

# Justifying Dissent\*

Leonardo Bursztyn<sup>†</sup>    Georgy Egorov<sup>‡</sup>    Ingar Haaland<sup>§</sup>  
Aakaash Rao<sup>¶</sup>    Christopher Roth<sup>||</sup>

July 2022

## Abstract

Dissent plays an important role in any society, but dissenters are often silenced through social sanctions. Beyond their persuasive effects, rationales providing arguments supporting dissenters' causes can increase the public expression of dissent by providing a "social cover" for voicing otherwise-stigmatized positions. Motivated by a simple theoretical framework, we experimentally show that liberals are more willing to post a Tweet opposing the movement to defund the police, are seen as less prejudiced, and face lower social sanctions when their Tweet implies they had first read credible scientific evidence supporting their position. Analogous experiments with conservatives demonstrate that the same mechanisms facilitate anti-immigrant expression. Our findings highlight both the power of rationales and their limitations in enabling dissent and shed light on phenomena such as social movements, political correctness, propaganda, and anti-minority behavior.

**Keywords:** Dissent; rationales; social image; social media

**JEL Classification:** D83, D91, P16, J15

---

\*We thank Davide Cantoni, Daniel Gottlieb, Ro'ee Levy, Pietro Ortoleva, Andrei Shleifer, Marco Tabellini, David Yang, Noam Yuchtman, and numerous seminar participants for very helpful suggestions. We thank Stelios Michalopoulos for a highly constructive discussion. We thank Danil Fedchenko, Takuma Habu, Hrishikesh Iyengar, Melisa Kurtis, and Stan Xie for outstanding research assistance. We gratefully acknowledge financial support from the Pearson Institute for the Study and Resolution of Global Conflicts and the UChicago Social Sciences Research Center. Roth acknowledges funding from the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany's Excellence Strategy EXC 2126/1-390838866. The research described in this article was approved by the University of Chicago Social and Behavioral Sciences Institutional Review Board and the Humanities and Social Sciences Research Ethics Committee at the University of Warwick.

<sup>†</sup>University of Chicago and NBER, [bursztyn@uchicago.edu](mailto:bursztyn@uchicago.edu)

<sup>‡</sup>Kellogg School of Management and NBER, [g-egorov@kellogg.northwestern.edu](mailto:g-egorov@kellogg.northwestern.edu)

<sup>§</sup>NHH Norwegian School of Economics and CESifo, [ingar.haaland@nhh.no](mailto:ingar.haaland@nhh.no)

<sup>¶</sup>Harvard University, [arao@g.harvard.edu](mailto:arao@g.harvard.edu)

<sup>||</sup>University of Cologne and CEPR, [roth@wiso.uni-koeln.de](mailto:roth@wiso.uni-koeln.de)

# 1 Introduction

From speaking out against injustice to victimizing protected groups, dissent can be a force for or against social change and therefore plays a consequential role in any society. Fundamental to dissent are *rationales* — narratives disseminated by political entrepreneurs, social movements, and media outlets — that provide arguments supporting dissenters’ causes. Some rationales spur dissent through persuasion: they change people’s views and, as a result, their public behavior. Yet dissent is often limited not because few people hold dissenting opinions, but rather because these people fear speaking their mind. Indeed, 62 percent of Americans agree that “The political climate these days prevents me from saying things I believe because others might find them offensive” (Ekins, 2020).

Consider Democrats who oppose the movement to defund the police. In many settings, publicly expressing this opposition generates social costs: opposition to police defunding may be seen as a signal of racial intolerance either by a majority or by a small but vocal minority. Suppose that a credible study is publicized suggesting that defunding the police would increase violent crime. This new study might increase an individual’s willingness to publicly oppose police defunding even if the study does not change her convictions, as long as she is able to *attribute* her views to the study. The key point is that the availability of this rationale opens up explanations other than racial intolerance for her position, reducing the social costs incurred by voicing it publicly and thus making her more willing to dissent.

In this paper, we present experiments exploring the power and potential limitations of rationales in facilitating the expression of dissent. Motivated by a simple theoretical framework, we experimentally examine the expression and interpretation of dissent in two contentious and policy-relevant domains: liberals’ opposition to defunding the police and conservatives’ support for deporting illegal immigrants. We focus on social media, where rationales from both mainstream and fringe sources proliferate and where people often face large social costs of expressing controversial opinions.

We begin by studying opposition to police reform among liberals. In a first experiment, respondents read a Washington Post article written by a Princeton criminologist arguing that “One of the most robust, most uncomfortable findings in criminology is that putting more officers on the street leads to less violent crime”.<sup>1</sup> Respondents then choose whether to join a campaign opposing the movement to defund the police and, conditional on doing so, decide whether to post a Tweet promoting the campaign. The experimental

---

<sup>1</sup>See “Why do we need the police?” Sharkey, Patrick. *The Washington Post*, June 12, 2020.

manipulation subtly varies the availability of a social cover in the Tweet while holding fixed other potential motives to post. In particular, in the *Cover* condition, respondents' Tweets indicate that they were shown the article *before* joining the campaign, while in the *No Cover* condition, respondents' Tweets indicate that they were shown the rationale *after* joining the campaign.<sup>2</sup> The implied timing in the *Cover* condition provides these respondents with a social cover — the (implicit) justification that they joined the campaign because they were persuaded by the article's claims — while the timing implied by the *No Cover* condition eliminates this social cover. Differences in the “willingness to Tweet” thus cannot be explained by the persuasiveness of the rationale — all respondents in both groups read the article — or by respondents' expectations that the rationale will persuade their followers — both versions of the Tweet contain an identical description of and link to the article.

The availability of a social cover strongly affects posting behavior: respondents are 12 percentage points more likely to post the Tweet in the *Cover* condition than in the *No Cover* condition. In a placebo experiment with an identical design, but with a Tweet expressing support for a non-stigmatized cause, we find no difference between posting rates in the *Cover* and *No Cover* conditions, suggesting effects are indeed driven by (anticipated) changes in the stigma associated with dissenting expression rather than some other independent effect of the treatment. An additional experiment in which respondents describe the considerations on their mind when posting potentially controversial content corroborates the importance of the social cover effect of rationales. However, lowering the credibility of the rationale by removing the references to the author's academic credentials and to the scientific evidence underlying the article's claims strongly reduces the treatment effects, highlighting the limits of rationales in facilitating dissent.

We conduct a second experiment, again with liberal respondents, to examine how the social cover shifts an audience's inferences about the motives underlying dissent and the resulting sanctions levied upon dissenters. Respondents are matched with a participant who posted the Tweet from the previous experiment — either a previous participant assigned to the *No Cover* condition or to the *Cover* condition — and are shown the anti-defunding Tweet their matched participant chose to post. They choose whether to deny a bonus to their matched participant, a measure of social sanctions. We also elicit respondents' inferences about their matched participant's underlying prejudice: respondents guess whether

---

<sup>2</sup>Both Tweets are factually correct, as respondents in both conditions were shown the article both before and after joining the campaign.

or not the participant authorized a donation to a pro-Black organization.

The results confirm that the availability of social cover shifts inference and resulting social sanctions. Respondents matched with a participant in the *Cover* condition are 7 percentage points more likely to think that their matched participant authorized the pro-Black donation (relative to a *No Cover* mean of 27 percent) and are 7 percentage points less likely to deny their matched participant the \$1 bonus (relative to a *No Cover* mean of 47 percent). As in the first experiment, slightly lowering the credibility of the rationale dramatically reduces these estimated treatment effects.

We next study the effects of rationales among a different sample, conservatives, and in a different policy context, anti-immigrant policies. Here, supporting the immediate deportation of all illegal immigrants from Mexico is a stigmatized opinion that people may be reluctant to publicly express, but a similar rationale as studied in the previous experiments — concerns about crime — may be effective in shifting inference about motives and thus decreasing social sanctions. In addition to speaking to the robustness of our previous findings and examining the use of rationales by a different population (conservative rather than liberal respondents), these experiments allow us to examine how rationales can generate social cover vis-a-vis different types of audience. In particular, opposition to police defunding is primarily stigmatized by liberals’ in-group (fellow liberals) rather than their out-group (conservative); in contrast, support for deportation is primarily stigmatized by conservatives’ out-group (liberals) rather than their in-group (fellow conservatives).

The experimental manipulation follows the logic in our first experiment: in the *Cover* condition, respondents’ Tweets indicate that they were exposed to a rationale — a clip of Fox News anchor Tucker Carlson arguing that illegal immigrants commit violent crimes at vastly higher rates than citizens — *before* joining the campaign, while in the *No Cover* condition, respondents’ Tweets indicate that they were exposed to the rationale *after* joining the campaign. Our findings corroborate the importance of rationales in facilitating the expression of dissent: respondents are 17 percentage points more likely to post the Tweet in the *Cover* condition than the *No Cover* condition, relative to a *No Cover* mean of 47 percent. A further experiment shows that this rationale once again has strong effects on inference: respondents matched with a participant who chose to post the *Cover* Tweet are 5 percentage points more likely to believe that this participant authorized the pro-immigrant donation (relative to a *No Cover* mean of 9 percent) and are 7 percentage points less likely to deny their matched participant the bonus (relative to a *No Cover* mean of 80 percent).

Taken together, our evidence highlights the importance of rationales in facilitating dis-



sent on both sides of the political spectrum; and it sheds light on the mechanisms by which individuals and institutions can influence public behavior by shaping the supply of rationales and perceptions of their social acceptability. Our findings have important implications for how the expression of dissent responds to the availability of new narratives. First, rationales are only effective to the extent to which observers believe that they genuinely change the dissenter’s beliefs: an obscure or non-credible rationale may fail to shift inference, and may even backfire, if it signals the dissenter’s underlying type. For example, if only intolerant people tend to read a particular source, citing a novel rationale provided by this source will fail to generate social cover. This implies that the endorsement of rationales by prominent figures such as politicians or celebrities may generate particularly large “social amplifiers”: such figures may not only be more credible and *directly* persuade more people, but also more able to generate *common knowledge* such that dissenters can claim they were exposed to the rationale without seeking it out directly from stigmatized sources.

Conversely, groups seeking to suppress dissent have strong incentives to silence or marginalize potential sources of rationales (for example, disinviting campus speakers or branding certain news sources as fringe), because these tactics reduce the perceived probability that people will be exposed to rationales “by chance.” If successful, these groups can create and sustain a “political correctness” culture — for better or for worse — in which certain rationales are ineffective because citing the stigmatized source undermines social cover. Indeed, at the time of our experiment, only 25% of Democrats privately supported decreasing police funding Parker and Hurst (2021). By challenging the credibility of rationales or explicitly linking them to stigmatized positions, a vocal group, even a vocal *minority*, can silence a majority.

Our paper contributes to an emerging literature on narratives as powerful drivers of economic and political behavior (Michalopoulos and Xue, 2021; Shiller, 2017). Related to work is Foerster and van der Weele (2021), which studies the communication of rationales for and against donating to prosocial causes, and Bénabou et al. (2020), which models the production and circulation of justifications for morally questionable actions. Our contribution to this literature is to characterize and experimentally identify an important channel — the “social cover” effect — through which narratives, or rationales, shape the expression and the interpretation of dissent. Our theoretical framework and experimental evidence suggest means by which individuals and institutions can exploit this channel to facilitate or suppress dissent.

Thus, our work relates to a literature examining how social norms influence public behavior (Kuran, 1997; Bénabou and Tirole, 2006; Ali and Lin, 2013; Lacetera and Macis, 2010; Perez-Truglia and Cruces, 2017), and to a theoretical literature on political correctness (Morris, 2001; Golman, 2021). Braghieri (2022) shows that publicly expressed views, which may be affected by political correctness norms, are not fully informative of private views. Like some of this previous work (Bursztyn et al., 2020a,b), our paper examines how previously-stigmatized public behavior can become socially acceptable, but it differs conceptually and in its implications for equilibrium expression. Conceptually, we show that rationales make public actions less informative about dissenters’ underlying type and increase the public expression of dissent by lowering its social cost. This enables moderates who previously would have been unwilling to express dissent for fear of being labeled an extremist to voice their opinions, further hindering inference about dissenters’ underlying type. In other words, our mechanism generates a “social amplifier” that magnifies rationales’ persuasive effects. We discuss how political entrepreneurs can strategically supply rationales to make the expression of unpopular views more mainstream.

This latter channel helps explain the mechanisms by which media and propaganda can promote socially undesirable behavior, such as anti-minority violence (e.g. Yanagizawa-Drott 2014; Adena et al. 2015; Enikolopov and Petrova 2015). Studies in this vein examining persuasion in field settings often find substantial effects (e.g. Caprettini et al. 2021) — in contrast to the relatively small effects of persuasion typically documented in a vast literature using information provision experiments (Haaland et al., 2021)). Among other plausible explanations for this discrepancy is the “social amplifier” channel: widespread propaganda creates common knowledge of rationales, generating greater social cover and magnifying the effect of rationales on public behavior. Thus, our work also connects to a literature on populist political movements (e.g. Acemoglu et al. 2013; Guriev and Papaioannou 2020; Patir et al. 2021) insofar as authoritarian populists are often highly skilled at producing and disseminating rationales normalizing the victimization of minority groups.

Finally, our paper relates to a lab experimental literature documenting that individuals seize upon even flimsy (self)-excuses for selfish behavior.<sup>3</sup> These findings can be understood through a behavioral model of self-signaling, as in Bénabou and Tirole (2011); similarly, Grossman and Van Der Weele (2017) formalize a mechanism by which individuals engage in willful ignorance as an excuse for selfish behavior. Our work holds this channel constant

---

<sup>3</sup>See, for example, Dana et al. (2007); Hamman et al. (2010); Cunningham and de Quidt (2015); Lazear et al. (2012); Exley (2016); Golman et al. (2017); Saccardo and Serra-Garcia (2020).

— all individuals in our experiments privately voice their agreement with the Tweet — and we instead examine signaling vis-a-vis *others*. Moreover, our work highlights the importance of the credibility of rationales: unlike in “self-excuse” experiments and in the classic “Xerox” experiment of Langer et al. (1978), our framework predicts, and our experiments demonstrate, that only credible rationales are effective in facilitating expression and shifting inference. Individuals and institutions can thus manipulate this credibility channel to enable or suppress dissent.

The remainder of this paper proceeds as follows. In Section 2, we present a simple model of the use and interpretation of rationales facilitating dissenting expression. In Section 3, we present experiments studying how the availability of a social cover shapes liberal respondents’ willingness to publicly oppose the movement to defund the police, and how this social cover shifts their audience’s beliefs about and behavior toward them. In Section 4, we present similar experiments focusing on conservative respondents in the context of anti-immigrant expression. Section 5 discusses implications of our findings and concludes. We list all main and auxiliary experiments in Appendix Table B.1.

## 2 Theoretical Framework

To organize these ideas and guide the experimental design, we start with a theoretical framework. All formal proofs are provided in Appendix A.

### 2.1 Setup

The society  $A$  consists of a continuum of citizens facing a binary policy decision between the status quo ( $Q$ ) and change ( $C$ ). There is some objective measure of social welfare from decision  $C$ , and we denote this value  $w$ . The welfare under the status quo  $Q$  is normalized to zero. From the citizens’ perspective, this value is distributed normally:  $w \sim \mathcal{N}(w_0, \sigma_w^2)$ . This social welfare may incorporate the expected economic payoff to each citizen from enacting decision  $C$ , but it may also include externalities to people outside the society or other factors inasmuch as citizens care about them.

Apart from the objective economic consequences captured by  $w$ , citizens have idiosyncratic tastes. Specifically, citizen  $i$  gets additional utility  $t_i$  if policy  $C$ , as opposed to  $Q$ , is enacted; we refer to  $t_i$  as  $i$ ’s type. We assume that  $t_i$  is distributed with c.d.f.  $H(\cdot)$  and p.d.f.  $h(\cdot)$ , and that it satisfies the monotone hazard rate property ( $\frac{h(x)}{1-H(x)}$  is increasing

in  $x$ , which is satisfied, e.g., for the normal and uniform distributions). To avoid corner cases, we assume that  $t_i$  has full support on the real line.

A citizen  $i \in A$  is given a chance to publicly state support for change (decision  $d_i = 1$ ) before an audience. Doing so results in expressive benefit  $B$  but social cost  $S$ , so  $U_i(d_i = 1) = B - S$ . We assume that

$$B = \beta (\mathbb{E}(w \mid *) + t_i);$$

in other words, the benefit is proportional to the sum of citizen  $i$ 's posterior belief about  $w$  using all available information and  $i$ 's own type. The social cost  $S$  is borne because action  $d_i = 1$  may be revealing about  $i$ 's type  $t_i$ , and having a high type is stigmatized by the audience. For simplicity, we assume that stigma is linear in the audience's posterior about citizen  $i$ 's type:

$$S = \gamma \mathbb{E}_{-i}(t_i \mid d_i = 1, *).$$

Lastly, the utility from inaction ( $d_i = 0$ ) is normalized to 0:  $U_i(d_i = 0) = 0$ .<sup>4</sup>

## 2.2 Analysis

In the absence of new information, the posterior of citizen  $i$  about  $w$  equals the prior  $w_0$ , and thus the benefit of action  $d_i = 1$  is  $B = \beta(w_0 + t_i)$ . Citizen  $i$  makes the decision holding his social cost  $S$  fixed. Therefore, he chooses  $d_i = 1$  if and only if

$$t_i \geq \frac{1}{\beta} S - w_0.$$

Thus, any equilibrium takes the threshold form, with the threshold  $\tau$  satisfying the condition

$$\tau = \frac{\gamma}{\beta} \mathbb{E}(t_i \mid t_i > \tau) - w_0. \quad (1)$$

Generally speaking, the threshold need not be unique due to strategic complementarity: if not only extreme right but also moderate types choose action  $d_i = 1$ , the social cost is lower, which increases citizens' propensity to choose  $d_i = 1$ . However, if the distribution of  $t_i$  satisfies the monotone hazard rate property, the equilibrium is unique.

---

<sup>4</sup>We implicitly assume that the audience does not observe that  $i$  had a chance to make the action, and thus if he chooses  $d_i = 0$  he is pooled with a continuum of citizens who are passive in this model. If the audience observes that inaction is by choice, there may be social consequences in this case as well. Nevertheless, all the results go through as stated.

**Proposition 1.** *Suppose that  $\gamma < \beta$ . Then there is a unique equilibrium that takes the form of a threshold: individuals with  $t_i > \tau$  choose  $d_i = 1$  and those with  $t_i < \tau$  choose  $d_i = 0$ .*

In other words, the equilibrium is unique provided that the citizen's choice is not driven solely by social image concerns and that the expressive benefit from their choice is sufficiently high.

## 2.3 Persuasive Rationales

Suppose that citizen  $i$ , prior to making the decision, received an informative signal  $s = w + \varepsilon$ , where  $\varepsilon \sim \mathcal{N}(0, \sigma_\varepsilon^2)$ . His posterior expectation about  $w$  is then equal to

$$w_1 = \mathbb{E}(w \mid s) = w_0 \frac{\sigma_\varepsilon^2}{\sigma_w^2 + \sigma_\varepsilon^2} + s \frac{\sigma_w^2}{\sigma_w^2 + \sigma_\varepsilon^2},$$

which exceeds  $w_0$  if and only if  $s > w_0$ . Now, if indeed the signal is positive ( $s > w_0$ ), then for a fixed social cost  $S$ , this would prompt more citizens to choose  $d_i = 1$  (specifically, all citizens with  $t_i \geq \frac{1}{\beta}S - w_1$  would do so). This corresponds to a *persuasion* mechanism. Now that more moderate people choose  $d_i = 1$ , the social cost of doing so is lower: intuitively, publicly supporting  $C$  is no longer a sign of extremism. Of course, a decrease in  $S$  will prompt even more people to choose  $d_i = 1$  (a “social amplifier”). In the end, we have the following characterization of the new equilibrium.

**Proposition 2.** *Suppose that citizen  $i$  makes his decision after receiving informative signal  $s > w_0$ . This citizen then has a higher posterior about  $w$  than the prior, and the ex ante probability that citizen  $i$  chooses  $d_i = 1$  is higher. The equilibrium social cost  $S$  is lower with signal  $s$  than without. An increase in  $\sigma_\varepsilon^2$  weakens all these effects.*

The last part of Proposition 2 highlights that all the effects are attenuated if the signal is noisier and therefore less informative. The citizens update less and are less likely to choose  $d_i = 1$ , and the associated social cost does not increase as much either. Practically, this means that if the same information is obtained from a more questionable or less credible source, the changes in behavior and social cost will be smaller, and in the limit, an uninformative signal will have no effect.

## 2.4 Polarizing Rationales

In reality, individuals are often presented with the same evidence, but the evidence has heterogeneous consequences (e.g. some individuals react favorably to news that a neighborhood is diversifying, while others react unfavorably) or is interpreted differently (e.g. due to differences in background knowledge, cognitive limitations, or behavioral biases). Can rationales still be effective even if they are not persuasive *on average* — that is, they “dissuade” as many people as they persuade?

To study this possibility, we assume that share  $\mu$  of citizens get a high signal  $s_h > w_0$  (with the corresponding posterior  $w_h > w_0$ ) and share  $1 - \mu$  get a low signal  $s_l < w_0$  (and their posterior is  $w_l < w_0$ ). We prove the following result.

**Proposition 3.** *Suppose that*

$$\mu (H(\tau) - H(\tau - (w_h - w_0))) \geq (1 - \mu) (H(\tau + (w_0 - w_l)) - H(\tau)), \quad (2)$$

*where  $\tau$  is the equilibrium threshold in the basic model (Proposition 1). Then the ex ante probability that citizen  $i$  chooses  $d_i = 1$  is higher than in the basic model, and the equilibrium social cost is lower.*

In other words, if the mass of people who are persuaded to choose  $d_i = 1$  by high signal  $s_h$  (holding the social cost fixed) is at least as large as the mass of people who are dissuaded from doing so by low signal  $s_l$ , then the social cost of choosing  $d_i = 1$  goes down in equilibrium, and more people do so in equilibrium. Intuitively, the audience now faces the inference problem: citizen  $i$  may have chosen  $d_i = 1$  either because  $t_i$  is high, or because he got a high signal  $s_h$ . More precisely, the set of citizens who would choose to support  $S$  now contains some types with  $t_i < \tau$  (moderates who got a high signal  $s_h$ ) and lacks some types with  $t_i > \tau$  (extremists who got a low signal  $s_l$ ). As long as the share of the former is not too small, the posterior of  $t_i$  conditional on choosing  $d_i = 1$  goes down. As a result, more citizens choose  $d_i = 1$  and face a lower social cost for doing so. This result is not knife-edge: it applies even if somewhat more people are dissuaded.

Taken together, Propositions 2 and 3 imply that while informative and persuasive evidence can reduce the social cost of a stigmatized public action and lead to more people doing it, evidence that dissuades as many people as it persuades can also be effective due to the social inference problem that such evidence creates. Put differently, for a rationale to be effective it does not have to be persuasive, so long as it hinders inference about the

motives for a public action.

### 3 Opposition to Defunding the Police

The experiments in this paper examine the expression of dissent on social media. Expression on social media is of direct interest: over 70 percent of Americans report using social media daily, many politicians and other prominent figures have turned to social media as a primary channel of communication with the public, and social media has been linked to a number of important real-world outcomes: protests (Enikolopov et al., 2020), hate crimes (Müller and Schwarz, 2018; Bursztyn et al., 2019), and social movements (Levy and Mattsson, 2021). Second, expressing dissent on social media — like doing so in real-world offline settings, and unlike doing so in more artificial lab settings — may have real social costs vis-a-vis a natural population about whose opinions respondents care — family members, friends, acquaintances, and current and/or future employers. Indeed, a substantial majority of hiring managers report using social media accounts as a screening tool (O’Brien, 2018).

Our first two experiments examine the use and interpretation of rationales for opposing the movement to defund the police. The slogan “defund the police” rose to national prominence after the murder of George Floyd in May 2020; advocates seek to decrease funding for police departments, and many favor restricting the responsibilities of law enforcement primarily to violent crime, redirecting resources to specialized response teams such as social workers and conflict-resolution specialists to deliver other services (Thompson, 2020). Popular opposition to police defunding is relatively high: as of an October 2021 Pew Research survey, only 15 percent of adults, 25 percent of Democrats, and 23 percent of Blacks support reducing spending on policing in their area (Parker and Hurst, 2021). Nonetheless, because the movement is closely linked to concerns about racial injustice — most advocates claim that the American law enforcement system is fundamentally racist and requires radical reform (or abolition) — it seems *a priori* plausible that many liberals would feel uncomfortable publicly voicing opposition to defunding. This is particularly true given that liberal Twitter users are more interested in social justice causes and are more likely to call out perceived injustice than liberals at large (Cohn and Quealy, 2019).

### 3.1 Experiment 1: Rationales and Anti-Defunding Expression

#### 3.1.1 Motivation for experimental design

Experiment 1 studies how the social cover provided by rationales affects respondents’ willingness to post a Tweet on their account opposing the movement to defund the police. Identifying this effect is challenging from both a design and ethical perspective. From a design perspective, we need to manipulate the availability of a social cover, ruling out other possible reasons for why a rationale might change posting behavior. For example, the rationale may affect posting behavior by changing respondents’ private beliefs (persuasion), or respondents might cite the rationale to persuade others (anticipated persuasion). Identifying the cover effect requires us to hold these other channels fixed across experimental conditions. At the same time, we wish to avoid a complicated or heavy-handed intervention in order to maximize the extent to which our results can speak to the expression of dissent in real-world contexts. From an ethical perspective, while we want to examine the most natural possible outcome — respondents’ willingness to Tweet — we prefer to avoid leading respondents to actually post political content on Twitter (a particular concern in our similarly-structured Experiment 3, which studies willingness to publicly support a campaign to deport all illegal Mexican immigrants). A related and conflicting goal is to avoid explicitly deceiving respondents. We address these design and ethical difficulties with an experiment that (1) holds the *persuasion* and *anticipated persuasion* effects constant while varying only the availability of a social cover; (2) measures respondents’ revealed-preference willingness to express dissent on their Twitter account; (3) avoids respondents actually posting these Tweets; and (4) avoids explicit deception.

#### 3.1.2 Sample and experimental design

We conducted our pre-registered Experiment 1 in October 2021 with a sample of 1,122 Democrats and Independents.<sup>5</sup> As explained below, this resulted in a final sample for analysis of 523 respondents. Given the need for respondents to (1) have an active Twitter account and (2) be willing to log into the survey using their Twitter account, as described below, recruiting respondents to participate in this experiment was more difficult than we anticipated. To reach our pre-registered minimum of 500 complete responses, we recruited respondents from both Luc.id and CloudResearch, two survey providers widely used in the

---

<sup>5</sup>Our experiment was pre-registered in the AEA RCT registry under ID AEARCTR-0008432. The full set of experimental instructions is included in Appendix E.1.



social sciences (Litman et al., 2017; Wood and Porter, 2019).<sup>6</sup> Our final sample is well-balanced on observables across treatment arms (Appendix Table B.2). Again to facilitate recruiting our pre-registered minimum number of respondents, we kept this experiment as short as possible; we probe underlying mechanisms in depth in Section 3.1.6.

Figure B.1 outlines the structure of Experiment 1. After completing a short attention check, we ask respondents to log in to our survey using their Twitter account through “Tweatability,” a Twitter application we created using Twitter’s Application Programming Interface (API) that allows us to schedule Tweets to be posted on the users’ accounts at a future date. To an observer, these Tweets look as though they were posted by the respondent him or herself. We automatically capture respondents’ Twitter handles after they log in. Respondents are assured that we will never use this application to access any private information from accounts, that all data will be securely stored until its deletion by no later than December 1, 2021, and that we will never schedule posts on their accounts without their explicit permission. Respondents then respond to a set of basic demographic and other background questions.

We then present respondents with an op-ed written in the Washington Post by Patrick Sharkey, a professor of public affairs and criminology at Princeton University.<sup>7</sup> In the article, Sharkey argues that a vast body of evidence shows that increasing policing decreases violent crime, that defunding the police is thus likely to increase violence, and that other solutions (e.g. granting communities more resources to maintain safety) will likely be more effective. After reading the article, respondents are asked if they would like to join a campaign to oppose the movement to defund the police. The survey terminates for respondents who do not join, leaving us with 529 remaining respondents. These respondents are presented with the article again and informed that they can spend as long as they wish reading it.

Once they continue, we inform respondents that the campaign involves circulating a petition on Twitter opposing the movement to defund the police. We show them a screenshot of the Tweet and ask if they are willing to schedule the Tweet to be posted on their account. We inform respondents that the Tweets of all respondents will be posted if and when we have surveyed people in all US counties (a strategy which, as we explain

---

<sup>6</sup>Our final analysis sample consists of 382 respondents from Lucid and 147 respondents from CloudResearch. The two estimates using the samples individually are very similar in size (12.6 p.p. on CloudResearch vs 11.3 p.p. on Lucid) and statistically indistinguishable.

<sup>7</sup>The article is available at <https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/>.

to respondents, is often used in social media campaigns to make certain topics “trend” on users’ timelines). In practice, because we target fewer respondents than the number of counties in the US, we ensure Tweets will never be posted.

Respondents in the *Cover* condition are asked whether they would like to schedule the following Tweet:

I have joined a campaign to oppose defunding the police: [LINK]. Before joining, I was shown this article written by a Princeton professor on the strong scientific evidence that defunding the police would increase violent crime: [LINK]

The Tweet is identical for respondents in the *No Cover* condition, with one exception: the second sentence begins “**After** I joined the campaign...”. Both Tweets are factually correct (all respondents were in fact shown the article both before and after joining the campaign), but this difference in wording suggests to potential readers of the Tweet that respondents in the *Cover* condition had been exposed to the scientific evidence against defunding the police before joining the campaign — and thus had a strong rationale for doing so. In contrast, the *No Cover* Tweet suggests that respondents had only been exposed to the evidence after joining, and thus that the evidence could not have led them to join the campaign. This design therefore isolates the cover effect of rationales while fixing the persuasion channel (all respondents are exposed to the same information) and the anticipated persuasion channel (all respondents know their Tweet’s readers will be exposed to the article, since it is linked in the Tweet) across conditions. By employing a one-word manipulation, we also hold other potential confounds, such as the length of the Tweet, fixed across conditions.

**Discussion of ethical considerations** Although our experiment avoids explicit deception — all statements subjects see are factually true — our design clearly misleads subjects: they believe that their Tweets might be posted (if we recruit respondents in every US county), when in fact we purposefully recruit fewer respondents than the number of counties such that there is no chance this condition will ever be met. In experimental economics, deceiving or misleading respondents is often considered problematic due to concerns that it will lead subjects to expect deception in future experiments, potentially changing their behavior. Because subjects do not know, and never learn, that we recruited fewer respondents than the the number of US counties, this concern does not apply to our

experiment.<sup>8</sup> More generally, we concluded that the benefits of protecting participants’ privacy and avoiding contributing to a political campaign outweighed the costs of misleading respondents. Moreover, our design ensures that the Twitter *followers* of the respondents in our survey will not be misled by respondents’ Tweets as to whether they read the article before or after joining the campaign — given that these Tweets are never posted. We discuss the ethical considerations underlying all experimental designs in greater detail in Appendix C.

### 3.1.3 Results

Figure 1 displays the results, which we also show in regression table form in Table 1. 57% of respondents authorize the Tweet in the *No Cover* condition compared to 69% of respondents in the *Cover* condition ( $p < 0.01$ ). These effects are stable to the inclusion of demographic and partisan controls; the effect size corresponds to 0.25 standard deviations, comparable to or larger than the effects on persuasion generally documented in information provision experiments (Haaland et al., 2021) and the effects of image concerns generally documented in experiments varying the observability of decisions (Bursztyn and Jensen, 2015).<sup>9</sup> This relatively large effect underscores the importance of the cover effect in driving the expression of dissent.

We next present the results of several experiments designed to rule out potential confounds and shed light on the underlying mechanisms. We summarize these experiments in Table 2.

### 3.1.4 Placebo experiment

One potential concern is that respondents are more willing to schedule the *Cover* Tweet (“Before I joined the campaign...”) than the *No Cover* Tweet (“After I joined the campaign...”) for reasons unrelated to the availability of the social cover. For example, respondents might think the “before” wording in the *Cover* Tweet sounds more natural

---

<sup>8</sup>Even if this concern did apply, it would be less relevant given that we recruit subjects from online survey platforms (which are widely used by psychologists and researchers from adjacent disciplines frequently using deception) rather than experimental economics labs. In Appendix C, we provide direct evidence that our intervention did not change respondents’ subsequent survey behavior.

<sup>9</sup>Indeed, in our pre-registered Auxiliary Experiment 1 with the same rationale, we estimate a persuasion effect on private attitudes of 0.12 standard deviations ( $p=0.059$ ). See Appendix B.1.2 for details, Appendix E.5 for experimental instructions, and Appendix D for balance and representativeness tables for all auxiliary experiments.

than the “after” wording in the *No Cover* Tweet. Alternatively, they may believe that the *No Cover* formulation would mislead their followers as to when they viewed the article, and thus may be more reluctant to post the Tweet.

To address this concern, we run a placebo experiment (Auxiliary Experiment 2)<sup>10</sup> with the same design and manipulation, but in a different, non-stigmatized context — conservation of the Amazon rainforest — and with a different rationale — an article reporting a new study which finds that over 10,000 species are at risk due to deforestation in the Amazon. Panel A of Table 2 shows no significant difference between posting rates in the *Cover* and *No Cover* conditions. The difference in effect sizes between the defunding experiment and the placebo experiment is large in magnitude and significant at the 1% level, suggesting effects are indeed driven by (anticipated) changes in the stigma associated with dissenting expression rather than some other independent effect of the before/after wording.

The placebo results also deliver additional insight into the effect sizes documented in the main experiment. The difference in the fraction of respondents authorizing the post in the *No Cover* treatment, *conditional on privately joining the campaign* — 83% in the placebo experiment, compared to 57% in the main experiment — constitutes suggestive evidence for the existence of (perceived) social sanctions for opposing police defunding and suggests that credible rationales may significantly reduce the extent to which these sanctions prevent the public expression of dissent.

### 3.1.5 Ruling out anticipated persuasion

While implausible, it remains possible that respondents anticipate that the *Cover* Tweet will be more persuasive to followers than the *No Cover* Tweet, and that this difference drives our estimated treatment effects. Alternatively, it could be that respondents believe their followers are more likely to read the article upon seeing one Tweet than the other, or that those who do not read the article themselves (which may constitute the vast majority of those who see the Tweet) will infer that the article is more convincing from one Tweet than the other.

To directly address this concern, we run an auxiliary experiment (Auxiliary Experiment 3) in which we present Democratic and Independent Twitter users with either the *Cover* or *No Cover* Tweet and then ask them to estimate the share of their followers who would join the campaign after seeing their Tweet, a summary statistic for the combined effects

---

<sup>10</sup>See Appendix B.1.3 for details and Appendix E.6 for experimental instructions.

of all channels above.<sup>11</sup> Panel B of Table 2 shows a small and insignificant 1.9 percentage point difference; we can rule out differences of greater than 4.2 percentage points with 95% confidence. This suggests that differences in posting rates are not driven by differences in the anticipated persuasiveness of the Tweets, as respondents’ posting decisions would need to be unrealistically elastic to their beliefs about their audience’s persuadability in order to generate the 12 percentage point treatment effect documented in Experiment 1. We provide further evidence against this mechanism below.

### 3.1.6 Direct evidence on social cover mechanism

We now provide direct evidence that our manipulation varies the perceived availability of social cover, and that this availability is an important consideration on respondents’ minds when considering the expression of dissent. We conduct Auxiliary Experiment 4 with a sample of 402 Democrats with Twitter accounts recruited from Prolific. This broader sample allows us to probe the external validity of our findings. In particular, respondents are not required to grant our “Tweatability” app permissions to schedule posts on their Twitter account, which may induce selection into Experiment 1.

**Experimental design** Respondents begin by reading the article presented in Experiment 1 describing the evidence that defunding the police would increase violent crime. We ask them to imagine that at this stage, they joined a campaign to oppose defunding the police. As in the main experiment, all respondents are then given the chance to read the article again.<sup>12</sup> Then, respondents randomized into the *Cover* condition are asked which of two Tweets they would *hypothetically* prefer to post: the Tweet from the *Cover* condition in Experiment 1, or a *Control* Tweet omitting any reference to a rationale:

I have joined a campaign to oppose defunding the police: [LINK].

Respondents randomized into the *No Cover* condition are instead asked about their hypothetical preference between posting the Tweet from the *No Cover* condition in Experiment 1 or the *Control* Tweet above. After respondents choose their preferred Tweet, we ask them to “Please explain why you chose this Tweet rather than the other Tweet.” Our object of interest is the difference in respondents’ explanations between conditions.

---

<sup>11</sup>See Appendix B.1.4 for details and Appendix E.7 for experimental instructions.

<sup>12</sup>See Appendix E.8 for experimental instructions.

A few comments about the experimental design are in order. First, we separately study preferences for the *Cover* Tweet over the *Control* Tweet and for the *No Cover* Tweet over the *Control* Tweet, rather than directly estimating preferences for the *Cover* Tweet over the *No Cover* Tweet. Our design thus avoids making the “Before/After” distinction between the Tweets salient, better capturing behavior both in our main experiment and in real-world settings and reducing the scope for experimenter demand effects. Similarly, our use of open-ended text to elicit motives, rather than structured questions, avoids priming respondents on particular motivations and better captures what naturally comes to mind when making their choice.

We hand-code open-ended responses across three categories.<sup>13</sup> “Social cover” responses mention that the respondent’s preferred Tweet indicates to followers that the article affected the respondent’s choice to join the campaign.<sup>14</sup> “Anticipated persuasion” responses mention that the article might persuade others.<sup>15</sup> Finally, “Information” responses mention that the article is informative or credible, or that it provides an explanation for why people might want to join the campaign, but do not explicitly relate the information to the respondent’s own views or other people’s views.<sup>16</sup> Many respondents classified as “Information” may have had the “Social cover” or “Anticipated persuasion” mechanisms in mind, but wrote responses that we could not unambiguously classify into either category. We chose a conservative coding scheme for “Social cover” and “Anticipated persuasion” in order to provide a plausible lower bound.

**Results** We begin by analyzing respondents’ preferences over which Tweet to post. 83% of respondents in the *No Cover* condition prefer the Tweet linking to the evidence over the *Control* Tweet without the evidence, compared to 87% of respondents in the *Cover* condition.<sup>17</sup> The high fraction choosing the Tweet with the rationale (whether the *Cover*

---

<sup>13</sup>Our categories themselves are mutually exclusive, but a response might fall under multiple categories if the respondent mentions multiple reasons for their choice. Our two coders were blind to treatment status.

<sup>14</sup>For example, one respondent writes: “I think the evidence provided in the article is an important catalyst in why I would have joined the campaign and without any context that first tweet could be misconstrued, or even cause me to be publicly shamed.”

<sup>15</sup>For instance, one respondent writes: “The tweet is meant to not only inform people of your decision, but to also advertise others to do the same. Having supporting evidence for your cause will increase the chance of others to side and agree with you. Tweet B does this, Tweet A doesn’t.”

<sup>16</sup>For example, one respondent writes: “I would want others to see this article and know that I have some evidence to back my tweet.”

<sup>17</sup>The treatment effect is not comparable with the effect estimated in Experiment 1: for example, we might observe zero treatment effect in this experiment and a strong treatment effect in Experiment 1 if most respondents prefer the *Cover* Tweet to the *No Cover* Tweet, but strongly prefer either Tweet to

or the *No Cover* version) over the *Control* Tweet suggests a widespread preference for citing evidence when engaging in dissenting expression, while the high fraction choosing the *No Cover* version constitutes further evidence that respondents do not avoid the “After” wording due to concerns about it being misleading or unnatural.

We next turn to the open-ended text. The perceived social costs of dissent in this setting are further evidenced by the substantial number of Tweets mentioning some form of social sanctions. A relatively large fraction of respondents (20 percent) explicitly mention the social cover mechanism, three times the number who mention the anticipated persuasion mechanism (7 percent). The majority of responses (53 percent) fall into the “Information” category, though many responses in this category likely meant to convey concerns relating to social cover. Focusing on treatment effects across conditions, reported in Panel C.1 of Table 2, the one-word manipulation indeed induces substantially more respondents to mention social cover (a 10 percentage point difference, or a 67 percent effect relative to the *No Cover* mean). Consistent with the results of Auxiliary Experiment 4, the manipulation appears to have no effect on the probability that respondents mention that their followers will find the article persuasive.

To gauge potential confounds, we also hand-code any responses suggesting potential confounds to our main mechanism of interest: “Unnatural” responses mention that one Tweets seems more unnatural or strangely-worded than another; “Misleading” responses mention that one Tweet seems more misleading or deceptive than another; “Signaling” responses mention that one Tweet suggests that the respondent supports the cause more strongly than the other; and “Experimenter demand” responses mention that the experimenter wants the respondent to choose one Tweet over another, or that the respondents’ followers will believe this is the case. As shown in Panel C.2 of Table 2, almost no Tweets fall into any of these categories.<sup>18</sup>

Together, the placebo experiment, the anticipated persuasion experiment, and this experiment eliciting participants’ reasoning establish that the treatment effects documented in Experiment 1 are indeed driven by differences in the availability of a social cover.

---

the *Control* Tweet (while a minority of respondents exhibit strong preferences for the shorter *Control* Tweet). Nonetheless, it is reassuring that the treatment effect is positive (though statistically insignificant,  $p = 0.311$ ).

<sup>18</sup>The small fraction of respondents who choose the *Control* Tweet without a rationale generally cite its shorter length as the reason for doing so. Given that the one-word manipulation in Experiment 1 holds the length of the Tweet fixed, preferences for shorter or longer Tweets will not affect our results.

### 3.1.7 The role of credibility

In Section 2, we show that for rationales to decrease the social cost of dissent, people must believe that they move at least some people’s opinions. In other words, the credibility of rationales matters: a rationale that is perceived to come from a questionable source, or whose credibility is otherwise undermined, is likely to be less effective. The wording of the Tweet in the main experiment emphasizes the credibility of the rationale, explicitly stating that the author is a Princeton professor and that the article is based on strong scientific evidence; our theory implies that reducing the credibility of the rationale will reduce its effect on posting behavior and increase the associated social sanctions.

We examine the role of credibility with Auxiliary Experiment 5, which investigates the effects of less credible rationales. We also use this experiment to probe another dimension of external validity. In particular, the sample of Experiment 1 consists of respondents who were willing to grant our app permissions to post on their Twitter account, and thus is likely unrepresentative of the population of social media users. To assess the importance of social cover in facilitating dissent among this broader population, we ask respondents whether they would have been willing to publish the post on their account if it was included as a campaign feature. We thus do not require them to grant our app permission to access their accounts.

**Experimental design** Auxiliary Experiment 5 is closely related to the design of Experiment 1. As explained above, all respondents who report actively using Facebook and Twitter are eligible to participate, and they are asked whether they would hypothetically be willing to make the post in question. To probe mechanisms, we also ask an incentivized (post-outcome) question eliciting perceived social sanctions: respondents estimate the share of Democrats who, upon seeing the post, chose to deny the poster a bonus. Finally, and most importantly, we cross-randomize a “credibility” manipulation with our previous manipulation of social cover, resulting in four conditions. In particular, to construct “lower-credibility” versions of the Tweets, we remove the references to Sharkey’s academic credentials and to the scientific evidence underlying the article’s claims. The revised lower credibility Tweets read:

I have joined a campaign to oppose defunding the police: [LINK]. [Before/After]  
joining, I was shown this article arguing that defunding the police would increase violent crime: [LINK]



Our framework predicts that this less credible rationale will generate less social cover and thus will be less effective in facilitating dissent.

**Results** We present results in Panel B of Figure 2 and in Panel D of Table 2. Restricting attention to the higher-credibility version of the post (i.e. the version used in Experiment 1), we find an almost identical treatment effect to that documented in Experiment 1, confirming that our results generalize to the broader sample of social media users. Turning to the lower-credibility version, we find a smaller and statistically insignificant treatment effect.

For perceived social punishment, we find a similar pattern when we instead examine respondents’ guesses as to the number of Democrats who would deny a person who made the post a bonus (our measure of perceived social sanctions): respondents believe that the social cover is effective in reducing social sanctions when the rationale is highly credible. Yet, when the rationale is less credible the effects on perceived social sanctions is smaller and statistically insignificant.

How accurate are respondents’ beliefs about the effects of social cover on actual social punishment? On average, respondents provided with the highly credible rationale expect that the fraction of Democrats engaging in social punishment is 5 percentage points lower than the expected fraction among respondents provided with the less credible rationale, a difference very similar to the 7 percentage point effect of the social cover on actual punishment. Similarly, respondents provided with the less credible rationale expect a reduction in the fraction of Democrats imposing social sanctions by 1 percentage point, virtually identical to the actual effects of the cover on punishment of 2 percentage points. On average, respondents are also fairly well-calibrated about the *levels* of punishment: pooling across all conditions, they expect around half of Democrats to deny the bonus, relative to the actual share of 43%. Thus, our mechanism does not require respondents to over- or under-estimate the share of their audience who would sanction them for expressing dissent, nor does it require this share to be a substantial majority.

**Discussion** These results are particularly striking given the subtle nature of the credibility manipulation. The article — published in the reputable *Washington Post* — remains constant, as does every other aspect of the post. The manipulation arguably generates a fairly modest reduction in credibility: far more modest than, for example, citing a right-leaning outlet or making such a claim without any supporting evidence. Nonetheless, even

this modest reduction in credibility halves the estimated effect of the rationale on posting.

Only 25% of Democrats privately support decreasing funding for police in their area, compared with 34% of Democrats who privately support increasing funding (Parker and Hurst, 2021). Thus, the results of Experiment 1 and Auxiliary Experiment 5 jointly illustrate how public dissent can be silenced by a vocal minority. Over 40% of respondents in the *No Cover* conditions who *privately oppose the movement* are unwilling to schedule a post (or hypothetically make a post) expressing this view. Some of these respondents undoubtedly refrain from posting for reasons other than perceived sanctions: for example, because they dislike posting about social causes in general. We can estimate the fraction of such respondents from the non-stigmatized rainforest setting, in which 20% of respondents who privately agree with the cause choose not to schedule the post. Thus, a reasonable estimate is that half of the 40% of respondents in the *No Cover* conditions who do not post refrain due to anticipated social sanctions. The availability of a highly credible rationale cuts this estimated fraction from 20% to 10%, but a very slightly less credible rationale only cuts the fraction to 16%. To the extent that this phenomenon generalizes, then, it suggests that for politically charged issues, only highly credible rationales may be effective in facilitating liberal dissent — potentially stifling dissent from the “politically correct” position on issues for which a strong scientific consensus does not yet exist.

## 3.2 Experiment 2: Interpretation of Anti-Defunding Rationale

Our theoretical framework implies that rationales lower the social cost of dissent by making the action less informative about type. As documented in Section 3.1, respondents are more willing to dissent when they can draw upon credible rationales because they *expect* such rationales to reduce the informativeness of dissent for prejudice and thus lower the associated social costs. In Experiment 2, we examine whether rationales indeed serve this purpose.

### 3.2.1 Sample and experimental design

We conducted our pre-registered Experiment 2 in November 2021 with a sample of Democrats and Independents recruited from Prolific.<sup>19</sup> Our final sample of 1,040 Democrats and Independents is mostly balanced on observables across treatment arms (Appendix Table B.5).

---

<sup>19</sup>Our experiment was pre-registered in the AEA RCT registry under ID AEARCTR-0005462. The full set of experimental instructions is included in Appendix E.2.

Figure B.2 outlines the structure of Experiment 2. After completing a battery of demographic and other background questions, respondents are informed that they have been matched with a previous survey participant who joined a campaign to oppose the movement to defund the police. They are then randomized into a *Cover* and a *No Cover* condition: respondents in the *Cover* condition are told that their matched participant authorized the Tweet corresponding to the *Cover* condition of Experiment 1 (“Before I joined the campaign...”) whereas respondents in the *No Cover* condition are told that their matched participant authorized the *No Cover* Tweet (“After I joined the campaign...”).

We begin by asking respondents to respond to the following open-ended question: “Why do you think your matched participant chose to join the campaign to oppose defunding the police?” This approach avoids priming respondents to think about particular dimensions and instead directly elicits “what comes to mind” (Gennaioli and Shleifer, 2010). As a more direct measure of inference about their matched participant’s prejudice, we subsequently tell them that their matched participant had the opportunity to authorize a \$5 donation to the National Association for the Advancement of Colored People (NAACP) and ask them to guess whether or not the participant donated. Finally, we also give respondents the opportunity to authorize a \$1 bonus to their matched respondent (at no cost to themselves): declining to do so is our measure of social sanction.

### 3.2.2 Results

We estimate statistically and economically significant treatment effects on all three measures of type inference. Sub-panel (a) in Figure 3 displays the fraction of participants in the *Cover* and *No Cover* condition who believe their matched participant donated to the NAACP (results reported in regression table form in Panel A, Columns 1–3 of Table 3). 27% of respondents in the *No Cover* condition believe their matched participant donated, compared to 35% of respondents in the *Cover* condition ( $p = 0.012$ ). Similarly, sub-panel (b) displays the fraction of participants who deny their matched participant a bonus (results reported in regression table form in Panel B, Columns 1–3 of Table 3). 47% of respondents in the *No Cover* condition deny their matched participant a bonus, compared to 40% of respondents in the *Cover* condition ( $p = 0.016$ ). As shown in Table 3, these estimates are stable to the inclusion of demographic and partisan controls. As implied by our framework, even respondents who *privately agree* with their matched participant’s opposition to defunding the police may choose to levy social sanctions if they believe that the only people who would be comfortable *publicly* expressing such an opinion are prejudiced.

To analyze the open-ended text, we look for the words or phrases of up to three words that are most characteristic of each condition. More precisely, we follow Gentzkow and Shapiro (2010) to calculate Pearson’s  $\chi^2$  statistic for each phrase.<sup>20</sup> This statistic is higher when the use of the phrase is more asymmetric across treatment conditions and lower for phrases that are used rarely across both conditions. Appendix Figure B.3 plots the top 20 most characteristic phrases of each condition. Consistent with our framework and the treatment effects on the structured measures of inference, we find that respondents in the *Cover* condition are more likely to describe their partner using phrases related to the article or the associated evidence — for example, “read an article,” “convincing,” “increase violent crime,” “study” — while respondents in the *No Cover* instead use phrases such as “Republican,” “racist,” and “probably white”.<sup>21</sup>

### 3.2.3 Credibility

To investigate the role of credibility, we run a slightly revised version of Experiment 2 (Auxiliary Experiment 6) with a sample of 506 Democrats and Independents: we instead show respondents the “lower-credibility” versions of the Tweets, as described in Section 3.1.7.<sup>22</sup> We display results in Panel B of Figure 3 and Columns 4–6 of Table 3. While the point estimate of the effect of the rationale on both structured measure of inference remains positive, it is substantially smaller: 30% of respondents in the *No Cover* condition believe their matched partner donated, compared to 33% in the *Cover* condition ( $p = 0.58$ ) and 44% of respondents in the *No Cover* condition deny their matched partner the donation, compared to 42% in the *Cover* condition.<sup>23</sup> While we are underpowered to conclude that these treatment effects are statistically significantly smaller than the treatment effects estimated using the more credible rationale, the evidence is qualitatively consistent with this slightly less credible rationale being substantially less effective.

Our revised experiment also speaks to one of the most common complaints surrounding

---

<sup>20</sup>This statistic is given by:  $\chi_p^2 = \frac{(n_p^R n_{\sim p}^{NR} - n_p^{NR} n_{\sim p}^R)^2}{(n_p^R + n_p^{NR})(n_p^R + n_{\sim p}^R)(n_p^{NR} + n_{\sim p}^{NR})(n_{\sim p}^R + n_{\sim p}^{NR})}$ , where  $n_p^R$ ,  $n_p^{NR}$  are the number of times  $p$  appears across all responses in the *Cover* condition and *No Cover* condition, respectively, and  $n_{\sim p}^i$  is the total number of times a phrase that is *not*  $p$  appears in condition  $i$ .

<sup>21</sup>These open-ended responses also allow us to mitigate concerns about other potential explanations for our findings: for example, that respondents in the *Cover* condition believed that their matched participant felt pressured by the experimenter to join the campaign and this pressure led them to do so. No respondents mention this or other related confounds.

<sup>22</sup>See Appendix E.10 for experimental instructions.

<sup>23</sup>As shown in Appendix D, our results are unchanged if we reweight responses to match the demographics of the sample in the higher-credibility variation.

“political correctness” culture: the alleged tendency of people to “take things out of context”. The article prominently lists both Sharkey’s academic credentials and, in the first few paragraphs, unequivocally states that “One of the most robust, most uncomfortable findings in criminology is that putting more officers on the street leads to less violent crime.” Nonetheless, the revised Tweet appears substantially less effective in shifting inference and reducing social sanctions (suggesting that most respondents do not read the article before deciding whether to sanction their partner). Requirements for dissenters to ensure that no part of their argument can be taken out of context and stripped of accompanying rationales may leave limited scope for expressing nuanced arguments. Conversely, evidence (such as scientific or media articles) may serve as a rationale even if few people actually examine it, so long as it appears compelling at first glance. We discuss implications for the spread of fake and misleading news and for political entrepreneurship in Section 5.

## 4 Support for Deporting Illegal Immigrants

Our next set of experiments examine the use and interpretation of rationales among a different population — conservatives — and to justify a different stigmatized position — support for a campaign to immediately deport all illegal Mexican immigrants. We examine our mechanism in this different context for three primary reasons. First, defunding the police is a highly salient but novel policy proposal, and it is thus unclear whether the power of rationales also extends to more “traditional” policy questions, for which there may be more common knowledge about a greater body of evidence and partisan talking points. Second, opposition to defunding the police is likely stigmatized by the in-group (Democrats) but not the out-group (Republicans); in contrast, supporting the immediate deportation of all illegal Mexican immigrants is less stigmatized by the in-group (Republicans), but is highly stigmatized by the out-group (Democrats). This setting thus allows us to examine whether rationales can be used to mitigate social sanctions levied by the out-group as well as from the in-group. Finally, understanding the drivers of anti-immigrant narratives on social media is of direct interest.

As in the previous experiment on the expression of dissent, we study the expression of xenophobia on social media. Given the widespread and growing importance of right-wing media as suppliers of anti-immigrant narratives, we examine a different form of rationale: a thirty-second clip from one of the most popular cable news shows in the US, *Tucker Carlson Tonight*. In the clip, Carlson draws upon statistics from the US Sentencing Commission

to argue that illegal immigrants commit violent crimes at substantially higher rates than citizens.<sup>24</sup>

## 4.1 Experiment 3: Rationales and Pro-Deportation Expression

### 4.1.1 Sample and experimental design

We conducted our pre-registered Experiment 3 in March 2021 with a sample of Republicans and Independents.<sup>25</sup> We recruited 1,130 participants through Luc.id. After screening out respondents who did not want to join the campaign (as described below), we are left with a final sample of 508 respondents. Our sample is balanced on observables across treatment arms (Appendix Table B.8).

Our experimental design is broadly similar to that of Experiment 1; we provide a diagram in Figure B.4. As in Experiment 1, respondents log into our survey using their Twitter account and respond to a set of demographic and other background questions. Respondents then view the clip from *Tucker Carlson Tonight*, which is embedded into the survey, and are randomized into the *Cover* condition or the *No Cover* condition. Respondents in the *Cover* condition, but not in the *No Cover* condition, are then provided with the URL to the video. We then ask all respondents whether they would like to join a campaign to immediately deport all illegal Mexican immigrants. The survey terminates for respondents who do not join the campaign, leaving us with 517 remaining respondents. Those respondents in the *No Cover* group who do join the campaign are provided the URL to the video. In other words, at this point in the survey, the only difference between conditions is whether respondents are provided with the video URL before (*Cover*) or after (*No Cover*) joining the campaign — though all respondents watch the clip before joining the campaign. As we discuss below, this difference in timing is key to avoiding explicit deception in our experimental manipulation.

Respondents who join the campaign are informed that one component of the campaign involves circulating a petition on Twitter calling for illegal Mexican immigrants to be deported. We show them a screenshot of the Tweet and ask them if they are willing to schedule it to be posted on their account. As in Experiment 1, we inform respondents that all Tweets will be posted all at once if and when we have surveyed people in all US

---

<sup>24</sup>The clip is available at [https://www.youtube.com/embed/SDdkkTLCUUQ?autoplay=1&controls=0&end=166&fs=0&modestbranding=1&start=113&iv\\_load\\_policy=3](https://www.youtube.com/embed/SDdkkTLCUUQ?autoplay=1&controls=0&end=166&fs=0&modestbranding=1&start=113&iv_load_policy=3).

<sup>25</sup>Our experiment was pre-registered in the AEA RCT registry under ID AEARCTR-0007379. The full set of experimental instructions is included in Appendix E.3.

counties, that this is a common tactic used to make campaigns trend on Twitter, and that we will delete all identifying information by no later than August 1, 2021. Again as in Experiment 1, because we target fewer respondents than the number of US counties, we ensure that Tweets will never be posted. The ethical considerations underlying our design are much the same as those of Experiment 1; we discuss these considerations in depth in Appendix C.

Respondents in the *Cover* condition are asked whether they would like to schedule the following Tweet:

I have joined a campaign to immediately deport all illegal Mexican immigrants.  
 Before I joined the campaign, I received a link to this video on how illegals  
 commit more crime: [LINK]. Sign this petition to immediately deport all illegal  
 Mexicans: [LINK]

The key experimental manipulation is similar to that of Experiment 1: respondents in the *No Cover* condition are presented with an identical Tweet, but with the “Before I joined the campaign...” replaced with “After I joined the campaign...”. Although all respondents in fact watched the video before joining the campaign, it is true that respondents in the *Cover* condition *received the link* to the video before joining, while those in the *No Cover* condition received the link after joining.<sup>26</sup> This difference in wording suggests to potential readers of the Tweet that respondents in the *Cover* group had been exposed to the video by Tucker Carlson before joining the campaign — and thus potentially joined because they were convinced by the clip’s evidence — while respondents in the *No Cover* group had *not* been exposed before joining the campaign, and thus could not have joined due to the clip. As in Experiment 1, then, this manipulation varies the availability of social cover while fixing the persuasion channel (all respondents are exposed to the same video) and the anticipated persuasion channel (all respondents know their Tweet’s readers will be exposed to the video, since it is linked in the Tweet).<sup>27</sup>

---

<sup>26</sup>One potential concern is that providing a link to respondents in the *Cover* condition, but not in the *No Cover* condition, induces differential selection into the campaign. Because we make the source of the clip obvious, we do not view this as a plausible confound. Indeed, we find no statistically significant difference in selection into the campaign between groups (a 2.6 percentage point difference,  $p = 0.474$ ), and our worst-case estimate under Lee (2009) bounds remains statistically significant at the 1% level.

<sup>27</sup>In principle, we could have used a similar design as Experiment 1: showing the video to respondents both before and after they join the campaign. We concluded that such a manipulation would be less natural for a 30-second video than for a longer article, as in Experiment 1.

### 4.1.2 Results

Figure 4 displays the results, which we also show in regression table form in Panel A of Table 4. We again find an economically and statistically significant cover effect: 47% of respondents in the *No Cover* condition authorize the Tweet, while 64% of respondents in the *Cover* condition authorize the Tweet ( $p < 0.01$ , a 0.35 standard deviation effect). This estimate is stable to the inclusion of demographic and partisan controls. The fact that the effect is larger than that estimated in Experiment 1 may reflect that Republicans feel greater stigma in joining a pro-deportation campaign than Democrats feel in joining an anti-defunding campaign (which is also consistent with the lower mean authorization rates in this experiment than in Experiment 1); or that Republicans perceive the *Tucker Carlson video* as a more compelling rationale vis-a-vis their Twitter followers than Democrats perceive the *Washington Post* article vis-a-vis their followers.<sup>28</sup>

## 4.2 Experiment 4: Interpretation of Pro-Deportation Rationale

We next examine how the availability of the social cover provided by the *Tucker Carlson Tonight* clip shapes an audience’s inference about a dissenter’s underlying motivations and the resulting social sanctions the dissenter faces.

### 4.2.1 Sample and experimental design

We conducted our pre-registered Experiment 4 in November 2021 with a sample of 1,082 Democrats and Independents recruited from Prolific.<sup>29</sup> We focus on Democrats and Independents, as anti-immigrant expression is less likely to be stigmatized among Republicans. Our sample is balanced on observables across treatment arms (Appendix Table B.11).

Experiment 4 follows the structure of Experiment 2; Figure B.2 outlines the structure of the experiments (with red text corresponding to Experiment 4). Respondents are informed that they have been matched with a previous survey participant who joined a campaign to deport all illegal Mexican immigrants. As in Experiment 2, they are then randomized into a *Cover* and a *No Cover* condition: respondents in the *Cover* condition are told that their matched participant authorized the Tweet corresponding to the *Cover* condition

---

<sup>28</sup>In our pre-registered Auxiliary Experiment 7 designed to measure the persuasiveness of the rationale, we find mixed evidence for persuasive effects on private opinions; see Appendix B.2.2 for details and Appendix E.11 for experimental instructions.

<sup>29</sup>Our experiment was pre-registered in the AEA RCT registry under ID AEARCTR-0005462. The full set of experimental instructions is included in Appendix E.4.



of Experiment 3 (“Before I joined the campaign...”) whereas respondents in the *No Cover* condition are told that their matched participant authorized the *No Cover* Tweet (“After I joined the campaign...”). Respondents then respond to the following open-ended question: “Why do you think your matched participant chose to donate to the campaign?”. Subsequently, they guess whether their matched participant authorized a \$5 donation to the US Border Crisis Children’s Relief Fund (an organization that seeks to provide care and basic hygiene items to children along the US–Mexico border) when given the opportunity to do so, and they choose whether or not to deny a \$1 bonus to their matched participant.<sup>30</sup>

#### 4.2.2 Results

Panel B of Figure 4 displays the fraction of participants in the *Cover* and *No Cover* condition who believe their matched participant donated to the pro-immigrant organization and the corresponding fractions of participants who deny their matched respondent a bonus. 8.5% of respondents in the *No Cover* condition believe their matched participant donated, compared to 13.4% of respondents in the *Cover* condition ( $p = 0.01$ ); 80% of respondents in the *No Cover* condition deny their matched participant a bonus, compared to 74% of respondents in the *Cover* condition ( $p = 0.011$ ). As shown in Panels B–C of Table 4, these estimates are stable to the inclusion of demographic and partisan controls.

We plot results from our analysis of open-ended text in Appendix Figure B.5 using the same procedure described in Section 3.2.2. As in Experiment 2, respondents in the *Cover* condition are substantially more likely to use words referencing the rationale — “watched a video,” “right wing media,” “link” — whereas respondents in the *No Cover* condition mention phrases such as “Republican,” “extremist,” and “biased”.

## 5 Discussion and Conclusion

This paper examines how rationales facilitate dissent by lowering the social cost of expressing controversial opinions. In our model, rationales change some people’s private views or beliefs about social welfare, but they can also be used to justify dissent, shifting an audience’s inference about the dissenter’s motivations. We explore these mechanisms among both liberal and conservative respondents, focusing primarily on a natural setting and outcome: willingness to express dissent on social media. First, we show that liberal

---

<sup>30</sup>We randomized the order of these two different outcomes and detect no significant order effects.

respondents are more likely to authorize a Tweet opposing the movement to defund the police when they can credibly ascribe their views to strong scientific evidence. Consistent with our framework, a credible rationale shifts an audience’s inference about the respondents and reduces resulting social sanctions. Similarly, conservative respondents are more likely to authorize a Tweet calling for the deportation of all illegal immigrants from Mexico — and are seen as less intolerant after doing so — when they can ascribe their views to a Fox News clip.

We now discuss some implications of our framework and empirical results, which may provide fruitful avenues for future research.

**Political correctness and the limitations of rationales** In a “political correctness” culture, certain arguments (rationales) cannot be voiced because they are seen as legitimizing dangerous or undesirable causes, and so anyone who voices such an argument is seen as supporting the cause itself. For example, people who argue for the presence of reverse discrimination against men in labor markets may be seen as sexists: that is, even scientific arguments such as correspondence studies — which are typically effective rationales — may fail to provide a social cover. In some cases, this may be socially desirable: for instance, equating the use of a rationale with sexism may prevent sexist individuals from citing rationales they do not believe or cherry-picking arguments to support their claims. In other cases, political correctness culture may stifle socially important forms of dissenting expression by stigmatizing rationales that would typically be seen as highly credible.

Individuals or institutions seeking to eliminate certain forms of public behavior — for better or for worse — may use multiple levers to silence dissenters. One lever, explored in Section 3.1.7, is to undermine the credibility of rationales directly. Another lever is to manipulate the real or perceived correlation between knowledge of a rationale and underlying type, tying the rationale directly to the stigmatized motive.<sup>31</sup> Indeed, in the limit in which only people with stigmatized motives are aware of a certain rationale — e.g. because only they consume the extreme news sources through which the rationale is broadcast, or because only they follow a fringe public figure who spreads the rationale — the rationale is completely ineffective, as to use it is to reveal one’s motives with certainty. Tactics to ma-

---

<sup>31</sup>For example, during the Second Red Scare, Joseph McCarthy and his allies explicitly tied several rationales for dissenting with government policy to Communist sympathies. Famously, physicist J. Robert Oppenheimer — credited as the “father of the atomic bomb” — was stripped of his security clearances when political opponents attributed his opposition to the development of the hydrogen bomb to alleged Soviet loyalties (Cassidy, 2019).

nipulate the real or perceived correlation between motive and rationale include disallowing controversial opinions a public platform (e.g. disinviting campus speakers or banning social media accounts), or branding particular media sources or speakers as fringe.<sup>32</sup> Further exploring the conditions under which rationales are most effective, and the unifying features of effective rationales, is an important direction for future research.

**Political entrepreneurship and populism** Successful politicians often base their campaigns on simple messages that resonate with the general public. Many populist politicians are particularly skilled at scapegoating minority groups.<sup>33</sup> Our framework can shed light on why some appeals are more effective than others. While the persuasive effects of propaganda are doubtless important (Adena et al., 2015), propaganda may also generate social cover, enabling supporters to speak their mind more openly and spread the message through their social circle (Satyanath et al., 2017; Caesmann et al., 2021). The strength of this “social amplifier” channel depends not only on the number of individuals who hold stigmatized views, but the number of individuals who *could not express these views* prior to the rationale becoming widespread. This distinction can provide one explanation for why the Nazis were able to leverage social networks and associations while other parties, including communists, could not: if antisemitism was stigmatized, but relatively common and persistent (Voigtländer and Voth, 2012; Cantoni et al., 2019), then Nazi rhetoric blaming Jews for Germany’s problems generated a large social amplifier, thereby furthering Nazi views. Blaming elites, on the other hand, was less stigmatized, and thus generated far smaller amplifiers.

**Fake and misleading news** Our findings speak to the debate about the influence of fake and misleading news on society. Some recent studies suggest that their persuasive effect is limited (Allcott and Gentzkow, 2017; Nyhan, 2018), while others suggest that they can be effective at changing behavior (Barrera et al., 2020) and that individuals may have trouble distinguishing between fake and real news (Angelucci and Prat, 2021) or between facts and opinions (Bursztyn et al., forthcoming). Our results point, however, to an alternative mechanism through which misleading news can affect public expression. Specifically, fake

---

<sup>32</sup>This can also help explain how censorship techniques such as China’s “Great Firewall” can be highly effective in repressing discourse unfriendly to the regime, even if citizens can bypass them relatively easily (Chen and Yang, 2019).

<sup>33</sup>See Guriev and Papaioannou (2020) for a review on the political economy of populism. Bursztyn et al. (2022) applies our framework to explore the scapegoating of minorities during economic crises.

news can generate a “social amplifier”: rationales that plausibly persuade a subset of the population can change public behavior among a much larger fraction of the population, increasing their willingness to express otherwise-stigmatized views. Interestingly, in Barrera et al. (2020), subjects exposed to fake news were not only more willing to support an extreme candidate (Marine Le Pen), but also were unlikely to change their opinion after being exposed to fact-checks — even though these fact-checks improved factual knowledge. This evidence is difficult to explain by the persuasive power of fake news alone, but it is consistent with the role of fake news as rationales: fake and misleading news can generate social cover for individuals to express extreme views, and debunking does not eliminate social cover as long as the fact-check can be plausibly dismissed.

This insight has implications for debunking fake news spread online and offline. Among other platforms, Facebook and Twitter have conducted small-scale experiments evaluating strategies to curtail the spread of misinformation, including warning users before they post an article flagged as fake news and flagging fake or misleading news when it appears on users’ timelines (e.g., because a friend shared it). The former initiative decreases the persuasive effect of fake news for a user who seeks to spread it, while the latter decreases the anticipated persuasiveness of the rationale. Yet because these experiments have occurred only among a small fraction of users, people have a ready-made social cover when sharing fake news: they can credibly claim that they were not warned the news was fake.<sup>34</sup>

Our results highlight the potential importance of eliminating social cover: ensuring that the audience *knows that the poster knew the news had been debunked* and nonetheless chose to post it. A simple path would be to scale the debunking experiments to the entire userbase, thus generating common knowledge that all users are warned before posting fake news. Because the general equilibrium results of such a change differ significantly from the partial equilibrium results, current estimates of the effects of debunking on users’ propensity to share fake news may substantially understate the true effects that would be realized if platforms were to fully scale up the feature.

---

<sup>34</sup>Indeed, both Twitter and Facebook’s fact-checking efforts have been widely criticized for a lack of transparency, and it is thus certain that most users lack information about how the platforms fight misinformation (Nyhan, 2017).

## References

- Acemoglu, Daron, Georgy Egorov, and Konstantin Sonin**, “A Political Theory of Populism,” *Quarterly Journal of Economics*, 2013, 128 (2), 771–805.
- Adena, Maja, Ruben Enikolopov, Maria Petrova, Veronica Santarosa, and Ekaterina Zhuravskaya**, “Radio and the Rise of The Nazis in Prewar Germany,” *Quarterly Journal of Economics*, 2015, 130 (4), 1885–1939.
- Ali, S Nageeb and Charles Lin**, “Why People Vote: Ethical Motives and Social Incentives,” *American Economic Journal: Microeconomics*, 2013, 5 (2), 73–98.
- Allcott, Hunt and Matthew A. Gentzkow**, “Social Media and Fake News in the 2016 Election,” *Journal of Economic Perspectives*, 2017, 31 (2), 211–36.
- Angelucci, Charles and Andrea Prat**, “Is Journalistic Truth Dead? Measuring How Informed Voters Are about Political News,” Working Paper 3593002, Social Science Research Network, 2021.
- Barrera, Oscar, Sergei Guriev, Emeric Henry, and Ekaterina Zhuravskaya**, “Facts, alternative facts, and fact checking in times of post-truth politics,” *Journal of Public Economics*, 2020, 182, 104123.
- Bénabou, Roland and Jean Tirole**, “Incentives and Prosocial Behavior,” *American Economic Review*, 2006, 96 (5), 1652–1678.
- and –, “Identity, Morals, and Taboos: Beliefs as Assets,” *Quarterly Journal of Economics*, 2011, 126 (2), 805–855.
- , **Armin Falk, and Jean Tirole**, “Narratives, Imperatives, and Moral Persuasion,” Working Paper 24798, National Bureau of Economic Research, 2020.
- Braghieri, Luca**, “Political Correctness, Social Image, and Information Transmission,” Working paper, 2022.
- Bursztyn, Leonardo, Aakaash Rao, Christopher Roth, and David Yanagizawa-Drott**, “Opinions as Facts,” *Review of Economic Studies*, forthcoming.
- , **Alessandra L. González, and David Yanagizawa-Drott**, “Misperceived Social Norms: Women Working Outside the Home in Saudi Arabia,” *American Economic Review*, 2020, 110 (10), 2997–3029.
- and **Robert Jensen**, “How Does Peer Pressure Affect Educational Investments?,” *Quarterly Journal of Economics*, 2015, 130 (3), 1329–1367.
- , **Georgy Egorov, and Stefano Fiorin**, “From Extreme to Mainstream: The Erosion of Social Norms,” *American Economic Review*, 2020, 110 (11), 3522–48.
- , –, **Ingar Haaland, Aakaash Rao, and Christopher Roth**, “Scapegoating during Crises,” *AEA Papers and Proceedings*, 2022, 112, 151–55.
- , –, **Ruben Enikolopov, and Maria Petrova**, “Social media and xenophobia: evidence from Russia,” Working Paper 26567, National Bureau of Economic Research, 2019.

- Caesmann, Marcel, Bruno Caprettini, Hans-Joachim Voth, David Yanagizawa-Drott et al.**, “Going Viral: Propaganda, Persuasion and Polarization in 1932 Hamburg,” Technical Report, Centre for Economic Policy Research 2021.
- Cantoni, Davide, Felix Hagemeister, and Mark Westcott**, “Persistence and activation of right-wing political ideology,” Working paper, 2019.
- Caprettini, Bruno, Marcel Caesmann, Hans-Joachim Voth, and David Yanagizawa-Drott**, “Going Viral: Propaganda, Persuasion and Polarization in 1932 Hamburg,” Working Paper 16356, Center for Economic and Policy Research, 2021.
- Cassidy, David C.**, *J. Robert Oppenheimer and the American century*, Plunkett Lake Press, 2019.
- Chen, Yuyu and David Y. Yang**, “The Impact of Media Censorship: 1984 or Brave New World?,” *American Economic Review*, 2019, 109 (6), 2294–2332.
- Cohn, Nate and Kevin Quealy**, “The Democratic Electorate on Twitter Is Not the Actual Democratic Electorate,” *The New York Times*, 2019.
- Cunningham, Tom and Jonathan de Quidt**, “Implicit Preferences Inferred from Choice,” Working Paper 2709914, Social Science Research Network, 2015.
- Dana, Jason, Roberto A. Weber, and Jason Xi Kuang**, “Exploiting Moral Wiggle Room: Experiments Demonstrating an Illusory Preference for Fairness,” *Economic Theory*, 2007, 33 (1), 67–80.
- Ekins, Emily E.**, “Poll: 62% of Americans Say They Have Political Views They’re Afraid to Share,” Working Paper 3659953, Social Science Research Network, 2020.
- Enikolopov, Ruben, Alexey Makarin, and Maria Petrova**, “Social media and protest participation: Evidence from Russia,” *Econometrica*, 2020, 88 (4), 1479–1514.
- **and Maria Petrova**, “Chapter 17 - Media Capture: Empirical Evidence,” in Simon P. Anderson, Joel Waldfogel, and David Strömberg, eds., *Handbook of Media Economics*, Vol. 1, North-Holland, 2015, pp. 687–700.
- Exley, Christine L.**, “Excusing Selfishness in Charitable Giving: The Role of Risk,” *Review of Economic Studies*, 2016, 83 (2), 587–628.
- Foerster, Manuel and Joel J van der Weele**, “Persuasion, Justification and the Communication of Social Impact,” *The Economic Journal*, 2021, 131, 2887–2919.
- Gennaioli, Nicola and Andrei Shleifer**, “What Comes to Mind,” *Quarterly Journal of Economics*, 2010, 125 (4), 1399–1433.
- Gentzkow, Matthew and Jesse M Shapiro**, “What drives media slant? Evidence from US daily newspapers,” *Econometrica*, 2010, 78 (1), 35–71.
- Golman, Russell**, “Acceptable Discourse: Social Norms of Beliefs and Opinions,” Working paper, 2021.
- **, David Hagmann, and George Loewenstein**, “Information Avoidance,” *Journal of Economic Literature*, 2017, 55 (1), 96–135.

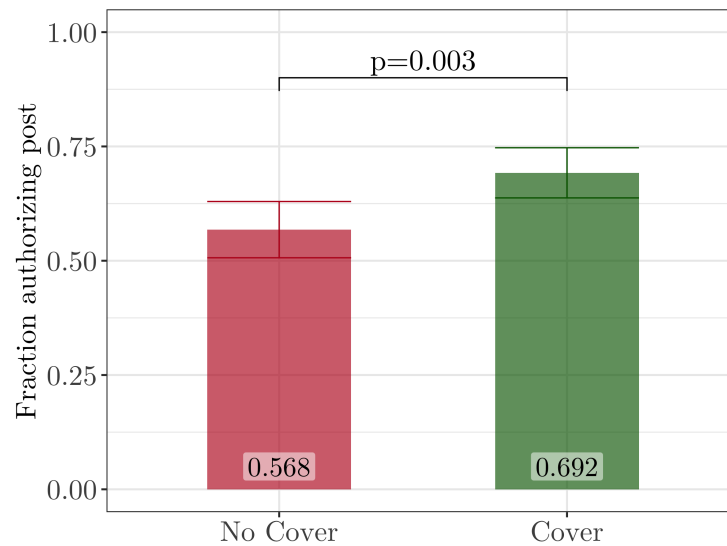
- Grossman, Zachary and Joel J Van Der Weele**, “Self-image and willful ignorance in social decisions,” *Journal of the European Economic Association*, 2017, 15 (1), 173–217.
- Guriey, Sergei and Elias Papaioannou**, “The Political Economy of Populism,” *Journal of Economic Literature*, 2020.
- Haaland, Ingar, Christopher Roth, and Johannes Wohlfart**, “Designing Information Provision Experiments,” *Journal of Economic Literature*, 2021.
- Hamman, John R., George Loewenstein, and Roberto A. Weber**, “Self-Interest through Delegation: An Additional Rationale for the Principal-Agent Relationship,” *American Economic Review*, 2010, 100 (4), 1826–1846.
- Kuran, Timur**, *Private Truths, Public Lies: The Social Consequences of Preference Falsification*, Cambridge, Massachusetts: Harvard University Press, 1997.
- Lacetera, Nicola and Mario Macis**, “Social Image Concerns and Prosocial Behavior: Field Evidence from a Nonlinear Incentive Scheme,” *Journal of Economic Behavior and Organization*, 2010, 76 (2), 225–237.
- Langer, Ellen J, Arthur Blank, and Ben Zion Chanowitz**, “The mindlessness of ostensibly thoughtful action: The role of “placebic” information in interpersonal interaction,” *Journal of Personality and Social Psychology*, 1978, 36 (6), 635.
- Lazear, Edward P., Ulrike Malmendier, and Roberto A. Weber**, “Sorting in Experiments with Application to Social Preferences,” *American Economic Journal: Applied Economics*, 2012, 4 (1), 136–163.
- Lee, David S.**, “Training, Wages, and Sample Selection: Estimating Sharp Bounds on Treatment Effects,” *Review of Economic Studies*, 2009, 76 (3), 1071–1102.
- Levy, Roee and Martin Mattsson**, “The Effects of Social Movements: Evidence from #MeToo,” Working Paper 3496903, Social Science Research Network, 2021.
- Litman, Leib, Jonathan Robinson, and Tzvi Abberbock**, “TurkPrime.com: A versatile crowdsourcing data acquisition platform for the behavioral sciences,” *Behavior Research Methods*, 2017, 49 (2), 433–442.
- Michalopoulos, Stelios and Melanie Meng Xue**, “Folklore,” *Quarterly Journal of Economics*, 2021, 136 (4), 1993–2046.
- Morris, Stephen**, “Political correctness,” *Journal of Political Economy*, 2001, 109 (2), 231–265.
- Müller, Karsten and Carlo Schwarz**, “Making America Hate Again? Twitter and Hate Crime Under Trump,” Working Paper 3149103, Social Science Research Network, 2018.
- Nyhan, Brendan**, “Why the Fact-Checking at Facebook Needs to Be Checked,” *The New York Times*, 2017.
- , “Fake News and Bots May Be Worrisome, but Their Political Power Is Overblown,” *The New York Times*, 2018.

- O'Brien, Sarah**, "Employers check your social media before hiring. Many then find reasons not to offer you a job," *CNBC*, 2018.
- Ousey, Graham C. and Charis E. Kubrin**, "Immigration and Crime: Assessing a Contentious Issue," *Annual Review of Criminology*, 2018, 1 (1), 63–84.
- Parker, Kim and Kiley Hurst**, "Growing share of Americans say they want more spending on police in their area," *Pew Research Center's Report*, 2021.
- Patir, Assaf, Bnaya Dreyfuss, and Moses Shayo**, "On the Workings of Tribal Politics," Working Paper 3797290, Social Science Research Network, 2021.
- Perez-Truglia, Ricardo and Guillermo Cruces**, "Partisan Interactions: Evidence from a Field Experiment in the United States," *Journal of Political Economy*, 2017, 125 (4), 1208–1243.
- Saccardo, Silvia and Marta Serra-Garcia**, "Cognitive Flexibility or Moral Commitment? Evidence of Anticipated Belief Distortion," Working paper 3676711, Social Science Research Network, 2020.
- Satyanath, Shanker, Nico Voigtländer, and Hans-Joachim Voth**, "Bowling for fascism: Social capital and the rise of the Nazi Party," *Journal of Political Economy*, 2017, 125 (2), 478–526.
- Science Panel for the Amazon**, "Amazon Assessment Report 2021," *Science Panel for the Amazon*, 2021.
- Shiller, Robert J**, "Narrative economics," *American Economic Review*, 2017, 107 (4), 967–1004.
- Thompson, Derek**, "Unbundle the Police," *The Atlantic*, 2020.
- Voigtländer, Nico and Hans-Joachim Voth**, "Persecution perpetuated: the medieval origins of anti-Semitic violence in Nazi Germany," *Quarterly Journal of Economics*, 2012, 127 (3), 1339–1392.
- Wood, Thomas and Ethan Porter**, "The Elusive Backfire Effect: Mass Attitudes' Steadfast Factual Adherence," *Political Behavior*, 2019, 41 (1), 135–163.
- Yanagizawa-Drott, David**, "Propaganda and Conflict: Evidence from the Rwandan Genocide," *Quarterly Journal of Economics*, 2014, 129 (4), 1947–1994.



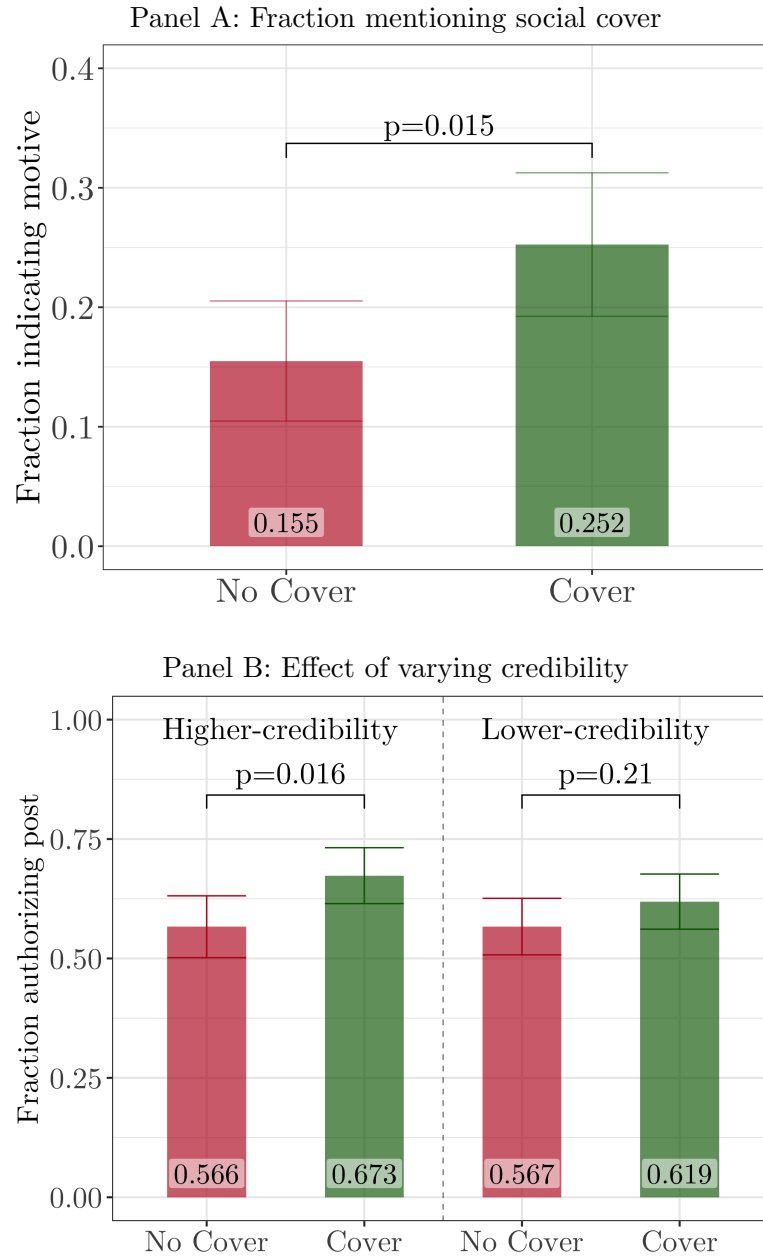
## Figures

**Figure 1:** Willingness to post anti-defunding Tweet



*Notes:* Figure presents results from Experiment 1 ( $n = 523$ ). We plot fraction of respondents authorizing the Tweet indicating their opposition to the movement to defund the police, separately by treatment. Error bars indicate 95% confidence intervals.  $p$ -values obtained from a two-sample  $t$ -test of equality of means.

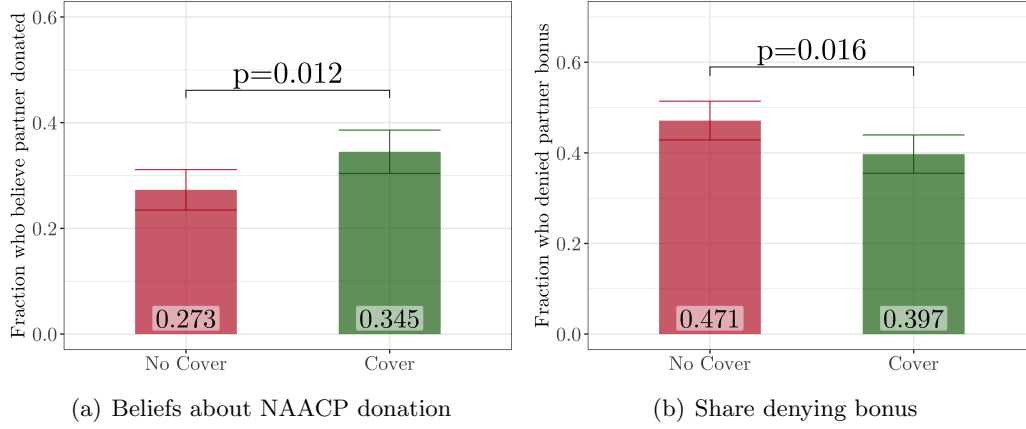
**Figure 2:** Willingness to post anti-defunding Tweet: investigating mechanisms



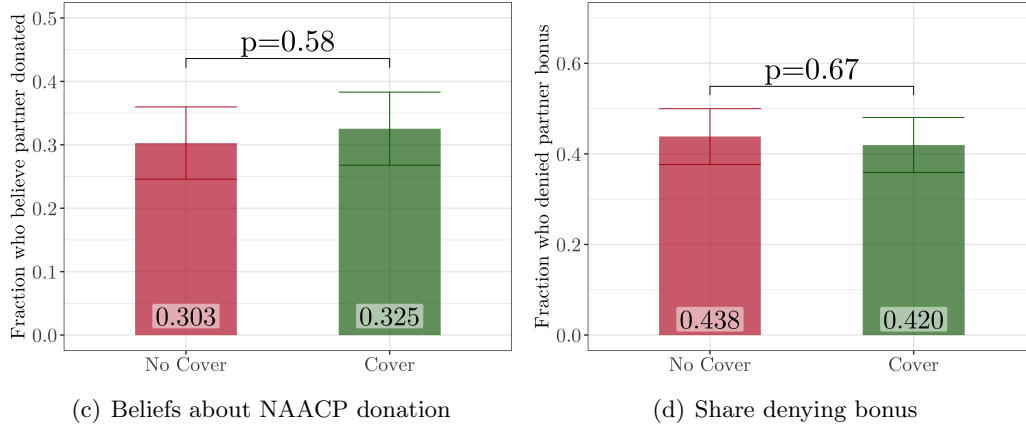
*Notes:* Figure presents results from Auxiliary Experiment 4 ( $n = 402$ ) and Auxiliary Experiment 5 ( $n = 1017$ ). Panel A displays the fraction of respondents who mention that “social cover”, as described in Section 3.1.6, was a motive underlying their choice of Tweet to post. Panel B displays the fraction of respondents who are willing to post the Tweet indicating their opposition to the movement to defund the police, separately by treatment group. Error bars indicate 95% confidence intervals.  $p$ -values obtained from a two-sample  $t$ -test of equality of means.

**Figure 3:** Interpretation of anti-defunding Tweet

Panel A: Higher-credibility rationale

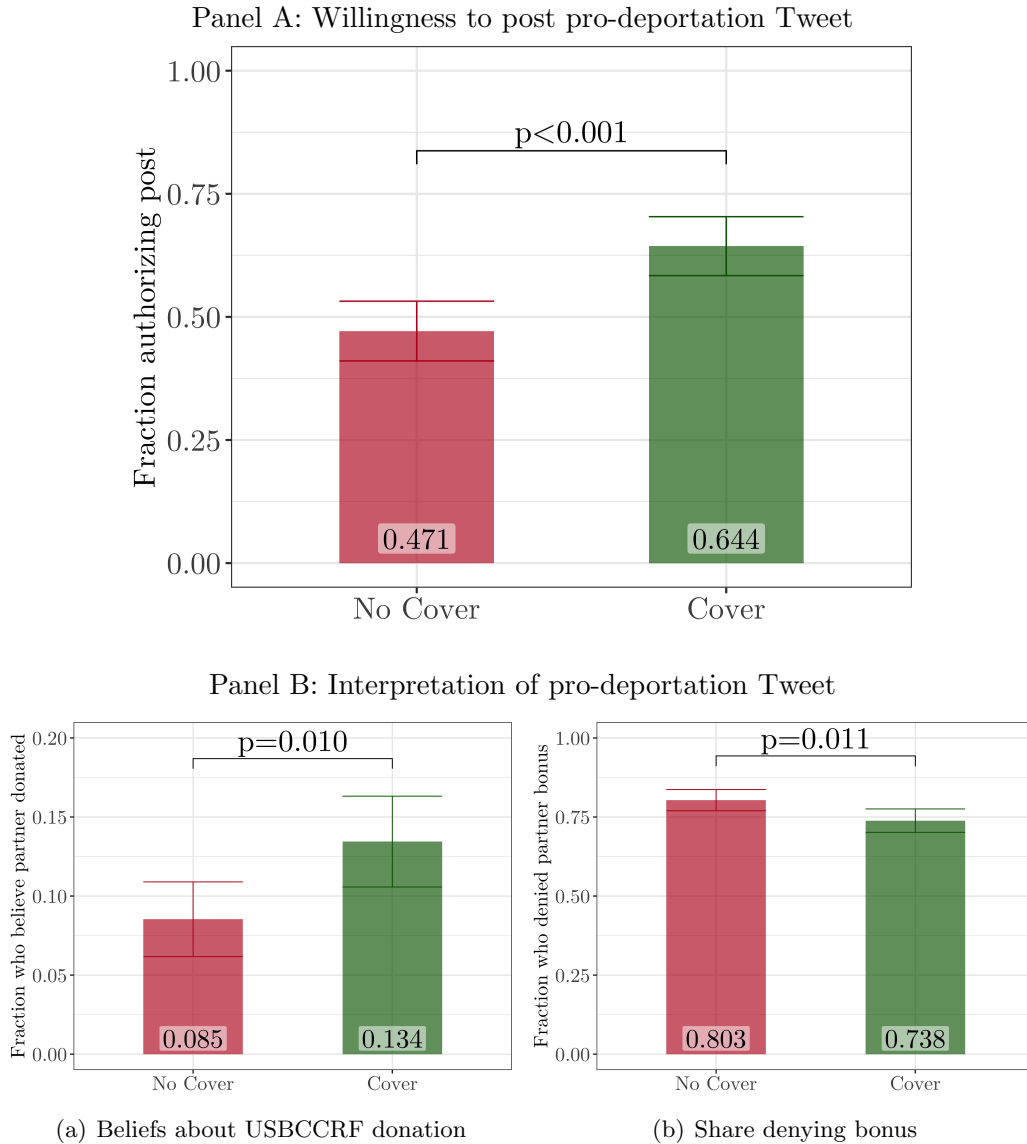


Panel B: Lower-credibility rationale



*Notes:* Panel A presents results from Experiment 2 ( $n = 1040$ ), in which respondents are shown the higher-credibility rationale; Panel B presents results from Auxiliary Experiment 6, in which respondents are shown the lower-credibility rationale ( $n = 506$ ). Sub-panels (a) and (c) present the fraction of respondents who believe their matched participant donated to the NAACP; sub-panels (b) and (d) present the fraction of respondents who deny their matched participant a \$1 bonus. Error bars indicate 95% confidence intervals.  $p$ -values obtained from a two-sample  $t$ -test of equality of means.

**Figure 4:** Expression and interpretation of pro-deportation Tweet



*Notes:* Figure presents results from Experiment 3 ( $n = 508$ ) and Experiment 4 ( $n = 1082$ ). Panel A displays the fraction of respondents authorizing the Tweet indicating their support for immediately deporting all illegal Mexican immigrants. Panel B presents the fraction of respondents who believe their matched participant donated to the US Border Crisis Children's Relief Fund (USBCCRF) on the left side and the fraction of respondents who deny their matched participant a \$1 bonus on the right side. Error bars indicate 95% confidence intervals.  $p$ -values obtained from a two-sample  $t$ -test of equality of means.

## Tables

**Table 1:** Willingness to post anti-defunding Tweet

	<i>Scheduled Tweet</i>		
	(1)	(2)	(3)
Cover	0.124*** (0.042)	0.119*** (0.042)	0.119*** (0.042)
No Cover mean	0.568	0.568	0.568
Observations	523	523	523
Demographic controls	No	Yes	Yes
Partisan controls	No	No	Yes

*Notes:* Table reports results from Experiment 1. The dependent variable is an indicator taking value 1 if the respondent chose to schedule the post. Demographic controls include age, age squared, a set of race indicators, a Hispanic indicator, a male indicator, a set of education indicators. Partisan controls include indicators for “Very conservative”, “Conservative”, “Neither liberal nor conservative” (omitted), “Liberal”, and “Very liberal”. Robust standard errors are reported.

**Table 2:** Interpreting effects of rationale on willingness to post anti-defunding Tweet

	Mean		Treatment effect	
	No Cover	Cover	Coefficient (s.e.)	<i>p</i> -value
<b>Panel A:</b> Rainforest placebo				
Scheduled post	0.83	0.79	-0.04 (0.04)	0.38
<b>Panel B:</b> Anticipated persuasion				
Estimated share persuaded	25.34	27.23	1.90 (2.12)	0.37
<b>Panel C:</b> Open-ended motive elicitation				
<i>C.1: Primary motives</i>				
<i>Respondent mentions...</i>				
Social cover	0.15	0.25	0.10 (0.04)	0.02
Anticipated persuasion	0.07	0.06	-0.01 (0.02)	0.67
Information	0.57	0.50	-0.07 (0.05)	0.13
<i>C.2: Potential confounds</i>				
<i>Respondent mentions...</i>				
Unnatural	0.01	0.01	0.01 (0.01)	0.32
Misleading	0.00	0.00	0.00 (0.00)	—
Signaling	0.00	0.00	0.00 (0.00)	—
Experimenter demand	0.00	0.00	0.00 (0.00)	—
<b>Panel D:</b> Credibility manipulation				
<i>D.1: Hypothetical willingness to post</i>				
Willing to post (high cred.)	0.57	0.67	0.11 (0.04)	0.02
Willing to post (low cred.)	0.57	0.62	0.05 (0.04)	0.21
<i>D.2: Beliefs about social sanctions</i>				
Share denying bonus (high cred.)	53.14	48.05	-5.09 (2.31)	0.03
Share denying bonus (low cred.)	53.99	53.00	-0.99 (2.06)	0.63

*Notes:* In Panel A, the DV is an indicator for whether the respondent chose to schedule the post (Auxiliary Experiment 2,  $n = 315$ ). In Panel B, the DV is the respondent's guess as to the percentage of their followers who would join the campaign if they saw the Tweet (Auxiliary Experiment 3,  $n = 501$ ). In Panel C, the DVs are indicators for whether the respondent's stated motive falls in each of the listed categories (Auxiliary Experiment 4,  $n = 402$ ). In Panel D, the DVs in lines 1–2 are indicators for whether the respondent was willing to post, while the DVs in lines 3–4 are stated beliefs about the share of Democrats who denied a bonus to the poster of the Tweet (Auxiliary Experiment 5,  $n = 1017$ ).

**Table 3:** Inference about and social sanctions toward matched anti-defunding respondent

	Higher-credibility			Lower-credibility		
	(1)	(2)	(3)	(4)	(5)	(6)
<b>Panel A:</b>	<i>Belief partner donated</i>					
Cover	0.072** (0.029)	0.072** (0.029)	0.067** (0.029)	0.023 (0.041)	0.023 (0.042)	0.019 (0.042)
No Cover mean	0.273	0.273	0.273	0.303	0.303	0.303
<b>Panel B:</b>	<i>Denied bonus to partner</i>					
Cover	-0.074** (0.031)	-0.074** (0.031)	-0.067** (0.030)	-0.019 (0.044)	-0.028 (0.044)	-0.015 (0.043)
No Cover mean	0.471	0.471	0.471	0.438	0.438	0.438
Observations	1,040	1,037	1,036	506	506	506
Demographic controls	No	Yes	Yes	No	Yes	Yes
Partisan controls	No	No	Yes	No	No	Yes

*Notes:* Table reports results from Experiment 2 (columns 1–3) and Auxiliary Experiment 6 (columns 4–6). The dependent variable in Panel A is an indicator taking value 1 if the respondent reports believing that his or her matched partner donated to the US Border Crisis Children’s Relief Fund. The dependent variable in Panel B is an indicator taking value 1 if the respondent denied his or her matched partner a \$1 bonus. Demographic controls include age, age squared, a set of race indicators, a Hispanic indicator, a male indicator, a set of education indicators. Partisan controls include indicators for “Very conservative”, “Conservative”, “Neither liberal nor conservative” (omitted), “Liberal”, and “Very liberal”. Robust standard errors are reported.

**Table 4:** Expression and interpretation of pro-deportation Tweet

Experiment 3			
<b>Panel A:</b>	<i>Scheduled Tweet</i>		
Cover	0.172*** (0.044)	0.179*** (0.044)	0.178*** (0.043)
No Cover mean	0.471	0.471	0.471
Observations	508	508	508
Experiment 4			
<b>Panel B:</b>	<i>Belief partner donated</i>		
Cover	0.049*** (0.019)	0.051*** (0.019)	0.048** (0.019)
No Cover mean	0.085	0.085	0.085
Observations	1,082	1,081	1,081
<b>Panel C:</b>	<i>Denied bonus to partner</i>		
Cover	-0.065** (0.026)	-0.065** (0.026)	-0.061** (0.026)
No Cover mean	0.803	0.803	0.803
Observations	1,082	1,081	1,081
Demographic controls	No	Yes	Yes
Partisan controls	No	No	Yes

*Notes:* Panel A presents the results of Experiment 3, in which the dependent variable is an indicator taking value 1 if the respondent chose to schedule the post. Panels B and C present the results of Experiment 4. The dependent variable in Panel B is an indicator taking value 1 if the respondent reports believing that his or her matched partner donated to the US Border Crisis Children’s Relief Fund (USBCCRF). The dependent variable in Panel C is an indicator taking value 1 if the respondent denied his or her matched partner a \$1 bonus. Demographic controls include age, age squared, a set of race indicators, a Hispanic indicator, a male indicator, a set of education indicators. Partisan controls include indicators for “Very conservative”, “Conservative”, “Neither liberal nor conservative” (omitted), “Liberal”, and “Very liberal”. Robust standard errors are reported.



# Online Appendix:

## Not for publication

Our supplementary material is structured as follows. Appendix A provides proofs of all theoretical results in Section 2. Appendix B.1 provides supporting material for the experiments presented in Section 3. Appendix B.2 provides supporting material for the experiments presented in Section 4. Appendix C discusses the ethical considerations underlying all experimental designs. Finally, Appendix E provides the instruments for all experiments described in the paper.

### A Theoretical Results

#### A.1 Proof of Proposition 1

We first prove that for random variable  $t$  distributed with c.d.f.  $H(\cdot)$  and p.d.f.  $h(\cdot)$ ,

$$\frac{d}{d\tau} \mathbb{E}(t \mid t > \tau) \leq 1.$$

Let  $z_\tau = t - \tau$  be a family of random variables indexed by  $\tau$ ; we need to show that

$$\mathbb{E}(z_\tau \mid z_\tau \geq 0)$$

is non-increasing in  $\tau$ . Denoting the c.d.f. of  $z_\tau$  by  $F_\tau(\cdot)$  and its p.d.f. by  $f_\tau(\cdot)$ , we have

$$\mathbb{E}(z_\tau \mid z_\tau \geq 0) = \frac{1}{1 - F_\tau(0)} \int_0^{+\infty} y f_\tau(y) dy.$$

The integral may be rewritten as

$$\begin{aligned} \int_0^{+\infty} y f_\tau(y) dy &= \int_0^{+\infty} f_\tau(y) \left( \int_0^y 1 dx \right) dy = \int_0^{+\infty} \int_0^y f_\tau(y) dx dy \\ &= \int_0^{+\infty} \int_x^{+\infty} f_\tau(y) dy dx = \int_0^{+\infty} (1 - F_\tau(x)) dx, \end{aligned}$$

where we used Fubini's theorem to change the order of integration.

Note that  $F_\tau(x) = \Pr(z_\tau \leq x) = \Pr(t \leq x + \tau) = H(x + \tau)$ . We therefore have

$$\mathbb{E}(z_\tau \mid z_\tau \geq 0) = \int_0^{+\infty} \frac{1 - F_\tau(x)}{1 - F_\tau(0)} dx = \int_0^{+\infty} \frac{1 - H(x + \tau)}{1 - H(\tau)} dx.$$

The integrand is non-increasing in  $\tau$  pointwisely (i.e., for any fixed  $x \geq 0$ ), because

$$\begin{aligned} \frac{d}{d\tau} \left( \frac{1 - H(x + \tau)}{1 - H(\tau)} \right) &= \frac{h(\tau)(1 - H(x + \tau)) - h(x + \tau)(1 - H(\tau))}{(1 - H(\tau))^2} \\ &= \frac{1 - H(x + \tau)}{1 - H(\tau)} \left( \frac{h(\tau)}{1 - H(\tau)} - \frac{h(x + \tau)}{1 - H(x + \tau)} \right) \leq 0, \end{aligned} \quad (3)$$

because the first term is positive and the second is nonpositive due to monotone hazard rate property. This proves that  $\mathbb{E}(z_\tau \mid z_\tau \geq 0)$  is non-increasing in  $\tau$ , and thus  $\frac{d}{d\tau} \mathbb{E}(t \mid t > \tau) \leq 1$ .

Now, for any fixed social cost  $S$ , type  $t_i$  would choose  $d_i = 1$  if  $t_i > \frac{1}{\beta}S - w_0$  and would choose  $d_i = 0$  if the opposite inequality holds. Thus, every equilibrium is characterized by a threshold  $\tau$ . This threshold  $\tau$  satisfies the condition

$$G(\tau) = -w_0, \quad (4)$$

where

$$G(\tau) = \tau - \frac{\gamma}{\beta} \mathbb{E}(t_i \mid t_i > \tau). \quad (5)$$

Since, as we proved,  $\frac{d}{d\tau} \mathbb{E}(t_i \mid t_i > \tau) \leq 1$  and  $\gamma < \beta$ , the  $G(\tau)$  is strictly increasing in  $\tau$ , and furthermore

$$\frac{d}{d\tau} G(\tau) \geq 1 - \frac{\gamma}{\beta} > 0.$$

This shows that equation (4) has a unique solution, which completes the proof. ■

## A.2 Proof of Proposition 2

Since the distributions are normal, the posterior of citizen  $i$  is given by the usual formula

$$w_1 = \mathbb{E}(w \mid s) = w_0 \frac{\sigma_\varepsilon^2}{\sigma_w^2 + \sigma_\varepsilon^2} + s \frac{\sigma_w^2}{\sigma_w^2 + \sigma_\varepsilon^2}.$$

We have

$$w_1 - w_0 = \frac{\sigma_w^2}{\sigma_w^2 + \sigma_\varepsilon^2} (s - w_0),$$

so  $w_1 > w_0$ . From the proof of Proposition 1, the new equilibrium again takes the form of a threshold  $\tau'$  that satisfies

$$G(\tau') = -w_1,$$

where  $G(\cdot)$  is defined in (5). Since  $\frac{d}{d\tau} G(\tau) > 0$  and  $-w_1 < -w_0$ , we have  $\tau' < \tau$  (and furthermore, since  $\frac{d}{d\tau} G(\tau) < 1$ , the difference  $\tau - \tau' > w_1 - w_0$ , so the decrease in threshold  $\tau$  is larger than the increase in  $w$ ). Now,  $\tau' < \tau$  implies that the share of citizens

choosing  $d_i = 1$  has increased  $(1 - H(\tau') > 1 - H(\tau))$ . Lastly, the social cost is now equal  $\gamma \mathbb{E}(t_i | t_i > \tau') < \gamma \mathbb{E}(t_i | t_i > \tau)$ , so it is lower than without the signal  $s$ .

Now consider an increase in  $\sigma_\varepsilon^2$ , to  $\tilde{\sigma}_\varepsilon^2$ . The new expectation of  $w$  will be  $\tilde{w}_1$  that satisfies

$$\tilde{w}_1 - w_0 = \frac{\sigma_w^2}{\sigma_w^2 + \tilde{\sigma}_\varepsilon^2} (s - w_0),$$

and since  $s > w_0$ , we have  $w_1 > \tilde{w}_1 > w_0$ , with  $\tilde{w}_1 \rightarrow w_0$  as  $\tilde{\sigma}_\varepsilon^2 \rightarrow \infty$ . By monotonicity, we have that the new equilibrium threshold  $\tilde{\tau}$  satisfies  $\tau' < \tilde{\tau} < \tau$ , which by the same argument implies that the share of citizens choosing  $d_i = 1$  increases by a smaller amount (vanishing if  $\tilde{\sigma}_\varepsilon^2 \rightarrow \infty$ ), and the same is true about the increase in the social cost. This completes the proof. ■

### A.3 Proof of Proposition 3

We start by establishing the uniqueness of equilibrium in this case.<sup>35</sup> Let  $\bar{S}$  be the social cost of choosing  $d_i = 1$  in a hypothetical equilibrium. Then the citizen would choose  $d_i = 1$  if  $t_i > \frac{1}{\beta} \bar{S} - w_h$  following signal  $s_h$  and if  $t_i > \frac{1}{\beta} \bar{S} - w_l$  following signal  $s_l$ . This implies that there are two thresholds,  $\tau_h$  and  $\tau_l$ , that satisfy  $\tau_l - \tau_h = w_h - w_l$ . Denote  $\bar{\tau} = \frac{1}{\beta} \bar{S} - w_0$ ; then  $\tau_h = \bar{\tau} + w_0 - w_h$  and  $\tau_l = \bar{\tau} + w_0 - w_l$ . From now on we describe the equilibrium in terms of  $\bar{\tau}$ .

In what follows, we use the following probabilities. We denote

$$p(x, y) = \mu(1 - H(x)) + (1 - \mu)(1 - H(y)),$$

so

$$p(\bar{\tau} + w_0 - w_h, \bar{\tau} + w_0 - w_l) = p\left(\frac{1}{\beta} \bar{S} - w_h, \frac{1}{\beta} \bar{S} - w_l\right)$$

is the probability of choosing  $d_i = 1$  if the citizen faces social cost  $\bar{S}$ . We also let

$$q(x, y) = \frac{\mu(1 - H(x))}{p(x, y)},$$

so  $q(\bar{\tau} + w_0 - w_h, \bar{\tau} + w_0 - w_l)$  is the equilibrium conditional probability that citizen  $i$  got signal  $s_h$  conditional on choosing  $d_i = 1$ .

---

<sup>35</sup>Notice first that our assumption of rational expectation of  $t_i$  conditional on  $d_i = 1$  allows us to bypass the discussion of whether members of the audience get signals  $s_l$ ,  $s_h$ , or both. Rational expectation can be formed in practice if people had prior interactions with those who choose  $d_i = 1$  and learned their type, which allows them to make a correct expectation in equilibrium about individuals who choose  $d_i = 1$  with a given piece of evidence. An alternative way is to assume that the audience is sophisticated, understands the whole signal structure, but does not know which signal citizen  $i$  got, and faces the signal decomposition problem as a result.

Define the function

$$\begin{aligned}\bar{S}(z) &= \gamma q(z + w_0 - w_h, z + w_0 - w_l) \mathbb{E}(t_i \mid t_i > z + w_0 - w_h) \\ &\quad + \gamma(1 - q(z + w_0 - w_h, z + w_0 - w_l)) \mathbb{E}(t_i \mid t_i > z + w_0 - w_l).\end{aligned}$$

In equilibrium characterized by  $\bar{\tau}$ , the social cost of choosing  $d_i = 1$  equals  $\bar{S}(\bar{\tau})$ . Given the above, thresholds  $\tau_h = \bar{\tau} + w_0 - w_h$  and  $\tau_l = \bar{\tau} + w_0 - w_l$  are equilibrium thresholds for choosing  $d_i = 1$  after getting signals  $s_h$  and  $s_l$ , respectively, if and only if  $\bar{\tau}$  solves the equation

$$\bar{\tau} - \frac{1}{\beta} \bar{S}(\bar{\tau}) = -w_0. \quad (6)$$

Let us show that  $\frac{d}{dz} \frac{1}{\gamma} \bar{S}(z) \leq 1$ . Indeed, from the proof of Proposition 1, we have

$$\begin{aligned}\frac{d}{dz} \mathbb{E}(t_i \mid t_i > z + w_0 - w_h) &\leq 1; \\ \frac{d}{dz} \mathbb{E}(t_i \mid t_i > z + w_0 - w_l) &\leq 1.\end{aligned}$$

Furthermore,

$$\mathbb{E}(t_i \mid t_i > z + w_0 - w_l) > \mathbb{E}(t_i \mid t_i > z + w_0 - w_h).$$

Lastly, we have

$$\begin{aligned}q(z + w_0 - w_h, z + w_0 - w_l) &= \frac{\mu(1 - H(z + w_0 - w_h))}{\mu(1 - H(z + w_0 - w_h)) + (1 - \mu)(1 - H(z + w_0 - w_l))} \\ &= \frac{1}{1 + \frac{1-\mu}{\mu} \frac{1-H(z+w_0-w_l)}{1-H(z+w_0-w_h)}}.\end{aligned}$$

Now,

$$\frac{d}{dz} \frac{1 - H(z + w_0 - w_l)}{1 - H(z + w_0 - w_h)} = \frac{d}{du} \frac{1 - H(u + (w_h - w_l))}{1 - H(u)} \leq 0,$$

where we denoted  $u = z + w_0 - w_h$  and used the calculation (3) from the proof of Proposition 1. This immediately implies that  $\frac{d}{dz} q(z + w_0 - w_h, z + w_0 - w_l) \geq 0$ . Summing up, we have

$$\begin{aligned}\frac{d}{dz} \frac{1}{\gamma} \bar{S}(z) &= q(z + w_0 - w_h, z + w_0 - w_l) \frac{d}{dz} \mathbb{E}(t_i \mid t_i > z + w_0 - w_h) \\ &\quad + (1 - q(z + w_0 - w_h, z + w_0 - w_l)) \frac{d}{dz} \mathbb{E}(t_i \mid t_i > z + w_0 - w_l) \\ &\quad + \left( \frac{d}{dz} q(z + w_0 - w_h, z + w_0 - w_l) \right) \\ &\quad \times (\mathbb{E}(t_i \mid t_i > z + w_0 - w_h) - \mathbb{E}(t_i \mid t_i > z + w_0 - w_l)).\end{aligned}$$

Notice that the sum of the first two lines does not exceed 1 (since both derivatives do not exceed 1), and term on the third line is positive and the one on the fourth is negative, so their product is negative. This proves that  $\frac{d}{dz} \frac{1}{\gamma} \bar{S}(z) \leq 1$ . Now, as in the proof of Proposition 1 this implies that the equation (6) has a unique solution  $\bar{\tau}$ , which proves the uniqueness of equilibrium in this case.

Let us now show that in this solution,  $\bar{\tau} < \tau$  and  $\bar{S}(\bar{\tau}) < S(\tau)$ , where  $S(\tau) = \frac{1}{\gamma} \mathbb{E}(t_i | t_i > \tau)$  is the equilibrium social cost in the absence of any signal, in the unique solution  $\tau$ . To do this, it is sufficient to show that  $\bar{S}(\tau) < S(\tau)$ . Indeed, this would imply that

$$\tau - \frac{1}{\beta} \bar{S}(\tau) > \tau - \frac{1}{\beta} S(\tau) = -w_0,$$

and since  $\bar{\tau}$  satisfies (6) and the function  $x - \frac{1}{\beta} \bar{S}(x)$  is increasing, we would get  $\bar{\tau} < \tau$ . Then we would get

$$\bar{S}(\bar{\tau}) = \beta(\bar{\tau} + w_0) < \beta(\tau + w_0) = S(\tau),$$

as required. So, to complete the proof, we need to show that  $\bar{S}(\tau) < S(\tau)$ .

In the light of condition (2) and by continuity of  $H(\cdot)$ , there exists  $\hat{w}_h \in (0, w_h)$  such that

$$\mu(H(\tau) - H(\tau - (\hat{w}_h - w_0))) = (1 - \mu)(H(\tau + (w_0 - w_l)) - H(\tau)).$$

Let  $\hat{S}$  denote the value

$$\begin{aligned} \hat{S} &= \gamma q(\tau + w_0 - \hat{w}_h, \tau + w_0 - w_l) \mathbb{E}(t_i | t_i > \tau + w_0 - \hat{w}_h) \\ &\quad + \gamma(1 - q(\tau + w_0 - \hat{w}_h, \tau + w_0 - w_l)) \mathbb{E}(t_i | t_i > \tau + w_0 - w_l); \end{aligned}$$

in other words, the expression for  $\hat{S}$  is analogous to  $\bar{S}(\tau)$ , except that  $w_h$  is replaced by  $\hat{w}_h$ .

We now show that  $\bar{S}(\tau) < \hat{S} < S(\tau)$ . To prove the first inequality, we use some algebra to establish that

$$\frac{1}{\gamma} \bar{S}(\tau) = (1 - \rho) \frac{1}{\gamma} \hat{S} + \rho \mathbb{E}(t_i | t_i \in (\tau + w_0 - w_h, \tau + w_0 - \hat{w}_h)),$$

where

$$\rho = q(\tau + w_0 - w_h, \tau + w_0 - w_l) \frac{H(\tau + w_0 - \hat{w}_h) - H(\tau + w_0 - w_h)}{1 - H(\tau + w_0 - w_h)}.$$

Since  $\rho > 0$  and  $\frac{1}{\gamma} \hat{S} < \mathbb{E}(t_i | t_i \in (\tau + w_0 - w_h, \tau + w_0 - \hat{w}_h))$  as the former is an expectation taken over values to the right of  $\tau + w_0 - \hat{w}_h$  while the latter expectation is taken over values to the left of that point, we get  $\bar{S}(\tau) < \hat{S}$ .

Let us now prove that  $\hat{S} < S(\tau)$ . Spelling out  $q(\tau + w_0 - \hat{w}_h, \tau + w_0 - w_l)$  and expec-

tations in the definition of  $\hat{S}$ , we have

$$\begin{aligned} \frac{1}{\gamma} \left( S(\tau) - \hat{S} \right) &= \frac{\int_{\tau}^{\infty} x h(x) dx}{1 - H(\tau)} \\ &\quad - \frac{\mu \int_{\tau+w_0-\hat{w}_h}^{\infty} x h(x) dx + (1-\mu) \int_{\tau+w_0-w_l}^{\infty} x h(x) dx}{\mu (1 - H(\tau + w_0 - \hat{w}_h)) + (1-\mu) (1 - H(\tau + w_0 - w_l))}. \end{aligned}$$

Notice that by the definition of  $\hat{w}_h$  the denominators in both terms are equal, hence  $S(\tau) - \hat{S}$  has the same sign as

$$\begin{aligned} &\int_{\tau}^{\infty} x h(x) dx - \left( \mu \int_{\tau+w_0-\hat{w}_h}^{\infty} x h(x) dx + (1-\mu) \int_{\tau+w_0-w_l}^{\infty} x h(x) dx \right) \\ &= (1-\mu) \int_{\tau}^{\tau+w_0-w_l} x h(x) dx - \mu \int_{\tau+w_0-\hat{w}_h}^{\tau} x h(x) dx \\ &= (1-\mu) (H(\tau + w_0 - w_l) - H(\tau)) \mathbb{E}(t_i \mid t_i \in (\tau, \tau + w_0 - w_l)) \\ &\quad - \mu (H(\tau) - H(\tau + w_0 - \hat{w}_h)) \mathbb{E}(t_i \mid t_i \in (\tau + w_0 - \hat{w}_h, \tau)). \end{aligned}$$

Since the coefficients in front of the expectations in the last two lines are the same (again, by the choice of  $\hat{w}_h$ ), the sign of this expression is the same as the sign of

$$\mathbb{E}(t_i \mid t_i \in (\tau, \tau + w_0 - w_l)) - \mathbb{E}(t_i \mid t_i \in (\tau + w_0 - \hat{w}_h, \tau)),$$

which is positive, because the first term is greater than  $\tau$  and the second is less than that. Therefore,  $\hat{S} < S(\tau)$ .

We have thus proved that  $\bar{S}(\tau) < \hat{S} < S(\tau)$  which, as we showed earlier, implies the results stated. This completes the proof. ■

## B Additional Details on Experiments

**Table B.1:** Overview of Data Collections

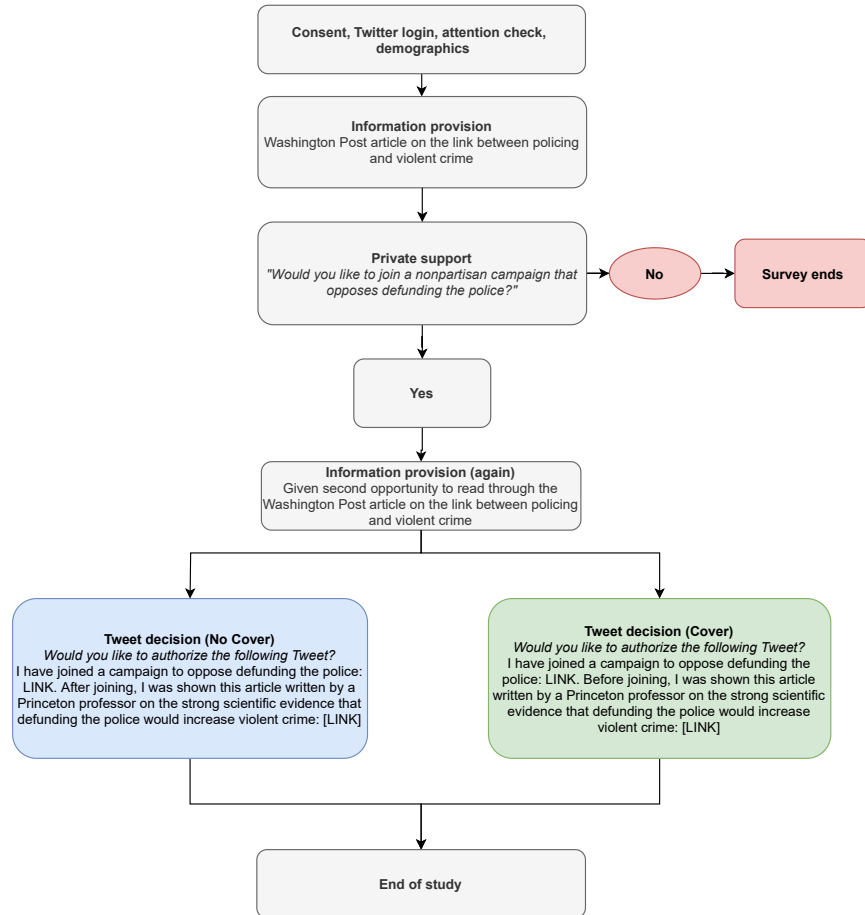
Experiment	Provider	Dates
<b>Panel A:</b> Main Experiments		
Experiment 1: Willingness to post anti-defunding Tweet – Democrats/Independents authorizing Twitter access (N=1,122)	Luc.id, Cloudresearch	October 2021
Experiment 2: Interpretation of anti-defunding Tweet – Democrats/Independents (N=1,040)	Prolific	November 2021
Experiment 3: Willingness to post pro-deportation Tweet – Republicans/Independents authorizing Twitter access (N=1,130)	Luc.id	March 2021
Experiment 4: Interpretation of pro-deportation Tweet – Democrats/Independents (N=1,082)	Prolific	November 2021
<b>Panel B:</b> Auxiliary Experiments		
Auxiliary Experiment 1: Persuasiveness of anti-defunding article – Democrats (N=1,008)	Prolific	December 2021
Auxiliary Experiment 2: Placebo: willingness to post pro-conservation Tweet – respondents authorizing Twitter access (N=483)	Luc.id, Cloudresearch	December 2021 and January 2022
Auxiliary Experiment 3: Anticipated persuasiveness of anti-defunding Tweet – Democrats (N=501)	Prolific	November 2021
Auxiliary Experiment 4: Motives underlying the choice – Democrats with Twitter account (N=402)	Prolific	January 2022
Auxiliary Experiment 5: Credibility and social cover – Democrats (N=1,017)	Luc.id	July 2022
Auxiliary Experiment 6: Interpretation of lower-credibility anti-defunding Tweet – Democrats/Independents (N=506)	Prolific	November 2021
Auxiliary Experiment 7: Persuasiveness of pro-deportation Tweet – Republicans (N=2,012)	Prolific, Lucid	December 2021

*Notes:* Reported sample sizes for Experiment 1, Experiment 3, and Auxiliary Experiment 2 include respondents who chose not to join the campaigns and therefore are not included in the sample we analyze.

## B.1 Anti-Defunding Experiments

### B.1.1 Experiment 1: Additional Figures and Tables

Figure B.1: Experiment 1: flow of dissent design





**Table B.2:** Experiment 1: Balance of covariates

	Overall		Cover	No Cover	$p$ -value
	mean	std. dev.	mean	mean	(C = NC)
Age	39.811	15.127	39.249	40.424	0.375
Black	0.214	0.411	0.231	0.196	0.334
Asian	0.076	0.266	0.073	0.080	0.773
White	0.671	0.470	0.667	0.676	0.821
Hispanic	0.185	0.389	0.161	0.212	0.136
Male	0.585	0.493	0.579	0.592	0.759
High school diploma	0.975	0.156	0.974	0.976	0.904
Bachelors degree	0.436	0.496	0.421	0.452	0.480

*Notes:*  $p$ -values based on robust standard errors reported.

**Table B.3:** Experiment 1: Sample representativeness

	Defund	Pew (Inds and Dems)
Age	39.81	45.86
Black	0.21	0.18
White	0.67	0.59
Asian	0.08	0.05
Hispanic	0.19	0.15
Male	0.59	0.46
High school diploma	0.98	0.89
Bachelors degree or higher	0.44	0.35
Observations	523	6,627

*Notes:* Table displays mean characteristics, comparing the experimental sample with the 2018 Pew Research Center’s American Trends Panel, Wave 39. Attriters are dropped from sample.

### B.1.2 Auxiliary Experiment 1: Persuasiveness of Defunding Rationale

We conducted this pre-registered experiment in December 2021 with a sample of 1,008 Democrats and Independents recruited from Prolific.<sup>36</sup> After completing a set of demographic questions, respondents assigned to the treatment group read Sharkey’s article in

<sup>36</sup>The pre-registration is available in the AEA RCT registry under ID AEARCTR-0008624.

the *Washington Post*, while respondents assigned to the control group did not read the article. They then respond to the following two questions: “Do you think that funding for the police should be increased, decreased, or stay the same?” and “How do you think increasing funding for the police would affect violent crime?”. We code both questions from -2 (“Decreased a lot” and “Strongly decrease violent crime”, respectively) to 2 (“Increased a lot” and “Strongly increase violent crime”, respectively).

Table B.4 displays results, with Columns 1–3 corresponding to the first measure and Columns 4–6 corresponding to the second measure. We find a significant effect on both measures, though effects are weaker for policy preferences and are no longer significant once we control for demographics and partisan affiliation. Effect of credibility on social cover

**Table B.4:** Persuasive effects of anti-defunding article

	<i>Belief</i>			<i>Policy preference</i>		
	(1)	(2)	(3)	(4)	(5)	(6)
Provided article	-0.236*** (0.056)	-0.245*** (0.055)	-0.223*** (0.054)	0.135* (0.071)	0.121* (0.068)	0.072 (0.060)
No Article mean	0.038	0.038	0.038	-0.636	-0.636	-0.636
Observations	1,008	1,007	1,004	1,008	1,007	1,004
Demographic controls	No	Yes	Yes	No	Yes	Yes
Partisan controls	No	No	Yes	No	No	Yes

*Notes:* Table reports results from Auxiliary Experiment 1. The dependent variable in Columns 1–3 is the respondent’s reported belief as to the effect of increasing funding for the police on violence crime, coded between -2 (“Strongly decrease violent crime”) and 2 (“Strongly increase violent crime”). The dependent variable in Columns 4–6 is the respondent’s reported preference for changing police funding, ranging from -2 (“Decreased a lot”) to 2 (“Increased a lot”). Demographic controls include age, age squared, a set of race indicators, a Hispanic indicator, a male indicator, a set of education indicators. Partisan controls include indicators for “Very conservative”, “Conservative”, “Neither liberal nor conservative” (omitted), “Liberal”, and “Very liberal”. Robust standard errors are reported.

### B.1.3 Auxiliary Experiment 2: Rainforest Placebo

We conducted this experiment in December 2021 and January 2022 with a sample of 483 Democrats and Independents recruited from Luc.id and CloudResearch. Respondents logged in to the survey with their Twitter accounts using the same procedure as in Experiment 1. The design is similar to that of Experiment 1, but examines a different (non-stigmatized) context: willingness to post a Tweet supporting efforts to conserve the Amazon rainforest. Rather than reading an article about the likely effects of defunding the police, respondents read a Reuters article reporting on a study conducted by the Science Panel for the Amazon which finds that over 10,000 species are at risk from deforestation

in the Amazon (Science Panel for the Amazon, 2021). The *Cover* Tweet reads:

I’ve joined a campaign to immediately stop the destruction of the Amazon rainforest! Before I joined the campaign, I was shown this article about how 10,000 species risk extinction in Amazon: [LINK]. Join the campaign and sign the petition: [LINK].

The *No Cover* Tweet is identical, but replaces “Before I joined the campaign...” with “After I joined the campaign...”.

As shown in Panel A of Table 2, we find no significant difference between posting rates in the *Cover* and *No Cover* conditions. The difference in effect sizes between the defunding experiment and the placebo experiment is large in magnitude (16 percentage points) and significant at the 5% level, suggesting effects are indeed driven by (anticipated) changes in the stigma associated with dissenting expression rather than some other independent effect of the wording.

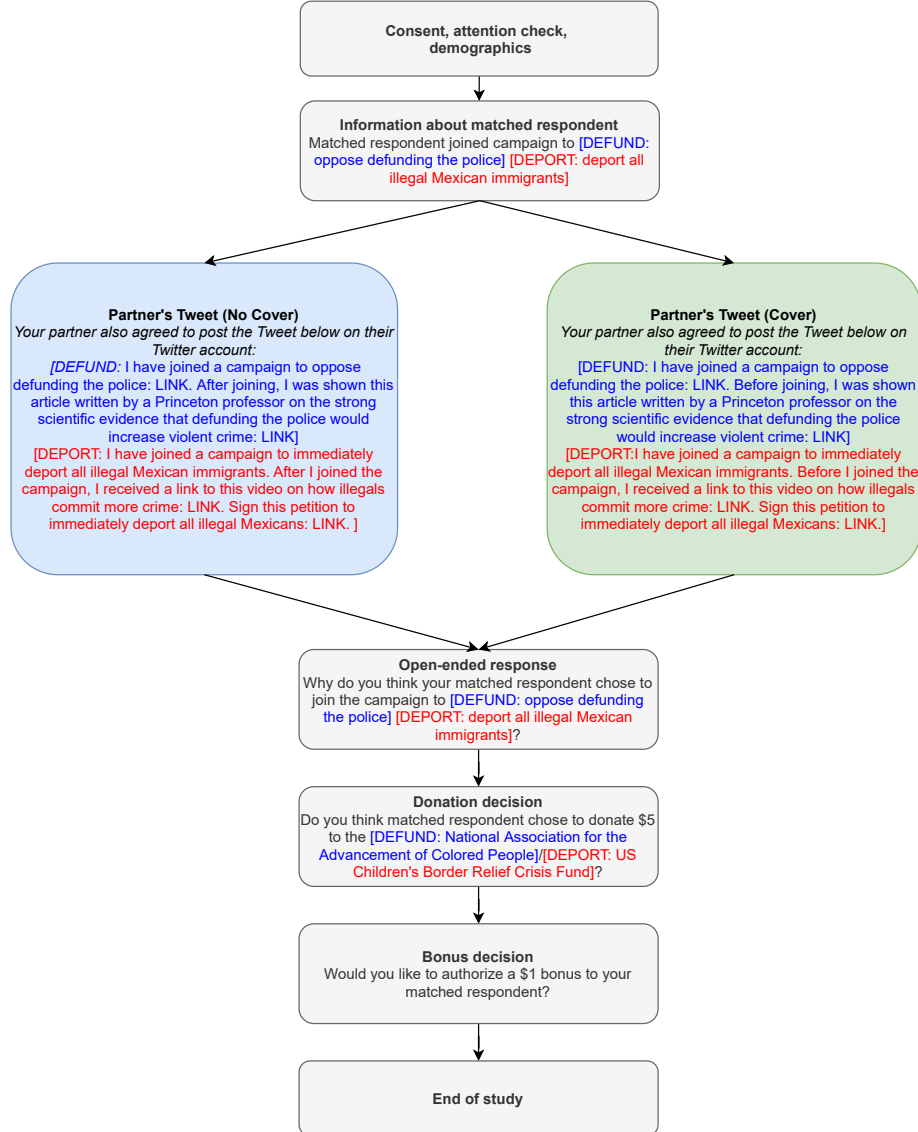
#### **B.1.4 Auxiliary Experiment 3: Anticipated Persuasion Experiment**

We conducted this experiment in November 2021 with a sample of 501 Democrats and Independents recruited from Prolific. Only Democrats and Independents with Twitter accounts were eligible to take the survey. After completing a set of demographic questions, respondents read Sharkey’s article in the *Washington Post*. As in Experiment 1, respondents are asked if they would like to join the campaign to oppose the movement to defund the police, only those who indicate that they would like to join the campaign proceed with the experiment, and those who do proceed are given a chance to re-read the article. They are then randomly shown either the *Cover* or the *No Cover* Tweet from Experiment 1 and are asked: “Suppose you posted the Tweet above on your account. If you had to guess, what percentage of people who saw your Tweet would choose to join the campaign to oppose defunding the police?”

Panel B of Table 2 displays results. Reassuringly, we find no significant difference between the anticipated persuasiveness of the Tweets, suggesting that differential posting rates are instead driven by changes in anticipated stigma.

### B.1.5 Experiment 2: Additional Figures and Tables

Figure B.2: Experiments 2 and 4: flow of inference design



*Notes:* Experiments 2 and 4 have identical structures, so we present both experiments jointly. Blue text corresponds to Experiment 2, studying opposition to the movement to defund the police; red text corresponds to Experiment 4, studying support for immediately deporting all illegal Mexican immigrants.

**Table B.5:** Experiment 2: Balance of covariates

	Overall		Cover	No Cover	$p$ -value
	mean	std. dev.	mean	mean	(C = NC)
Age	30.725	11.258	30.686	30.763	0.912
Black	0.070	0.256	0.086	0.055	0.057
Asian	0.085	0.279	0.089	0.080	0.589
White	0.773	0.419	0.766	0.781	0.563
Hispanic	0.112	0.315	0.093	0.130	0.060
Male	0.374	0.484	0.384	0.365	0.522
High school diploma	0.997	0.054	0.996	0.998	0.552
Bachelors degree	0.572	0.495	0.562	0.582	0.520

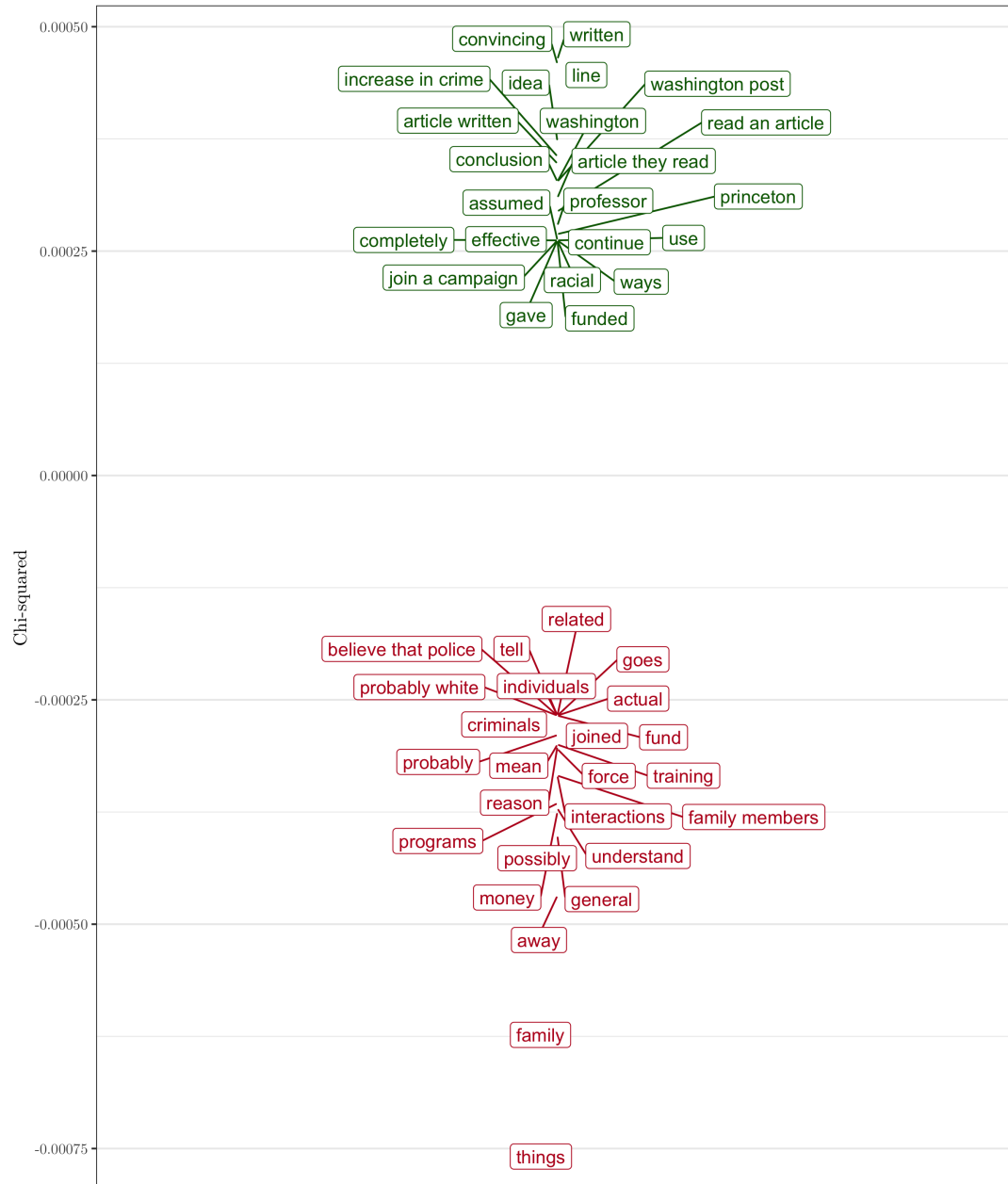
*Notes:*  $p$ -values based on robust standard errors reported.

**Table B.6:** Experiment 2: Sample representativeness

	Defund	Pew (Inds and Dems)
Age	30.73	45.86
Black	0.07	0.18
White	0.77	0.59
Asian	0.08	0.05
Hispanic	0.11	0.15
Male	0.37	0.46
High school diploma	1.00	0.89
Bachelors degree or higher	0.57	0.35
Observations	1,040	6,627

*Notes:* Table displays mean characteristics, comparing the experimental sample with the 2018 Pew Research Center's American Trends Panel, Wave 39. Attriters are dropped from sample.

**Figure B.3:** Experiment 2: most distinctive phrases in each condition



*Notes:* Appendix Figure B.3 plots phrases by their associated  $\chi^2$  statistic, limiting to the top 50 phrases and multiplying the  $\chi^2$  of phrases more characteristic of the “No Cover” condition by -1. The words “article” and “read” have  $\chi^2$  values of greater than 0.001 and have been suppressed to facilitate visualization of the remaining phrases.

**Table B.7:** Experiments 2 and 4: Heterogeneity by partisan affiliation

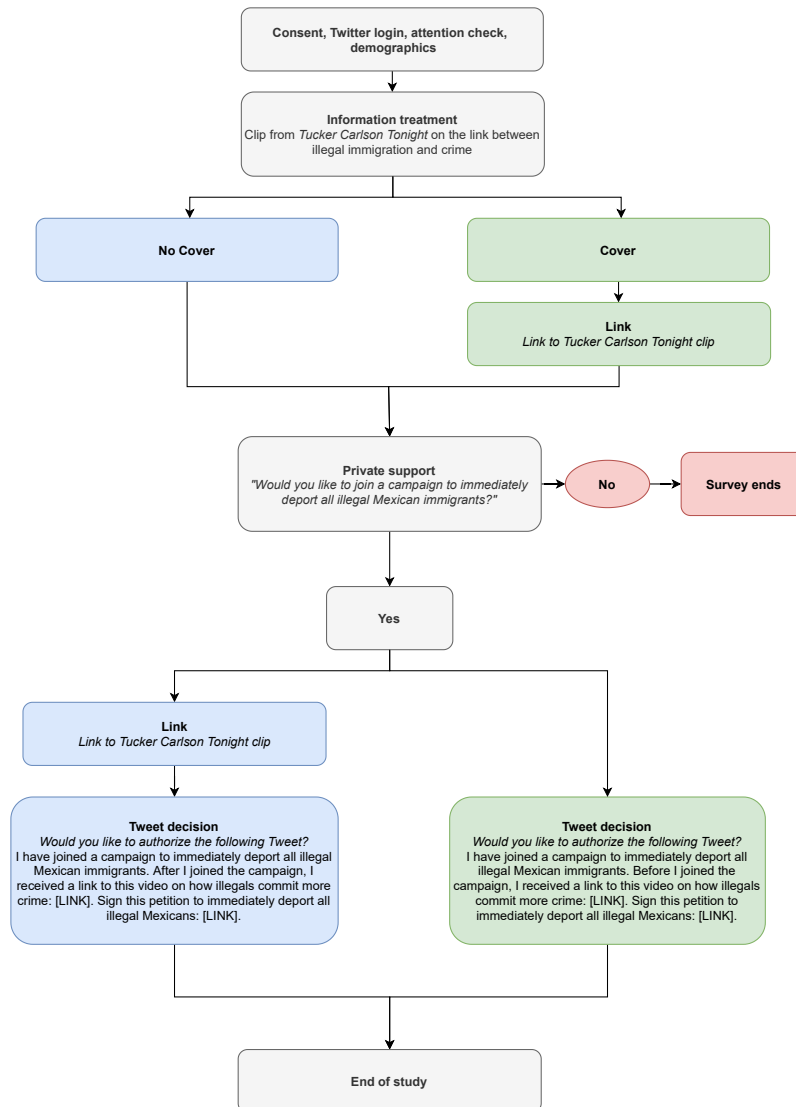
	<i>Belief partner donated</i>		<i>Denied bonus to partner</i>	
	(1)	(2)	(3)	(4)
<b>Panel A:</b>	Experiment 2			
Cover	0.058 (0.067)	0.068 (0.067)	−0.054 (0.071)	−0.056 (0.071)
Cover × Strong partisan	0.013 (0.074)	0.001 (0.074)	−0.021 (0.079)	−0.019 (0.079)
Strong partisan	−0.155*** (0.054)	−0.139** (0.055)	0.207*** (0.057)	0.172*** (0.058)
No Cover mean	0.273	0.273	0.471	0.471
Observations	1,037	1,036	1,037	1,036
<b>Panel B:</b>	Experiment 4			
Cover	0.012 (0.046)	0.014 (0.046)	−0.036 (0.062)	−0.015 (0.062)
Cover × Strong partisan	0.042 (0.050)	0.042 (0.051)	−0.032 (0.068)	−0.057 (0.069)
Strong partisan	−0.125*** (0.037)	−0.105*** (0.037)	0.120** (0.050)	0.109** (0.050)
No Cover mean	0.085	0.085	0.803	0.803
Observations	1,081	1,081	1,081	1,081
Demographic controls	No	Yes	No	Yes

*Notes:* The dependent variable in Columns 1–2 is an indicator taking value 1 if the respondent reports believing that his or her matched partner donated to the US Border Crisis Children’s Relief Fund; The dependent variable in Columns 3–4 taking value 1 if the respondent denied his or her matched partner a \$1 bonus. Demographic controls include age, age squared, a set of race indicators, a Hispanic indicator, a male indicator, a set of education indicators. Strong partisan is an indicator taking value 1 if the respondent is “Liberal” or “Very liberal”. Robust standard errors are reported. Panel A presents the results of Experiment 2; Panel B presents the results of Experiment 4.

## B.2 Anti-Immigrant Experiments

### B.2.1 Experiment 3: Additional Figures and Tables

Figure B.4: Experiment 3: design





**Table B.8:** Experiment 3: Balance of covariates

	Overall		Cover	No Cover	$p$ -value
	mean	std. dev.	mean	mean	(C = NC)
Age	49.226	13.550	48.510	49.904	0.247
Black	0.012	0.108	0.012	0.011	0.946
Asian	0.016	0.125	0.016	0.015	0.938
White	0.955	0.208	0.951	0.958	0.728
Hispanic	0.065	0.247	0.053	0.077	0.274
Male	0.504	0.500	0.490	0.517	0.538
High school diploma	0.994	0.077	0.996	0.992	0.596
Bachelors degree	0.380	0.486	0.340	0.418	0.072

Notes:  $p$ -values based on robust standard errors reported.

**Table B.9:** Experiment 3: Sample representativeness

	Deport	Pew (Inds and Reps)
Age	49.23	47.17
Black	0.01	0.05
White	0.95	0.77
Asian	0.02	0.03
Hispanic	0.06	0.11
Male	0.50	0.52
High school diploma	0.99	0.93
Bachelors degree or higher	0.38	0.31
Observations	508	5,501

Notes: Table displays mean characteristics, comparing the experimental sample with the 2018 Pew Research Center's American Trends Panel, Wave 39. Attriters are dropped from sample.

## B.2.2 Auxiliary Experiment 7: Persuasiveness of Deportation Rationale

We conducted a first pre-registered experiment in December 2021 with a sample of 1,008 Republicans recruited from Prolific.<sup>37</sup> After completing a set of demographic questions, respondents assigned to the treatment group viewed the clip from *Tucker Carlson Tonight*, while respondents assigned to the control group did not view the clip. They then indicated

<sup>37</sup>The pre-registration is available in the AEA RCT registry under ID AEARCTR-0008624.

their agreement with the following two statements: “Illegal immigrants are not much more likely to commit serious crimes than U.S. citizens” (beliefs) and “The US should immediately deport all illegal Mexican immigrants” (policy preference). We code both questions from -2 (“Strongly disagree”) to 2 (“Strongly agree”).

Panel A of Table B.10 displays results. While we found significant effects on the beliefs outcome, we found no treatment effects on the policy preference outcome. Two logistical problems complicate interpretation of this result. First, when setting up the survey, we forgot to exclude respondents from some previous experiments which included the video. Thus, some respondents in the *Control* condition had seen the video in previous experiments. Second, there was a highly limited sample of Republicans available on Prolific (fewer than 2000 who met our screening criteria), and we had to pay a higher than usual rate in order to meet our pre-registered sample size. This potentially induced selection into the survey.

We thus ran the same experiment on Lucid, with the same sample restrictions. , with Columns 1–3 corresponding to the first measure and Columns 4–6 corresponding to the second measure. We find a significant effect on both measures, with an effect size of around 0.12 standard deviations for the first outcome and 0.18 standard deviations for the second outcome.

Overall, we take the evidence for the effects of the clip on persuasion as mixed.

**Table B.10:** Persuasive effects of *Tucker Carlson Tonight* video

	<i>Belief</i>			<i>Policy preference</i>		
	(1)	(2)	(3)	(4)	(5)	(6)
<b>Panel A:</b>	<i>Prolific Sample</i>					
Provided article	0.533*** (0.063)	0.538*** (0.063)	0.544*** (0.061)	−0.005 (0.074)	−0.025 (0.073)	−0.023 (0.070)
No Article mean	0.300	0.300	0.300	0.541	0.541	0.541
Observations	1,008	1,008	1,008	1,008	1,008	1,008
<b>Panel B:</b>	<i>Lucid Sample</i>					
Provided article	0.751*** (0.066)	0.743*** (0.066)	0.761*** (0.064)	0.177** (0.074)	0.179** (0.073)	0.192*** (0.071)
No Article mean	0.251	0.251	0.251	0.652	0.652	0.652
Observations	1,004	1,002	1,002	1,004	1,002	1,002
Demographic controls	No	Yes	Yes	No	Yes	Yes
Partisan controls	No	No	Yes	No	No	Yes

*Notes:* Table reports results from Auxiliary Experiment 7. The dependent variable in Columns 1–3 is the respondent’s reported agreement with the statement “Illegal immigrants are more likely to commit serious crimes than US citizens,” coded between -2 (“Strongly disagree”) and 2 (“Strongly agree”). The dependent variable in Columns 4–6 is the respondent’s reported agreement with the statement “The US should immediately deport all illegal Mexican immigrants,” ranging from -2 (“Strongly disagree”) to 2 (“Strongly agree”). Demographic controls include age, age squared, a set of race indicators, a Hispanic indicator, a male indicator, a set of education indicators. Partisan controls include indicators for “Very conservative”, “Conservative”, “Neither liberal nor conservative” (omitted), “Liberal”, and “Very liberal”. Robust standard errors are reported. Panel A uses the sample from Prolific, and Panel B uses the sample from Lucid.

### B.2.3 Experiment 4: Additional Figures and Tables

**Table B.11:** Experiment 4: Balance of covariates

	Overall		Cover	No Cover	$p$ -value
	mean	std. dev.	mean	mean	(C = NC)
Age	31.729	12.256	32.408	31.046	0.068
Black	0.069	0.254	0.063	0.076	0.389
Asian	0.100	0.300	0.090	0.109	0.297
White	0.767	0.423	0.785	0.750	0.174
Hispanic	0.118	0.323	0.109	0.128	0.325
Male	0.479	0.500	0.492	0.466	0.392
High school diploma	0.995	0.068	0.994	0.996	0.659
Bachelors degree	0.589	0.492	0.590	0.588	0.939

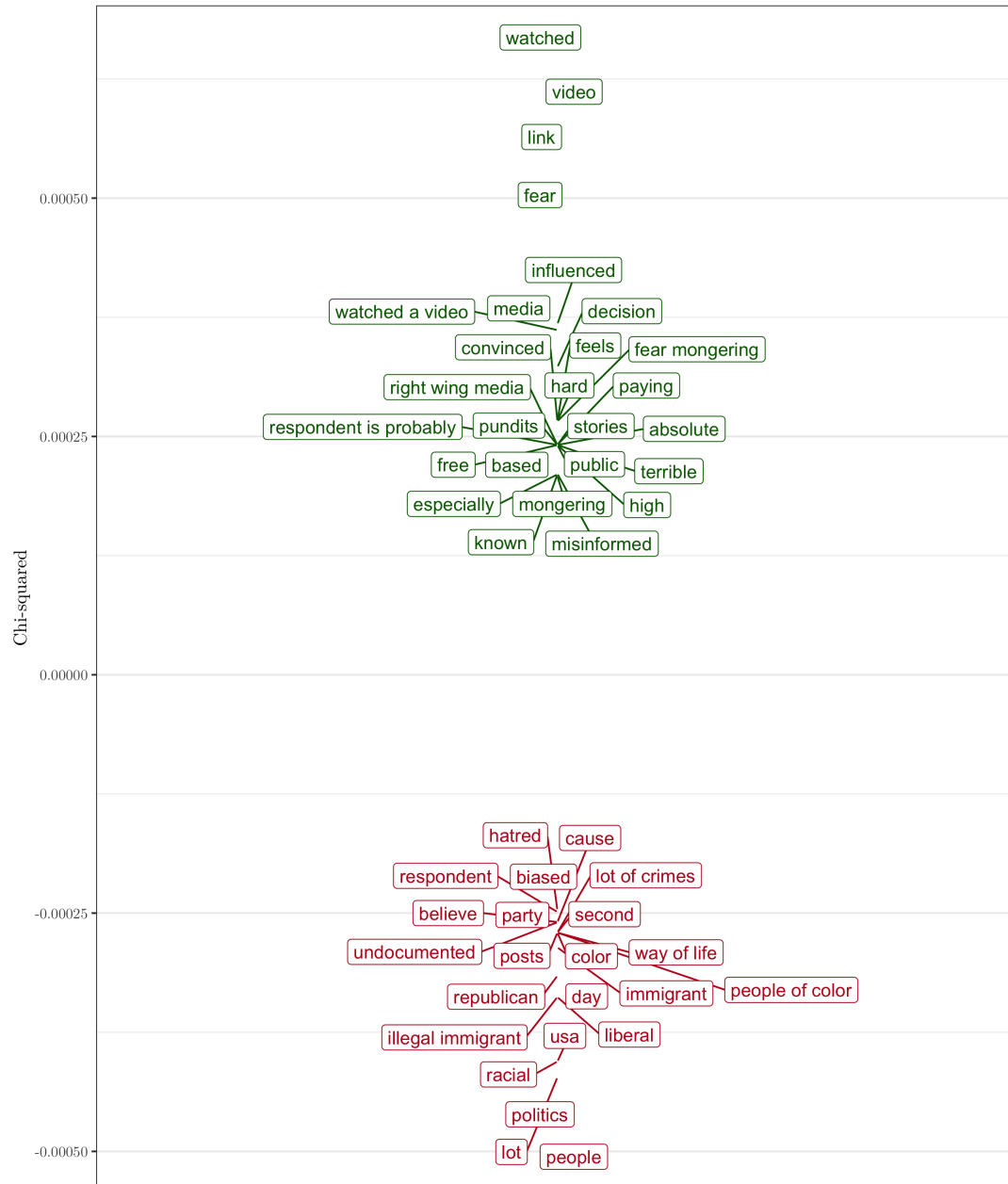
*Notes:*  $p$ -values based on robust standard errors reported.

**Table B.12:** Experiment 4: Sample representativeness

	Deport	Pew (Inds and Dems)
Age	31.73	45.86
Black	0.07	0.18
White	0.77	0.59
Asian	0.10	0.05
Hispanic	0.12	0.15
Male	0.48	0.46
High school diploma	1.00	0.89
Bachelors degree or higher	0.59	0.35
Observations	1,082	6,627

*Notes:* Table displays mean characteristics, comparing the experimental sample with the 2018 Pew Research Center's American Trends Panel, Wave 39. Attriters are dropped from sample.

**Figure B.5:** Experiment 4: most distinctive phrases in each condition



*Notes:* Appendix Figure B.5 plots phrases by their associated  $\chi^2$  statistic, limiting to the top 50 phrases and multiplying the  $\chi^2$  of phrases more characteristic of the “No Cover” condition by -1.

## C Ethical Considerations

Understanding dissenting expression is of great social importance. Identifying the drivers of xenophobic expression is crucial in designing policies best-suited to curbing it, while understanding barriers to dissenting expression in situations where such expression is desirable — for example, speaking out against unjust practices or systems — may help design contexts with lower such barriers.

Nonetheless, ethically conducting revealed-preference experiments on dissenting expression — particularly xenophobic expression — requires balancing three often contradictory objectives: avoiding explicitly deceiving respondents, avoiding compromising respondents’ privacy, and avoiding increasing public xenophobic expression. In this section, we provide a more detailed explanation of how our experimental designs balance these objectives. Of course, all experiments obtained approval from multiple Institutional Review Boards.

### C.1 Considerations related to information provision (Experiments 3–4)

The raw numbers pertaining to violent crime cited in the *Tucker Carlson Tonight* clip that we provide to respondents in Experiments 3–4 are taken from the U.S. Sentencing Commission and are factually correct. Nonetheless, the clip paints an incomplete picture of the academic literature, which generally finds null or negative effects of illegal immigration on violent crime. Although we do not endorse this evidence, we nonetheless debrief all respondents at the end of the study, providing them with an accessible academic overview of the link between illegal immigration and violent crime (Ousey and Kubrin, 2018) and a list of further readings. (The debriefing is strictly speaking unnecessary, as the numbers cited in the video clip are not factually wrong.)

### C.2 Considerations related to privacy and deception (Experiments 1 and 3)

Given that our mechanism examines the effect of perceived social stigma on behavior, it is crucial that respondents in Experiments 1 and 3 believe that their decisions will be visible to others. Although our experiments avoid explicit deception, protecting participants’ privacy and avoiding starting a political campaign in these contexts required us to mislead respondents. We distinguish between the ethical and practical problems associated with deception (the latter relating to concerns about subject pool contamination), addressing the first concern in this section and the second in Section C.3.

**Twitter login** All respondents were required to log in via their Twitter accounts to the “Tweetability” app we created. This app is governed by the Twitter API’s terms of service and has the second most restrictive set of permissions among the three application scopes Twitter provides (“Read” and “Write”). That is, the app does not have access to

users’ passwords, messages, or account settings, but it is able to post Tweets from the users’ accounts. We do not use this functionality in any way, and no information that could compromise users’ accounts is ever accessed or downloaded. We explicitly inform respondents of the app’s permissions in transparent language and give them the option to end the survey if they are uncomfortable granting the app these permissions. We also inform respondents that the app’s data, including the tokens that give us access to post on their accounts, will be deleted by no later than August 1, 2021 (Experiment 3) and December 1, 2021 (Experiment 1). Tokens were indeed deleted immediately after collection.

**Twitter posts** Our key outcome is whether respondents are willing to post a Tweet including a link to a petition to immediately deport all illegal Mexican immigrants. We were not willing to consider designs that asked respondents to actually post such Tweets. We thus asked them to “schedule” their Tweet for the future (using the Tweetability app), to be posted “if/when we have finished surveying people in all US counties”. Because we targeted fewer total respondents than the total number of US counties, these posts will never be published. This formulation is therefore misleading, even if it is not explicitly deceptive. Given our desire to avoid leading respondents to publicly post political content (particularly xenophobic content, as in Experiment 3) as part of our survey, we and our Institutional Review Board felt comfortable with this formulation.

### **C.3 Considerations related to subject pool contamination (Experiments 1 and 3)**

An important concern with deceptive or misleading experiments is that they can contaminate the subject pool by lowering trust in scientists and making respondents less likely to participate in future research studies. Of course, this can only happen if respondents know that they are being misled.

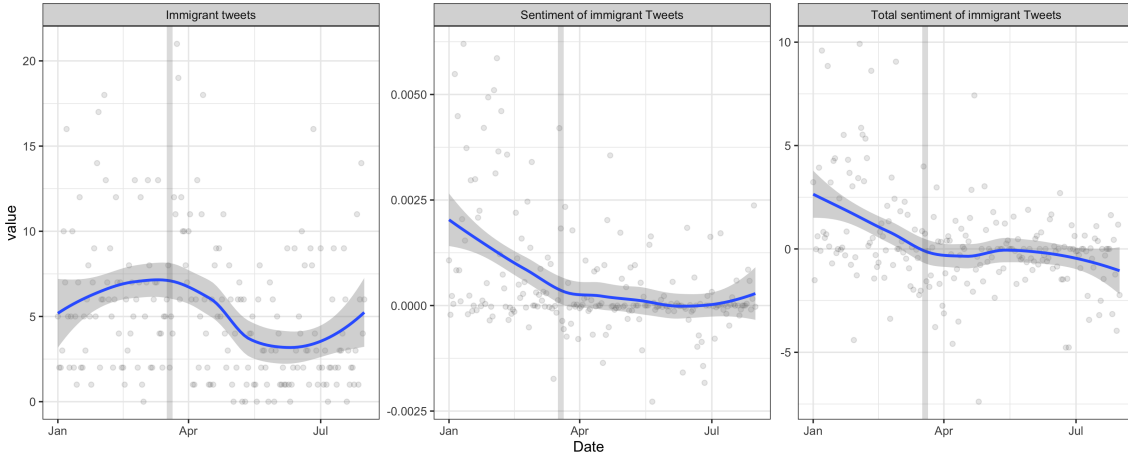
In Experiments 1 and 3, subjects are told we will post their Tweets when and if we reach survey respondents on all US counties before August 1, 2021 (Experiment 3) or December 1, 2021 (Experiment 1). Although we privately targeted fewer respondents than the number of US counties, ensuring that this condition would not be met, subjects do not know (and never learn) this is the case. In other words, it is not possible for respondents to know that they have been misled about the implementation of the main outcomes (unless they independently find our working paper). Furthermore, concerns about contaminating the experimental subject pool are most important in an economic lab with clear rules against deception. In online survey marketplaces, where survey participants are expected to regularly participate in studies by psychologists in which explicit deception is common, considerations about contaminating the subject pool are less relevant.

#### C.4 Considerations related to starting political Twitter campaigns (Experiments 1 and 3)

As discussed in Appendix C.2, we designed our experiment to ensure that none of the Tweets would ever be posted. It is of course possible that respondents independently posted political content on Twitter as a result of our experiment. This is a concern for Experiment 3, in which respondents were exposed to a clip presenting a misleading narrative about the link between illegal immigration and crime.

To examine whether this was the case, we accessed all Twitter posts made by respondents between the date of experimental collection and August 1, 2021 (the date by which we promised respondents that our access to their accounts and any Twitter-related data would be deleted). We used simple text analysis techniques to identify which posts concern immigrants and quantify the sentiment and content of these posts. The results of this analysis are presented in Figure C.1 and Table C.1. We find no evidence that respondents in our experiment begin posting more immigrant-related Tweets or more negative content about immigrants after participating (Figure C.1). Restricting to the period after the experiment, we find no evidence that respondents in the *Cover* condition post more or fewer Tweets in general, more or fewer Tweets specifically about immigrants, or more or less negative Tweets about immigrants than respondents in our *No Cover* condition (Table C.1). This evidence further strengthens our confidence that our experiment did not contribute to anti-immigrant discourse on social media.

**Figure C.1:** Twitter activity of respondents before and after experiment



*Notes:* Figure C.1 presents various measures of the Twitter activity of respondents before and after Experiment 3, conducted between March 17 and March 22, 2021 (shaded in a gray rectangle). The left panel of the figure presents the average number of immigrant-related Tweets; the middle panel the average sentiment of immigrant-related Tweets; and the right panel the total expressed sentiment of immigrant-related Tweets.



**Table C.1:** Subsequent Twitter behavior of respondents

	<i>Dependent variable:</i>					
	Tw. (1)	Tw. (w) (2)	Imm. Tw. (3)	Imm. Tw. (w) (4)	Imm. sent. (5)	Tot. imm. sent. (6)
Cover	-44.414 (29.941)	-9.298 (9.462)	-0.583 (0.416)	-0.152 (0.117)	0.005 (0.012)	0.024 (0.062)
Constant	80.075*** (20.862)	35.951*** (6.593)	0.970*** (0.290)	0.383*** (0.082)	0.003 (0.008)	-0.052 (0.043)
Observations	517	517	517	517	517	517

*Notes:* Table C.1 presents the results of our analysis of the subsequent Twitter behavior of the respondents in Experiment 3 between the end of our experiment and August 1, 2021. Table presents regressions of various measures of behavior on an indicator for whether the respondent was in the *Cover* condition: Columns 1 and 2 consider the total number of Tweets, Columns 3 and 4 the total number of immigrant-related Tweets, Column 5 the sentiment of immigrant-related Tweets, and Column 6 the sentiment of immigrant-related Tweets multiplied by the number of Tweets. Columns 2 and 4 winsorize the dependent variable at the 0.98 quantile.

## D Additional Exhibits for Auxiliary Experiments

This appendix reports balance and representativeness tables for all auxiliary experiments.

**Table D.1:** Auxiliary Experiment 1: Balance of covariates

	Overall		Article	No Article	<i>p</i> -value
	mean	std. dev.	mean	mean	(A = NA)
Age	36.955	14.084	37.315	36.598	0.420
Black	0.080	0.272	0.070	0.091	0.213
Asian	0.092	0.290	0.076	0.109	0.069
White	0.774	0.419	0.793	0.754	0.141
Hispanic	0.099	0.299	0.111	0.087	0.199
Male	0.494	0.500	0.489	0.499	0.752
High school diploma	0.996	0.063	0.992	1.000	0.044
Bachelors degree	0.632	0.483	0.618	0.646	0.357

*Notes:* *p*-values based on robust standard errors reported.

**Table D.2:** Auxiliary Experiment 1: Sample representativeness

	Defund	Pew (Inds and Dems)
Age	36.96	45.86
Black	0.08	0.18
White	0.77	0.59
Asian	0.09	0.05
Hispanic	0.10	0.15
Male	0.49	0.46
High school diploma	1.00	0.89
Bachelors degree or higher	0.63	0.35
Observations	1,008	6,627

*Notes:* Table displays mean characteristics, comparing the experimental sample with the 2018 Pew Research Center's American Trends Panel, Wave 39. Attriters are dropped from sample.

**Table D.3:** Auxiliary Experiment 2: Balance of covariates

	Overall		Cover	No Cover	$p$ -value
	mean	std. dev.	mean	mean	(C = NC)
Age	39.137	13.234	38.146	40.120	0.186
Black	0.137	0.344	0.172	0.101	0.068
Asian	0.044	0.206	0.025	0.063	0.104
White	0.768	0.423	0.739	0.797	0.219
Hispanic	0.159	0.366	0.166	0.152	0.740
Male	0.473	0.500	0.529	0.418	0.049
High school diploma	0.984	0.125	0.981	0.987	0.648
Bachelors degree	0.387	0.488	0.376	0.399	0.677

*Notes:*  $p$ -values based on robust standard errors reported.

**Table D.4:** Auxiliary Experiment 2: Sample representativeness

	Placebo	Pew (Inds, Dems and Reps)
Age	39.14	46.99
Black	0.14	0.13
White	0.77	0.67
Asian	0.04	0.04
Hispanic	0.16	0.13
Male	0.47	0.47
High school diploma	0.98	0.90
Bachelors degree or higher	0.39	0.33
Observations	315	9,506

*Notes:* Table displays mean characteristics, comparing the experimental sample with the 2018 Pew Research Center's American Trends Panel, Wave 39. Attriters are dropped from sample.

**Table D.5:** Auxiliary Experiment 3: Balance of covariates

	Overall		Cover	No Cover	$p$ -value
	mean	std. dev.	mean	mean	(C = NC)
Age	34.838	14.330	34.595	35.094	0.697
Black	0.076	0.265	0.070	0.082	0.615
Asian	0.112	0.315	0.117	0.107	0.719
White	0.762	0.426	0.755	0.770	0.682
Hispanic	0.102	0.303	0.105	0.098	0.805
Male	0.463	0.499	0.471	0.455	0.722
High school diploma	0.994	0.077	0.988	1.000	0.091
Bachelors degree	0.483	0.500	0.447	0.520	0.102

*Notes:*  $p$ -values based on robust standard errors reported.

**Table D.6:** Auxiliary Experiment 3: Sample representativeness

	Anticipated persuasion	Pew (Inds and Dems)
Age	34.84	45.86
Black	0.08	0.18
White	0.76	0.59
Asian	0.11	0.05
Hispanic	0.10	0.15
Male	0.46	0.46
High school diploma	0.99	0.89
Bachelors degree or higher	0.48	0.35
Observations	501	6,627

*Notes:* Table displays mean characteristics, comparing the experimental sample with the 2018 Pew Research Center's American Trends Panel, Wave 39. Attriters are dropped from sample.

**Table D.7:** Auxiliary Experiment 4: Balance of covariates

	Overall		Cover	No Cover	$p$ -value
	mean	std. dev.	mean	mean	(C = NC)
Age	33.299	12.532	33.361	33.235	0.920
Black	0.077	0.267	0.079	0.075	0.875
Asian	0.132	0.339	0.139	0.125	0.688
White	0.751	0.433	0.752	0.750	0.954
Hispanic	0.117	0.322	0.139	0.095	0.174
Male	0.493	0.501	0.500	0.485	0.764
High school diploma	0.998	0.050	1.000	0.995	0.316
Bachelors degree	0.617	0.487	0.589	0.645	0.250

*Notes:*  $p$ -values based on robust standard errors reported.

**Table D.8:** Auxiliary Experiment 4: Sample representativeness

	Motives	Pew (Inds and Dems)
Age	33.30	45.86
Black	0.08	0.18
White	0.75	0.59
Asian	0.13	0.05
Hispanic	0.12	0.15
Male	0.49	0.46
High school diploma	1.00	0.89
Bachelors degree or higher	0.62	0.35
Observations	402	6,627

*Notes:* Table displays mean characteristics, comparing the experimental sample with the 2018 Pew Research Center's American Trends Panel, Wave 39. Attriters are dropped from sample.

**Table D.9:** Auxiliary Experiment 5: Balance of covariates

	Overall		Cover	No Cover	$p$ -value
	mean	std. dev.	mean	mean	(C = NC)
Higher-credibility					
Age	41.549	14.019	41.008	42.142	0.381
Black	0.209	0.407	0.198	0.221	0.528
Asian	0.046	0.211	0.052	0.040	0.516
White	0.696	0.460	0.698	0.695	0.946
Hispanic	0.152	0.359	0.149	0.155	0.864
Male	0.418	0.494	0.415	0.420	0.912
High school diploma	0.981	0.137	0.980	0.982	0.845
Bachelors degree	0.483	0.500	0.480	0.487	0.881
Lower-credibility					
Age	40.802	14.575	40.504	41.104	0.632
Black	0.193	0.395	0.198	0.189	0.793
Asian	0.037	0.189	0.033	0.041	0.631
White	0.715	0.452	0.736	0.693	0.261
Hispanic	0.166	0.372	0.154	0.178	0.454
Male	0.455	0.498	0.473	0.437	0.407
High school diploma	0.987	0.113	0.985	0.989	0.715
Bachelors degree	0.449	0.498	0.440	0.459	0.645

Notes:  $p$ -values based on robust standard errors reported.

**Table D.10:** Auxiliary Experiment 5: Sample representativeness

	Defund	Pew (Inds and Dems)
	Higher-credibility	
Age	41.55	45.86
Black	0.21	0.18
White	0.70	0.59
Asian	0.05	0.05
Hispanic	0.15	0.15
Male	0.42	0.46
High school diploma	0.98	0.89
Bachelors degree or higher	0.48	0.35
Observations	474	6,627
	Lower-credibility	
Age	40.80	45.86
Black	0.19	0.18
White	0.71	0.59
Asian	0.04	0.05
Hispanic	0.17	0.15
Male	0.45	0.46
High school diploma	0.99	0.89
Bachelors degree or higher	0.45	0.35
Observations	543	6,627

*Notes:* Table displays mean characteristics, comparing the experimental sample with the 2018 Pew Research Center's American Trends Panel, Wave 39. Attriters are dropped from sample.

**Table D.11:** Auxiliary Experiment 6: lower-credibility variation, reweighted to match higher-credibility sample in Experiment 2 on demographics

	Auxiliary Experiment 6	
	(1)	(2)
<b>Panel A:</b>	<i>Belief partner donated</i>	
Cover	0.010 (0.041)	0.016 (0.041)
No Cover mean	0.310	0.310
<b>Panel B:</b>	<i>Denied bonus to partner</i>	
Cover	0.016 (0.045)	0.007 (0.045)
No Cover mean	0.429	0.429
Observations	494	494
Demographic controls	No	Yes

*Notes:* The dependent variable in Panel A is an indicator taking value 1 if the respondent reports believing that his or her matched partner donated to the US Border Crisis Children’s Relief Fund. The dependent variable in Panel B is an indicator taking value 1 if the respondent denied his or her matched partner a \$1 bonus. Columns 1–2 report results for the lower-credibility experiment. Demographic controls include age, age squared, a set of race indicators, a Hispanic indicator, a male indicator, a set of education indicators. This table reweights observations to match the higher-credibility sample (Experiment 2) on the demographics reported in Table D.12. 12 observations are dropped due to the region of common support of the demographics covariates.



**Table D.12:** Auxiliary Experiment 6: Balance of covariates

	Overall		Cover	No Cover	$p$ -value
	mean	std. dev.	mean	mean	(C = NC)
Age	35.366	14.585	35.275	35.458	0.888
Black	0.053	0.225	0.043	0.064	0.303
Asian	0.132	0.339	0.137	0.127	0.747
White	0.771	0.421	0.773	0.769	0.923
Hispanic	0.107	0.309	0.141	0.072	0.011
Male	0.496	0.500	0.502	0.490	0.789
High school diploma	0.996	0.063	0.992	1.000	0.160
Bachelors degree	0.597	0.491	0.600	0.594	0.884

*Notes:*  $p$ -values based on robust standard errors reported.

**Table D.13:** Auxiliary Experiment 6: Sample representativeness

	Placebo	Pew (Inds and Dems)
Age	35.37	45.86
Black	0.05	0.18
White	0.77	0.59
Asian	0.13	0.05
Hispanic	0.11	0.15
Male	0.50	0.46
High school diploma	1.00	0.89
Bachelors degree or higher	0.60	0.35
Observations	506	6,627

*Notes:* Table displays mean characteristics, comparing the experimental sample with the 2018 Pew Research Center's American Trends Panel, Wave 39. Attriters are dropped from sample.

**Table D.14:** Auxiliary Experiment 7: Balance of covariates

	Overall		Article	No Article	$p$ -value
	mean	std. dev.	mean	mean	(A = NA)
Age	44.566	15.446	44.762	44.373	0.572
Black	0.026	0.159	0.027	0.025	0.724
Asian	0.028	0.166	0.026	0.031	0.552
White	0.917	0.275	0.916	0.919	0.768
Hispanic	0.061	0.239	0.059	0.062	0.795
Male	0.434	0.496	0.425	0.444	0.385
High school diploma	0.983	0.129	0.978	0.988	0.074
Bachelors degree	0.412	0.492	0.404	0.421	0.448

*Notes:*  $p$ -values based on robust standard errors reported.

**Table D.15:** Auxiliary Experiment 7: Sample representativeness

	Deport	Pew (Inds and Reps)
Age	44.57	47.17
Black	0.03	0.05
White	0.92	0.77
Asian	0.03	0.03
Hispanic	0.06	0.11
Male	0.43	0.52
High school diploma	0.98	0.93
Bachelors degree or higher	0.41	0.31
Observations	2,012	5,501

*Notes:* Table displays mean characteristics, comparing the experimental sample with the 2018 Pew Research Center's American Trends Panel, Wave 39. Attriters are dropped from sample.

## E Experimental Instructions

### E.1 Experiment 1: Expression of dissent – Democrats

#### E.1.1 Attention screener

The next question is about the following problem. In questionnaires like ours, sometimes there are participants who do not carefully read the questions and just quickly click through the survey. This means that there are a lot of random answers which compromise the results of research studies. To show that you read our questions carefully, please choose **both** “Extremely interested” and “Not at all interested” on the question below.

**Given the text above,** how interested are you in sports?

☐ Extremely interested

☐ Very interested

☐ A little bit interested

☐ Very little interested

☐ Not at all interested

>>

## E.1.2 Twitter information and login

Since our survey is about Twitter and current events, it requires you to grant the system access to your Twitter account through the "Tweetability" app.

Please note that we are **bound by agreement** with the Social and Behavioral Sciences Institutional Review Board at the University of Chicago to adhere to the following terms (in addition to the Twitter terms of service):

- We will **never** use the app to access non-public information from your account (including your posts)
- We will **never** use the app to make posts on your account without your **explicit consent**
- The app **does not give us access to your direct messages or email address**
- All identifying information will be stored on **password-protected directories** secured with **two-factor authentication**, and only **authorized research personnel** will have access
- All identifying information, **including your Twitter handle**, will be deleted by no later than December 1, 2021. Therefore, **the app will lose all access to your account** after this date (if not earlier)

If you have any questions for the researchers, you can contact the researchers at: [twitter.study@uchicago.edu](mailto:twitter.study@uchicago.edu)

If you have any questions or complaints, you can contact the Social and Behavioral Sciences Institutional Review Board at the University of Chicago at:

The Social & Behavioral Sciences Institutional Review Board,  
University of Chicago  
Phone: (773) 834-7835  
E-mail: [sbs-irb@uchicago.edu](mailto:sbs-irb@uchicago.edu)

If you are uncomfortable with these terms in any way, please end the survey now. Otherwise, please click the button below to proceed by signing into Twitter.

Sign in with Twitter

### E.1.3 Background questions

Are you Spanish, Hispanic, or Latino or none of these?

☐ Yes

☐ None of these

What is your year of birth?

What is your sex?

☐ Male

☐ Female

In politics, as of today, do you consider yourself a Republican, a Democrat, or an Independent?

☐ Republican

☐ Democrat

☐ Independent

>>

What is the highest level of school you have completed or the highest degree you have received?

- ☐ Less than high school degree
- ☐ High school graduate (high school diploma or equivalent including GED)
- ☐ Some college but no degree
- ☐ Associate degree in college (2-year)
- ☐ Bachelor's degree in college (4-year)
- ☐ Master's degree
- ☐ Doctoral degree
- ☐ Professional degree (JD, MD)

Which of the following best describes your race or ethnicity?

- ☐ African American/Black
- ☐ Asian/Asian American
- ☐ Caucasian/White
- ☐ Native American, Inuit or Aleut
- ☐ Native Hawaiian/Pacific Islander
- ☐ Other

Who did you vote for in the 2020 presidential election?

- ☐ Donald Trump
- ☐ Joe Biden
- ☐ Other
- ☐ Did not vote

Are you liberal or conservative?

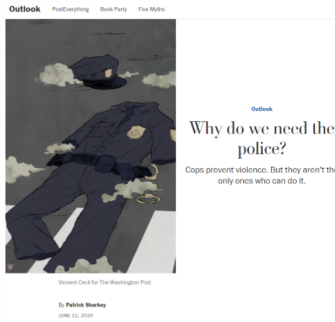
- ☐ Very liberal
- ☐ Liberal
- ☐ Neither liberal nor conservative
- ☐ Conservative
- ☐ Very conservative



#### E.1.4 Pre-treatment outcomes

On the next page, you will be provided with a recent Washington Post column written by **Princeton Professor of Criminology Patrick Sharkey**, in which he discusses evidence showing that more policing leads to less violent crime.

>>



© Patrick Sharkey  
JUNE 10, 2020

He calls to end policing as we know it contains a sort of trap. The best evidence we have makes clear that police are effective in reducing violence, and without designating some group to combat this problem, efforts to weaken them through budget cuts — “defund the police” — are likely to have unanticipated consequences and to destabilize communities. In many cities this is likely to lead to a rise in violence. And research shows that, when violence increases, Americans of all races become more punitive, supporting harsher policing and criminal justice policies. That’s how we got to this point.



Yet none of this means that the police, which have served as an institution of racialized control throughout our nation’s history, are the only group capable of reducing violence.

Community leaders and residents have proved adept at overseeing their neighborhoods, caring for their populations and maintaining safe streets. Studies show that this work lowers crime, sometimes dramatically. What happens if we put those people in charge of containing violence, too?

Over the past 10 years, an expanding body of research has shown just how damaging violence is to community life, children’s academic trajectories and healthy child development. We have rigorous, causal evidence that every shooting in a neighborhood affects children’s sleep and their ability to focus and learn. When a neighborhood becomes violent, it begins to fall apart, as public spaces empty, businesses close, parks and playgrounds turn dangerous, and families try to move elsewhere. Violence is the fundamental challenge for cities: Nothing works if public space is unsafe.

Those who argue that the police have no role in maintaining safe streets are arguing against lots of strong evidence. One of the most robust, most uncomfortable findings in criminology is that putting more officers on the street leads to less violent crime. We know this from randomized experiments involving “hot spots policing” and natural experiments in which more officers were brought to the streets because of something other than crime — a shift in the terror alert level or the timing of a federal grant — and violent crime fell. After the unrest around the deaths of Freddie Gray in Baltimore and Michael Brown in Ferguson, Mo., police officers stepped back from their duty to protect and serve; arrests for all kinds of low-level offenses dropped, and violence rose. This shouldn’t be interpreted to mean that protests against violent policing lead to more violence; rather, it means that when police don’t do their jobs, violence often results.

Considered alongside the brutal response to protests over the past few weeks, this evidence forces us to hold two incongruent ideas: Police are effective at reducing violence, the most damaging feature of urban inequality. And yet one can argue that law enforcement is an authoritarian institution that historically has inflicted violence on black people and continues to do so today.



Would you like to join a nonpartisan campaign that opposes defunding the police?

☐ Yes

☐ No

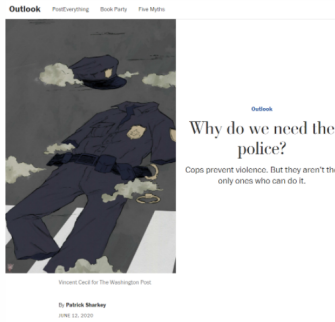
>>

### **You have successfully joined the campaign.**

Since you chose to join the campaign, we wanted to give you more time reading the Washington Post column written by **Princeton Professor of Criminology Patrick Sharkey**, where he discusses evidence showing that more policing leads to less violent crime.

The article is available on the next page, and you can spend as much time as you want reading it before you continue with the remaining part of the survey.

>>



© Patrick Sharkey  
JUNE 10, 2020

He calls to end policing as we know it contains a sort of trap. The best evidence we have makes clear that police are effective in reducing violence, and without designating some group to combat this problem, efforts to weaken them through budget cuts — “defund the police” — are likely to have unanticipated consequences and to destabilize communities. In many cities this is likely to lead to a rise in violence. And research shows that, when violence increases, Americans of all races become more punitive, supporting harsher policing and criminal justice policies. That’s how we got to this point.



Yet none of this means that the police, which have served as an institution of racialized control throughout our nation’s history, are the only group capable of reducing violence.

Community leaders and residents have proved adept at overseeing their neighborhoods, caring for their populations and maintaining safe streets. Studies show that this work lowers crime, sometimes dramatically. What happens if we put those people in charge of containing violence, too?

Over the past 10 years, an expanding body of research has shown just how damaging violence is to community life, children’s academic trajectories and healthy child development. We have rigorous, causal evidence that every shooting in a neighborhood affects children’s sleep and their ability to focus and learn. When a neighborhood becomes violent, it begins to fall apart, as public spaces empty, businesses close, parks and playgrounds turn dangerous, and families try to move elsewhere. Violence is the fundamental challenge for cities: Nothing works if public space is unsafe.

Those who argue that the police have no role in maintaining safe streets are arguing against lots of strong evidence. One of the most robust, most uncomfortable findings in criminology is that putting more officers on the street leads to less violent crime. We know this from randomized experiments involving “hot spots policing” and natural experiments in which more officers were brought to the streets because of something other than crime — a shift in the terror alert level or the timing of a federal grant — and violent crime fell. After the unrest around the deaths of Freddie Gray in Baltimore and Michael Brown in Ferguson, Mo., police officers stepped back from their duty to protect and serve; arrests for all kinds of low-level offenses dropped, and violence rose. This shouldn’t be interpreted to mean that protests against violent policing lead to more violence; rather, it means that when police don’t do their jobs, violence often results.

Considered alongside the brutal response to protests over the past few weeks, this evidence forces us to hold two incongruent ideas: Police are effective at reducing violence, the most damaging feature of urban inequality. And yet one can argue that law enforcement is an authoritarian institution that historically has inflicted violence on black people and continues to do so today.

### E.1.5 Treatment: “Before” wording (rationale)

This nonpartisan campaign involves signing up people on Twitter **to make a post encouraging their friends and followers to sign a petition** opposing the movement to defund the police.

The posts will be made public if/when we have finished surveying people in all U.S. counties. This strategy is often used to make campaigns “trend” on Twitter. To coordinate these efforts, we will use the *Tweetability* app you signed into earlier to schedule the posts.

I have joined a campaign to oppose defunding the police: [bit.ly/3DK3UEr](https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/). Before joining, I was shown this article written by a Princeton professor on the strong scientific evidence that defunding the police would increase violent crime:  
<https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/>



washingtonpost.com  
Perspective | Cops prevent violence. But they aren't the only ones wh...  
Communities already know how to police their own. Now put them in charge of it.

Do you authorize the *Tweetability* app to schedule the post above to be posted on your account? (If you choose “no,” then nothing will be posted on your account.)

☐ Yes

☐ No

>>

### E.1.6 Treatment: “After” wording (no rationale)

This nonpartisan campaign involves signing up people on Twitter **to make a post encouraging their friends and followers to sign a petition** opposing the movement to defund the police.

The posts will be made public if/when we have finished surveying people in all U.S. counties. This strategy is often used to make campaigns “trend” on Twitter. To coordinate these efforts, we will use the *Tweetability* app you signed into earlier to schedule the posts.

I have joined a campaign to oppose defunding the police: [bit.ly/3DK3UEr](https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/). After joining, I was shown this article written by a Princeton professor on the strong scientific evidence that defunding the police would increase violent crime:  
<https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/>



washingtonpost.com

Perspective | Cops prevent violence. But they aren't the only ones wh...  
Communities already know how to police their own. Now put them in charge of it.

Do you authorize the *Tweetability* app to schedule the post above to be posted on your account? (If you choose “no,” then nothing will be posted on your account.)

☐ Yes

☐ No

>>

## E.2 Experiment 2: Interpretation of dissent – Democrats

### E.2.1 Attention screener and background questions

The next question is about the following problem. In questionnaires like ours, sometimes there are participants who do not carefully read the questions and just quickly click through the survey. This means that there are a lot of random answers which compromise the results of research studies. To show that you read our questions carefully, please choose both **“Extremely interested”** and **“Not at all interested”** as your answer in the below question.

Given the above, how interested are you in sports?

☐ Extremely interested

☐ Very interested

☐ A little bit interested

☐ Almost not interested

☐ Not at all interested



What is your sex?

- ☐ Male
- ☐ Female

What is your year of birth?

In politics, as of today, do you consider yourself a Republican, a Democrat, or an Independent?

- ☐ Republican
- ☐ Democrat
- ☐ Independent

What is the highest level of school you have completed or the highest degree you have received?

- ☐ Less than high school degree
- ☐ High school graduate (high school diploma or equivalent including GED)
- ☐ Some college but no degree
- ☐ Associate degree in college (2-year)
- ☐ Bachelor's degree in college (4-year)
- ☐ Master's degree
- ☐ Doctoral degree
- ☐ Professional degree (JD, MD)

Are you Spanish, Hispanic, or Latino or none of these?

- ☐ Yes
- ☐ None of these

Which of the following best describes your race or ethnicity?

- ☐ African American/Black
- ☐ Asian/Asian American
- ☐ Caucasian/White
- ☐ Native American, Inuit or Aleut
- ☐ Native Hawaiian/Pacific Islander
- ☐ Other



Do you lean toward the Republican Party or the Democratic Party?

☐ Lean toward the Republican Party

☐ Lean toward the Democratic Party

Who did you vote for in the 2020 presidential election?

☐ Donald Trump

☐ Joe Biden

☐ Other

☐ Did not vote

Are you liberal or conservative?

☐ Very liberal

☐ Liberal

☐ Neither liberal nor conservative

☐ Conservative

☐ Very conservative



## E.2.2 Treatment: “Before” wording (rationale)

We conducted a survey about political and social attitudes in the United States earlier this year. You have been matched with one of the participants from that survey.

Your matched participant **joined a campaign to oppose defunding the police.**

Your matched participant **also agreed to post the Tweet below** on their Twitter profile to promote the campaign.

I have joined a campaign to oppose defunding the police: [bit.ly/3DK3UEr](https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/). Before joining, I was shown this article written by a Princeton professor on the strong scientific evidence that defunding the police would increase violent crime:  
<https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/>



washingtonpost.com  
Perspective | Cops prevent violence. But they aren't the only ones wh...  
Communities already know how to police their own. Now put them in charge of it.

Why do you think your matched respondent chose to join the campaign to oppose defunding the police?





### Matched Respondent's Donation Decision

We gave your matched respondent the opportunity to donate \$10 to the **National Association for the Advancement of Colored People (NAACP)**, America's oldest and largest civil rights organization.

Below, we will ask you to guess whether or not your matched respondent donated \$10 to the National Association for the Advancement of Colored People (NAACP).

**Reminder:** Your matched participant agreed to post the Tweet below on their Twitter account.

I have joined a campaign to oppose defunding the police: [bit.ly/3DK3UEr](https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/). Before joining, I was shown this article written by a Princeton professor on the strong scientific evidence that defunding the police would increase violent crime:  
<https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/>



washingtonpost.com

Perspective | Cops prevent violence. But they aren't the only ones wh...  
Communities already know how to police their own. Now put them in charge of it.

Do you think that your matched participant chose to donate \$5 to the National Association for the Advancement of Colored People (NAACP)?

☐ Yes, I think my matched respondent chose to donate

☐ No, I think my matched respondent **did not** choose to donate



You now have the opportunity to authorize a \$1 bonus payment to your matched respondent. **The bonus payment will not be deducted from your payment.** Your matched respondent did not know you would have the opportunity to decide their bonus.

**Reminder:** Your matched participant agreed to post the Tweet below on their Twitter account.

I have joined a campaign to oppose defunding the police: [bit.ly/3DK3UEr](https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/). Before joining, I was shown this article written by a Princeton professor on the strong scientific evidence that defunding the police would increase violent crime:  
<https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/>



Do you want to authorize a \$1 bonus to your matched respondent?

- ☐ Yes, I would like to authorize a \$1 bonus
- ☐ No, I would not like to authorize a \$1 bonus



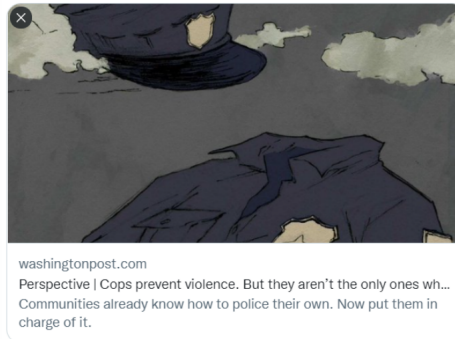
### E.2.3 Treatment: “After” wording (no rationale)

We conducted a survey about political and social attitudes in the United States earlier this year. You have been matched with one of the participants from that survey.

Your matched participant **joined a campaign to oppose defunding the police.**

Your matched participant **also agreed to post the Tweet below** on their Twitter profile to promote the campaign.

I have joined a campaign to oppose defunding the police: [bit.ly/3DK3UEr](https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/). After joining, I was shown this article written by a Princeton professor on the strong scientific evidence that defunding the police would increase violent crime: <https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/>



Why do you think your matched respondent chose to join the campaign to oppose defunding the police?



### Matched Respondent's Donation Decision

We gave your matched respondent the opportunity to donate \$10 to the **National Association for the Advancement of Colored People (NAACP)**, America's oldest and largest civil rights organization.

Below, we will ask you to guess whether or not your matched respondent donated \$10 to the National Association for the Advancement of Colored People (NAACP).

**Reminder:** Your matched participant agreed to post the Tweet below on their Twitter account.

I have joined a campaign to oppose defunding the police: [bit.ly/3DK3UEr](https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/). After joining, I was shown this article written by a Princeton professor on the strong scientific evidence that defunding the police would increase violent crime:  
<https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/>



washingtonpost.com  
Perspective | Cops prevent violence. But they aren't the only ones wh...  
Communities already know how to police their own. Now put them in charge of it.

Do you think that your matched participant chose to donate \$5 to the National Association for the Advancement of Colored People (NAACP)?

☐ Yes, I think my matched respondent chose to donate

☐ No, I think my matched respondent **did not** choose to donate



You now have the opportunity to authorize a \$1 bonus payment to your matched respondent. **The bonus payment will not be deducted from your payment.** Your matched respondent did not know you would have the opportunity to decide their bonus.

**Reminder:** Your matched participant agreed to post the Tweet below on their Twitter account.

I have joined a campaign to oppose defunding the police: [bit.ly/3DK3UEr](https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/). After joining, I was shown this article written by a Princeton professor on the strong scientific evidence that defunding the police would increase violent crime:  
<https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/>



washingtonpost.com  
Perspective | Cops prevent violence. But they aren't the only ones wh...  
Communities already know how to police their own. Now put them in charge of it.

Do you want to authorize a \$1 bonus to your matched respondent?

- ☐ Yes, I would like to authorize a \$1 bonus
- ☐ No, I would not like to authorize a \$1 bonus



## E.3 Experiment 3: Expression of dissent – Republicans

### E.3.1 Attention screener

The next question is about the following problem. In questionnaires like ours, sometimes there are participants who do not carefully read the questions and just quickly click through the survey. This means that there are a lot of random answers which compromise the results of research studies. To show that you read our questions carefully, please choose **both** "Extremely interested" and "Not at all interested" on the question below.

**Given the text above,** how interested are you in sports?

☐ Extremely interested

☐ Very interested

☐ A little bit interested

☐ Very little interested

☐ Not at all interested

>>

### E.3.2 Twitter information and login

Since our survey is about Twitter and current events, it requires you to grant the system access to your Twitter account through the "Tweatability" app.

Please note that we are **bound by agreement** with the Social and Behavioral Sciences Institutional Review Board at the University of Chicago to adhere to the following terms (in addition to the Twitter terms of service):

- We will **never** use the app to access non-public information from your account (including your posts)
- We will **never** use the app to make posts on your account without your **explicit consent**
- The app **does not give us access to your direct messages or email address**
- All identifying information will be stored on **password-protected directories** secured with **two-factor authentication**, and only **authorized research personnel** will have access
- All identifying information, **including your Twitter handle**, will be deleted by no later than August 1, 2021. Therefore, **the app will lose all access to your account** after this date (if not earlier)

If you have any questions for the researchers, you can contact the researchers at: [twitter.study@uchicago.edu](mailto:twitter.study@uchicago.edu)

If you have any questions or complaints, you can contact the Social and Behavioral Sciences Institutional Review Board at the University of Chicago at:

The Social & Behavioral Sciences Institutional Review Board,  
University of Chicago  
Phone: (773) 834-7835  
E-mail: [sbs-irb@uchicago.edu](mailto:sbs-irb@uchicago.edu)

If you are uncomfortable with these terms in any way, please end the survey now. Otherwise, please click the button below to proceed by signing into Twitter.

Sign in with Twitter

## Authorize Tweetability: Schedule Tweets to access your account?



### Tweetability: Schedule Tweets

This app was created to use the Twitter API.

☐ Remember me · [Forgot password?](#)

Sign In

Cancel

#### This application will be able to:

- See Tweets from your timeline (including protected Tweets) as well as your Lists and collections.
- See your Twitter profile information and account settings.
- See accounts you follow, mute, and block.
- Follow and unfollow accounts for you.
- Update your profile and account settings.
- Post and delete Tweets for you, and engage with Tweets posted by others (Like, un-Like, or reply to a Tweet, Retweet, etc.) for you.
- Create, manage, and delete Lists and collections for you.
- Mute, block, and report accounts for you.

Learn more about third-party app permissions in the [Help Center](#).



### E.3.3 Demographics

Are you Spanish, Hispanic, or Latino or none of these?

☐ Yes

☐ None of these

What is your year of birth?

What is your sex?

☐ Male

☐ Female

In politics, as of today, do you consider yourself a Republican, a Democrat, or an Independent?

☐ Republican

☐ Democrat

☐ Independent

>>

What is the highest level of school you have completed or the highest degree you have received?

- ☐ Less than high school degree
- ☐ High school graduate (high school diploma or equivalent including GED)
- ☐ Some college but no degree
- ☐ Associate degree in college (2-year)
- ☐ Bachelor's degree in college (4-year)
- ☐ Master's degree
- ☐ Doctoral degree
- ☐ Professional degree (JD, MD)

Which of the following best describes your race or ethnicity?

- ☐ African American/Black
- ☐ Asian/Asian American
- ☐ Caucasian/White
- ☐ Native American, Inuit or Aleut
- ☐ Native Hawaiian/Pacific Islander
- ☐ Other

Who did you vote for in the 2020 presidential election?

- ☐ Donald Trump
- ☐ Joe Biden
- ☐ Other
- ☐ Did not vote

Are you liberal or conservative?

- ☐ Very liberal
- ☐ Liberal
- ☐ Neither liberal nor conservative
- ☐ Conservative
- ☐ Very conservative



### E.3.4 Video clip

Please see the short video below where Fox News host **Tucker Carlson** presents evidence on whether illegal immigrants commit more crime.



### E.3.5 Treatment: “After” wording

Would you like to join a campaign to immediately deport all illegal Mexican immigrants?

☐ Yes

☐ No

>>

In case you want save the video with Tucker Carlson talking about immigration and crime, here is the link: <https://www.youtube.com/watch?v=SDdkkTLCUUQ>

>>

**You have successfully joined the campaign.** This campaign involves signing up people on Twitter **to make a post encouraging their friends and followers to sign a petition** in favor of immediately deporting all illegal Mexican immigrants.

The posts will be made public if/when we have finished surveying people in all U.S. counties. This strategy is often used to make campaigns "trend" on Twitter. To coordinate these efforts, we will use the *Tweetability* app you signed into earlier to schedule the posts.



Below, we will ask you if you want to authorize the following Tweet to be posted on your account:

I have joined a campaign to immediately deport all illegal Mexicans! After I joined the campaign, I received a link to this video about how illegals commit more crime: [y2u.be/SDdkkTLCUUQ](https://y2u.be/SDdkkTLCUUQ). Sign this petition to immediately deport all illegal Mexicans: [bit.ly/2OttyJ6](https://bit.ly/2OttyJ6)



Tucker: Why didn't we know truth about illegals and crime?  
Tucker's Thoughts: For years, we were told illegal immigrants were more law-abiding than American citizens. In fact, the ...  
[youtube.com](https://youtube.com)

Do you authorize the *Tweetability* app to schedule the post above to be posted on your account? (If you choose "no," then nothing will be posted on your account.)

☐ Yes

☐ No

>>

### E.3.6 Treatment: “Before” wording

In case you want save the video with Tucker Carlson talking about immigration and crime, here is the link: <https://www.youtube.com/watch?v=SDdkkTLCUUQ>

>>

Would you like to join a campaign to immediately deport all illegal Mexican immigrants?

☐ Yes

☐ No

>>

**You have successfully joined the campaign.** This campaign involves signing up people on Twitter **to make a post encouraging their friends and followers to sign a petition** in favor of immediately deporting all illegal Mexican immigrants.

The posts will be made public if/when we have finished surveying people in all U.S. counties. This strategy is often used to make campaigns "trend" on Twitter. To coordinate these efforts, we will use the *Tweetability* app you signed into earlier to schedule the posts.





Below, we will ask you if you want to authorize the following Tweet to be posted on your account:

I have joined a campaign to immediately deport all illegal Mexicans! Before I joined the campaign, I received a link to this video about how illegals commit more crime: [y2u.be/SDdkkTLCUUQ](https://www.youtube.com/watch?v=SDdkkTLCUUQ). Sign this petition to immediately deport all illegal Mexicans: [bit.ly/2OttyJ6](https://bit.ly/2OttyJ6)



Tucker: Why didn't we know truth about illegals and crime?  
Tucker's Thoughts: For years, we were told illegal immigrants were more law-abiding than American citizens. In fact, the ...  
[🔗 youtube.com](https://www.youtube.com/watch?v=SDdkkTLCUUQ)

Do you authorize the *Tweetability* app to schedule the post above to be posted on your account? (If you choose "no," then nothing will be posted on your account.)

☐ Yes

☐ No

>>

## E.4 Experiment 4: Interpretation of dissent – Republicans

### E.4.1 Attention screener and background questions

The next question is about the following problem. In questionnaires like ours, sometimes there are participants who do not carefully read the questions and just quickly click through the survey. This means that there are a lot of random answers which compromise the results of research studies. To show that you read our questions carefully, please choose both **“Extremely interested”** and **“Not at all interested”** as your answer in the below question.

Given the above, how interested are you in sports?

☐ Extremely interested

☐ Very interested

☐ A little bit interested

☐ Almost not interested

☐ Not at all interested



What is your sex?

☐ Male

☐ Female

What is your year of birth?

In politics, as of today, do you consider yourself a Republican, a Democrat, or an independent?

☐ Republican

☐ Democrat

☐ Independent

What is the highest level of school you have completed or the highest degree you have received?

☐ Less than high school degree

☐ High school graduate (high school diploma or equivalent including GED)

☐ Some college but no degree

☐ Associate degree in college (2-year)

☐ Bachelor's degree in college (4-year)

☐ Master's degree

☐ Doctoral degree

☐ Professional degree (JD, MD)

Are you Spanish, Hispanic, or Latino or none of these?

☐ Yes

☐ None of these

Which of the following best describes your race or ethnicity?

☐ African American/Black

☐ Asian/Asian American

☐ Caucasian/White

☐ Native American, Inuit or Aleut

☐ Native Hawaiian/Pacific Islander

☐ Other



Who did you vote for in the 2020 presidential election?

☐ Donald Trump

☐ Joe Biden

☐ Other

☐ Did not vote

Are you liberal or conservative?

☐ Very liberal

☐ Liberal

☐ Neither liberal nor conservative

☐ Conservative

☐ Very conservative



#### E.4.2 Treatment: “Before” condition (rationale)

We conducted a survey about political and social attitudes in the United States earlier this year. You have been matched with one of the participants from that survey.

Your matched participant **joined a campaign to immediately deport all illegal Mexican immigrants.**

Your matched participant **also agreed to post the Tweet below** on their Twitter profile to promote the campaign.

I have joined a campaign to immediately deport all illegal Mexicans! Before I joined the campaign, I received a link to this video about how illegals commit more crime: [y2u.be/SDdkkTLCUUQ](https://y2u.be/SDdkkTLCUUQ). Sign this petition to immediately deport all illegal Mexicans: [bit.ly/2OttyJ6](https://bit.ly/2OttyJ6)



Tucker: Why didn't we know truth about illegals and crime?  
Tucker's Thoughts: For years, we were told illegal immigrants were more law-abiding than American citizens. In fact, the ...  
[youtube.com](https://www.youtube.com)

Why do you think your matched respondent chose to join the campaign to immediately deport all illegal Mexican immigrants?



### Matched Respondent's Donation Decision

We gave your matched respondent the opportunity to authorize a \$5 donation to the US Border Crisis Children's Relief Fund, which delivers humanitarian aid to migrant children and families at the US-Mexico border. The organization is working with local partners to ensure that children and families have necessities such as hygiene kits, diapers and clothing. We told your matched respondent that we would make the donation on their behalf, so the donation did not affect their payment.

Below, we will ask you to guess whether or not your matched respondent authorized the \$5 donation to the US Border Crisis Children's Relief Fund.

**Reminder:** Your matched participant agreed to post the Tweet below on their Twitter profile to promote the campaign.

I have joined a campaign to immediately deport all illegal Mexicans! Before I joined the campaign, I received a link to this video about how illegals commit more crime: [youtu.be/SDdkkTLCUUQ](https://youtu.be/SDdkkTLCUUQ). Sign this petition to immediately deport all illegal Mexicans: [bit.ly/2OttyJ6](https://bit.ly/2OttyJ6)



Tucker: Why didn't we know truth about illegals and crime?  
Tucker's Thoughts: For years, we were told illegal immigrants were more law-abiding than American citizens. In fact, the ...  
[youtube.com](https://www.youtube.com)

Do you think that your matched participant chose to authorize the \$5 donation to the US Border Crisis Children's Relief Fund?

☐ Yes, I think my matched respondent chose to authorize the donation

☐ No, I think my matched respondent **did not** choose to authorize the donation



You now have the opportunity to authorize a \$1 bonus payment to your matched respondent. **The bonus payment will not be deducted from your payment.** Your matched respondent did not know that you would have the opportunity to decide on their bonus.

**Reminder:** Your matched participant agreed to post the Tweet below on their Twitter profile to promote the campaign.

I have joined a campaign to immediately deport all illegal Mexicans! Before I joined the campaign, I received a link to this video about how illegals commit more crime: [y2u.be/SDdkkTLCUUQ](https://y2u.be/SDdkkTLCUUQ). Sign this petition to immediately deport all illegal Mexicans: [bit.ly/2OttyJ6](https://bit.ly/2OttyJ6)



Tucker: Why didn't we know truth about illegals and crime?  
Tucker's Thoughts: For years, we were told illegal immigrants were more law-abiding than American citizens. In fact, the ...  
[youtube.com](https://youtube.com)

Do you want to authorize a \$1 bonus to your matched respondent?

☐ Yes, I would like to authorize a \$1 bonus

☐ No, I would not like to authorize a \$1 bonus



### E.4.3 Treatment: “After” condition (no rationale)

We conducted a survey about political and social attitudes in the United States earlier this year. You have been matched with one of the participants from that survey.

Your matched participant **joined a campaign to immediately deport all illegal Mexican immigrants.**

Your matched participant **also agreed to post the Tweet below** on their Twitter profile to promote the campaign.

I have joined a campaign to immediately deport all illegal Mexicans! After I joined the campaign, I received a link to this video about how illegals commit more crime: [y2u.be/SDdkkTLCUUQ](https://y2u.be/SDdkkTLCUUQ). Sign this petition to immediately deport all illegal Mexicans: [bit.ly/2OttyJ6](https://bit.ly/2OttyJ6)



Tucker: Why didn't we know truth about illegals and crime?  
Tucker's Thoughts: For years, we were told illegal immigrants were more law-abiding than American citizens. In fact, the ...  
[youtube.com](https://youtube.com)

Why do you think your matched respondent chose to join the campaign to immediately deport all illegal Mexican immigrants?





You now have the opportunity to authorize a \$1 bonus payment to your matched respondent. **The bonus payment will not be deducted from your payment.** Your matched respondent did not know that you would have the opportunity to decide on their bonus.

**Reminder:** Your matched participant agreed to post the Tweet below on their Twitter profile to promote the campaign.

I have joined a campaign to immediately deport all illegal Mexicans! After I joined the campaign, I received a link to this video about how illegals commit more crime: [y2u.be/SDdkkTLCUUQ](https://youtu.be/SDdkkTLCUUQ). Sign this petition to immediately deport all illegal Mexicans: [bit.ly/2OttyJ6](https://bit.ly/2OttyJ6)



Tucker: Why didn't we know truth about illegals and crime?  
Tucker's Thoughts: For years, we were told illegal immigrants were more law-abiding than American citizens. In fact, the ...  
[youtube.com](https://www.youtube.com)

Do you want to authorize a \$1 bonus to your matched respondent?

- ☐ Yes, I would like to authorize a \$1 bonus
- ☐ No, I would not like to authorize a \$1 bonus



### Matched Respondent's Donation Decision

We gave your matched respondent the opportunity to authorize a \$5 donation to the US Border Crisis Children's Relief Fund, which delivers humanitarian aid to migrant children and families at the US-Mexico border. The organization is working with local partners to ensure that children and families have necessities such as hygiene kits, diapers and clothing. We told your matched respondent that we would make the donation on their behalf, so the donation did not affect their payment.

Below, we will ask you to guess whether or not your matched respondent authorized the \$5 donation to the US Border Crisis Children's Relief Fund.

**Reminder:** Your matched participant agreed to post the Tweet below on their Twitter profile to promote the campaign.

I have joined a campaign to immediately deport all illegal Mexicans! After I joined the campaign, I received a link to this video about how illegals commit more crime: [youtu.be/SDdkkTLCUUQ](https://youtu.be/SDdkkTLCUUQ). Sign this petition to immediately deport all illegal Mexicans: [bit.ly/2OttyJ6](https://bit.ly/2OttyJ6)



Tucker: Why didn't we know truth about illegals and crime?  
Tucker's Thoughts: For years, we were told illegal immigrants were more law-abiding than American citizens. In fact, the ...  
[youtube.com](https://www.youtube.com)

Do you think that your matched participant chose to authorize the \$5 donation to the US Border Crisis Children's Relief Fund?

☐ Yes, I think my matched respondent chose to authorize the donation

☐ No, I think my matched respondent **did not** choose to authorize the donation



## E.5 Auxiliary Experiment 1: Persuasion experiment – Democrats

### E.5.1 Pre-treatment beliefs

How do you think decreasing funding for the police, commonly referred to as "defunding the police," would affect violent crime?

- ☐ Strongly increase violent crime
- ☐ Somewhat increase violent crime
- ☐ Neither increase nor decrease violent crime
- ☐ Somewhat decrease violent crime
- ☐ Strongly decrease violent crime



### E.5.2 Information treatment (treatment group only)

According to a recent article in the Washington Post written by Princeton Professor of Criminology Patrick Sharkey, **one of the most robust findings in criminology is that putting more police officers on the street leads to less violent crime.**

If you want to learn more, you can read the article here: <https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/>



### E.5.3 Post-treatment outcomes

Do you think that funding for the police should be increased, decreased, or stay the same?

☐ Increased a lot

☐ Increased a little

☐ Stay about the same

☐ Decreased a little

☐ Decreased a lot



How do you think **increasing** funding for the police would affect violent crime?

- ☐ Strongly increase violent crime
- ☐ Somewhat increase violent crime
- ☐ Neither increase nor decrease violent crime
- ☐ Somewhat decrease violent crime
- ☐ Strongly decrease violent crime



## E.6 Auxiliary Experiment 2: Rainforest placebo

### E.6.1 Pre-treatment questions

On the next page, you will be provided with a recent Reuters article reporting about a new landmark study showing that more than 10,000 species are at high risk of extinction due to the destruction of the Amazon rainforest.



## Environment

# Over 10,000 species risk extinction in Amazon, says landmark report

By Stephen Eisenhammer and Oliver Griffin

SAO PAULO/BOGOTA, July 14 (Reuters) - More than 10,000 species of plants and animals are at high risk of extinction due to the destruction of the Amazon rainforest - 35% of which has already been deforested or degraded, according to the draft of a landmark scientific report published on Wednesday.

Produced by the Science Panel for the Amazon (SPA), the 33-chapter report brings together research on the world's largest rainforest from 200 scientists from across the globe. It is the most detailed assessment of the state of the forest to date and both makes clear the vital role the Amazon plays in global climate and the profound risks it is facing.

Cutting deforestation and forest degradation to zero in less than a decade "is critical," the report said, also calling for massive restoration of already destroyed areas.

The rainforest is a vital bulwark against climate change both for the carbon it absorbs and what it stores.

Would you like to join a nonpartisan campaign to immediately stop the destruction of the Amazon rainforest?

☐ Yes

☐ No

>>

**You have successfully joined the campaign.**

Since you chose to join the campaign, we wanted to give you more time reading the Reuters article covering the landmark study showing that more than 10,000 species are at high risk of extinction due to the destruction of the Amazon rainforest.

The article is available on the next page, and you can spend as much time as you want reading it before you continue with the remaining part of the survey.

>>

## Environment

# Over 10,000 species risk extinction in Amazon, says landmark report

By Stephen Eisenhammer and Oliver Griffin

SAO PAULO/BOGOTA, July 14 (Reuters) - More than 10,000 species of plants and animals are at high risk of extinction due to the destruction of the Amazon rainforest - 35% of which has already been deforested or degraded, according to the draft of a landmark scientific report published on Wednesday.

Produced by the Science Panel for the Amazon (SPA), the 33-chapter report brings together research on the world's largest rainforest from 200 scientists from across the globe. It is the most detailed assessment of the state of the forest to date and both makes clear the vital role the Amazon plays in global climate and the profound risks it is facing.

Cutting deforestation and forest degradation to zero in less than a decade "is critical," the report said, also calling for massive restoration of already destroyed areas.

The rainforest is a vital bulwark against climate change both for the carbon it absorbs and what it stores.



## E.6.2 Treatment: “Before” wording (rationale)

This nonpartisan campaign involves signing up people on Twitter **to make a post encouraging their friends and followers to sign a petition** to immediately stop the destruction of the Amazon rainforest.

The posts will be made public if /when we have finished surveying people in all U.S. counties. This strategy is often used to make campaigns “trend” on Twitter. To coordinate these efforts, we will use the *Tweetability* app you signed into earlier to schedule the posts.

Below, we will ask you if you want to authorize the following Tweet to be posted on your account:

I've joined a campaign to immediately stop the destruction of the Amazon rainforest! Before I joined the campaign, I was shown this article about how 10,000 species risk extinction in Amazon: <https://www.reuters.com/business/environment/over-10000-species-risk-extinction-amazon-says-landmark-report-2021-07-14/> Join the campaign and sign the petition: [bit.ly/3whrwxT](https://bit.ly/3whrwxT)



reuters.com  
Over 10,000 species risk extinction in Amazon, says landmark report  
More than 10,000 species of plants and animals are at high risk of extinction due to the destruction of the Amazon rainforest - 35% of ...

Do you authorize the *Tweetability* app to schedule the post above to be posted on your account? (If you choose “no,” then nothing will be posted on your account.)

☐ Yes

☐ No



### E.6.3 Treatment: “After” wording (no rationale)

This nonpartisan campaign involves signing up people on Twitter **to make a post encouraging their friends and followers to sign a petition** to immediately stop the destruction of the Amazon rainforest.

The posts will be made public if/when we have finished surveying people in all U.S. counties. This strategy is often used to make campaigns “trend” on Twitter. To coordinate these efforts, we will use the *Tweetability* app you signed into earlier to schedule the posts.

Below, we will ask you if you want to authorize the following Tweet to be posted on your account:

I've joined a campaign to immediately stop the destruction of the Amazon rainforest! After I joined the campaign, I was shown this article about how 10,000 species risk extinction in Amazon:  
<https://www.reuters.com/business/environment/over-10000-species-risk-extinction-amazon-says-landmark-report-2021-07-14/> Join the campaign and sign the petition: [bit.ly/3whrwxt](https://bit.ly/3whrwxt)



reuters.com

Over 10,000 species risk extinction in Amazon, says landmark report  
More than 10,000 species of plants and animals are at high risk of extinction due to the destruction of the Amazon rainforest - 35% of ...

Do you authorize the *Tweetability* app to schedule the post above to be posted on your account? (If you choose “no,” then nothing will be posted on your account.)

☐ Yes

☐ No



## E.7 Auxiliary Experiment 3: Anticipated persuasion – Democrats

### E.7.1 Treatment: “Before” wording (rationale)

This nonpartisan campaign involves signing up people on Twitter **to make a post encouraging their friends and followers to sign a petition** opposing the movement to defund the police.

I have joined a campaign to oppose defunding the police: [bit.ly/3DK3UEr](https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/). Before joining, I was shown this article written by a Princeton professor on the strong scientific evidence that defunding the police would increase violent crime: <https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/>



Suppose you posted the Tweet above on your account. If you had to guess, what percentage of people who saw your Tweet would choose to join the campaign to oppose defunding the police?

0 10 20 30 40 50 60 70 80 90 100

Percentage of people who join



>>

### E.7.2 Treatment: “After” wording (no rationale)

This nonpartisan campaign involves signing up people on Twitter **to make a post encouraging their friends and followers to sign a petition** opposing the movement to defund the police.

I have joined a campaign to oppose defunding the police: [bit.ly/3DK3UEr](https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/). After joining, I was shown this article written by a Princeton professor on the strong scientific evidence that defunding the police would increase violent crime:



Suppose you posted the Tweet above on your account. If you had to guess, what percentage of people who saw your Tweet would choose to join the campaign to oppose defunding the police?

0 10 20 30 40 50 60 70 80 90 100  
Percentage of people who join

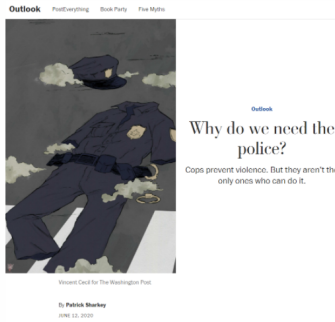


## E.8 Auxiliary Experiment 4: Open-ended explanations of preferred anti-defunding Tweet – Democrats

### E.8.1 Pre-treatment questions

On the next page, you will be provided with a recent Washington Post column written by **Princeton Professor of Criminology Patrick Sharkey**, in which he discusses evidence showing that more policing leads to less violent crime.

>>



He calls to end policing as we know it contains a sort of trap. The best evidence we have makes clear that police are effective in reducing violence, and without designating some group to combat this problem, efforts to weaken them through budget cuts — “defund the police” — are likely to have unanticipated consequences and to destabilize communities. In many cities this is likely to lead to a rise in violence. And research shows that, when violence increases, Americans of all races become more punitive, supporting harsher policing and criminal justice policies. That’s how we got to this point.



**Patrick Sharkey**  
Associate professor in a  
department of sociology and  
public affairs at Princeton  
University’s Woodrow Wilson School of  
Public and International Affairs. His most  
recent books are *Threats, Politics, The Street*,  
*Crime Decline, The Renewal of City Life, and*  
*The Best Way to Violence?*

Yet none of this means that the police, which have served as an institution of racialized control throughout our nation’s history, are the only group capable of reducing violence.

Community leaders and residents have proved adept at overseeing their neighborhoods, caring for their populations and maintaining safe streets. Studies show that this work lowers crime, sometimes dramatically. What happens if we put those people *in charge* of containing violence, too?

Over the past 10 years, an expanding body of research has shown just how damaging violence is to community life, children’s academic trajectories and healthy child development. We have rigorous, causal evidence that every shooting in a neighborhood affects children’s sleep and their ability to focus and learn. When a neighborhood becomes violent, it begins to fall apart, as public spaces empty, businesses close, parks and playgrounds turn dangerous, and families try to move elsewhere. Violence is the *fundamental challenge* for cities: Nothing works if public space is unsafe.

Those who argue that the police have no role in maintaining safe streets are arguing against lots of strong evidence. One of the most robust, most uncomfortable findings in criminology is that putting more officers on the street leads to less violent crime. We know this from *randomized experiments* involving “hot spots policing” and natural experiments in which more officers were brought to the streets because of something other than crime — a shift in the terror alert level or the timing of a federal grant — and violent crime fell. After the unrest around the deaths of Freddie Gray in Baltimore and Michael Brown in Ferguson, Mo., police officers stepped back from their duty to protect and serve; arrests for all kinds of low-level offenses dropped, and violence rose. This shouldn’t be interpreted to mean that protests against violent policing lead to more violence; rather, it means that when police don’t do their jobs, violence often results.

Considered alongside the brutal response to protests over the past few weeks, this evidence forces us to hold two incongruent ideas: Police are effective at *reducing violence*, the most damaging feature of urban inequality. And yet one can argue that law enforcement is an authoritarian institution that historically has inflicted violence on black people and continues to do so today.

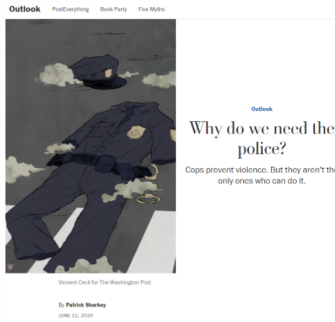
**Imagine** that at this point in the study, you indicated that you wanted to join a campaign that opposes the movement to defund the police.

>>

**Imagine that you successfully have joined the campaign.**

Since you joined the campaign, we wanted to give you more time reading the Washington Post column written by **Princeton Professor of Criminology Patrick Sharkey**, where he discusses evidence showing that more policing leads to less violent crime.

The article is available on the next page, and you can spend as much time (or as little time) as you want reading it before you continue with the remaining part of the survey.



He calls to end policing as we know it contains a sort of trap. The best evidence we have makes clear that police are effective in reducing violence, and without designating some group to combat this problem, efforts to weaken them through budget cuts — “defund the police” — are likely to have unanticipated consequences and to destabilize communities. In many cities this is likely to lead to a rise in violence. And research shows that, when violence increases, Americans of all races become more punitive, supporting harsher policing and criminal justice policies. That’s how we got to this point.



**Patrick Sharkey**  
Associate professor in a  
department of sociology and  
public affairs at Princeton  
University’s Woodrow Wilson School of  
Public and International Affairs. His most  
recent books are “Threats, Politics, The Street  
Crime Decline, The Renewal of City Life, and  
the Street War on Violence.”

Yet none of this means that the police, which have served as an institution of racialized control throughout our nation’s history, are the only group capable of reducing violence.

Community leaders and residents have proved adept at overseeing their neighborhoods, caring for their populations and maintaining safe streets. Studies show that this work lowers crime, sometimes dramatically. What happens if we put those people in charge of containing violence, too?

Over the past 10 years, an expanding body of research has shown just how damaging violence is to community life, children’s academic trajectories and healthy child development. We have rigorous, causal evidence that every shooting in a neighborhood affects children’s sleep and their ability to focus and learn. When a neighborhood becomes violent, it begins to fall apart, as public spaces empty, businesses close, parks and playgrounds turn dangerous, and families try to move elsewhere. Violence is the fundamental challenge for cities: Nothing works if public space is unsafe.

Those who argue that the police have no role in maintaining safe streets are arguing against lots of strong evidence. One of the most robust, most uncomfortable findings in criminology is that putting more officers on the street leads to less violent crime. We know this from randomized experiments involving “hot spots policing” and natural experiments in which more officers were brought to the streets because of something other than crime — a shift in the terror alert level or the timing of a federal grant — and violent crime fell. After the unrest around the deaths of Freddie Gray in Baltimore and Michael Brown in Ferguson, Mo., police officers stepped back from their duty to protect and serve; arrests for all kinds of low-level offenses dropped, and violence rose. This shouldn’t be interpreted to mean that protests against violent policing lead to more violence; rather, it means that when police don’t do their jobs, violence often results.

Considered alongside the brutal response to protests over the past few weeks, this evidence forces us to hold two incongruent ideas: Police are effective at reducing violence, the most damaging feature of urban inequality. And yet one can argue that law enforcement is an authoritarian institution that historically has inflicted violence on black people and continues to do so today.



## E.8.2 Treatment: “Before” wording (rationale)

As part of the campaign, we plan to ask people **to make a post encouraging their friends and followers to sign a petition** opposing the movement to defund the police.

Imagine that you had joined the campaign. If you were going to post **one** of the following two Tweets on your Twitter account, which would you prefer to post?

### Tweet A

I have joined a campaign to oppose defunding the police: <https://bit.ly/3DK3UEr>.

### Tweet B

I have joined a campaign to oppose defunding the police: <https://bit.ly/3DK3UEr>. Before joining, I was shown this article written by a Princeton professor on the strong scientific evidence that defunding the police would increase violent crime: <https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/>

Which of the above Tweets would you have preferred to post on your account?

☐ Tweet A

☐ Tweet B

Please explain why you chose this Tweet rather than the other Tweet.



### E.8.3 Treatment: “After” wording (no rationale)

As part of the campaign, we plan to ask people **to make a post encouraging their friends and followers to sign a petition** opposing the movement to defund the police.

Imagine that you had joined the campaign. If you were going to post **one** of the following two Tweets on your Twitter account, which would you prefer to post?

#### Tweet A

I have joined a campaign to oppose defunding the police: <https://bit.ly/3DK3UEr>.

#### Tweet B

I have joined a campaign to oppose defunding the police: <https://bit.ly/3DK3UEr>. After joining, I was shown this article written by a Princeton professor on the strong scientific evidence that defunding the police would increase violent crime: <https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/>

Which of the above Tweets would you have preferred to post on your account?

☐ Tweet A

☐ Tweet B

Please explain why you chose this Tweet rather than the other Tweet.



1 of 2

## E.9 Auxiliary Experiment 5: Credibility and social cover – Democrats

### E.9.1 Attention screener

The next question is about the following problem. In questionnaires like ours, sometimes there are participants who do not carefully read the questions and just quickly click through the survey. This means that there are a lot of random answers which compromise the results of research studies. To show that you read our questions carefully, please choose **both** “Extremely interested” and “Not at all interested” on the question below.

**Given the text above,** how interested are you in sports?

☐ Extremely interested

☐ Very interested

☐ A little bit interested

☐ Very little interested

☐ Not at all interested

>>

### E.9.2 Background questions

Are you Spanish, Hispanic, or Latino or none of these?

☐ Yes

☐ None of these

What is your year of birth?

What is your sex?

☐ Male

☐ Female

In politics, as of today, do you consider yourself a Republican, a Democrat, or an Independent?

☐ Republican

☐ Democrat

☐ Independent

>>

What is the highest level of school you have completed or the highest degree you have received?

- ☐ Less than high school degree
- ☐ High school graduate (high school diploma or equivalent including GED)
- ☐ Some college but no degree
- ☐ Associate degree in college (2-year)
- ☐ Bachelor's degree in college (4-year)
- ☐ Master's degree
- ☐ Doctoral degree
- ☐ Professional degree (JD, MD)

Which of the following best describes your race or ethnicity?

- ☐ African American/Black
- ☐ Asian/Asian American
- ☐ Caucasian/White
- ☐ Native American, Inuit or Aleut
- ☐ Native Hawaiian/Pacific Islander
- ☐ Other

Who did you vote for in the 2020 presidential election?

- ☐ Donald Trump
- ☐ Joe Biden
- ☐ Other
- ☐ Did not vote

Are you liberal or conservative?

- ☐ Very liberal
- ☐ Liberal
- ☐ Neither liberal nor conservative
- ☐ Conservative
- ☐ Very conservative



Which social media platform do you use the most?

☐ Twitter

☐ Facebook

☐ I do not use Twitter or Facebook

>>

### E.9.3 Pre-treatment outcomes

On the next page, you will be provided with a recent Washington Post column written by **Princeton Professor of Criminology Patrick Sharkey**, in which he discusses evidence showing that more policing leads to less violent crime.

>>



Illustration by Vincent Gadd for The Washington Post

By Patrick Sharkey  
JUNE 12, 2015

## Why do we need the police?

Cops prevent violence. But they aren't the only ones who can do it.

He calls to end policing as we know it contain a sort of trap. The best evidence we have makes clear that police are effective in reducing violence, and without designating some group to combat this problem, efforts to weaken them through budget cuts — “defund the police” — are likely to have unanticipated consequences and to destabilize communities. In many cities this is likely to lead to a rise in violence. And research shows that, when violence increases, Americans of all races become more punitive, supporting harsher policing and criminal justice policies. That’s how we got to this point.



**Patrick Sharkey**  
Patrick Sharkey is a professor of sociology and public affairs at Princeton University’s Woodrow Wilson School of Public and International Affairs. His most recent books are “Unsettled Justice: The Great Crime Decline, the Renewal of City Life, and the Next War on Violence.”

Yet none of this means that the police, which have served as an institution of racialized control throughout our nation’s history, are the only group capable of reducing violence.

Community leaders and residents have proved adept at overseeing their neighborhoods, caring for their populations and maintaining safe streets. Studies show that this work lowers crime, sometimes dramatically. What happens if we put those people in charge of containing violence, too?

Over the past 10 years, an expanding body of research has shown just how damaging violence is to community life, children’s academic trajectories and healthy child development. We have rigorous, causal evidence that every shooting in a neighborhood affects children’s sleep and their ability to focus and learn. When a neighborhood becomes violent, it begins to fall apart, as public spaces empty, businesses close, parks and playgrounds turn dangerous, and families try to move elsewhere. Violence is the fundamental challenge for cities: Nothing works if public space is unsafe.

Those who argue that the police have no role in maintaining safe streets are arguing against lots of strong evidence. One of the most robust, most uncomfortable findings in criminology is that putting more officers on the street leads to less violent crime. We know this from randomized experiments involving “hot spots policing” and natural experiments in which more officers were brought to the streets because of something other than crime — a shift in the terror alert level or the timing of a federal grant — and violent crime fell. After the unrest around the deaths of Freddie Gray in Baltimore and Michael Brown in Ferguson, Mo., police officers stepped back from their duty to protect and serve; arrests for all kinds of low-level offenses dropped, and violence rose. This shouldn’t be interpreted to mean that protests against violent policing lead to more violence; rather, it means that when police don’t do their jobs, violence often results.

Considered alongside the brutal response to protests over the past few weeks, this evidence forces us to hold two incongruent ideas: Police are effective at reducing violence, the most damaging feature of urban inequality. And yet one can argue that law enforcement is an authoritarian institution that historically has inflicted violence on black people and continues to do so today.

Would you like to join a nonpartisan campaign that opposes defunding the police?

☐ Yes

☐ No

>>

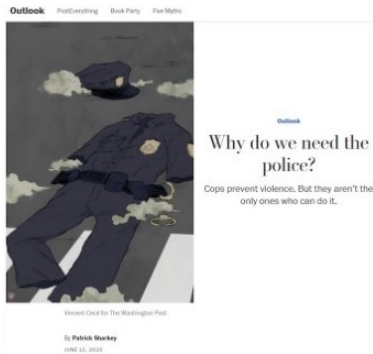
**You have successfully joined the campaign.**

Since you chose to join the campaign, we wanted to give you more time reading the Washington Post column written by **Princeton Professor of Criminology Patrick Sharkey**, where he discusses evidence showing that more policing leads to less violent crime.

The article is available on the next page, and you can spend as much time as you want reading it before you continue with the remaining part of the survey.

>>





He calls to end policing as we know it contain a sort of trap. The best evidence we have makes clear that police are effective in reducing violence, and without designating some group to combat this problem, efforts to weaken them through budget cuts — “defund the police” — are likely to have unanticipated consequences and to destabilize communities. In many cities this is likely to lead to a rise in violence. And research shows that, when violence increases, Americans of all races become more punitive, supporting harsher policing and criminal justice policies. That’s how we got to this point.



**Patrick Sharkey**  
Sharkey is a professor of sociology and public affairs at Princeton University’s Woodrow Wilson School of Public and International Affairs. His most recent book is “Uneasy Peace: The Great Crime Decline, the Renewal of City Life, and the Next War on Violence.”

Yet none of this means that the police, which have served as an institution of racialized control throughout our nation’s history, are the only group capable of reducing violence.

Community leaders and residents have proved adept at overseeing their neighborhoods, caring for their populations and maintaining safe streets. Studies show that this work lowers crime, sometimes dramatically. What happens if we put those people in charge of containing violence, too?

Over the past 10 years, an expanding body of research has shown just how damaging violence is to community life, children’s academic trajectories and healthy child development. We have rigorous, causal evidence that every shooting in a neighborhood affects children’s sleep and their ability to focus and learn. When a neighborhood becomes violent, it begins to fall apart, as public spaces empty, businesses close, parks and playgrounds turn dangerous, and families try to move elsewhere. Violence is the fundamental challenge for cities: Nothing works if public space is unsafe.

Those who argue that the police have no role in maintaining safe streets are arguing against lots of strong evidence. One of the most robust, most uncomfortable findings in criminology is that putting more officers on the street leads to less violent crime. We know this from randomized experiments involving “hot spots policing” and natural experiments in which more officers were brought to the streets because of something other than crime — a shift in the terror alert level or the timing of a federal grant — and violent crime fell. After the unrest around the deaths of Freddie Gray in Baltimore and Michael Brown in Ferguson, Mo., police officers stepped back from their duty to protect and serve; arrests for all kinds of low-level offenses dropped, and violence rose. This shouldn’t be interpreted to mean that protests against violent policing lead to more violence; rather, it means that when police don’t do their jobs, violence often results.

Considered alongside the brutal response to protests over the past few weeks, this evidence forces us to hold two incongruent ideas: Police are effective at reducing violence, the most damaging feature of urban inequality. And yet one can argue that law enforcement is an authoritarian institution that historically has inflicted violence on black people and continues to do so today.

#### E.9.4 Treatment (higher-credibility): “Before” wording (rationale)

As part of this nonpartisan campaign, we consider asking people to **publish a post on their Twitter profile encouraging their friends and followers to sign a petition** opposing the movement to defund the police.

We are therefore interested in whether you would have been willing to publish the post below on your Twitter profile if it was included as a campaign feature.

I have joined a campaign to oppose defunding the police: [bit.ly/3DK3UEr](https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/). Before joining, I was shown this article written by a Princeton professor on the strong scientific evidence that defunding the police would increase violent crime:  
<https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/>



washingtonpost.com

Perspective | Cops prevent violence. But they aren't the only ones wh...  
Communities already know how to police their own. Now put them in charge of it.

Would you have been willing to publish the post above on your Twitter profile?

☐ Yes

☐ No

>>

E.9.5 Treatment (higher-credibility): “After” wording (no rationale)

As part of this nonpartisan campaign, we consider asking people to **publish a post on their Twitter profile encouraging their friends and followers to sign a petition** opposing the movement to defund the police.

We are therefore interested in whether you would have been willing to publish the post below on your Twitter profile if it was included as a campaign feature.

I have joined a campaign to oppose defunding the police: [bit.ly/3DK3UEr](https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/). After joining, I was shown this article written by a Princeton professor on the strong scientific evidence that defunding the police would increase violent crime:  
<https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/>



washingtonpost.com

Perspective | Cops prevent violence. But they aren't the only ones wh...  
Communities already know how to police their own. Now put them in charge of it.

Would you have been willing to publish the post above on your Twitter profile?

☐ Yes

☐ No

>>

### E.9.6 Treatment (lower-credibility): “Before” wording (rationale)

As part of this nonpartisan campaign, we consider asking people to **publish a post on their Twitter profile encouraging their friends and followers to sign a petition** opposing the movement to defund the police.

We are therefore interested in whether you would have been willing to publish the post below on your Twitter profile if it was included as a campaign feature.

I have joined a campaign to oppose defunding the police: [bit.ly/3DK3UEr](https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/). Before joining, I was shown this article arguing that defunding the police would increase violent crime:

<https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/>



washingtonpost.com

Perspective | Cops prevent violence. But they aren't the only ones wh...  
Communities already know how to police their own. Now put them in charge of it.

Would you have been willing to publish the post above on your Twitter profile?

☐ Yes

☐ No

>>

### E.9.7 Treatment (lower-credibility): “After” wording (no rationale)

As part of this nonpartisan campaign, we consider asking people to **publish a post on their Twitter profile encouraging their friends and followers to sign a petition** opposing the movement to defund the police.

We are therefore interested in whether you would have been willing to publish the post below on your Twitter profile if it was included as a campaign feature.

I have joined a campaign to oppose defunding the police: [bit.ly/3DK3UEr](https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/). After joining, I was shown this article arguing that defunding the police would increase violent crime:  
<https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/>



washingtonpost.com

Perspective | Cops prevent violence. But they aren't the only ones wh...  
Communities already know how to police their own. Now put them in charge of it.

Would you have been willing to publish the post above on your Twitter profile?

☐ Yes

☐ No

>>

### E.9.8 Perceived social punishment

A few weeks ago, we asked a sample of Democrats whether they would approve or deny a \$1 bonus (at no cost to themselves) to a matched survey participant.

**They were told that their matched participant had been willing to publish the post on the previous page** on their Twitter profile. They were not told anything else about their matched participant.

How many percent of Democrats do you think chose to deny a \$1 bonus to their matched participant?

0 10 20 30 40 50 60 70 80 90 100

Percent



>>

## E.10 Auxiliary Experiments 6: Interpretation of dissent with low-credibility rationale – Democrats

### E.10.1 Treatment: “Before” condition (rationale)

We conducted a survey about political and social attitudes in the United States earlier this year. You have been matched with one of the participants from that survey.

Your matched participant **joined a campaign to oppose defunding the police.**

Your matched participant **also agreed to post the Tweet below** on their Twitter profile to promote the campaign.

I have joined a campaign to oppose defunding the police: [bit.ly/3DK3UEr](https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/). Before joining, I was shown this article that argues that defunding the police would increase violent crime:  
<https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/>



washingtonpost.com

Perspective | Cops prevent violence. But they aren't the only ones who can do it. Communities already know how to police their own. Now put them in charge of it.

Why do you think your matched respondent chose to join the campaign to oppose defunding the police?





### Matched Respondent's Donation Decision

We gave your matched respondent the opportunity to donate \$5 to the **National Association for the Advancement of Colored People (NAACP)**, America's oldest and largest civil rights organization.

Below, we will ask you to guess whether or not your matched respondent donated \$5 to the National Association for the Advancement of Colored People (NAACP).

**Reminder:** Your matched participant agreed to post the Tweet below on their Twitter account.

I have joined a campaign to oppose defunding the police: [bit.ly/3DK3UEr](https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/). Before joining, I was shown this article that argues that defunding the police would increase violent crime:  
<https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/>



Do you think that your matched participant chose to donate \$5 to the National Association for the Advancement of Colored People (NAACP)?

☐ Yes, I think my matched respondent chose to donate

☐ No, I think my matched respondent **did not** choose to donate





You now have the opportunity to authorize a \$1 bonus payment to your matched respondent. **The bonus payment will not be deducted from your payment.** Your matched respondent did not know you would have the opportunity to decide their bonus.

**Reminder:** Your matched participant agreed to post the Tweet below on their Twitter account.

I have joined a campaign to oppose defunding the police: [bit.ly/3DK3UEr](https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/). Before joining, I was shown this article that argues that defunding the police would increase violent crime:  
<https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/>



washingtonpost.com

Perspective | Cops prevent violence. But they aren't the only ones who can do it. Communities already know how to police their own. Now put them in charge of it.

Do you want to authorize a \$1 bonus to your matched respondent?

- ☐ Yes, I would like to authorize a \$1 bonus
- ☐ No, I would not like to authorize a \$1 bonus



## E.10.2 Treatment: “After” condition (no rationale)

We conducted a survey about political and social attitudes in the United States earlier this year. You have been matched with one of the participants from that survey.

Your matched participant **joined a campaign to oppose defunding the police.**

Your matched participant **also agreed to post the Tweet below** on their Twitter profile to promote the campaign.

I have joined a campaign to oppose defunding the police: [bit.ly/3DK3UEr](https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/). After joining, I was shown this article that argues that defunding the police would increase violent crime:  
<https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/>



Why do you think your matched respondent chose to join the campaign to oppose defunding the police?



### Matched Respondent's Donation Decision

We gave your matched respondent the opportunity to donate \$5 to the **National Association for the Advancement of Colored People (NAACP)**, America's oldest and largest civil rights organization.

Below, we will ask you to guess whether or not your matched respondent donated \$5 to the National Association for the Advancement of Colored People (NAACP).

**Reminder:** Your matched participant agreed to post the Tweet below on their Twitter account.

I have joined a campaign to oppose defunding the police: [bit.ly/3DK3UEr](https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/). After joining, I was shown this article that argues that defunding the police would increase violent crime:  
<https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/>



washingtonpost.com  
Perspective | Cops prevent violence. But they aren't the only ones who can do it. Communities already know how to police their own. Now put them in charge of it.

Do you think that your matched participant chose to donate \$5 to the National Association for the Advancement of Colored People (NAACP)?

- ☐ Yes, I think my matched respondent chose to donate
- ☐ No, I think my matched respondent **did not** choose to donate



You now have the opportunity to authorize a \$1 bonus payment to your matched respondent. **The bonus payment will not be deducted from your payment.** Your matched respondent did not know you would have the opportunity to decide their bonus.

**Reminder:** Your matched participant agreed to post the Tweet below on their Twitter account.

I have joined a campaign to oppose defunding the police: [bit.ly/3DK3UEr](https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/). After joining, I was shown this article that argues that defunding the police would increase violent crime:

<https://www.washingtonpost.com/outlook/2020/06/12/defund-police-violent-crime/>



washingtonpost.com

**Perspective | Cops prevent violence. But they aren't the only ones who can do it.**  
Communities already know how to police their own. Now put them in charge of it.

Do you want to authorize a \$1 bonus to your matched respondent?

☐ Yes, I would like to authorize a \$1 bonus

☐ No, I would not like to authorize a \$1 bonus



## E.11 Auxiliary Experiment 7: Persuasion experiment – Republicans

### E.11.1 Pre-treatment beliefs

Please see the short video below where Fox News host **Tucker Carlson presents evidence on whether illegal immigrants commit more crime.**



>>

**E.11.2 Information treatment (only shown to respondents in the treatment group)**

To what extent do you agree with the following statement: "The United States should immediately deport all illegal Mexican immigrants."

☐ Strongly agree

☐ Agree

☐ Neither agree nor disagree

☐ Disagree

☐ Strongly disagree



### E.11.3 Post-treatment outcomes

To what extent do you agree with the following statement: "Illegal immigrants are not much more likely to commit serious crimes than U.S. citizens."

☐ Strongly agree

☐ Agree

☐ Neither agree nor disagree

☐ Disagree

☐ Strongly disagree

