# Can the Government Deter Discrimination?
## Evidence from a Randomized Intervention in New York City[*]

Albert H. Fang[†]     Andrew M. Guess[‡]     Macartan Humphreys[§]

March 11, 2017

**Abstract**

Racial discrimination persists despite established anti-discrimination laws. A common government strategy to deter discrimination is to publicize the law and communicate potential penalties for violations. We study this strategy by coupling an audit experiment with a randomized intervention involving nearly 700 landlords in New York City and report the first causal estimates of the effect on rental discrimination against Blacks and Hispanics of a targeted government messaging campaign. We uncover discrimination levels higher than prior estimates indicate, especially against Hispanics, who are approximately six percentage points less likely to receive callbacks and offers than whites. We find suggestive evidence that government messaging can reduce discrimination against Hispanics, but not against Blacks. The findings confirm discrimination's persistence and suggest that government messaging can address it in some settings, but more work is needed to understand the contexts under which such appeals are most effective.

Abstract Word Count: 143

Article Word Count: 9064

Keywords: government communication; discrimination; political economy of race; behavioral policy compliance; field experiment

---

[†]Postdoctoral Associate, Institution for Social and Policy Studies, Yale University (albert.fang@yale.edu)
[‡]Postdoctoral Fellow, Social Media and Political Participation Lab, New York University (guess@nyu.edu)
[§]Professor, Department of Political Science, Columbia University (mh2245@columbia.edu)

The passage of civil rights laws in the 1960s marked a turning point in the development of racial politics in the United States. Against the historical backdrop of slavery, Jim Crow, longstanding racial inequalities, and the marginalization of non-whites, the legal prohibition of discrimination on the basis of race, color or national origin signaled an important political shift toward a more egalitarian racial order (King and Smith 2005; Massey and Denton 1993). In addition to examining the political conditions leading to the passage of civil rights laws, scholars of the politics surrounding the formal end of Jim Crow have focused considerable attention on understanding its aftermath. In particular, scholars have documented how racial inequalities have persisted despite established civil rights laws, and have provided a constellation of policy-centered, political economy explanations for their persistence (e.g., Alexander 2012; Forman 2012). Central to these explanations is the argument that racial disparities persist because racial discrimination is both pervasive and persistent in the contemporary neoliberal economy (Dawson 2016; Dawson and Francis 2015; Bobo, Kluegel and Smith 1997). That racial discrimination persists despite the existence of anti-discrimination laws calls into question the various roles that governments play in either actively curbing or permissively allowing discriminatory economic behavior, as well as the downstream consequences of government efforts to enforce anti-discrimination laws for racial inequalities, the political economy of race, and the development of racial orders.

Understanding this puzzle requires further developing a body of research organized around two related lines of inquiry. The first asks: *To what extent does racial discrimination persist and why?*[1] The first-order challenge of measuring discrimination levels is not trivial. For its potential targets, discrimination is difficult to observe because in many cases it is impossible observe how a counterfactual individual would be treated in a market interaction. Going beyond first-hand accounts, third-party observations of discrimination by lay bystanders are imperfect because these observations only occur if the third party is present and because, as has been documented by prior psychological research, many individuals are unwilling to "see" or make attributions to discrimi-

---

[1]In this article we focus on understanding the incidence of discrimination. We bracket questions about the causes of discrimination, which have been widely studied elsewhere by scholars in economics and sociology, as one that is beyond the scope of this project. See, for example, Heckman (1998); Heckman and Siegelman (1993); Ewens, Tomlin and Wang (2014); Neumark (2011); Bertrand, Chugh and Mullainathan (2005); Pager and Karafin (2009).

nation even when confronted with direct evidence of discrimination.[2] To address this challenge, governments use enforcement audits to collect evidence of discrimination on a case-by-case basis, and both scholars and governments have relied on audit studies to measure the aggregate-level incidence of discrimination in employment and housing[3]. But despite their proliferation, high-quality audits are both costly and relatively rare, and most audit studies focus on discrimination in early-stage market interactions and do not measure end-line outcomes (such as being offered a job or a housing unit), which are important to measure because they map onto levels of downstream racial inequalities and can provide evidence of whether disparate impacts exist. The second, related line of inquiry asks: *Which government strategies are effective at reducing discrimination and why?* Despite an abundance of theoretical work on policy enforcement in general and much historical research on the primacy of legal strategies in efforts to enforce anti-discrimination laws since their initial passage,[4] to our knowledge no work exists that theorizes the conditions under which governmental strategies to reduce discrimination are effective and that experimentally tests hypothesized expectations. Answering this second question is doubly difficult because it requires both measure of a hidden behavior and causal inferences around those behaviors.

In this article, we address both of these major questions and provide the first attempt in the literature to experimentally evaluate the effects on racial discrimination levels of any government strategy to reduce discrimination in the United States. We study the effectiveness of pre-emptive strategies to deter racial discrimination using official communication campaigns encouraging compliance with anti-discrimination law, and we focus on two commonly employed but understudied appeals: making the law itself salient and making the costs of violating the law salient. We work at a large scale which gives us a good handle on discrimination levels although, given fundamental measurement and inferential challenges, our estimates of treatment effects are still measured with considerable uncertainty.

---

[2]For example Eyer (2012); Dixon, Storen and Van Horn (2002); Nielsen (2004); Swim et al. (2003); Foster and Dion (2004); Kappen and Branscombe (2001); Major and Dover (2016)

[3]For example Bertrand and Duflo (2017); Pager and Shepherd (2008); Turner et al. (2013, 2001); Yinger (1986, 1995); Roychoudhury and Goodman (1992, 1996)

[4]For example Rosenberg (1991); Cover (1995); Epp (1998); Frymer (2003); Ritter (2007); Nielsen, Nelson and Lancaster (2010); Skrentny (2002).

Our analysis examines discrimination and interventions in field conditions, in the context of government efforts to enforce fair housing law in the New York City rental market. We focus on the rental market as it is the segment of the housing market in which reported discrimination is the most pervasive.[5] We partnered with the New York City municipal government to implement a large-scale field experiment that began in 2012 and lasted 20 months. We assess the effects on racial discrimination levels against Black and Hispanic rental applicants (as compared to whites) of government appeals that are delivered via targeted and personalized phone calls to landlords and brokers[6] who interact with confederates posing as rental housing applicants. The city randomly assigned to nearly 700 landlords either (1) a targeted live phone call from the city that drew attention to fair housing law and implicitly signaled increased government monitoring of housing agents (a "monitoring" condition), (2) a targeted live phone call that contained the contents of the monitoring message and additionally communicated information about the potential pecuniary punishments for discriminatory practices (a "punitive" condition), or (3) no call from the city (a pure control condition). The city then sent matched triples of confederates who varied by race (Black, Hispanic, and white) to visit the same advertised unit—meeting in person with the same landlord—and to collect detailed qualitative notes on their interactions over the course of the housing search process. The enhanced audit design embedded in the experiment provides leverage above and beyond existing audit methodologies to measure both subtle and direct forms of racial bias that pervade the housing market. Additionally, in contrast to prior housing audit studies focusing on differential treatment in landlords' responses to initial inquiries about rental listings, this study focuses on two outcomes occurring toward the end of the housing search process that are consequential for downstream inequalities in who gets to live where: discrimination in receiving a callback following an in-person meeting with the landlord to view the listed unit and discrimination in receiving a rental offer from the landlord after the appointment.

---

[5]According to the National Fair Housing Alliance (2014), of the 18,978 housing discrimination complaints reported by private fair housing groups in 2014, 16,694 (or about 88%) involved discrimination in the rental market.

[6]For ease of exposition, we refer to experimental subjects, defined as landlords or brokers associated with a sampled rental housing advertisement and who interact with a matched team of testers, simply as "landlords" in the article.

The theoretical framework motivating the experiment is relatively simple. Drawing on theoretical and empirical results from existing research, we expect that the monitoring condition should decrease racial discrimination levels relative to control, because receiving a personally targeted monitoring signal from a government known to be active in enforcing fair housing law is expected to reduce landlords' perceived benefits from discriminating and to increase their perceptions of the probability of punishment if they discriminate. In contrast, existing debates in the literature offer mixed expectations about the conditions under which coercive sanctions are effective at inducing pro-social behavior, with some arguments proposing possibly adverse effects; thus we are agnostic about the expected direction of the effect of punitive messaging.

Analyzing data from the field experiment, we find evidence of substantial baseline levels of discrimination in the New York City rental housing market, particularly against Hispanics: they are 28% less likely to receive a callback (6.1 percentage points, 95% C.I. = [1.00, 11.19]) and 49% less likely to receive an offer for an apartment than whites (5.7 percentage points, 95% C.I. = [1.32, 10.16]). This result is striking when compared to discrimination levels reported by the 2012 US Department of Housing and Urban Development (HUD) audit of racial rental discrimination in the New York metropolitan area, which found a lack of discrimination against minorities in receiving callbacks (discrimination in receiving offers was not studied). Using a larger sample where estimates are less likely to be affected by sampling variability, we show that racial discrimination in fact persists in the New York City rental market. We also find suggestive evidence that treatment messages reduce levels of discrimination in receiving callbacks, though these effects are statistically weak and are observed only for discrimination against Hispanics. There is no evidence that discrimination against Blacks is affected by these government messaging interventions. We also find no evidence that the punitive messages are more effective than simple monitoring messages. Our results suggest that a bundled governmental messaging strategy that invokes fair housing law and makes salient the costs of punishment among lawbreakers can be effective at reducing racial housing discrimination, but not under all circumstances. Accordingly, we argue that additional work is needed to theorize and experimentally test the conditions under which such appeals by the

4

government are effective at reducing racial discrimination.

In the next section we present the theoretical framework and expectations. In doing so, we situate the research questions and motivate the field experiment in the context of existing debates in the literature. We then describe the experiment's context, design, and our estimation strategy. Next, we present the main results. We close by discussing the implications of our findings and avenues for future research.

## THEORETICAL FRAMEWORK AND EXPECTATIONS

### Government Strategies to Reduce Racial Discrimination in the United States

Government strategies to enforce anti-discrimination law may be classified into one of two broad types: reactive or pre-emptive. Reactive strategies such as "fire-alarm" and "police patrol" models of policy enforcement (McCubbins and Schwartz 1984) are those in which governments investigate and litigate potential discriminators as an *ex post* response conditional on citizens and governments, respectively, observing and reporting discrimination in the first place. Given the challenges associated with observing discrimination (as previously discussed) and given the expectation that pursuing litigation may be disproportionately costly – both financially and psychologically – for the targets of discrimination relative to defendants (Feagin 1991; Bumiller 1988; Donohue and Siegelman 1991; Galanter 1974), reactive strategies are limited in their ability to identify, much less punish, most individuals who discriminate.

In contrast, pre-emptive strategies are those in which official messaging campaigns attempt to deter discrimination *ex ante* by advertising existing anti-discrimination laws and making various appeals that attempt to induce citizens to comply with the law and not discriminate. Such campaigns complement reactive strategies and are attractive in part because, if effective, they could offer a lower-cost means by which governments can induce behavioral change (as compared to deploying an intensive and constant auditing program). However, existing government communication campaigns typically involve bundled messages and appeals that aim to reach multiple potential audiences that include the targets of discrimination, discriminators, and third-party observers.

5

Disentangling the effects of specific appeals on its intended recipients is important in order to test whether certain appeals are more effective at reducing discrimination than others and whether the effectiveness of an appeal may be attributed to a theorized psychological process (such as by making salient certain considerations that alter the perceived returns to a behavior). In particular, there is a need to empirically test whether specific appeals that actually affect potential discriminators are effective at reducing discrimination by affecting their perceptions of the expected returns to discriminating.

**When Do Messaging Campaigns Reduce Discriminatory Behavior?**

Under what conditions might government messaging campaigns be effective at reducing discriminatory behavior? Existing research undertheorizes why official appeals would be effective at reducing discrimination. We develop a simple decision-theoretic model to help clarify the basic logic. In the context of a housing market interaction, the individual of interest is the landlord, broker, or housing agent who interacts with housing applicants and may potentially engage in discriminatory behavior. Here, discrimination (as an individual-level behavior) simply refers to the differential treatment of two individuals who vary on a relevant attribute such as race. This definition remains agnostic as to why people discriminate, which may arise from group-specific prejudices and a taste for discrimination (Becker 1957) or for strategic reasons, where people make inferences about a person using stereotypes about the average member of a group to which that person is perceived to belong (e.g., Phelps 1972; Arrow 1973; Aigner and Cain 1977).

Let individual utility from discrimination given government policy $x$ be expressed as:

$$u(x) = b(x) - p(x)c(x) + e_i$$

where $b$ denotes marginal benefits from discriminating (which could be driven by prejudice or statistical discrimination, for example), $p$ the subject's expectation of punishment when she discriminates, $c$ the expected severity of the punishment, and $e_i$ individual heterogeneity, which is distributed according to $F$. In this case the expected amount of discrimination is simply $1 -$

$F(-b(x) + p(x)c(x))$. We make explicit the fact that in principle each of these three parameters can be affected by government policy — appeals. Most obviously, the beliefs about $p$ are plausibly altered upwards, as government appeals communicate the fact that government is active on the issue; $c$ can also be affected, either upwards or downwards depending on the content of the appeal and the subjects' priors; less obviously $b$ — the subjective benefits of discrimination — could also be affected by an appeal, either upwards or downwards, following logics outlined below.

Building on this framework we aim to estimate both the baseline level of housing discrimination against Black and Hispanic renters as well as the causal effects on racial housing discrimination of governmental messaging that encourages landlords and brokers to comply with fair housing law. We do so in the context of the New York City rental housing market and in the City of New York's jurisdiction. We examine two distinct types of policy in this study: a monitoring signal that invokes the law and implicitly signals increased government monitoring of housing agents, and a punitive appeal that additionally primes the costs of discrimination. For practical reasons—in particular to satisfy our implementing partner's (i.e., the government's) preference for treatment realism and to avoid depleting statistical power—the experiment is designed to test realistic, bundled treatments, but does so at the expense of being able to test specific parameters. We are cognizant that this is a major limitation of this study, and we recommend future work to design and test treatments that target each parameter separately. Nevertheless, the framework we specify above remains useful to develop theoretical expectations about why the appeals tested in this experiment would affect discriminatory behavior. As we explain further below, both treatments plausibly alter $p(x)$ while the second is intended to alter $c(x)$ differentially. Moreover, we are cognizant of the fact that $b$ may also be responsive to $x$.

We test the hypotheses that sending monitoring messages alone decreases rental housing discrimination rates against Blacks and Hispanics (hypothesis family *H1*). There are two possible reasons for why it might. First, simply contacting landlords and making the law salient may raise expectations of sanctioning, $p$, even if this is not explicitly invoked and no information on sanctions is provided. Given the context of this study, the mode of treatment delivery, and the particular

messenger of the treatment appeal, we argue that this is plausible in this experimental setting. Because the treatment calls in this study are personalized and targeted messages for each landlord, because they are sent by a city government with strong financial and administrative resources to monitor and enforce fair housing law, and because the city has a known history of actively monitoring and enforcing fair housing law, receiving such a call provides a strong and credible signal to the landlord that the city is already tracking the subject and monitoring their behavior. Thus receiving any targeted call from the city may increase $p$ above zero.

It is also possible however that the law increases the salience of injunctive norms about fair housing in the landlord's thinking and thus reduce the value of $b$. Simply invoking the law may activate compliance norms and may also render compliance focal (McAdams and Nadler 2005; Tyler 2006). For discrimination in particular, Mendelberg (2001) and others have highlighted how *explicit* priming of discrimination considerations can result in less discrimination (relative to implicit priming) due to an invoking of social nondiscrimination norms. In addition, if discriminatory behavior is perceived as a descriptive norm among housing agents,[7] descriptive norm perceptions may drive beliefs that the law itself is illegitimate and thus result in lower levels of compliance (Tyler 2004, 2006). In settings where these "peer effects" exists, appeals that communicative injunctive norms have been theorized to be effective at crowding out perceptions of descriptive norms that are used to justify noncompliance (Weaver 2015). Our estimates of the monitoring effect thus capture joint effects on $p$ and, to the extent that preferences are not fixed, $b$.

In addition we seek to assess is the effect of altering $c$; we do so by assessing the effect of sending "punitive" messages, relative both to the control condition (thus the combined effect of monitoring and punitive content) and to the monitoring condition (thus the additional effect of punitive content as compared to the monitoring condition). Thus in practice we assess the hypotheses that sending punitive messages decreases rental housing discrimination rates against Blacks and Hispanics relative to the baseline condition and to a monitoring condition (hypothesis family *H2*).[8]

---

[7]For example, if landlords think that it is common practice to discriminate against non-white tenants because there are shared negative stereotypes about the characteristics of minority tenants and shared beliefs that renting to white tenants instead of non-white tenants generates greater returns in terms of expected future rental income.

[8]For the latter comparison, the monitoring condition can be conceptualized as a kind of placebo condition, such

Although we have clear expectations about the marginal effect of costs, $c$, we are cognizant that our punitive treatment may also have the effect of altering $b$, relative to the monitoring treatment.

Existing research suggests mixed expectations about whether amplifying the salience of the severity of punishment will be effective at reducing discrimination. On the one hand, theoretical work on the politics of inducing behavioral compliance with public policies argues that when individuals face insufficient incentives to comply with the law or perceive complying with the law as being costly (either in terms of direct costs or perceived opportunity costs), a potentially effective strategy to induce compliant behavior involves increasing and making salient the costs of noncompliance such that they are perceived to be greater than the costs of compliance (e.g., Dixit 2006; Hadfield and Weingast 2014; Weaver 2014, 2015). Prior research evaluating the effectiveness of similar official communication campaigns has shown that appeals highlighting the costs of non-compliance can, under certain circumstances, induce greater compliance with the law in domains such as paying taxes (Blumenthal et al. 2001; Slemrod, Blumenthal and Christian 2001; Iyer, Reckers and Sanders 2010, but see Dunning et al. (2015)) and paying delinquent fines (Haynes et al. 2013). On the other hand, a growing line of social psychological research challenges this expectation and argues that coercive sanctions to induce behavioral change may backfire and instead crowd out social norms needed to motivate pro-social behavior (e.g., Bowles 2016; Gneezy and Rustichini 2000; Cardenas, Stranlund and Willis 2000). Thus while we gather direct evidence on the effect of the punitive treatment, this can be interpreted as the effect of costs, $c$, only under the assumption of no *differential* effects, relative to the monitoring treatment, on $b$ and $p$.

## EXPERIMENTAL DESIGN

Partnering with the City of New York, we designed and analyze data from a field experiment where the city randomly sent targeted messages to landlords associated with a specific advertised rental unit, which is pursued by a trio of matched testers who vary by race. The matched audit provides leverage to measure racial discrimination, which we operationalize as differential landlord

---

that comparing discrimination levels between the punitive and monitoring conditions allows us to assess the pure effect of explicit priming of punishment on responses by landlords.

behavior toward testers by the tester's race for the same advertised unit. This section describes the experiment's setting, design, implementation, and analysis.

## Study Context

New York City is a useful political setting for testing the effectiveness of official communication campaigns encouraging compliance with anti-discrimination and fair housing law for two reasons. It is well known for having one of the strongest anti-discrimination laws in the country.[9] In addition, the agency in charge of enforcing the city's anti-discrimination law, the New York City Commission on Human Rights, is well known among real estate and legal professionals for having the administrative and financial capacity to enforce fair housing law and for applying that capacity to those ends. Thus, the study context allows us to adduce that receiving any targeted and personalized treatment phone message from the city government is arguably perceived as a credible signal of increased government effort to monitor landlords and enforce fair housing law.

We focus in particular on studying the behavior of landlords and brokers who post rental listings on the online classified advertising site Craigslist. While Craigslist is one of multiple sources of classified rental advertisements that housing agents and housing seekers could use in the tenant and housing search processes, respectively, it is the primary forum used to advertise and pursue rental listings in general and in New York City.[10]

## Audit Design

We briefly summarize how rental listings are sampled from Craigslist and how they are pursued by matched auditors.[11] On each day of the study's implementation, a set of rental housing ads from the current day was selected from Craigslist using an automated script. First, using keywords, a list of "likely discrimination" (LD) ads were identified and selected with 100% probability with the

---

[9]Title 8 of the Administrative Code of the City of New York. See also Gurian (2005).

[10]As a prior rental market discrimination audit study by Ewens, Tomlin and Wang (2014) reports, "Craigslist receives 95% of visits to online classified sites" and "approximately 2.5% of all U.S. Internet visits are to Craigslist... [as compared to]... only 0.14% of U.S. Internet visits [for] other classified websites combined" (126).

[11]Complete details of the sampling and audit procedures are provided in the Supplemental Information (SI) online.

goal of increasing statistical power if baseline discrimination levels were low.[12]  Second, among

the remainder of ads posted on Craigslist that day (excluding those identified in the LD search), a

sample of ads was randomly drawn in a way that was representative of the distribution of advertised

vacant rental housing stock by New York City borough.[13]  Only advertisements that invited housing

seekers to reply by phone were pursued.[14]

Testers responded to sampled ads posing as individuals interested in renting the listed apart-

ment.  A project manager pre-screened all sampled ads against a master database of prior audits

to ensure that there were no duplicate landlords or brokers in the sample of ads ever pursued by

testers. Each ad was pursued in a randomized order by a matched team of three testers of the same

gender who varied by race: one white, one Black, and one Hispanic.[15]  By extensively controlling

for aspects of testers' assumed biographies within each trio, we employ a design-based approach

to reduce the likelihood that observed racial discrimination levels are driven by statistical discrim-

ination.[16]  Assigned biographies logically consistent with the rental price and size of the advertised

unit were automatically generated at the time of ad sampling.

Upon reaching an individual when replying to an ad, testers were instructed to provide *limited*

information about themselves in this initial stage and were instructed to schedule an appointment to

view the unit and meet with the person on the same day at the earliest convenience of the individual

with whom they were speaking. What was revealed over the phone typically included the tester's

---

[12]All LD ads were pursued (n=156) and of these, 44 were associated with landlords admitted into the experiment.
The inclusion of the 44 (of 653) subjects associated with LD ads may increase the mean level of discrimination in
the control group and decrease the estimated effect if the landlords associated with these ads are people who are
more likely (than landlords associated with non-LD ads) to discriminate in the absence of treatment and who will
discriminate regardless of which treatment arm they're assigned to. We re-estimate the ITT effects excluding the LD
subsample and find that the main results are not materially affected when subjects from LD ads are excluded. Results
are shown in the SI.

[13]35% Manhattan, 30% Brooklyn, 20% Queens, 10% Bronx, and 5% Staten Island.  These are based on the
distribution of ads by borough as identified in a pilot study.

[14]In a pilot study, we found that the contact rate and appointment scheduling rate are significantly higher when
replying to ads by phone instead of email, and that response rates for emailed inquiries were extremely low.

[15]The features of testers' assumed biographies that were matched included their credit score range, income level,
household composition, occupation, perceived stability of employment and income, gender, interpersonal skills, and
perceived age.

[16]While scholars have argued that conclusions about market discrimination cannot be inferred from information on
differential treatment in paired audits (Heckman 1998; Heckman and Siegelman 1993), these objections do not threaten
the identification of treatment effects directly because our design uses a matched audit as an outcome measurement
strategy in an experimental setting and we maintain measurement symmetry across treatment groups.

assumed name,[17] their interest in pursuing and renting the unit, their availability for an in-person appointment to view the unit, and their financial qualifications to rent. If testers were asked about other aspects of their biographies, they volunteered this information accordingly. During this early stage of the housing search process, landlords are provided minimal but adequate information about the identity of housing seekers (such as their name, voice, and information about their assumed biographies) that can affect how they interact with testers prior to random assignment.[18]

When all three testers in a matched team successfully secured an appointment to view the same advertised housing unit with the same landlord, the landlord associated with that unit was admitted into the experiment and randomly assigned to a treatment condition. We refer to the landlord-housing unit (and thus the unit of randomization) as a "case." Testers then made their individually scheduled appointments, viewed the unit, interacted with the subject, and recorded their interactions afterwards. Testers were blind to treatment, but they were not blind to the city's interest in assessing discrimination as one of a laundry list of characteristics of the rental housing market. As such, they were extensively instructed and coached not to fish for particular reactions, not to let their personal opinions about landlord behavior interfere with their ability to continue interactions during the audit, and not to raise suspicions of an audit.

**Definition of Treatments and Random Assignment**

Landlords admitted into the experiment were randomly assigned to one of three conditions: a control condition where no message was sent, a monitoring messaging condition, or punitive messaging condition.

Treatment scripts were designed by the city and includes official language that the city uses to communicate aspects of fair housing law to the public.[19] In both the monitoring and punitive con-

---

[17]Testers mentioned their first name when introducing themselves on the call; if landlords or brokers asked for the tester's full name, this was provided. Testers had an assumed first name and an assumed last name for the duration of the study, each of which was randomly drawn from a database of names tagged with racial and ethnic identifiers. See the SI for details on how assumed names were generated.

[18]While testers were screened at the time of hire such that their manner of speaking did not strongly signal a particular racial identity over the phone, testers' assumed names did contain a signal of their race (Butler and Homola Forthcoming; Fryer and Levitt 2004).

[19]The full text of the treatment scripts is provided in the SI.

ditions, a city employee delivered the assigned appeal via a personalized, targeted phone message to the landlord. Calls were sent about two hours after testers successfully scheduled appointments over the phone and about two hours before the first scheduled appointment.

In the monitoring condition, the treatment script verified that the subject was on the line, informed the subject that the call was from the New York City Commission on Human Rights *"as part of an ongoing informational campaign to remind landlords and brokers of their obligations under fair housing law,"* and provided the Commission's web address for more information. While the monitoring message does not provide specific information about what exactly is illegal and is in fact designed to avoid priming subjects to think about racial discrimination in particular, our prior is that most of the subjects already knew about fair housing law in the absence of treatment.[20] Thus the monitoring message can be interpreted an intensive intervention that makes salient injunctive norms surrounding fair housing law simply by invoking it. In the punitive condition, the treatment script was the same as the monitoring treatment script but includes, prior to providing the Commission's web address, the following advisement emphasizing the punitive power of the state and the potential pecuniary costs of violating the city's fair housing law: *"It is illegal to discriminate against a person seeking housing due to their membership in a protected class. If you are found to have broken the law, you may be ordered to pay damages, provide reasonable accommodation, or incur civil penalties of up to $250,000."*

A block randomization procedure was used where blocks were defined by the sampling frame associated with the ad[21] and by treatment regime.[22] Table A3 in the SI summarizes the distribution of cases across blocks and treatment assignments.[23] Using randomization inference, we conduct a

---

[20]This is because real estate brokers, who comprise nearly 85% of the subjects in the experiment, are required to demonstrate knowledge of fair housing law to be licensed in New York, and because Craigslist shows advisory language stating that discrimination is illegal to landlords and brokers who are posting rental ads.

[21]These include: the Bronx, Brooklyn, Manhattan, Queens, Staten Island, and the LD oversample.

[22]We define a treatment regime as a distinct design and randomization procedure. There are three in this experiment. In Regime 1, the design had 5 arms (including 2 not examined in this paper), included the LD oversample, and employed equal assignment probabilities by sampling strata. In Regime 2, the design was reduced to 3 arms and the probability of assignment to control (to either punitive or monitoring conditions) was 50% (25%) by sampling stratum. In Regime 3, the LD oversample was discontinued and equal assignment probabilities by sampling strata were restored. Figure A2 in the SI shows the duration of each regime.

[23]Additional details on the randomization procedure are provided in the SI.

randomization check and infer that the randomization was valid.[24]

**Data and Measurement**

We combine the following data: (1) scraped data on advertised rental listings, (2) automated assumed biographical and treatment assignment data, (3) data on treatment delivery and compliance, and (4) detailed field reports testers recorded about their interactions with landlords *prior to the visit* (pre-treatment); *during a housing unit visit* (post-treatment); and *after the visit* (post-treatment).[25] We construct case-level measures of net discrimination, which is defined as the difference in favorable treatment between the majority and minority group tester in any matched pair.[26] We apply this measurement strategy to construct two objective measures of discrimination, which we pre-specified as our primary post-treatment outcome measures: differential treatment by race in landlord or broker efforts to call back and follow up with testers after the appointment; and to offer the unit to the tester. These measures are computed for each majority-minority pairing (white-Black, white-Hispanic, and Black-Hispanic) and can take three values at the case level: -1 if only the minority tester is treated favorably, 0 if both the minority and majority testers are treated equally, and 1 if only the majority tester is treated favorably.[27] When examining average levels of the net measure, 1 means 100% net discrimination against the minority group, -1 means 100% net discrimination against the majority group, and 0 means that the two groups are treated equally.[28]

We also construct pre-treatment measures of discrimination that occurs over the phone between the time testers initially contact the landlord when inquiring about the ad and the time of

---

[24]The probability of obtaining a log-likelihood statistic (from a multinomial logistic regression model regressing treatment assignment on 122 pre-treatment covariates) at least as large as the observed statistic is $p=0.15$. A balance table is presented in the SI.

[25]Details on data collection procedures are provided in the SI.

[26]We provide a detailed explanation of how the net measure of discrimination is constructed and why we prefer it to alternative measures of discrimination in the SI.

[27]If the landlord didn't honor the appointment, the net discrimination measure is coded 0 (i.e., both in any pair were treated the same as neither received a callback/an offer).

[28]A potential source of measurement error is that estimates of discrimination levels are driven by the particular composition of testers in each racial group. We address this concern by estimating a non-nested hierarchical model regressing a landlord-tester level indicator for favorable treatment on pre-treatment covariates as well as tester, tester, race, and tester team gender random effects. We estimate this model for each of the objective outcome indicators (receiving a callback and receiving an offer) among the control group and among the experimental sample. We find that the estimated variance of the varying tester intercepts is negligible and infer that discrimination levels are not driven by the particular composition of testers in each racial group. See the SI for details.

14

random assignment. We measure early-stage discrimination for all cases pursued by testers (not just those admitted into the experiment).[29] These measures are used to describe baseline levels of discrimination and to improve precision as covariates in the experimental analysis.

**Sample Definition**

The study was in the field from April 13, 2012 to December 20, 2013. We focus on two key samples for the analysis: the audit sample, which contains 2,711 cases, and the experimental sample, which contains 653 cases. Figure A1 in the SI presents a flow diagram summarizing the sample filtering procedure and Figure A2 in the SI summarizes the cumulative number of cases admitted into the experimental sample over this period and by treatment regime. The characteristics and geographic distribution of the rental housing stock in our samples appear to be broadly similar to the rental housing stock in the New York City rental housing market during this period.[30]

**Estimation and Inference**

We use both design-based and model-based approaches to estimation and inference, with both approaches producing near-identical results.[31] For design-based (non-parametric) inference, we estimate sample intent-to-treat (ITT) effects[32] of government messaging on discrimination levels as the weighted average across blocks of differences in net discrimination levels between treatment groups within blocks, where weights are equal to the inverse of the probability of assignment to the condition to which the case was actually assigned. For this non-parametric analysis, standard errors are calculated using the conservative weighted Neyman estimator.

---

[29]Early stage discrimination indicators include: differences in the ability of matched testers to make contact with the landlord; the ability to schedule an appointment; the number of biographical attributes landlords inquire about over the phone while screening testers; the number and percentage of attributes for which testers receive positive, negative, neutral, and skeptical responses from landlords after providing relevant information; and whether testers receive any negative or any skeptical feedback on any aspect of their biography communicated over the phone.

[30]See the SI for summary statistics of the characteristics of the housing stock in each sample.

[31]Additional analyses are described and reported in the Discussion and in the SI.

[32]We observe subject noncompliance with treatment assignment in a few cases because landlords assigned to a treatment arm hang up mid-message or because the city staffer administering the treatment calls could not reach the landlord prior to tester appointments. Because the rate of treatment compliance is relatively high (between 71-81%), our estimates of Complier Average Causal Effects (CACEs)—the effect of receiving governmental messages—are qualitatively similar to our ITT estimates. We therefore focus our presentation of main results on the ITT estimates. See the SI for CACE estimates.

In addition, we estimate effects using the following linear model:

$$Y_{ib} \;=\; \alpha_0 + \beta_1 T_{ib} + \gamma_b + u_{ib} \tag{1}$$

where $i$ indexes landlords and $b$ indexes experimental block, $Y$ is the net discrimination outcome measure, and $T$ is treatment assignment, variously defined when comparing mean outcomes between monitoring versus control, punitive versus control, or punitive versus monitoring. $\gamma$ is a full set of block fixed effects and $u$ is a disturbance term. We estimate Equation 1 on the subset of the data assigned to each pair of treatment arms being compared; $\beta_1$ is the effect of the treatment group relative to the comparison group. For this analysis we again use inverse propensity weights. We calculate $p$-values corresponding to a one-sided test of the null hypothesis of equality of means for the monitoring-control and punitive-control comparisons. We calculate $p$-values corresponding to a two-sided test of the null hypothesis of equality of means for the punitive-monitoring comparison and for all analyses involving net discrimination against Hispanic (vs. Black) testers because we have no strong priors about the expected direction of racial bias and treatment effects on racial bias when comparing the experiences of Black and Hispanic testers.

## RESULTS

Figure 1 summarizes the main findings and shows, separately for each outcome measure: (1) levels of favorable treatment for different racial groups (top left); (2) differences in favorable treatment rates between groups (i.e., net discrimination levels) by treatment assignment (lower left); (3) differences in favorable treatment rates across treatment conditions for the same group (top right); and (4) the effects of treatment assignment on net discrimination levels relative to the control or monitoring comparison group (lower right).[33] We focus our discussion on the lower-left and lower-right quadrants in particular, which present estimates of baseline net discrimination levels and of ITT effects of messaging on net discrimination levels, respectively. We also draw attention to our most important findings in the figure by highlighting key estimates.

---

[33]Estimates corresponding to the information presented in Figure 1 are shown in the SI.

## Baseline Discrimination Levels

First we assess baseline levels of discrimination in the outcome variables by examining the control group mean levels of net discrimination, which are defined as the control group mean difference in favorable treatment rates between groups, as shown the lower-left quadrants in Figure 1. Here we find statistically significant differences between Hispanic and white testers. We estimate that Hispanic testers were less likely than white testers to receive a callback from a landlord or broker—in 15.4% of cases compared to 21.5%, a difference of 6.1 percentage points ($p = 0.019$). They were also less likely to receive an offer for an apartment—in 6.1% of cases compared to 11.8%, a difference of 5.7 percentage points ($p = 0.011$). We find smaller, and statistically weaker, differences between Black and white testers. We estimate that Black testers were less likely than white testers to receive a callback from a landlord or broker—in 16.8% of cases compared to 21.5%, a difference of 4.7 percentage points ($p = 0.107$). Black testers were also less likely to receive an offer for the unit–in 9% of cases compared to 11.8%, a difference of nearly 2.9 percentage points ($p = 0.239$).

We can compare baseline levels of discrimination in callbacks to estimates of differential treatment in receiving a follow-up from the 2012 HUD audit of ethnic and racial discrimination in rental markets.[34] Our results run counter to the New York metropolitan area estimates from the 2012 HUD audit, which found that Hispanic testers received a follow-up from agents more frequently than white testers (in 6.9% and 5.4% of cases, respectively, a difference of $-1.5$ percentage points, $p = 0.804$) and that Black testers received a follow-up from agents more frequently than white testers (in 10.2% and 3.9% of cases, respectively, or a difference of $-6.3$ points, $p = 0.096$) (Turner et al. 2013, pp. 157-158). These differences in apparent baseline levels of discrimination might be due to two factors. First, the two estimates may be comparing discrimination baseline levels for two different populations in the New York area rental market. The HUD estimates for the New York City area correspond to the New York/Northeast New Jersey metropolitan area, and not to New York City specifically. Our estimates are specific to New York City, and in particular

---

[34]No comparable measure of discrimination in offers exists in the HUD audit study.

to landlords and agents who post ads on Craigslist and who schedule appointments with all three testers pursuing the listing. Second, there is less variability in our estimates than in the New York area estimates from the HUD audit since the sample size of our control group is nearly twice the size of the entire New York HUD audit. When comparing our estimates to national estimates of discrimination in follow-ups from the 2012 HUD study, the estimates are similar in that Black and Hispanic testers were less likely to receive follow-ups from agents than white testers: 10.5% for Black testers versus 11.0% for white testers, a difference of 0.6 percentage points that is not significant at the 0.05 level; and 7.9% for Hispanic testers versus 11.2% for white testers, a difference of 3.3 percentage points that is significant at the 0.1 level (Turner et al. 2013, pp. 44, 47).

**Does Government Messaging Reduce Racial Discrimination?**

Next, we report estimates of the effect of messaging on net discrimination levels, as shown in the lower-right quadrant of Figure 1. This section of the figure presents weighted nonparametric estimates (open markers) and regression estimates adjusted using block fixed effects and inverse probability weighting (filled markers; our preferred specification) with 95% confidence intervals.[35]

When compared to a pure control condition, sending a monitoring message decreases net discrimination against both Black and Hispanic testers (versus white testers) across the objective outcome measures, although the magnitudes of the average effects are substantively negligible and not statistically distinguishable from zero. Sending a monitoring signal decreases net discrimination against Black testers in receiving a callback ($-0.002; p = 0.486$) and in receiving an offer for the unit ($-0.003; p = 0.464$). For Hispanic testers, sending the monitoring message decreases net discrimination in receiving a callback ($-0.036; p = 0.201$) and receiving an offer for the unit ($-0.017; p = 0.31$).

The results on punitive messaging point generally to reductions in discrimination against Hispanics but increases in discrimination against Blacks. Sending a punitive message increases net discrimination against Black testers (versus white testers) across the outcome measures when com-

---

[35]We focus on the latter in our exposition of results, because the two estimators produce qualitatively similar estimates.

pared to control, but none of these effects are statistically significant at the 5% level. Sending the punitive message reduces the likelihood of receiving a post-visit follow-up callback from the landlord (0.02; $p = 0.675$) and reduces the likelihood of receiving a post-visit offer (0.016; $p = 0.679$). Punitive messaging instead decreases net discrimination against Hispanic testers in receiving a post-visit callback ($-0.066$; $p = 0.056$) and in receiving a post-visit offer for the unit ($-0.021$; $p = 0.268$). The effect of punitive messaging on net discrimination in receiving a post-visit callback is just shy of significance at 5% ($p = 0.056$). Notably, sending a punitive message worsens outcomes for Black compared to Hispanic testers across the outcome measures. Punitive messaging decreases net discrimination against Hispanic testers relative to Black testers in receiving callbacks (-0.085; $p = 0.037$), and in receiving offers for units (-0.039; $p = 0.242$). Note that of these two estimated treatment effects, only the effect on callbacks is statistically significant at the 5% level ($p = 0.037$).

We next examine the comparative effectiveness of punitive versus monitoring messaging. It is clear from Figure 1 that there is no statistically distinguishable difference in the effectiveness of the two treatments. We find mixed results when comparing outcomes between the two treatment groups—Black testers do slightly worse when interacting with landlords in the punitive condition as compared to the monitoring condition, while Hispanic testers do slightly better. All $p$-values for these objective outcomes range from 0.268 to 0.962.

**Hidden Discrimination: Estimates using Subjective Measures**

Our primary analysis focuses on the causal effects of government messaging on two objective measures of racial discrimination in the housing market. We believe that objective measures are important in part because discrimination in the housing market can be relatively invisible.[36] In addition to our objective measures of discrimination, we also have access to a set of subjective measures of discrimination that capture, among other aspects of the housing search process, sub-

---

[36]As described by the Fair Housing Justice Center (2013), in this sector, "discrimination takes place with a handshake and a smile." One expert in the area has argued that the "image of housing discrimination as a slammed door must be replaced with an invisible revolving door where people are courteously escorted in, out of, and ultimately away from the desired housing" (Freiberg 2015).

jects' perceptions of steering. As such, they reflect both the behavior of the landlords—which may be a function of treatment conditions—and the dispositions of confederates to notice and report indicators of discrimination—which is not subject to treatment conditions but may be associated with racial group membership. The measures, while rich, are also flawed in that they suffer from substantial missingness (on the order of 25-30%).[37] While this missingness is not correlated with racial group or with treatment condition, it is still plausibly related to potential outcomes. That is, it is possible that individuals were less likely to report if they felt little discrimination, or perhaps, if they felt substantial discrimination.

With this caveat in mind, in Figure A3 in the SI, we report messaging effects on an index measure of testers' subjective perceptions of favorable treatment and a net discrimination measure using the subjective index.[38] Although on the objective measures there is strong evidence for discrimination against Hispanics, this does not appear at all in the subjective measures. Moreover, while there is some evidence for discrimination against Blacks from the objective measures, for the subjective measure, Blacks report if anything better treatment than whites. Our results are thus consistent with the dominant viewpoint in the field, but they also have implications for the generality of our findings. The difficulty of identifying discrimination in individual cases or based on perceptions of treatment may be one factor that makes the punitive threat of law relatively weak, even if there are normative benefits. Different effects may arise in areas in which identification of violations is easier.

## DISCUSSION

Research on the political economy of race in the United States has shown that racial discrimination in market behavior plays a central role in perpetuating racial inequalities, but has long overlooked questions about whether governmental efforts to reduce discrimination—in particular strategies to implement and enforce anti-discrimination laws since the 1960s—are effective at reducing the

---

[37]We refer the reader to the SI for an analysis of missingness on subjective measures.

[38]We use an index measure to avoid a multiple comparisons problem. We conduct and report this analysis to maintain fidelity to our pre-analysis plan and to situate our main results in a broader policy context. Details about the data collection and measure construction procedures related to this analysis are presented in the SI.

incidence of discriminatory market behavior. To address this need in the literature, we conduct the first experimental test of two common governmental strategies to pre-emptively induce behavioral compliance with the law in the domain of fair housing law in the United States. Partnering with the City of New York, we assess whether government campaigns that invoke fair housing law and make salient the costs of violating the law reduce racial discrimination levels in the New York City rental market when these appeals are delivered as targeted live phone calls from the city. The design offers a strong test because it (1) individually targets government appeals at potential discriminators, (2) measures outcomes soon after the treatment is delivered (thereby avoiding concerns about whether observed effects might have decayed), and (3) occurs in a policy context where government signals to monitor landlords and enforce the law are credible.

First, in the absence of any government intervention, we find strong evidence of discrimination in the New York City rental market on outcomes that are consequential for racial housing inequality. In the control group, Hispanics are less likely than Blacks, who are themselves less likely than whites, to receive favorable treatment from landlords and brokers. In terms of the objective outcome measures, we find strong evidence of discrimination in the New York City rental market, particularly against Hispanics. They are 6.1 percentage points less likely to receive a callback and 5.7 percentage points less likely to receive an offer for an apartment than whites. These findings stand in stark contrast to estimates of discrimination levels from a comparable audit study of the New York metropolitan area rental market in 2012 which reports no discrimination against Blacks or Latinos in receiving a callback but relies on a small sample for which sampling variability is a concern. We also find suggestive, but not strong, evidence that punitive government messaging campaigns that both invoke the law and make salient the costs of violating the law can reduce net discrimination in receiving a callback from a landlord following an appointment, but only for Hispanics. We find no statistically distinguishable difference between the effects of punitive and monitoring messaging relative to control.[39] We assess the robustness of our main results (which are based on our pre-analysis plan) to alternative estimation strategies that, at least in principle,

---

[39]Employing a conservative Bonferroni approach to multiple comparisons across treatments, outcomes, and groups, we are unable to reject the null hypothesis that these treatments make no difference.

improve efficiency; our main results are not materially affected by the choice of estimator.[40]  In addition, our finding that the effectiveness of government messaging strategies varies by the race of the target of discrimination is also an important contribution to the literature. We speculate on potential explanations for these mixed findings and related implications for future research below.

Despite the lack of statistical significance, the magnitudes of the estimated messaging effects we observe in this study appear to be substantively significant both in relative terms and in terms of the mean level of discrimination under treatment. In relative terms, the effects are very large, albeit from a small baseline.[41]  Punitive messaging reduces discrimination against Hispanics relative to whites by 109% and relative to Blacks by 126%. Not only is there a large relative decrease in net discrimination against Hispanics, the sign of the predicted group mean changes from positive to negative suggesting more favorable treatment for Hispanics relative to whites and Blacks under the punitive messaging condition. Alternatively we can adduce the substantive meaning of effects by interpreting the mean discrimination level under treatment.  Our estimates of discrimination levels on key measures drop to zero under the punitive condition. This creates an inferential puzzle for policymakers as well as a need for further experimental replication, which we discuss below as well.  In the remainder of the article, we address concerns regarding the internal and external validity of our findings, implications for policymakers, and directions for future research.

**Addressing Internal Validity Concerns about Spillovers**

The main threat to internal validity arises if the non-interference assumption is violated, which could occur if landlords assigned to a messaging condition communicate with other landlords

---

[40]See the SI for a full description of the alternative analyses and results. One strategy pools observations across all arms and estimates models using binary indicators for each treatment condition, block fixed effects, and inverse probability weights.  This allows for possible gains in efficiency from tighter estimation of block effects.  We see however that estimates and inferences are materially unaffected by the choice of using a two- or three-group parametric estimator. Second, the introduction of additional covariates can in principle improve efficiency, though this introduces obvious risks of specification selection. To introduce covariates in a principled way, we employ a lasso regression as a machine learning approach to select covariates that are prognostic of outcomes while avoiding biases arising from ex post covariate selection (e.g., Bloniarz et al. 2016) and semi-parametrically estimate covariate adjusted ITT effects following Yuan, Zhang and Davidian (2012). We find again that the covariate adjusted ITT estimates and the substantive results are not materially different from our main results.

[41]Interpreting treatment effects as relative differences offers an alternative substantive interpretation and is customary in experimental analyses. We present the main ITT estimates in terms of percent differences in the SI.

(who are assigned to control or to the other messaging condition) about the treatment message they received. We argue that such spillovers are unlikely to occur. Given the context that is the New York City rental housing market and the sampling procedures used, we argue that there is a very low probability that subjects in the experiment interact with each other. To adduce this, we first estimate bounds on the probability that a landlord or broker (who posts rental ads on Craigslist) enters the audit and experiment samples and show that these probabilities are low.[42] We estimate that the probability a landlord or broker enters the audit sample is between 3.2% and 15%, and that the probability a landlord or broker enters the experiment sample is between 0.77% and 3.6%. For the experimental sample especially, we think this makes a strong prima facie case for minimal interference. Second, we characterize the experiment sample as a very small random sample of landlords and brokers in the New York City rental market. The actual number of active landlords and brokers in the New York City rental market is unknown. We therefore estimate the denominator to adduce a very conservative upper bound on this quantity. We infer that the experiment sample must be far less than 2% of the estimated population of landlords and brokers in the New York City rental market which, when interpreted as a very small random sample of the population, suggests that interactions among subjects in the experiment are highly unlikely.[43]

**Addressing External Validity Concerns**

Next, we address concerns related to external validity deriving from the ways study samples are constructed. Although care was taken to sample ads randomly from the universe of Craigslist rental listings, there are several scope conditions that limit the generalizability of our findings. First, our results are limited to ads placed via public listings on Craigslist, which is a subset of advertised market-rate rental units, and we may miss out on discrimination that may arise when prospective tenants are identified through other channels (especially private channels) and discrimination in the subsidized rental market. This would likely lead to underestimates of discrimination levels under the assumption that landlords who post rental ads on Craigslist are likely to be less discriminatory

---

[42]Details on the method and calculations are provided in the SI.
[43]See the SI for details on the method used to estimate this quantity and alternative approaches explored.

than those who select into non-public modes of advertising. This expectation is consistent with economic and sociological research on the role of segregated informal social networks that shape how people search for and get connected to employment and housing opportunities (Loury 2001; Pager and Shepherd 2008).

Second, our findings are also limited to landlords and brokers who require housing seekers to reply to rental ads by phone. We use a capture-recapturing sampling procedure and estimate that 55.8% of New York City rental listings on Craigslist require contact to be done by phone.[44] We acknowledge that this is a scope condition on our findings and that future research should explicitly assess whether findings differ by the mode of initial contact.

Third, a possibly relevant factor is that Craigslist employs its own anti-discrimination messaging, which is shown to users before posting and acknowledges federal, state, and city prohibitions on discrimination.[45] If this messaging deters would-be discriminators from advertising on Craigslist, then this would bias our estimates of baseline discrimination levels downward. If it simply changes the ways in which landlords and brokers advertise, then this makes the evidence for discrimination all the more striking as the discrimination we uncover cannot be readily detected from simple text analysis.[46] It is also possible that Craiglist's messaging, by having effects of its own on discrimination and reducing the baseline level of discrimination, reduces the magnitude of the effect of government messaging. If this is the case, then our results should be interpreted as evidence of the effects of government messaging *in a context in which non-governmental messaging also exists*.

Fourth, our experiment applied treatment only after appointments were scheduled (in order to avoid differential attrition). This might have the effect of limiting the generalizability of our analysis to encounters in which early-stage discrimination was absent. We can address this concern by assessing discrimination toward testers occurring over the phone *prior* to randomization—that

---

[44]Details on the method and calculations are provided in the SI.

[45]The message states, in part: "It is illegal to discriminate in the sale, rental or leasing of housing because of a person's race, color, creed, national origin, sexual orientation, marital status, familial status, or religion."

[46]However, other work suggests that some discriminatory text survives. Mirroring our own attempts to target "likely discrimination" in postings, Oliveri (2010) found that most discriminatory language is used by people seeking roommates.

is, across the full audit sample. Table A1 in the SI summarizes these findings and shows an overall lack of difference in early-stage discrimination. We find no statistically significant differences in the likelihood of making any contact when first replying to an advertisement by phone. And differences in success in setting up appointments were very small: the only significant difference from three comparisons is between white and Black testers, with Black testers somewhat *more* successful at scheduling appointments (36.1% versus 34.8%; $p = 0.035$). These results provide greater confidence that our results are not driven by sample selection resulting from all three testers successfully scheduling appointments.

Finally, given the potential rarity of encountering both a set of auditors and of receiving a targeted phone call from the government, there may be concerns that landlords are "spooked" by the intervention. Being spooked by the matched auditors is, by itself, does not threaten the internal validity of the experimental results because the audit-based measurement strategy is consistent across all treatment arms in the experiment.[47] Being spooked by the targeted call from the government is not a concern because it is part of the net effect of the bundled treatment that we are interested in understanding.[48]

**Policy Implications**

What are the policy implications of our findings? The results on levels of discrimination are relatively clear for policy makers: discrimination is a greater concern in New York than previously believed, and this may be true elsewhere also. The evidence is especially strong for Hispanics. The implications for interventions are less clear but just as important. The challenge here is that there is suggestive evidence of strong effects but many of the estimated effects do not hit conventional levels of significance, despite the large scale of this experiment. Power issues loom large when the behavior of interest—discrimination—is relatively rare, and is itself estimated with uncertainty. In such cases relying on conventional standards of significance may be overly conservative—an issue

---

[47]This could mean our estimates of discrimination levels in the control group are downwardly biased and could be interpreted as conservative estimates.

[48]There may also be an interaction effect between the two (encountering auditors and getting a phone call), but we remain agnostic about that effect. Understanding that interaction effect and understanding whether alternative messaging modes yield different results are beyond the scope of this study and should be tested in future research.

that has been of concern also in the study of rare diseases, leading to calls to alter standards in such cases (Gobburu and Pastoor 2016). Here, rather than altering standards, we calculate appropriate beliefs that the intervention is effective using a Bayesian framework by estimating the probability of a hypothesis (that each messaging strategy is effective) given the data. This is a fundamentally distinct quantity than the probability of the data given a hypothesis, which is the quantity of interest for determining statistical significance (Gill 1999). Understanding this alternative question is plausibly more relevant to policymakers who care about the probability that a given intervention works, rather than the probability of rejecting a specified null hypothesis. Importantly, we clarify that this is not a method to fish for statistically significant results, but is instead a method to answer a substantively distinct question that can be used following future replications (given a cumulation of experimental data) to improve how policymakers form posterior beliefs about the effectiveness of different interventions.

[FIGURE 2 HERE]

Figure 2 plots the posterior densities along with the estimated probabilities of effectiveness.[49] In the bottom row, we see that for each distribution, approximately half or less of the probability mass lies below zero: The data are not highly informative for posterior beliefs that either monitoring or punitive messages reduce discrimination against Blacks. In the case of punitive messages on callbacks, roughly three-quarters of probability mass lies *above* zero, providing some indication that the treatment may have even increased discrimination against Black testers. The top row shows that the story is different for Hispanics. For callbacks especially, the bulk of probability mass—at least 89%—lies beneath zero, implying posterior beliefs that the messages reduced discrimination against that group. We can be most confident in posterior beliefs about the effectiveness of the punitive treatment in reducing discrimination against Hispanics in callbacks: More than 99% of the mass lies below zero, centered around the posterior mean reduction in discrimination of 6.6 percentage points. For reference, that estimate is roughly twice the corresponding posterior mean reduction due to the monitoring treatment.

---

[49]See the SI for further details on the method.

26

**Directions for Future Research**

Lastly, we outline several directions for future research. First, an important limitation of the present field experiment is that despite the unusually large scale of our study (in comparison to prior in-person rental market audits), there is not enough statistical power in the design to detect treatment effects on the order of about 6 points. This is in part because the control group mean discrimination level is relatively low and because the estimation of treatment effects requires estimating an interaction between tester race and treatment (since discrimination is not directly observable at the case-tester level). As this is the first field experiment of its kind, it was not possible to conduct ex ante power analyses using previously reported effect sizes as a benchmark. Instead, the study was designed with the reasonable ex ante expectation that it was adequately powered to detect and reduce racial discrimination levels reported in rental market audit studies published prior to the start of the study.[50] Accordingly, additional experimental replications that use larger sample sizes are needed to more precisely estimate both discrimination levels and messaging effects.

Second, the finding of differential impacts for Hispanic and Black housing seekers is surprising given the lack of an expectation of heterogeneous effects in the literature. Existing work on the politics of policy enforcement and compliance does not explicitly develop explanations for why attempts to induce compliance with anti-discrimination law might vary by the group membership of the target of discrimination. We argue that this result points to a need for future research to develop and test hypotheses about the conditions under which we might expect variation in baseline discrimination levels and heterogeneous effects of governmental anti-discrimination appeals by the group membership of the targets of discrimination and by the group membership of potential discriminators.[51] We briefly speculate on several possible interpretations for the mixed findings in

---

[50]Hanson and Hawley (2011) report differences in response rates to initial inquiries of about 6 points between white and Black testers in the 2009 New York City rental market. Studying racial discrimination in the Los Angeles County rental market in 2003, Carpusor and Loges (2006) report differences in response rates to initial inquiries of between 16 and 44 percentage points, depending on the listed monthly rent of the unit, between white and Black testers. In a more recent study published after the completion of our field experiment, Ewens, Tomlin and Wang (2014) report a 9 percentage point difference in response rates to initial inquiries between white and Black testers.

[51]To begin to address this question, we conduct an exploratory analysis of heterogeneous messaging effects by the perceived race of landlords using the data from this experiment. Details and results are shown in the SI.

order to motivate directions for future research. One is that landlords may associate selected groups as the actual targets of housing discrimination and adjust their behavior with respect to only those groups when sent a punitive message. In this case it may be the case that landlords are actively discriminating against Hispanics at baseline, are acutely aware of this, and only adjust their behavior toward Hispanics in the punitive condition.[52] A second possible interpretation is that anti-Black prejudice is more entrenched than anti-Hispanic prejudice and stronger prejudices may be more difficult to change (Broockman and Kalla 2016). A third possibility—consistent with our evidence on discrimination levels—is that landlords associate Blacks as the main intended beneficiaries of contemporary campaigns for fair housing, given their roots in the civil rights movement, and already take some conscious steps to avoid treating them unfavorably. The punitive messages, then, might have made landlords' behavior toward other groups comparatively more salient. This would suggest that existing norms of fairness have not caught up to the demographic realities of New York City, where according to the 2010 Census more Hispanics live (27.5%) than Blacks (25.1%). We offer other possible explanations for the mixed findings in the SI due to space constraints.

Third, the treatments we tested in this study were, given practical considerations and constraints, bundled. They also did not test the full range of potential strategies government could deploy in pre-emptive messaging campaigns to reduce discrimination. Future studies could more finely operationalize treatments to test specific mechanisms (e.g., factors that affect specific parameters in our aforementioned decision-theoretic model) that build on existing theoretical frameworks from literatures on prejudice reduction (Paluck and Green 2009) and on behavioral policy compliance (Weaver 2014, 2015).[53] Future studies could also test whether effects vary by the mode of treatment delivery or by the messenger of the appeal (governmental versus non-governmental as

---

[52]Our inability to detect statistically significant discrimination levels against Blacks may be consistent with this explanation, but additional replication using an adequately powered sample is necessary to infer whether that is in fact the case.

[53]For example, future studies could randomly vary, among other features, whether the appeals provide new information about the law, in particular what exactly is illegal, which could change behavior subjects did not already know that information; whether the appeals remind landlords about the law among subjects who know about law in the first place; whether the appeals vary signals of the probability of monitoring and enforcement to affect landlords' perceived probability of detection and perceived probability of punishment. Future studies could also attempt to distinguish between responses that are driven by fear of punishment versus appeals to social norms by manipulating the informational and normative content of appeals.

well as variations in source credibility).

Fourth, future studies should be designed to validate why subjects do not comply with the law and to experimentally test theories about how different interventions could increase compliance with the law among different types of noncompliers. It may be the case that interventions that directly address reasons for noncompliance may be effective at increasing compliance. However, it may also be the case that the interventions that increase compliance may be effective for reasons orthogonal to the subjects' motivations for defying the law in the first place.

Finally, replication in other housing markets and political jurisdictions would be valuable to understand how messaging effects vary across political and economic contexts. There may be different baseline discrimination levels and reasons why market actors discriminate that moderate their behavioral response to governmental appeals. There may also be different expectations among market actors about the credibility of government efforts to enforce anti-discrimination laws. How citizens perceive the capacity of the state to enforce compliance (Weaver 2014) as well as the degree of forbearance in government in the domain of civil rights (Holland 2016)—in particular the preferences for discrimination among those in government and the types of legal, organizational, and partisan incentives and constraints that motivate bureaucrats to curb or allow discrimination (White, Nathan and Faller 2015; Einstein and Glick 2016; Mummolo 2017)—may affect how they perceive the credibility of government signals to monitor market behavior and enforce the law, the expected intensity of government efforts to do so, and ultimately the costs of violating the law.

## REFERENCES

Aigner, Dennis J and Glen G Cain. 1977. "Statistical theories of discrimination in labor markets." *Industrial and Labor relations review* pp. 175–187.

Alexander, Michelle. 2012. *The New Jim Crow: Mass Incarceration in the Age of Color-blindness*. New York: New Press.

Arrow, Kenneth J. 1973. The Theory of Discrimination. In *Discrimination in Labor Markets*, ed. Orley Ashenfelter and Albert Rees. Princeton: Princeton University Press pp. 3–33.

Becker, Gary S. 1957. *The Economics of Discrimination*. University of Chicago Press.

Bertrand, Marianne, Dolly Chugh and Sendhil Mullainathan. 2005. "Implicit Discrimination." *American Economic Review* 95(2):94–98.

Bertrand, Marianne and Esther Duflo. 2017. Field Experiments on Discrimination. In *Handbook of Field Experiments*, ed. Esther Duflo and Abhijit Banerjee. Vol. 1 Amsterdam: Elsevier.

Bloniarz, Adam, Hanzhong Liu, Cun-Hui Zhang, Jasjeet Sekhon and Bin Yu. 2016. "Lasso adjustments of treatment effect estimates in randomized experiments." *Proceedings of the National Academy of Sciences of the United States of America* 113(27):7383–7390.

Blumenthal, Marsha, Charles Christian, Joel Slemrod and Matthew G. Smith. 2001. "Do Normative Appeals Affect Tax Compliance? Evidence from Controlled Experiment in Minnesota." *National Tax Journal* 54(1):125–138.

Bobo, Lawrence, James R. Kluegel and Ryan A. Smith. 1997. Laissez-Faire Racism: The Crystallization of a Kinder, Gentler, Antiblack Ideology. In *Racial Attitudes in the 1990s: Continuity and Change*, ed. Steven A. Tuch and Jack K. Martin. Westport, CT: Praeger pp. 15–42.

Bowles, Samuel. 2016. *The Moral Economy: Why Good Incentives are No Substitute for Good Citizens*. New Haven: Yale University Press.

Broockman, David and Joshua Kalla. 2016. "Durably reducing transphobia: A field experiment on door-to-door canvassing." *Science* 352(6282):220–224.

Bumiller, Kristen. 1988. *The Civil Rights Society: The Social Construction of Victims*. Baltimore: Johns Hopkins University Press.

Butler, Daniel M. and Jonathan Homola. Forthcoming. "An Empirical Justification for the Use of Racially Distinctive Names to Signal Race in Experiments." *Political Analysis* .

Cardenas, Juan Camilo, John Stranlund and Cleve Willis. 2000. "Local Environmental Control and Institutional Crowding-Out." *World Development* 28(10):1719–1733.

Carpusor, Adrian G. and William E. Loges. 2006. "Rental Discrimination and Ethnicity in Names." *Journal of Applied Social Psychology* 36(4):934–952.

Cover, Robert. 1995. The Origins of Judicial Activism in the Protection of Minorities. In *Narrative, Violence, and the Law*, ed. Martha Minow, Michael Ryan and Austin Sarat. University of Michigan Press.

Dawson, Michael C. 2016. "Hidden in Plain Sight: A Note on Legitimation Crises and the Racial Order." *Critical Historical Studies* 3(1):143–161.

Dawson, Michael C. and Megan Ming Francis. 2015. "Black Politics and the Neoliberal Racial Order." *Public Culture* 28(1):23–62.

Dixit, Avinash K. 2006. *Lawlessness and Economics: Alternative Modes of Governance*. New York: Oxford University Press.

Dixon, K. A., Duke Storen and Carl E. Van Horn. 2002. *A Workplace Divided: How Americans View Discrimination and Race on the Job*. New Brunswick, NJ: Rutgers University, State University of New York, and John J Heldrich Center for Workplace Development.

Donohue, John J. and Peter Siegelman. 1991. "The Changing Nature of Employment Discrimination Litigation." *Stanford Law Review* 43(983-1033).

Dunning, Thad, Felipe Monestier, Rafeal Pinero, Fernando Rosenblatt and Guadalupe Tunon. 2015. "Positive vs. Negative Incentives for Compliance: Evaluating a Randomized Tax Holiday in Uruguay." Working Paper. Available at `http://dx.doi.org/10.2139/ssrn.2650105`.

Einstein, Katherine Levine and David M. Glick. 2016. "Does Race Affect Access to Government Services? An Experiment Exploring Street-Level Bureaucrats and Access to Public Housing." *American Journal of Political Science* 61(1):100–116.

Epp, Charles. 1998. *The Rights Revolution: Lawyers, Activists, and Supreme Courts in Comparative Perspective*. Chicago: University of Chicago Press.

Ewens, Michael, Bryan Tomlin and Liang Choon Wang. 2014. "Statistical Discrimination or Prejudice? A Large Sample Field Experiment." *The Review of Economics and Statistics* 96(1):119–134.

Eyer, Katie R. 2012. "That's Not Discrimination: American Beliefs and the Limits of Anti-Discrimination Law." *Minnesota Law Review* 96:1275–1361.

Fair Housing Justice Center. 2013. "A Matter of Place." Available at `http://www.fairhousingjustice.org/resources/film/`.

Feagin, Joseph R. 1991. "The Continuing Significance of Race: Antiblack Discrimination in Public Places." *American Sociological Review* 56(1):101–116.

Forman, James Jr. 2012. "Racial Critiques of Mass Incarceration: Beyond the New Jim Crow." *New York University Law Review* 87(1):101–146.

Foster, Mindi D. and Kenneth L. Dion. 2004. "The Role of Hardiness in Modering the Relationship between Global/Specific Attributions and Actions Against Discrimination." *Sex Roles* 51(3):161–169.

Freiberg, Fred. 2015. "Housing Discrimination Doesn't Begin or End with 'Poor Doors'." *City Limits*. Available at `http://citylimits.org/2015/09/08/housing-discrimination-doesnt-begin-or-end-with-poor-doors`.

Fryer, Roland G and Steven D Levitt. 2004. "The causes and consequences of distinctively black names." *The Quarterly Journal of Economics* 119(3):767–805.

Frymer, Paul. 2003. "Acting When Elected Officials Won't: Federal Courts and Civil Rights Enforcement in U.S. Labor Unions, 1935-1985." *American Political Science Review* 97(3):483–499.

Galanter, Marc. 1974. "Why the Haves Come Out Ahead: Speculations on the Limits of Legal Change." *Law and Society Review* 9(1):95–160.

Gill, Jeff. 1999. "The Insignificance of Null Hypothesis Significance Testing." *Political Research Quarterly* 52(3):647–674.

Gneezy, Uri and Aldo Rustichini. 2000. "A Fine is a Price." *Journal of Legal Studies* 29(1):1–18.

Gobburu, J and D Pastoor. 2016. "Drugs Against Rare Diseases: Are The Regulatory Standards Higher?" *Clinical Pharmacology & Therapeutics* 100(4):322–323.

Gurian, Craig. 2005. "A Return to Eyes on the Prize: Litigating Under the Restored New York City Human Rights Law." *Fordham Urban Law Journal* 33(2).

Hadfield, Gillian K. and Barry R. Weingast. 2014. "Microfoundations of the Rule of Law." *Annual Review of Political Science* 17:21–42.

Hanson, Andrew and Zackary Hawley. 2011. "Do Landlords Discriminate in the Rental Housing Market? Evidence from an Internet Field Experiment in U.S. Cities." *Journal of Urban Economics* 70(2-3):99–114.

Haynes, Laura C., Donald P. Green, Rory Gallagher, Peter John and David J. Torgerson. 2013. "Collection of Delinquent Fines: An Adaptive Randomized Trial to Assess the Effectiveness of Alternative Text Messages." *Journal of Policy Analysis and Management* 32(4):718–730.

Heckman, James J. 1998. "Detecting discrimination." *The Journal of Economic Perspectives* pp. 101–116.

Heckman, James J and Peter Siegelman. 1993. The Urban Institute Audit Studies: Their Methods and Findings. In *Clear and convincing Evidence: Measurement of Discrimination in America*, ed. M Fix and R Struyk. Urban Institute Press.

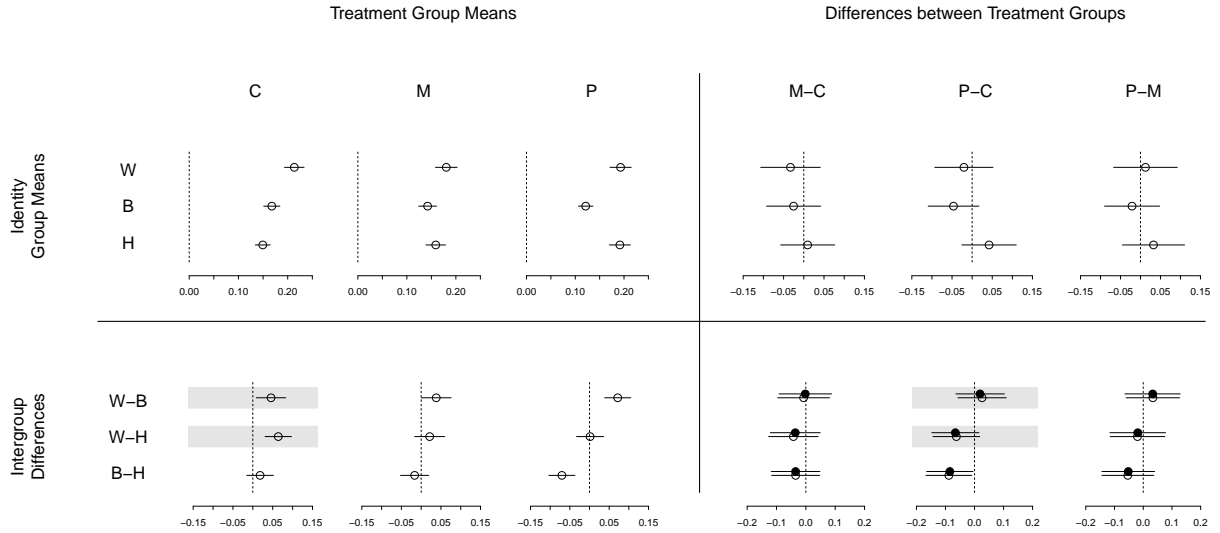Holland, Alisha C. 2016. "Forbearance." *American Political Science Review* 110(2).

Iyer, Govind S., Philip M. J. Reckers and Debra L. Sanders. 2010. "Increasing Tax Compliance in Washington State: A Field Experiment." *National Tax Journal* 63(1):7–32.

Kappen, Diane M. and Nyla R. Branscombe. 2001. "The effects of reasons given for ineligibility on perceived gender discrimination and feelings of injustice." *British Journal of Social Psychology* 40(2):295–313.

King, Desmond S. and Rogers M. Smith. 2005. "Racial Orders in American Political Development." *American Political Science Review* 99(1):75–92.

Loury, Glenn C. 2001. Politics, race, and poverty research. In *Understanding Poverty*, ed. Sheldon H. Danziger and Robert H. Haveman. Cambridge: Harvard University Press.

Major, Brenda and Tessa L. Dover. 2016. Attributions to Discrimination: Antecedents and Consequences. In *Handbook of Prejudice, Stereotyping, and Discrimination*, ed. Todd D. Nelson. 2nd ed. New York: Taylor and Francis.

Massey, Douglas S and Nancy A Denton. 1993. *American Apartheid: Segregation and the Making of the Underclass*. Cambridge: Harvard University Press.

McAdams, Richard H and Janice Nadler. 2005. "Testing the Focal Point Theory of Legal Compliance: The Effect of Third-Party Expression in an Experimental Hawk/Dove Game." *Journal of Empirical Legal Studies* 2(1):87–123.

McCubbins, Mathew D. and Thomas Schwartz. 1984. "Congressional oversight overlooked: Police patrols versus fire alarms." *American Journal of Political Science* pp. 165–179.

Mendelberg, Tali. 2001. *The race card: Campaign strategy, implicit messages, and the norm of equality*. Princeton University Press.

Mummolo, Jonathan. 2017. "Modern Police Tactics, Police-Citizen Interactions and the Prospects for Reform." Working Paper. Stanford University.

National Fair Housing Alliance. 2014. "Fair Housing Trends Report 2014. Expanding Opportunity: Systemic Approaches to Fair Housing." Available at `http://www.nationalfairhousing.org/LinkClick.aspx?fileticket=MqO6AE6loGY%3D&tabid=3917&mid=5321`.

Neumark, David. 2011. "Detecting Discrimination in Audit and Correspondence Studies." *Journal of Human Resources* 47(4):1128–1157.

Nielsen, Laura Beth. 2004. *License to Harass: Law, Hierarchy, and Offensive Public Speech*. Princeton: Princeton University Press.

Nielsen, Laura Beth, Robert L. Nelson and Ryon Lancaster. 2010. "Individual Justice or Collective Legal Mobilization? Employment Discrimination Litigation in the Post Civil Rights United States." *Journal of Empirical Legal Studies* 7(2):175–201.

Oliveri, Rigel Christine. 2010. "Discriminatory Housing Advertisements On-Line: Lessons from Craigslist." *Indiana Law Review* 43:1125.

Pager, Devah and Diana Karafin. 2009. "Bayesian Bigot? Statistical Discrimination, Stereotypes, and Employer Decision Making." *Annals of the American Academy of Political and Social Science* 621:70–93.

Pager, Devah and Hana Shepherd. 2008. "The Sociology of Discrimination: Racial Discrimination in Employment, Housing, Credit, and Consumer Markets." *Annual Review of Sociology* 34:181–209.

Paluck, Elizabeth Levy and Donald P Green. 2009. "Prejudice reduction: What works? A review and assessment of research and practice." *Annual review of psychology* 60:339–367.

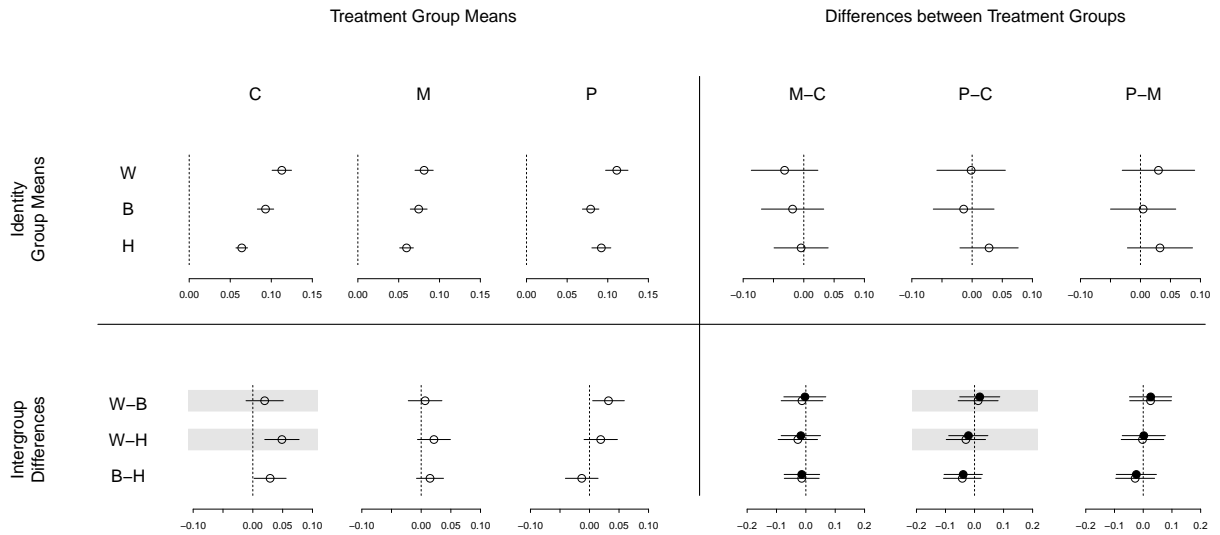Phelps, Edmund S. 1972. "The Statistical Theory of Racism and Sexism." *American Economic Review* 62(4):659–661.

Ritter, Gretchen. 2007. *The Constitution as Social Design: Gender and Civic Membership in the American Constitutional Order*. Stanford: Stanford University Press.

Rosenberg, Gerald N. 1991. *The Hollow Hope: Can Courts Bring About Social Change?* Chicago: University of Chicago Press.

Roychoudhury, Canopy and Allen C Goodman. 1992. "An Ordered Probit Model for Estimating Racial Discrimination through Fair Housing Audits." *Journal of Housing Economics* 2:358–373.

Roychoudhury, Canopy and Allen C Goodman. 1996. "Evidence of Racial Discrimination in Different Dimensions of Owner-Occupied Housing Search." *Real Estate Economics* 24:161–178.

Skrentny, John D. 2002. *The Minority Rights Revolution*. Cambridge: Belknap Press of Harvard University Press.

Slemrod, Joel, Marsha Blumenthal and Charles Christian. 2001. "Taxpayer Response to an Increased Probability of Audit: Evidence from a Controlled Experiment in Minnesota." *Journal of Public Economics* 79:455–483.

Swim, Janet K., Elizabeth D. Scott, Gretchen B. Sechrist, Bernadette Campbell and Charles Stangor. 2003. "The role of intent and harm in judgments of prejudice and discrimination." *Journal of Personality and Social Psychology* 84(5):944–959.

Turner, Margery A., Claudia Aranda, Diane K. Levy, Rob Pitingolo, Rob Santos and Doug Wissoker. 2013. "Housing Discrimination Against Racial And Ethnic Minorities 2012." Washington, D.C.: U.S. Department of Housing and Urban Development, Office of Policy Development and Research.

Turner, Margery Austin, Stephen L Ross, George C Galster and John Yinger. 2001. "Discrimination in Metropolitan Housing Markets: National Results from Phase 1 of the HDS 2000." Washington, D.C.: Urban Institute and US Department of Housing and Urban Development.

Tyler, Tom R. 2004. "Enhancing Police Legitimacy." *Annals of the American Academy of Political and Social Science* 693:84–99.

Tyler, Tom R. 2006. *Why people obey the law*. Princeton University Press.

Weaver, R. Kent. 2014. "Compliance regimes and barriers to behavioral change." *Governance* 27(2):243–265.

Weaver, R. Kent. 2015. "Getting People to Behave: Research Lessons for Policy Makers." *Public Administration Review* 75(6):806–816.

White, Ariel R., Noah L. Nathan and Julie K. Faller. 2015. "What Do I Need to Vote? Bureaucratic Discretion and Discrimination by Local Election Officials." *American Political Science Review* 109(1):129–142.

Yinger, John. 1986. "Measuring Discrimination with Fair Housing Audits: Caught in the Act." *American Economic Review* 76:881–893.

Yinger, John. 1995. *Closed doors, opportunities lost: The continuing costs of housing discrimination*. Russell Sage Foundation.

Yuan, Shuai, Hao Helen Zhang and Marie Davidian. 2012. "Variable selection for covariate-adjusted semiparametric inference in randomized clinical trials." *Statistics in Medicine* 31(29):3789–3804.

# FIGURES

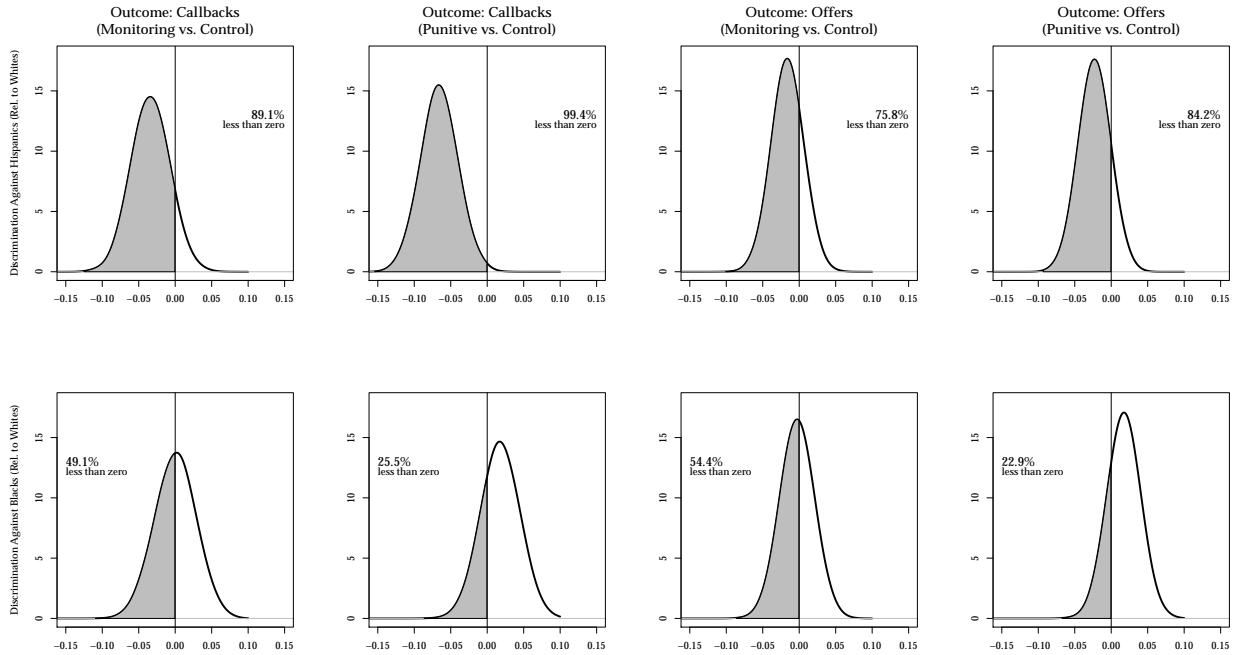## Main results on discrimination levels and treatment effects



**(a)** Outcome 1: Net Discrimination in Receiving Callbacks



**(b)** Outcome 2: Net Discrimination in Receiving Offers

**Figure 1:** Each panel shows, for each of the two objective outcome measures, levels of favorable treatment for different racial groups (top left quadrant), differences in favorable treatment rates between groups (i.e., net discrimination levels) by treatment assignment (lower left quadrant), differences in favorable treatment rates across treatment conditions for the same group (top right quadrant), and the effects of treatment assignment on net discrimination levels relative to the control or monitoring comparison group (lower right quadrant) with weighted nonparametric estimates shown using open markers and regression estimates adjusted using block fixed effects and inverse probability weighting shown using filled markers. Lines indicate 95% confidence intervals. Our main quantities of interest are highlighted in light gray.

**Figure 2:** Posterior densities of treatment effects on discrimination against Hispanics (top) and African Americans (bottom) relative to whites. Columns correspond to combinations of outcome measure (callbacks, offers) and treatment message (punitive, monitoring) relative to control. The area under the curve below zero is shaded in each plot, and the mean of each posterior density is shown with a dotted line. Each posterior is estimated via 10,000 Monte Carlo draws from the marginal distribution of $\beta_1$ conditional on $\sigma^2$ and $y$ given improper uniform priors, $\beta_1|\sigma^2,y \sim N(\hat{\beta}_1, V_{\beta_1}\sigma^2)$.