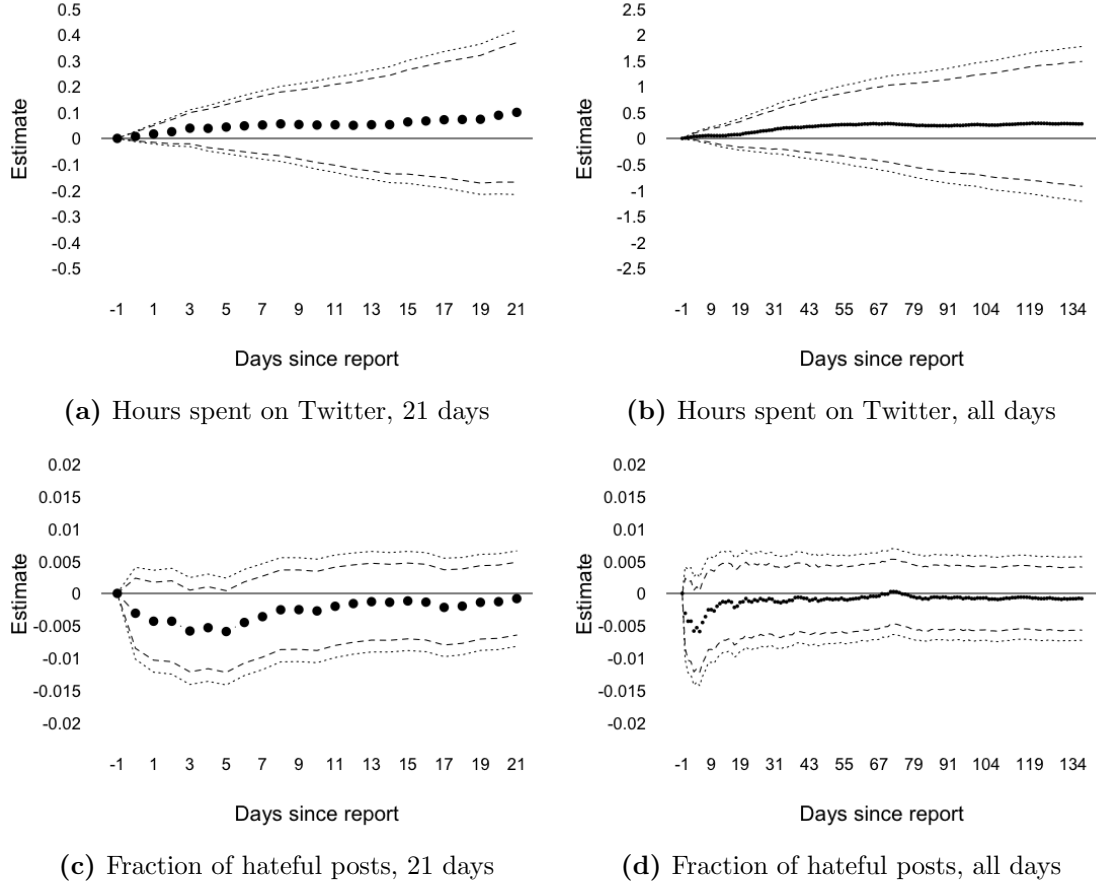
	Update	No update
Observable sanction	153	814
No observable sanction	117	1,990

Notes: This table presents the number of reported accounts that received an update and/or observable sanctions (deletions, suspensions, shadowbans) within 3 weeks of the reports.

**Table C.4:** Effects of reporting on other measures of activity

<i>Panel A: Tweets and Likes</i>						
	Tweets			Likes		
	(1)	(2)	(3)	(4)	(5)	(6)
Treatment	35.882* (20.694)	35.643* (20.753)	22.946 (15.419)	26.416 (41.243)	27.580 (41.467)	4.863 (29.651)
<i>y</i> Mean	405.47	405.89	405.89	846.49	847.86	847.86
<i>y</i> SD	782.50	783.58	783.58	1559.14	1560.95	1560.95
$R^2$	0.00	0.02	0.45	0.00	0.01	0.49
Obs.	5,717	5,697	5,697	5,717	5,697	5,697
<i>Panel B: other activity measures</i>						
	Winsorized time			Fraction of active days		
	(1)	(2)	(3)	(4)	(5)	(6)
Treatment	0.192 (0.133)	0.192 (0.134)	0.126 (0.109)	-0.001 (0.001)	-0.001 (0.001)	-0.001 (0.001)
<i>y</i> Mean	3.36	3.37	3.37	1.09	1.09	1.09
<i>y</i> SD	5.05	5.05	5.05	0.04	0.04	0.04
$R^2$	0.00	0.02	0.34	0.00	0.01	0.01
Obs.	5,717	5,697	5,697	5,727	5,708	5,708
Strata FE	No	Yes	Yes	No	Yes	Yes
Controls	No	No	Yes	No	No	Yes

Notes: This table reports estimates from OLS regressions on treatment assignment. Robust standard errors are parenthesized. Controls are selected using the double-lasso method of Belloni et al. (2014) recommended by Urminsky et al. (2016). \*, \*\*, and \*\*\* denote significance at the 10%, 5%, and 1% levels, respectively.



**Figure C.7:** Cumulative dynamic treatment effects on activity and hatefulness

Notes: This figure presents cumulative dynamic treatment effects, pointwise confidence intervals (dashed), and sup- $t$  simultaneous confidence bands (dotted).

**Table C.5:** Effects of reporting on other measures of hatefulness

<i>Panel A: extensive margin</i>						
	Posting toxicity $\geq 0.8$			Repeating the slur		
	(1)	(2)	(3)	(4)	(5)	(6)
Treatment	0.004 (0.008)	0.004 (0.008)	0.003 (0.008)	0.001 (0.013)	0.001 (0.013)	0.001 (0.011)
<i>y</i> Mean	0.90	0.90	0.90	0.62	0.61	0.61
<i>y</i> SD	0.30	0.30	0.30	0.49	0.49	0.49
$R^2$	0.00	0.01	0.03	0.00	0.02	0.34
Obs.	5,727	5,708	5,708	5,727	5,708	5,708
<i>Panel B: average scores</i>						
	Average toxicity			Average severe toxicity		
	(1)	(2)	(3)	(4)	(5)	(6)
Treatment	-0.001 (0.003)	-0.001 (0.003)	0.000 (0.003)	-0.001 (0.002)	-0.001 (0.002)	-0.000 (0.002)
<i>y</i> Mean	0.30	0.30	0.30	0.18	0.18	0.18
<i>y</i> SD	0.11	0.11	0.11	0.09	0.09	0.09
$R^2$	0.00	0.01	0.09	0.00	0.01	0.08
Obs.	5,631	5,616	5,616	5,631	5,616	5,616
Strata FE	No	Yes	Yes	No	Yes	Yes
Controls	No	No	Yes	No	No	Yes

Notes: This table reports estimates from OLS regressions on treatment assignment. Robust standard errors are parenthesized. Controls are selected using the double-lasso method of Belloni et al. (2014) recommended by Urminsky et al. (2016). \*, \*\*, and \*\*\* denote significance at the 10%, 5%, and 1% levels, respectively.

**Table C.6:** Effects of reporting on other measures of replied users' activity

<i>Panel A: Tweets and Likes</i>						
	Tweets			Likes		
	(1)	(2)	(3)	(4)	(5)	(6)
Treatment	59.739* (30.756)	58.073* (30.903)	50.849* (30.824)	152.538*** (56.932)	148.888*** (57.615)	141.722** (57.611)
$y$ Mean	656.20	657.39	657.39	1151.51	1151.96	1151.96
$y$ SD	1060.33	1062.04	1062.04	1963.42	1964.77	1964.77
$R^2$	0.00	0.02	0.03	0.00	0.02	0.03
Obs.	4,752	4,733	4,733	4,752	4,733	4,733
<i>Panel B: other activity measures</i>						
	Winsorized time			Fraction of active days		
	(1)	(2)	(3)	(4)	(5)	(6)
Treatment	0.420** (0.210)	0.408* (0.211)	0.362* (0.210)	-0.019 (0.113)	-0.032 (0.113)	-0.033 (0.113)
$y$ Mean	5.21	5.21	5.21	20.42	20.42	20.42
$y$ SD	7.23	7.24	7.24	3.91	3.91	3.91
$R^2$	0.00	0.02	0.04	0.00	0.02	0.02
Obs.	4,752	4,733	4,733	4,761	4,742	4,742
Strata FE	No	Yes	Yes	No	Yes	Yes
Controls	No	No	Yes	No	No	Yes

Notes: This table reports estimates from OLS regressions on treatment assignment. Robust standard errors are parenthesized. Controls are selected using the double-lasso method of Belloni et al. (2014) recommended by Urminsky et al. (2016). \*, \*\*, and \*\*\* denote significance at the 10%, 5%, and 1% levels, respectively.

**Table C.7:** Effects of reporting on other measures of replied users' activity, sample of attacks

<i>Panel A: Tweets and Likes</i>						
	Tweets			Likes		
	(1)	(2)	(3)	(4)	(5)	(6)
Treatment	77.136** (32.018)	74.865** (32.165)	71.894** (32.037)	157.203** (61.304)	156.095** (62.240)	152.089** (62.025)
<i>y</i> Mean	635.01	635.59	635.59	1140.39	1142.13	1142.13
<i>y</i> SD	1035.93	1037.27	1037.27	1983.49	1986.15	1986.15
$R^2$	0.00	0.02	0.03	0.00	0.02	0.03
Obs.	4,171	4,155	4,155	4,171	4,155	4,155
<i>Panel B: other activity measures</i>						
	Winsorized time			Fraction of active days		
	(1)	(2)	(3)	(4)	(5)	(6)
Treatment	0.512** (0.221)	0.499** (0.223)	0.478** (0.222)	-0.012 (0.123)	-0.028 (0.123)	-0.028 (0.123)
<i>y</i> Mean	5.07	5.08	5.08	20.35	20.35	20.35
<i>y</i> SD	7.14	7.15	7.15	3.97	3.97	3.97
$R^2$	0.00	0.02	0.03	0.00	0.02	0.02
Obs.	4,171	4,155	4,155	4,178	4,162	4,162
Strata FE	No	Yes	Yes	No	Yes	Yes
Controls	No	No	Yes	No	No	Yes

Notes: This table reports estimates from OLS regressions on treatment assignment. Robust standard errors are parenthesized. Controls are selected using the double-lasso method of Belloni et al. (2014) recommended by Urminsky et al. (2016). \*, \*\*, and \*\*\* denote significance at the 10%, 5%, and 1% levels, respectively.

**Table C.8:** Effects of reporting on the time spent of followers, follower and following counts

	Followers' time spent			Follower count			Following count		
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
Treatment	0.041 (0.169)	-0.067 (0.144)	-0.028 (0.123)	9.543 (8.247)	9.096 (7.812)	8.808 (7.820)	2.487 (2.061)	2.391 (2.067)	2.208 (2.048)
<i>y</i> Mean	3.68	3.68	3.68	25.58	25.48	25.48	14.13	14.14	14.14
<i>y</i> SD	12.70	12.70	12.70	311.83	312.19	312.19	77.93	78.05	78.05
$R^2$	0.00	0.03	0.05	0.00	0.02	0.02	0.00	0.02	0.03
Obs.	824,391	824,390	824,390	5,717	5,697	5,697	5,717	5,697	5,697
Strata FE	No	Yes	Yes	No	Yes	Yes	No	Yes	Yes
Controls	No	No	Yes	No	No	Yes	No	No	Yes

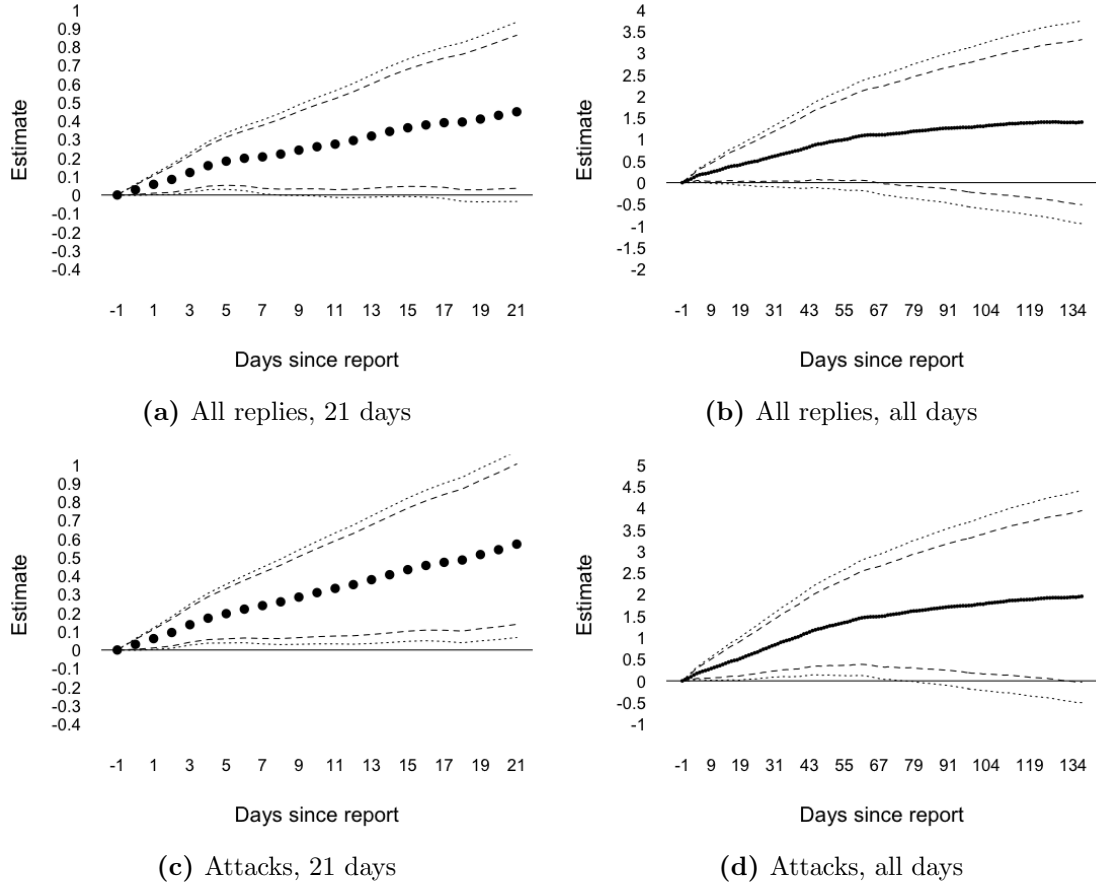
Notes: This table reports estimates from OLS regressions on treatment assignment. Clustered standard errors at the sampled user level for followers' time spent and robust standard errors for follower and following counts are parenthesized. Controls are selected using the double-lasso method of Belloni et al. (2014) recommended by Urminsky et al. (2016). \*, \*\*, and \*\*\* denote significance at the 10%, 5%, and 1% levels, respectively.

**Table C.9:** Effects of reporting on attrition

	Attrition on day 21			Attrition on day $\leq 21$		
	(1)	(2)	(3)	(4)	(5)	(6)
Treatment	0.002 (0.006)	0.001 (0.006)	0.001 (0.006)	0.002 (0.007)	0.002 (0.007)	0.001 (0.007)
$y$ Mean	0.07	0.07	0.07	0.08	0.08	0.08
$y$ SD	0.25	0.25	0.25	0.28	0.28	0.28
$R^2$	0.00	0.02	0.03	0.00	0.02	0.04
Obs.	6,148	6,134	6,134	6,148	6,134	6,134
Strata FE	No	Yes	Yes	No	Yes	Yes
Controls	No	No	Yes	No	No	Yes

Notes: This table reports estimates from OLS regressions on treatment assignment. Robust standard errors are parenthesized. Controls are selected using the double-lasso method of Belloni et al. (2014) recommended by Urminsky et al. (2016). \*, \*\*, and \*\*\* denote significance at the 10%, 5%, and 1% levels, respectively.





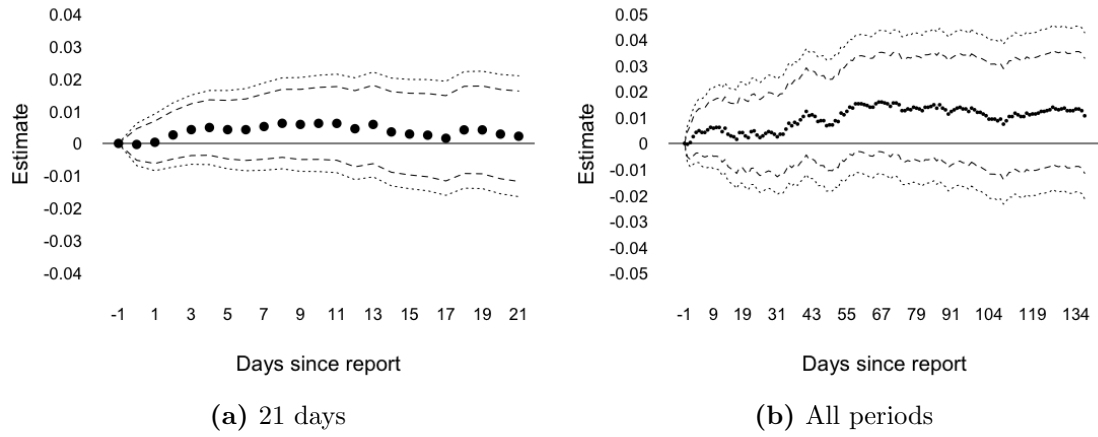
**Figure C.8:** Cumulative dynamic treatment effect on replied users' activity

Notes: This figure presents cumulative dynamic treatment effects, pointwise confidence intervals (dashed), and sup- $t$  simultaneous confidence bands (dotted). The outcome variable is a measure of time spent of the users that the posts in the sample reply to. It is a linear combination of Tweets and Likes.

**Table C.10:** Correlations between reporting account characteristics and first stage success

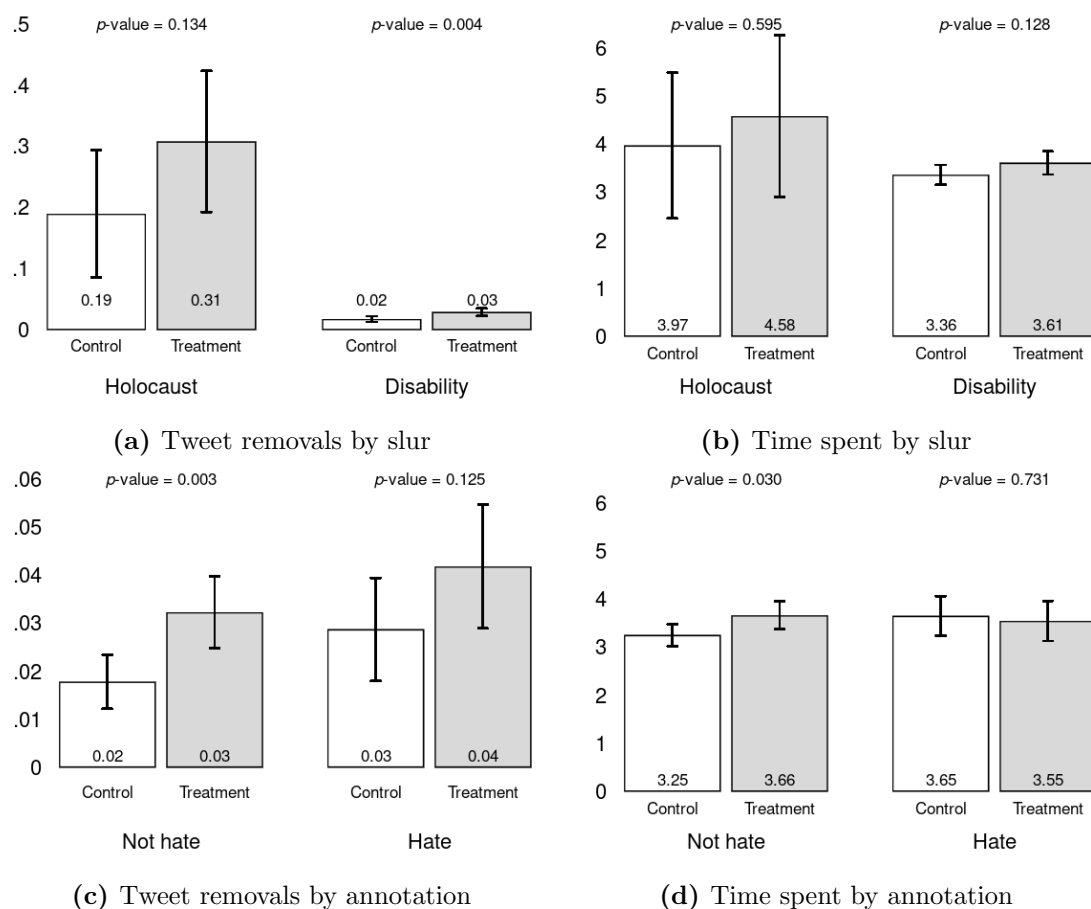
	(1)
Registered email	0.0058 (0.010)
Mobile device	0.0328** (0.015)
Firefox browser	-0.0233* (0.013)
Opera browser	0.0081 (0.010)
Safari browser	0.0199 (0.016)
Tweets/month	0.0018 (0.007)
Likes/month	-0.0062 (0.005)
Followed	0.0009 (0.001)
Followers	-0.0006 (0.000)
Has location	0.0211** (0.009)
$y$ Mean	0.04
$y$ Std. Dev.	0.18
$R^2$	0.01
$F$	2.01
Obs.	3,074

Notes: This table reports estimates from OLS regressions of an indicator variable of whether Tweets get deleted on reporting account characteristics, on the subsample of users assigned to the reporting treatment arm. Robust standard errors are parenthesized. \*, \*\*, and \*\*\* denote significance at the 10%, 5%, and 1% levels, respectively.



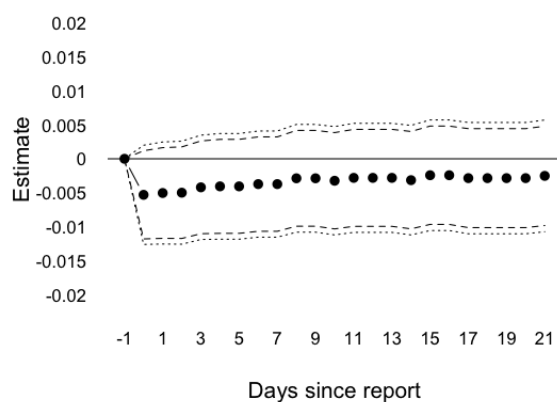
**Figure C.9:** Cumulative dynamic treatment effect on attrition

Notes: These figures presents dynamic treatment effects on an indicator of whether users drop from the sample at or before every day after reporting. Pointwise confidence intervals are dashed and sup- $t$  confidence bands are dotted.



**Figure C.10:** Heterogeneity by slur and hate annotation

Notes: This figure reports estimates of reporting on Tweet removals and users' time spent posting and liking by slur and hate annotation.



**Figure C.11:** Cumulative effect on the likelihood of mentioning the replied user

Notes: This figure presents dynamic treatment effects on an indicator of whether the users in the sample mention the replied users again. Pointwise confidence intervals are dashed and sup-t confidence bands are dotted.

**Table C.11:** Inference robustness and Lee Bounds

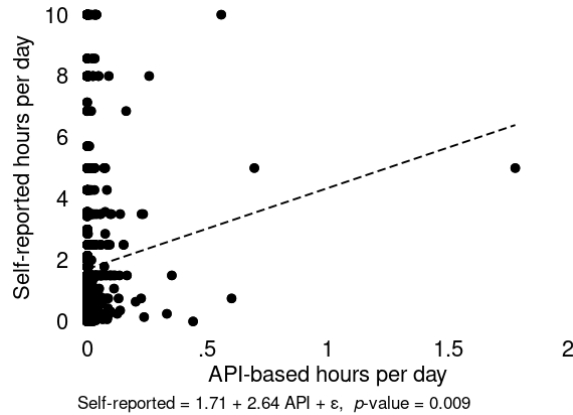
	Tweet removals (1)	Suspensions (2)	Shadowbans (3)	Hrs. spent (4)	Toxicity (5)	Replied hrs. (6)	Attacked hrs. (7)
Treatment	0.014*** (0.004)	0.000 (0.006)	0.001 (0.012)	0.254 (0.159)	-0.002 (0.003)	0.513** (0.232)	0.650*** (0.244)
<i>Inference robustness (p-values)</i>							
Robust	0.001	1.000	0.934	0.112	0.472	0.027	0.008
Wild bootstrap	0.001	1.000	0.927	0.104	0.478	0.025	0.008
Permutation test	0.001	1.000	0.935	0.111	0.478	0.026	0.007
MHT	0.000	1.000	0.998	0.361	0.849	0.133	0.039
<i>Lee bounds</i>							
Lower	-	-	-0.004 (0.013)	0.209 (0.748)	-0.003 (0.003)	0.508 (0.243)	0.022 (0.641)
Upper	-	-	0.003 (0.012)	0.255 (0.161)	0.002 (0.008)	0.564 (0.840)	0.722 (0.265)
$y$ Mean	0.03	0.05	0.26	3.50	0.13	5.33	5.19
$y$ SD	0.17	0.22	0.44	6.03	0.12	8.01	7.90
$R^2$	0.00	0.00	0.00	0.00	0.00	0.00	0.00
Obs.	6,148	6,148	5,692	5,717	5,631	4,752	4,171

Notes: This table reports estimates from OLS regressions on treatment assignment. Robust standard errors are parenthesized. \*, \*\*, and \*\*\* denote significance at the 10%, 5%, and 1% levels, respectively. MHT correction refers to the multiple-hypothesis testing procedure presented in List et al. (2019). Reported p-values for wild bootstrap, permutation tests and MHT are derived from running 5,000 replications.

**Table C.12:** Effects on sanctions among Tweets with replied and attacked users

<i>Panel A: Sample of replies</i>									
	Tweet removals			Suspensions			Shadow-bans		
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
Treatment	0.008*	0.008*	0.008*	0.003	0.002	0.002	-0.003	-0.001	-0.004
	(0.004)	(0.004)	(0.004)	(0.007)	(0.006)	(0.006)	(0.013)	(0.013)	(0.013)
<i>y</i> Mean	0.02	0.02	0.02	0.05	0.05	0.05	0.24	0.24	0.24
<i>y</i> SD	0.15	0.15	0.15	0.22	0.22	0.22	0.43	0.43	0.43
<i>R</i> <sup>2</sup>	0.00	0.08	0.08	0.00	0.03	0.03	0.00	0.02	0.08
Obs.	4,752	4,734	4,734	4,752	4,734	4,734	4,404	4,388	4,388
<i>Panel B: Sample of attacks</i>									
	Tweet removals			Suspensions			Shadow-bans		
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
Treatment	0.006	0.004	0.004	0.002	0.001	0.001	-0.010	-0.008	-0.010
	(0.005)	(0.005)	(0.005)	(0.007)	(0.007)	(0.007)	(0.013)	(0.013)	(0.013)
<i>y</i> Mean	0.02	0.02	0.02	0.05	0.05	0.05	0.22	0.22	0.22
<i>y</i> SD	0.15	0.15	0.15	0.22	0.22	0.22	0.42	0.42	0.42
<i>R</i> <sup>2</sup>	0.00	0.04	0.04	0.00	0.02	0.02	0.00	0.02	0.06
Obs.	4,165	4,149	4,149	4,165	4,149	4,149	3,860	3,845	3,845
Strata FE	No	Yes	Yes	No	Yes	Yes	No	Yes	Yes
Controls	No	No	Yes	No	No	Yes	No	No	Yes

Notes: This table reports estimates from OLS regressions on treatment assignment. Robust standard errors are parenthesized. Controls are selected using the double-lasso method of Belloni et al. (2014) recommended by Urminsky et al. (2016). \*, \*\*, and \*\*\* denote significance at the 10%, 5%, and 1% levels, respectively.

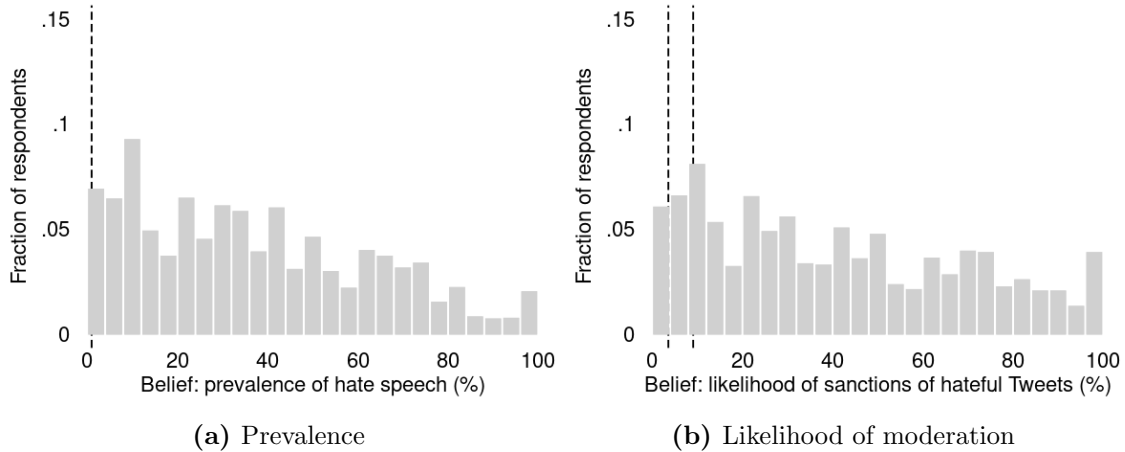
**Figure C.12:** Self-reported and API-based time spent on Twitter

Notes: This figure presents a comparison between the self-reported hours that participants spend on Twitter with the hours implied by their statuses and likes per day obtained through Twitter's API. The dashed line comes from a linear regression of self-reported hours on API-based hours.

**Table C.13:** Harassment and moderation experience by subsample

	Means			Difference <i>t</i> -stat.
	Survey	Minority	Not minority	Min.-Not min.
<i>Observations</i>	3,027	1,440	1,587	
Has been harassed	25.2	28.8	20.8	4.07
Prevalence of hate in feed	18.1	20.5	15.1	5.52
Has reported content	32.2	35.7	27.8	3.62
Has been sanctioned or reported	18.5	19.9	16.6	1.98
Tweet removal	9.6	10.4	8.8	1.33
Suspension	5.0	6.0	3.7	2.65
Shadow-ban	6.3	6.2	6.4	-0.16
Account locked	9.8	10.9	8.5	1.86
Has been reported	9.0	9.5	8.3	1.01

Notes: This table presents mean values of variables across different subsamples. It also presents *t*-statistics from tests of difference in means between minorities and not minorities. Observations are weighted to match representative Twitter users. Minority status Minority status depends on religion, sexual preference, gender, and race.

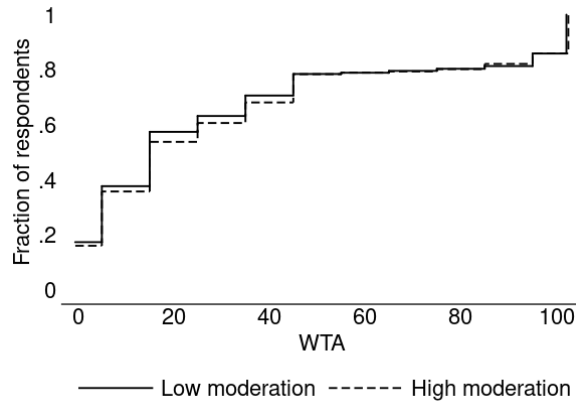
**Figure C.13:** Beliefs about prevalence and moderation of hate speech

Notes: These figures present histograms of beliefs about prevalence and moderation of hate speech among survey respondents. Prevalence is the fraction of Tweets that are classified as hate speech. Likelihood of moderation is the fraction of hate speech Tweets or users that get removed or de-platformed after 1 month of posting. The dashed lines indicate the observed values of prevalence and moderation in my sample of Tweets. One line in panel (b) corresponds to the majority rule and one to the consensus rule for classifying hate speech.

**Table C.14:** Effects of information on other measures of socia-media valuation

<i>Panel A</i>						
	WTA uniform distribution			WTA upper endpoint		
	(1)	(2)	(3)	(4)	(5)	(6)
High moderation	-0.150 (1.802)	0.024 (1.778)	0.024 (1.778)	-0.123 (1.879)	0.066 (1.853)	0.030 (1.853)
$y$ Mean	33.59	33.59	33.59	38.36	38.36	38.36
$y$ SD	36.33	36.33	36.33	37.91	37.91	37.91
$R^2$	0.00	0.02	0.02	0.00	0.02	0.02
Obs.	2,998	2,998	2,998	2,998	2,998	2,998
<i>Panel B</i>						
	WTA heuristic			TIOLI		
	(1)	(2)	(3)	(4)	(5)	(6)
High moderation	0.276 (2.902)	0.630 (2.855)	0.382 (2.857)	0.011 (0.020)	0.009 (0.020)	0.010 (0.020)
$y$ Mean	36.40	36.40	36.40	0.78	0.78	0.78
$y$ SD	58.81	58.81	58.81	0.42	0.42	0.42
$R^2$	0.00	0.02	0.03	0.00	0.01	0.02
Obs.	2,998	2,998	2,998	2,998	2,998	2,998
Strata FE	No	Yes	Yes	No	Yes	Yes
Controls	No	No	Yes	No	No	Yes

Notes: This table reports estimates from OLS regressions on treatment assignment. Robust standard errors are parenthesized. Controls are selected using the double-lasso method of Belloni et al. (2014) recommended by Urminsky et al. (2016). \*, \*\*, and \*\*\* denote significance at the 10%, 5%, and 1% levels, respectively. Observations are reweighted to match Twitter users from the ATP on observables.

**Figure C.14:** CDF of the WTA to stop using social media

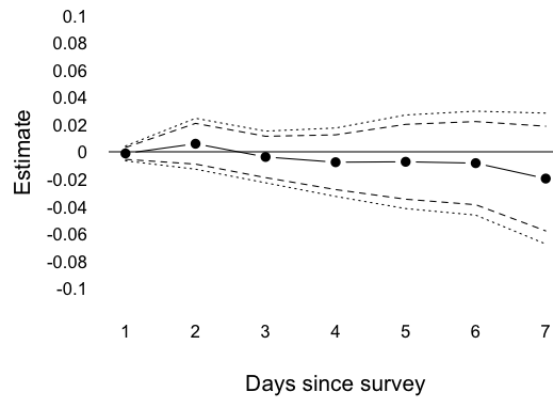
Notes: This figure displays the CDF of the WTA to stop using social media during one week, by treatment arm. Observations are reweighted to match Twitter users from the ATP on observables.



**Table C.15:** Effects of information on other measures of activity

<i>Panel A: Tweets and Likes</i>						
	Tweets			Likes		
	(1)	(2)	(3)	(4)	(5)	(6)
High moderation	3.488 (3.733)	3.359 (3.621)	1.567 (2.623)	13.480* (7.636)	13.752* (7.356)	8.927 (6.140)
$y$ Mean	9.64	9.64	9.64	27.97	27.97	27.97
$y$ SD	70.91	70.91	70.91	140.25	140.25	140.25
$R^2$	0.00	0.02	0.67	0.00	0.04	0.40
Obs.	1,427	1,427	1,427	1,427	1,427	1,427
<i>Panel B: other activity measures</i>						
	Winsorized time			Fraction of active days		
	(1)	(2)	(3)	(4)	(5)	(6)
High moderation	0.022 (0.017)	0.022 (0.016)	0.014 (0.014)	0.017 (0.027)	0.015 (0.026)	0.014 (0.023)
$y$ Mean	0.07	0.07	0.07	0.28	0.28	0.28
$y$ SD	0.27	0.27	0.27	0.37	0.37	0.37
$R^2$	0.00	0.05	0.34	0.00	0.04	0.21
Obs.	1,427	1,427	1,427	1,427	1,427	1,427
Strata FE	No	Yes	Yes	No	Yes	Yes
Controls	No	No	Yes	No	No	Yes

Notes: This table reports estimates from OLS regressions on treatment assignment. Robust standard errors are parenthesized. Controls are selected using the double-lasso method of Belloni et al. (2014) recommended by Urminsky et al. (2016). \*, \*\*, and \*\*\* denote significance at the 10%, 5%, and 1% levels, respectively. Observations are reweighted to match Twitter users from the ATP on observables.

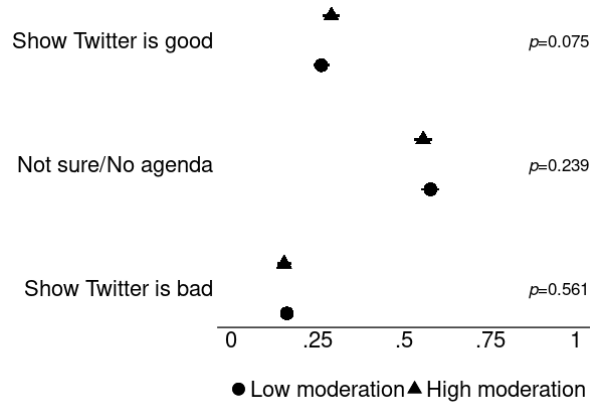
**Figure C.15:** Cumulative dynamic treatment effects on hours spent on Twitter

Notes: This figure presents dynamic treatment effects of hours spent one week after the survey, pointwise confidence intervals (dashed), and sup- $t$  simultaneous confidence bands (dotted).

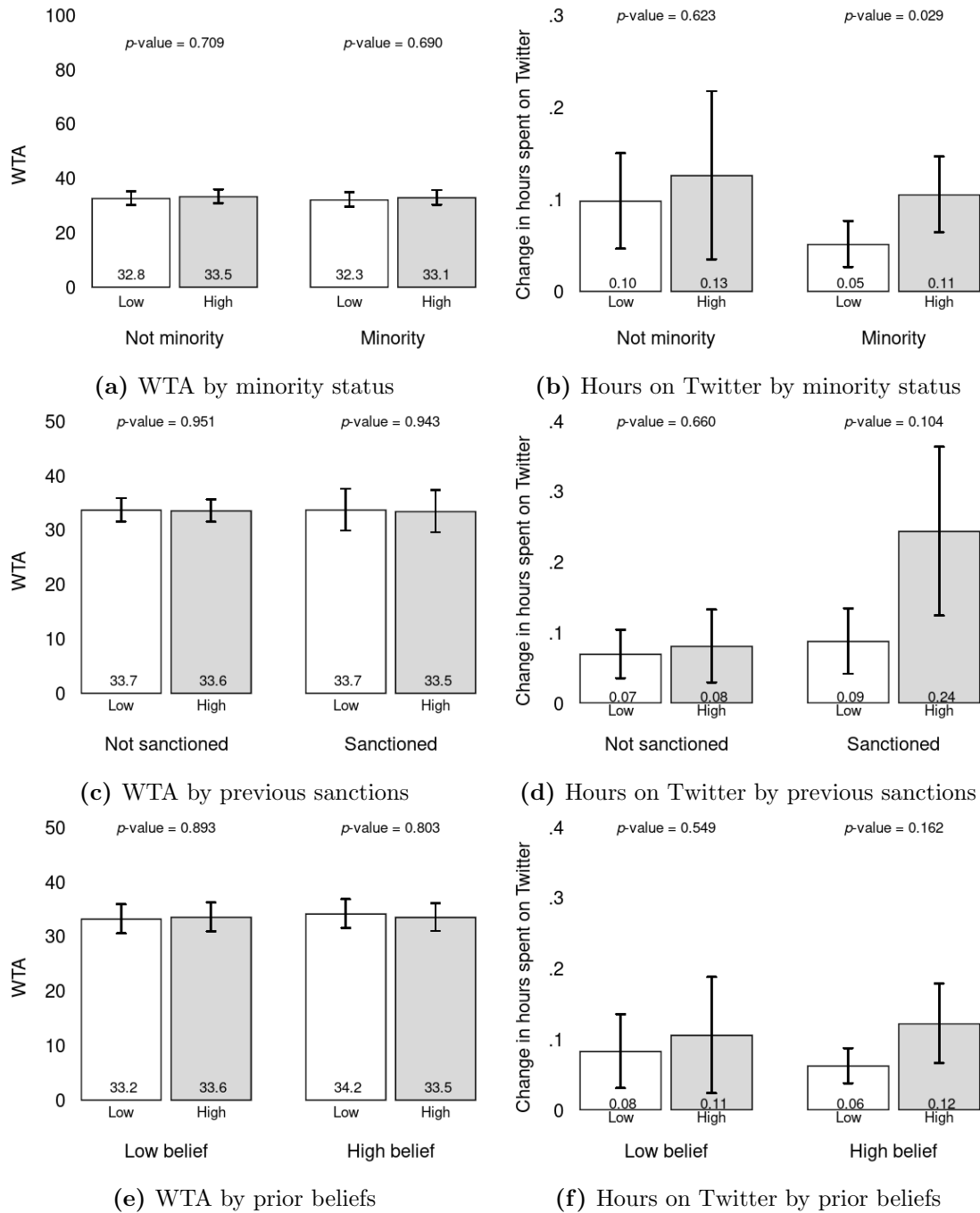
**Table C.16:** Effects of information on inattention and attrition

<i>Panel A: Tweets and Likes</i>						
	Inattention:  recollection-info.			Attrition		
	(1)	(2)	(3)	(4)	(5)	(6)
High moderation	1.252 (0.898)	1.013 (0.856)	1.127 (0.820)	0.004 (0.004)	0.005 (0.004)	0.005 (0.004)
$y$ Mean	8.90	8.90	8.90	0.01	0.01	0.01
$y$ SD	19.82	19.82	19.82	0.10	0.10	0.10
$R^2$	0.00	0.12	0.24	0.00	0.01	0.01
Obs.	2,997	2,997	2,997	3,027	3,027	3,027
Strata FE	No	Yes	Yes	No	Yes	Yes
Controls	No	No	Yes	No	No	Yes

Notes: This table reports estimates from OLS regressions on treatment assignment. Attrition indicates whether participants who finished the prescreening questions did not finish the survey. Robust standard errors are parenthesized. Controls are selected using the double-lasso method of Belloni et al. (2014) recommended by Urminsky et al. (2016). \*, \*\*, and \*\*\* denote significance at the 10%, 5%, and 1% levels, respectively. Observations are reweighted to match Twitter users from the ATP on observables.

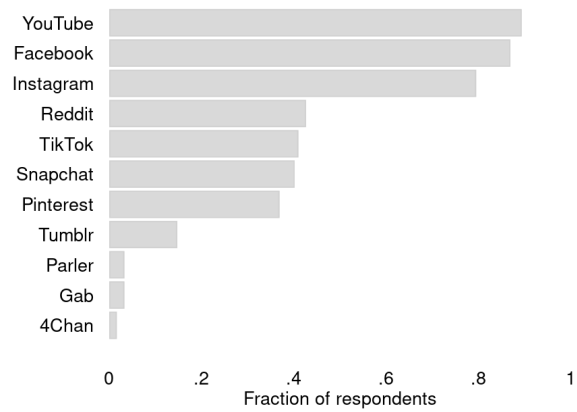
**Figure C.16:** Treatment effect on perceived experimenter's agenda

Notes: This figure presents means and 95% confidence intervals by treatment arm. The dependent variables are answers to the question "Do you think that the researchers in this study had an agenda?". The  $p$ -values come from independent OLS regressions.



**Figure C.17:** Heterogeneity of WTA and hours on Twitter by minority status, previous sanctions, and priors

Notes: These figures present means and 95% confidence intervals by treatment arm and minority status. The  $p$ -values come from OLS regressions. Observations are reweighted to match Twitter users from the ATP on observables.



**Figure C.18:** Other platforms frequented by Twitter users

Notes: This figure presents the fraction of respondents who use other platforms besides Twitter.

**Table C.17:** Effects of information on WTA and time spent on Twitter

<i>Panel A: Weighted (Twitter ATP)</i>						
	WTA			Time spent		
	(1)	(2)	(3)	(4)	(5)	(6)
High moderation	-0.3061 (2.234)	0.1035 (2.056)	0.1035 (2.056)	0.0542 (0.033)	0.0576* (0.034)	0.0576* (0.034)
$y$ Mean	33.57	33.57	33.57	0.10	0.10	0.10
$y$ Std. Dev.	36.75	36.75	36.75	0.57	0.57	0.57
$R^2$	0.00	0.03	0.03	0.00	0.03	0.03
<i>Panel B: Weighted (Social Media ATP)</i>						
	WTA			Time spent		
	(1)	(2)	(3)	(4)	(5)	(6)
High moderation	-0.0796 (2.126)	-0.0746 (2.093)	-0.0746 (2.093)	0.0400 (0.042)	0.0334 (0.036)	0.0334 (0.036)
$y$ Mean	34.98	34.98	34.98	0.10	0.10	0.10
$y$ Std. Dev.	37.26	37.26	37.26	0.61	0.61	0.61
$R^2$	0.00	0.02	0.02	0.00	0.05	0.05
N	2998.00	2998.00	2998.00	1427.00	1427.00	1427.00
<i>Panel C: Unweighted</i>						
	WTA			Time spent		
	(1)	(2)	(3)	(4)	(5)	(6)
High moderation	0.7230 (1.328)	0.7241 (1.320)	0.7241 (1.320)	0.0456 (0.036)	0.0461 (0.036)	0.0461 (0.036)
$y$ Mean	32.94	32.94	32.94	0.10	0.10	0.10
$y$ Std. Dev.	36.35	36.35	36.35	0.68	0.68	0.68
$R^2$	0.00	0.02	0.02	0.00	0.01	0.01
Observations	2,998	2,998	2,998	1,427	1,427	1,427
Strata FE	No	Yes	Yes	No	Yes	Yes
Controls	No	No	Yes	No	No	Yes

Notes: This table reports estimates from OLS regressions of the WTA and the change in time spent on Twitter on a treatment indicator. Panel A reweights observations to match a representative sample of Twitter users on observables. Panel B reweights observations to match a representative sample of social-media users on observables. Panel C includes unweighted estimates. Robust standard errors are parenthesized. Controls are selected using the double-lasso method of Belloni et al. (2014) recommended by Urminsky et al. (2016).

## D Survey Instruments

### D.1 Classification of random posts

By accepting this HIT, you confirm that you are at least 18 years old, have read and understood this [consent form](#) and are willing to participate in this classification exercise. Identifying information will **not** be shared (your MTurk ID will be replaced with an arbitrary alphanumeric code).

[Instructions](#) [Shortcuts](#) Please classify the following post:

Text of the Tweet

**Select an option**

Hate speech	1
Offensive but not hate speech	2
Neither offensive nor hate speech	3

(a) Task screen

**Hate speech classification instructions** ×

**Definition**

For this task, hate speech is defined as posts that would be censored in social media platforms. Please read [Twitter's definition](#) for clarification and some examples.

**Bonus payment**

I will give a bonus of \$20 to the 5 most accurate workers, among those who complete at least 100 HITs. Performance will be measured comparing responses to other workers' responses.

**Rejections**

I included some attention check posts. They will be easy to identify as long as you are reading the posts. Failing these attention checks will result in rejecting your HITs.

Close

(b) Instructions

**Figure D.1:** MTurk task to classify posts as hate speech

## D.2 Welfare survey

**Welcome and thank you for participating in this social media survey!**

Which of the following **social media** platforms did you use in the past month? Please select all that apply

Facebook	<input type="checkbox"/>
Instagram	<input type="checkbox"/>
Twitter	<input type="checkbox"/>
Snapchat	<input type="checkbox"/>
YouTube	<input type="checkbox"/>
TikTok	<input type="checkbox"/>
Gab	<input type="checkbox"/>
Parler	<input type="checkbox"/>
Reddit	<input type="checkbox"/>
Pinterest	<input type="checkbox"/>
Tumblr	<input type="checkbox"/>
4Chan	<input type="checkbox"/>
None	<input type="checkbox"/>

The next question is about your interest in sports. In reality, this is an attention check. **If you are reading carefully, please select "A little bit interested" and "Not at all interested"**. How interested are you in sports?

Extremely interested	<input type="checkbox"/>
Very interested	<input type="checkbox"/>
A little bit interested	<input type="checkbox"/>
Almost not interested	<input type="checkbox"/>
Not at all interested	<input type="checkbox"/>

After this survey, we will invite some respondents to a **follow-up study**. Below we will describe it in detail

We will also send **\$50 Amazon gift card bonuses** to some participants depending on survey answers and **attention checks**

To proceed with this survey, please confirm if you are willing to **provide a valid email** address. We will use it for the follow-up and bonuses

Yes	<input checked="" type="radio"/>
No: you exit this survey on the next screen	<input type="radio"/>

Thank you. Please enter a valid email address:

example@email.com
-------------------

Now some **demographic** questions

What is the highest level of **education** you have completed?

Less than high school	<input type="radio"/>
High school graduate	<input type="radio"/>
Some college but no degree	<input type="radio"/>
Associate's degree	<input type="radio"/>
Bachelor's degree or some postgraduate	<input type="radio"/>
Graduate degree (for example: MA, MBA, JD, PhD)	<input type="radio"/>



What was your 2020 annual household **income** before taxes?

Less than \$10,000	<input type="radio"/>
\$10,000 - \$19,999	<input type="radio"/>
\$20,000 - \$29,999	<input type="radio"/>
\$30,000 - \$39,999	<input type="radio"/>
\$40,000 - \$49,999	<input type="radio"/>
\$50,000 - \$59,999	<input type="radio"/>
\$60,000 - \$69,999	<input type="radio"/>
\$70,000 - \$79,999	<input type="radio"/>
\$80,000 - \$89,999	<input type="radio"/>
\$90,000 - \$99,999	<input type="radio"/>
\$100,000 - \$149,999	<input type="radio"/>
More than \$150,000	<input type="radio"/>
Prefer not to say	<input type="radio"/>

What is your sexual **orientation**?

Heterosexual	<input type="radio"/>
Homosexual	<input type="radio"/>
Bisexual	<input type="radio"/>
Other	<input type="radio"/>
Prefer not to say	<input type="radio"/>

What is your present **religion**, if any?

Christian	<input type="radio"/>
Jewish	<input type="radio"/>
Muslim	<input type="radio"/>
Buddhist	<input type="radio"/>
Hindu	<input type="radio"/>
Atheist	<input type="radio"/>
Agnostic	<input type="radio"/>
Other	<input type="radio"/>
Nothing in particular	<input type="radio"/>
Prefer not to say	<input type="radio"/>

Now some questions about your **social media** use

What is your **Twitter handle / user name**? Optional: you may leave it blank

If you choose to share it, we will only use it to get the account creation date and number of tweets and likes. This data improves the quality of the study

Please include the "@"

@username
-----------

How many days did you **use Twitter** last week?

0	<input type="radio"/>
1	<input type="radio"/>
2	<input type="radio"/>
3	<input type="radio"/>
4	<input type="radio"/>
5	<input type="radio"/>
6	<input type="radio"/>
7	<input type="radio"/>

On an average day that you used Twitter, how much **time** did you spend on it?

Less than 30 minutes	<input type="radio"/>
30 minutes to 1 hour	<input type="radio"/>
1 to 2 hours	<input type="radio"/>
2 to 3 hours	<input type="radio"/>
3 to 4 hours	<input type="radio"/>
4 to 6 hours	<input type="radio"/>
6 to 10 hours	<input type="radio"/>
More than 10 hours	<input type="radio"/>

Have you received one or more of the following **sanctions** on Twitter?

	Yes	No	Don't know
Account suspended permanently / de-platformed	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Account locked / suspended temporarily	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Tweet removed	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Shadow banned / limited visibility	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

Have **you reported** Tweets or accounts for violating the Twitter Rules?

Yes

No

Have you received a notification that someone **reported you** for violating the Twitter Rules?

Yes

No

Don't know

Have you been **harassed, threatened or attacked** online?

Yes

No

Out of every 100 Tweets that you read, how many would you say are **hate speech**—not just offensive?

0

10

20

30

40

50

60

70

80

90

100

Slide to select your guess

A crowd-sourced team of annotators helped us identify hate speech using 10,000 random Tweets in August 2020

We call a Tweet "hate speech" if most annotators label its text as hateful—not just offensive

**What percentage of Tweets do you think were classified as hate speech?**

We will send a **\$50** Amazon gift card to one random respondent among those with the **closest guess**

0 10 20 30 40 50 60 70 80 90 100

Slide to select your guess



How **confident** are you about your guess?

Very confident ☐

Somewhat confident ☐

Neutral ☐

Not too confident ☐

Not at all confident ☐

The next question is about the **removal, banning, suspension, or de-platforming of hate speech**

We checked if Twitter removed (de-platformed) the hateful Tweets in the sample or the accounts that posted them

What percentage of hateful Tweets or accounts do you think  
Twitter **removed (de-platformed)** within 1 month?

We will send a **\$50** Amazon gift card to one random respondent  
among those with the **closest guess**


0 10 20 30 40 50 60 70 80 90 100  
Slide to select your guess



How **confident** are you about your guess?

Very confident	<input type="radio"/>
Somewhat confident	<input type="radio"/>
Neutral	<input type="radio"/>
Not too confident	<input type="radio"/>
Not at all confident	<input type="radio"/>

Twitter **removed (de-platformed)** **3.6%** of hate speech Tweets or the  
accounts that posted them, within 1 month

**Less than 1%** of Tweets in the sample were classified as hate speech.  
Other popular platforms (Youtube, Facebook, and Reddit) have a  
similar prevalence of hate 

We will conduct a small study that **compensates** some participants to  
**stop using Twitter, Facebook, Instagram, YouTube, Snapchat,**  
**TikTok, and Reddit for 1 week**

Similar studies have been conducted in the past (Hunt et al. 2018,  
Mosquera et al., 2019, or Allcott et al., 2020)

To establish your compensation we ask below if you want to stop  
using social media **for different Amazon gift card amounts**

A computer will randomly choose some eligible participants whom we  
will contact for the follow-up

If you are selected, the computer will also choose one of your answers  
below at random

- If your answer is "yes", we will **ask you to stop using social media one week and pay you the offered amount**
- If your answer is "no", we will **not ask you to stop using social media**

**Each question could be "the one that counts" to determine your compensation**

Would you accept a **\$50** compensation to stop using Twitter, Facebook, Instagram, YouTube, Snapchat, TikTok, and Reddit for 1 week?

Yes	<input type="radio"/>
No	<input type="radio"/>

Thank you. 4 final questions:

What was the fraction of hate that Twitter **removes (de-platforms)** that **we told you**? Please give the closest number

We will send a **\$50** Amazon gift card to one random respondent among those with the **closest guess**

0 10 20 30 40 50 60 70 80 90 100  
Slide to select



Do you think the researchers in this study had an **agenda**?

I don't think they had a particular agenda	<input type="radio"/>
Yes, they wanted to show that Twitter is good	<input type="radio"/>
Yes, they wanted to show that Twitter is bad	<input type="radio"/>
Not sure	<input type="radio"/>

What do you think about Twitter's **removal of Tweets or de-platforming of users**?

Twitter removes too much	<input type="radio"/>
Neutral / Twitter removes enough	<input type="radio"/>
Twitter removes too little	<input type="radio"/>

What percentage of hate speech posts or accounts would you guess **Facebook removes (de-platforms)** within one month of posting?

0 10 20 30 40 50 60 70 80 90 100  
Slide to select

